# LMQ-Sketch: Lagom Multi-Query Sketch for High-Rate Online Analytics

# Martin Hilgendorf □ □

Chalmers University of Technology, Gothenburg, Sweden University of Gothenburg, Sweden

#### Marina Papatriantafilou □

Chalmers University of Technology, Gothenburg, Sweden University of Gothenburg, Sweden

#### Abstract

Data sketches balance resource efficiency with controllable approximations for extracting features in high-volume, high-rate data. Two important points of interest are highlighted separately in recent works; namely, to (1) answer multiple types of queries from a single data structure built in one pass over the data, and (2) perform both queries and updates concurrently. In this work, we now tackle the new challenges arising when combining these useful directions together.

We investigate the trade-offs around efficiency, consistency, and accuracy to be balanced and synthesize key ideas into LMQ-Sketch, a single, composite data sketch supporting concurrent updates and multiple queries (frequency point queries, frequency moments  $F_1$ , and  $F_2$  as representative selection). Our method "LAGOM" is a cornerstone of LMQ-Sketch for low-latency global querying (<100 µs), combining freshness, timeliness, and accuracy with a low memory footprint and high throughput (>2B updates/s). We analyze and evaluate the accuracy of LAGOM, which builds on a simple geometric argument and efficiently combines work distribution with synchronization for proper concurrency semantics – monotonicity of operations and intermediate value linearizability. Comparing with state-of-the-art methods, which, as mentioned, provide either mixed queries or concurrency separately, LMQ-Sketch shows highly competitive throughput, with additional accuracy guarantees and concurrency semantics, while also reducing the required memory budget by an order of magnitude. We expect the methodology to have broader impact on concurrent multi-query sketches.

**2012 ACM Subject Classification** Computing methodologies  $\rightarrow$  Shared memory algorithms; Theory of computation  $\rightarrow$  Concurrent algorithms; Information systems  $\rightarrow$  Data structures

Keywords and phrases Concurrent Data Structures, Data Sketches, IVL, Freshness, Synchronization

Digital Object Identifier 10.4230/LIPIcs.DISC.2025.36

Related Version Previous Version: https://arxiv.org/abs/2506.16928

**Supplementary Material** Software: https://gitlab.com/marhilg/lmq-sketch [26] archived at swh:1:dir:d33d3a06fac57f29b697bc57c19fa414b6f7c8cb

Funding Work supported by the Swedish Research Council project EPITOME 2021-05424, by the Marie Skłodowska-Curie Doctoral Network RELAX-DN, funded by the European Union under Horizon Europe 2021-2027 Framework Programme (grant nr. 101072456) and by Chalmers University AoA frameworks Energy and Production, WP. INDEED, and "Scalability, Big Data and AI".

# 1 Introduction

Data sketches provide approximate summaries of data streams and can answer questions of interest efficiently, with bounded memory requirements. Examples include estimation of element frequencies, set/multi-set size, frequency moments (norms), frequent elements, distributions, quantiles, and more [1, 5, 8, 10, 13, 15]. Summarizations involve suitable

hash functions, and the results are controllable approximations of the targeted aggregate – commonly for skewed data – in form of  $(\epsilon, \delta)$  guarantees (the estimate deviates at most  $\epsilon$ , with probability at least  $1-\delta$ ). The approximation reflects a trade-off between accuracy and sketch size, and, in consequence, the operations' time cost. Due to their usefulness in data analytics and feature extraction, sketches get a major role in data processing platforms [3, 2, 23, 28, 36].

Mixed query sketches. Data analysis requires more information than a single metric. Typically, this would require one sketch for each metric to be tracked [9], requiring additional memory and processing – a prohibitive overhead stressed as early as in [13] – and their contents may be inconsistent. Can it be done differently? Important questions regarding having a single sketch for answering mixed queries have been asked in works such as [16, 32, 38], showing potential to improve approximation guarantees and quality of results compared to using separate sketches. A notable sketch supporting, with auxiliary data structures, multiple queries is the Count-Min Sketch (CMS) [13, 15]. Besides, at the core of many statistics are frequency moments (the i-th such defined as the sum of the i-th power of the frequency of each element). E.g.,  $F_2$ , also called "surprise number", has uses in linear algebra [45], calculating the Gini index [7, 22], evaluating skewness of distributions, quantiles [20], range queries, heavy-hitters, histograms, or wavelet synopses helping for compression [11, 12, 19, 21]. In conjunction with other metrics,  $F_2$  enables to detect anomalies [10, 29, 33, 43] and associates with universal sketching [4, 31].  $F_2$  can be estimated via, e.g., CMS [14] or Fast-AGMS [10].

**Concurrency.** Given the rates of real-world, continuous streams, multi-threaded processing of updates and concurrent querying is both a *necessity* and a *challenge*, since, on one hand, the perceived order of updates by queries influences accuracy, and, on the other hand, synchronization overhead can dramatically influence throughput, operation timeliness and result freshness. Work to address concurrency among queries and updates on sketching was initiated in recent years, for single-query sketches [18, 27, 34, 42].

Our Targets. We integrate these directions, opening up the arising challenges not tackled by existing approaches in balancing concurrency and accuracy with resource footprint and timeliness, through workload distribution, data structure design, and synchronization for proper consistency guarantees. We target a single, low memory-footprint composite data structure with concurrent updates and queries for a representative selection of queries:  $F_1$ ,  $F_2$ , and point queries (individual element frequencies), with extensibility to other queries in mind.

**Challenges.** Concurrency adds new factors and trade-offs; the influence is complex: stronger consistency can imply better compliance with sequential sketch accuracy bounds but also higher synchronization overhead, increasing query latency and staleness, in turn aggravating accuracy and throughput. Further, mixed queries should not diverge excessively in observed input to ensure consistency. In short: various trade-offs impact the metrics of interest: set of queries, result quality, resource footprint, operation timeliness, and scalability.

Idea and Contributions. We study these trade-offs in detail and construct LMQ-Sketch ( $Lagom^1$  Multi-Query Sketch), a low memory-footprint sketch supporting concurrent and mixed frequency queries. The input domain and data structure are partitioned to support multiple threads processing updates, a design also present in [42] albeit for point queries only, showing benefits in accuracy and memory-efficiency. However, queries spanning multiple

partitions, global queries, face higher operation complexity, synchronization overhead and latency. To this end, our algorithmic design LAGOM employs Intermediate Value Linearizability (IVL) [39], a modular consistency property suitable for sketches, admitting more efficient implementations compared to requiring linearizability (unnecessarily strong for sketches which approximate by definition) and which can preserve  $(\epsilon, \delta)$  bounds of the sequential counterparts. LAGOM demonstrates significant gains in efficiency, timeliness, and result accuracy. It gathers "just enough" information (shown using a geometric argument) from partitions while maintaining accuracy in line with sequential methods. We also show monotonicity of scans [17] for cross-query consistency of LMQ-Sketch.

The analysis is complemented by a detailed empirical study using real-world CAIDA data [6] and synthetic stress-tests for extended insights. We compare to state-of-the-art methods: SW-SKT [9] employs separate sketches for multiple, albeit non-concurrent, queries; Delegation Sketch [42] supports concurrent updates and queries (for point queries only). Further, we define elementary baselines based on the competing factors in the trade-offs. Our evaluation shows how LMQ-Sketch efficiently balances memory, accuracy, and concurrency, scaling beyond 2B updates/sec concurrently with high-rate queries. We also propose a methodology to experimentally study the impact of concurrency on accuracy.  $F_2$  query latency between 1–100 µs implies freshness with error below 0.01 %, with a sketch just 4 MiB in size.

**Roadmap.** After preliminaries (§ 2) and problem analysis (§ 3), we present and analyze the methods constituting LMQ-Sketch (§ 4, 5), detail our empirical evaluation (§ 6, with open-source implementation [26]) and other related work (§ 7), concluding in § 8. Throughout, claims are backed by the main ideas; proof sketches are provided in Appendix A.

# 2 Background

Count-Min Sketch (CMS) and Frequency Moments. CMS [13] estimates the frequency of keys in the input stream, via point queries, with potential overestimation by a factor  $\epsilon$  with probability  $1 - \delta$ . It uses a matrix of counters, CMS $[H \times K]$ , alongside H hash functions, each mapped to a row. H and K balance memory against accuracy, by  $H = \lceil e/\epsilon \rceil$  and  $K = \lceil \ln(1/\delta) \rceil$ . A CMS update for a key a increments CMS $[j, h_j(a)]$  for each row j. To answer a point query for a, its frequency f(a) is estimated as  $\hat{f}(a) = \min_{i} \text{CMS}[j, h_i(a)]$ .

For a stream S of keys from a domain U, the n-th frequency moment is [46]:  $F_n = \sum_{a \in U} f(a)^n$ . CMS provides an exact  $F_1$  as  $\sum_k^K \mathsf{CMS}[j,k]$  for any row j. To estimate  $F_2$ , which gives an indication of the skewness of the frequency distribution, [14] introduces  $CM^-$  and  $CM^+$  for different levels of skewness, encoded by the z parameter of a Zipfian, as many real-world phenomena follow Zipf's Law;  $CM^-$  for  $z \leq 1$ , while  $CM^+$  yields better accuracy for skewed data (z > 1). Similar to CMS point queries,  $CM^+$  can only overestimate.

▶ **Definition 1** ( $CM^+$  [14]).  $CM^+$  estimates  $F_2$  from a  $CMS[H \times K]$  by  $\min_j \sum_k^K CMS[j,k]^2$  with relative error  $1 + \epsilon$  with probability  $1 - \delta = 1 - \frac{3}{4}^{-H}$ , where  $\epsilon = O\left(K^{\frac{-(1+z)}{2}}\right)$ .

**Augmented Sketch (ASketch).** ASketch [41] is a CMS extension for skewed streams, where a small subset of keys account for the majority of updates. It uses a *filter*, denoted AF, behaving as a fixed-size map of keys and counters, to count occurrences of the heaviest

Lagom (link) describes "just the right amount", indicating appropriateness rather than suggesting lack.

keys separately from the CMS matrix, improving throughput and accuracy. Alg. 1 shows its operations; updateAndPQ (line 1.10) is a simple extension of the CMS update to also get a point query estimate for the key by returning the min of the updated counters. AF can be efficiently searched for a key (line 1.2), e.g., using SIMD; if the key is found, the update increments just a single counter in AF (line 1.3), bypassing calculation of H hash functions and counter increments, and avoiding collisions with light keys in CMS. Keys are swapped between AF and CMS as needed to keep the most frequent keys in the filter (line 1.11). Note that  $\Delta$  (line 1.13) is the exact number of occurrences of key a while resident in the filter.

Algorithm 1 ASketch – with highlighted enhancements for LMQ-Sketch, detailed in § 4.

```
Function ASketchUpdateEnhanced(key a, value v)
  1.1
                if a in AF then
                                                                                                                                                                             Dincrement count in filter
  1.2
                      \begin{array}{l} \mathsf{AF}[a] \mathrel{+}= v; \\ \mathsf{AF}^{\mathrm{mavg}}[a] \mathrel{\longleftarrow} wv + (1\!-\!w) \mathsf{AF}^{\mathrm{mavg}}[a]; \end{array} 
  1.3
                                                                                                                                    \triangleright w = 0.8, bias towards recent changes in skew
  1.4
                else if AF not full then
                                                                                                                                                                                              ▷Add key to filter
  1.5
                     \begin{array}{l} \mathsf{AF}[a] \leftarrow v; \\ \mathsf{AF}^{\mathrm{old}}[a] \leftarrow 0; \\ \mathsf{AF}^{\mathrm{mavg}}[a] \leftarrow v; \end{array}
  1.6
  1.7
  1.8
  1.9
                else
                                                                                                                                                                                                      ▷ Update sketch
                      estimatedFreq \leftarrow CMS.updateAndPQEnhanced(a, v);
1.10
1.11
                      if estimatedFreq > min_{key in AF}(AF[key]) then
                            \begin{aligned} & \min \text{Key} \leftarrow & \operatorname{argmin}_{\text{key in AF}}(\text{AF[key]}); \\ & \Delta \leftarrow & \text{AF[minKey]} - \text{AF}^{\text{old}}[\text{minKey}]; \end{aligned}
1.12
1.13
                            CMS.updateEnhanced(minKey, \Delta);
1.14
                            \begin{array}{l} \mathsf{AF}[a] \leftarrow \mathsf{estimatedFreq}; \\ \mathsf{AF}^{\mathrm{old}}[a] \leftarrow \mathsf{estimatedFreq}; \\ \mathsf{AF}^{\mathrm{mavg}}[a] \leftarrow v; \end{array}
1.15
1.16
1.17
1.18
          Function ASketchQuery(key a)
               if a in AF then return AF[a]:
1.19
               else return CMS.pointQuery(a);
1.20
```

# **Problem Description and Analysis**

We target to estimate item frequencies and frequency moments using a single sketch, built from one pass over a high-rate input data stream. The stream consists of tuples (a, v), representing v occurrences of key a; we consider the cash register model, so v > 0. We target multi-core shared memory systems supporting atomic primitives including fetch-and-add and compare-and-swap. Threads communicate via a coherent memory model, and do not fail or crash. Operations on the sketch object are updates, processing input tuples, and queries of various kinds, and may all be invoked concurrently by different threads. A solution needs to balance a delicate multi-way trade-off that we analyze in the following.

**Accuracy and Consistency.** Strong consistency for queries requires atomic views of the sketch state, which may involve significant overhead; weakening these semantics arbitrarily may reduce accuracy in unclear ways. We aim for a balance: fresh query results with low synchronization overhead, yet with clear semantics, avoiding the accuracy pitfalls of both. To facilitate this, we build on the idea of Intermediate Value Linearizability (IVL) [39], that extends weak regularity [35], requiring that a concurrent query return a value in the interval between the minimum and maximum ones it could return in any linearization of the execution. In the cash register model, the considered query values increase monotonically with updates; hence:

▶ **Observation 2.** For a monotonically increasing function, the minimum return value of query Q, concurrent with updates, is given by an idealized return value of the query (e.g. by a perfect external observer) that observes updates up to the ones completed before Q's start (i.e. no overlapping ones). Similarly, the maximum return value of Q is given by an idealized return value that observes also all updates overlapping Q.

Q observes all updates completed before it started exactly once; the range of permitted return values, effectively the "inaccuracy" accepted as a consequence of IVL, is due to updates concurrent with Q. As we target high data rates, a short query latency is critical for minimizing this range and preserving freshness. To this end, performance can be improved by further relaxing consistency where tolerable: [39] combines r-relaxation [24] with IVL—allowing Q to also miss up to r updates preceding it—paraphrased here:

▶ **Definition 3.** A query Q is r-relaxed IVL if it returns a value between the min and max values it could return in any linearization that may reorder Q with up to r operations.

**Multiple Query Types.** The targeted queries include *local* and *global* ones. The former involve part of the data, e.g., a point query for a certain key; the latter consider complete contents of the sketch and associate with semantics for bulk operations [17, 35, 37]. Supporting multiple queries also concerns their relative consistency; we identify *operation monotonicity* as a useful notion to tell us how aligned the views of different queries are; paraphrasing from [17, 35]:

▶ **Definition 4** (Monotonicity of scans). For queries  $Q_1$  and  $Q_2$ , where  $Q_1$  precedes  $Q_2$  (denoted  $\longrightarrow$ ), all updates observed (i.e., accounted for) by  $Q_1$  must also be observed by  $Q_2$ .

Accuracy and Memory. The  $(\epsilon, \delta)$  bounds of sketches allow trading more memory for improved accuracy; e.g., in CMS, a query can overestimate due to hash collisions – more counters (or filters as in ASketch) reduce collisions. In the sequential setting, bounds are well-studied, and ideally should be preserved under concurrency. Similarly important are memory layout and management. Partitioning the input domain can aid parallelism as partitions can safely be summarized concurrently, enhancing accuracy and efficiency, shown for local queries in [42]. However, aforementioned accuracy-consistency challenges arise for global queries.

**Summary:** Goals and Challenges. (G1) *IVL* and associated relaxations are identified as consistency targets, due to explainability and efficiency advantages for synchronization. Low query latency becomes a significant goal, allying with freshness. (G2) Monotonicity of scans is desirable for multiple query types. (G3) Maintaining a single data structure for answering multiple queries can enable better use of the memory budget if the accuracy for each query can benefit from the full memory. We identify partitioning as an ancillary approach to improving accuracy and memory utilization, also aiding parallelization and hence throughput.

# 4 LMQ-Sketch – Design and Coordination

We describe the design space and coordination in LMQ-Sketch to enable concurrent estimation of multiple queries from a single, composite data structure (Fig. 1) in conjunction with partitioning, as motivated in § 3, for consistency, memory efficiency, and accuracy-friendliness.

The input domain is split into P partitions, as is the data structure (horizontal slices in Fig. 1). Updates are performed by P threads, each "owning" one partition (ASketch and other components) with exclusive update access to it. When thread  $T_i$  processes an update for a key owned by another partition, the update is buffered in a corresponding delegation filter DF (Alg. 2). Similar to ASketch filters, delegation filters behave as maps of C keys, and are small enough to be searched efficiently with SIMD operations. Periodically (upon a triggering condition, line 2.13) filters are handed over to the owning thread (using, e.g., a Treiber stack [44]), and contents are transferred in bulk to its local ASketch (Alg. 3).

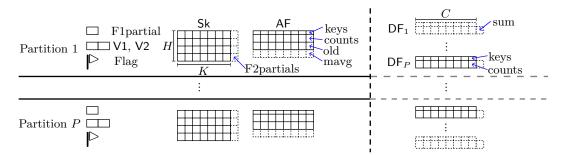


Figure 1 LMQ-Sketch's composite data structure. Partitions, horizontal slices of memory, are modifiable only by the owner thread. All updates eventually reach partition-local ASketches (Sk and AF, managing owned heavy keys), to the left of the vertical boundary. Delegation filters, DF, on the right, buffer updates owned by other partitions and are periodically handed over to the owning thread. C is a parameter, e.g. 16. Dotted components are for optimizations.

**Table 1** Notation used throughout.

Symbol	Description	Symbol	Description
P	No. of partitions & threads	Sk.F2partials	Sum for $CM^+$ per row (Def. 1)
K, H	CMS rows and columns	$T_i.AF$	ASketch filter for $T_i$
C	No. of slots in ASketch and delegation filters	$AF[a], \ AF^{\mathrm{old}}[a]$	Count & Old count of heavy key $a$ in $AF$
B	Max. no. of updates buffered in delegation filter	$AF^{\mathrm{mavg}}[a]$	Projected count of buffered occurrences of heavy key $\boldsymbol{a}$
$\overline{Owner(a)}$	Get partition owning key a	$T_i.DF_j$	Delegation filter of $T_i$ for up-
$T_i$	Updater thread for partition $i$		dates owned by partition $j$
	1	$DF_j[a]$	Count of key $a$ buffered in $DF_j$
$T_i.Sk$	Local sketch for $T_i$	$DF_i.\mathrm{size},$	Number of keys and updates
$T_i$ .F1partial	No. of updates completed by $T_i$	$DF_j.sum$	buffered in $DF_j$

Algorithm 2 Delegation Sketch update on  $T_i$ . Enhancements for LAGOM in green.

2.2

2.3

2.8

```
Algorithm 3 Processing delegated up-
dates on T_i. Enhancements for LAGOM in
green.
```

```
2.1 Function update(key a, value v)
          filter \leftarrow T_i.\mathsf{DF}_{\mathsf{Owner}(a)};
                                                                               3.1
           while filter.size = C or filter.sum \geq B do
                                                                               3.2
           processDelegatedUpdates();
 2.4
                                                                               3.3
          \overline{\mathbf{if}} a in filter then
 2.5
                                                                               3.4
              filter[a] += v;
 2.6
                                                                               3.5
 2.7
              filter.sum +=v;
                                                                               3.6
                                                                               3.7
              \mathrm{filter}[a] \leftarrow v;
 2.9
                                                                               3.8
2.10
              filter.size += 1
                                                                               3.9
              filter.sum +=v:
2.11
                                                                             3.10
          T_i.F1partial +=v;
if filter.size =C or filter.sum \geq B then
2.12
2.13
2.14
              T_{\mathsf{Owner}(a)}.pendingFilters.push(filter);
```

```
{\bf Function} \ {\tt processDelegatedUpdates}
    while T_i pending Filters is not empty do
       Wait until T_i.beingScanned = false;
       T_i.V1++; filter \leftarrow T_i.pendingFilters.pop();
       foreach a in filter do
          T_i.ASketchUpdateEnhanced(a, filter[a])
       Clear filter contents;
       filter.size \leftarrow 0:
       filter.sum \leftarrow 0;
       T_i.V2++;
```

**Point Queries.** The idea of domain-partitioning and delegation is also the basis of the Delegation Sketch [42] for point queries, summarized here for self-containment. Moreover, we here show their consistency properties as a query type in LMQ-Sketch. A point query (Alg. 4) for key a is answered by  $T_{\mathsf{Owner}(a)}$ , the thread owning a, estimating  $\hat{f}(a)$  as the sum of occurrences of a in the thread-local ASketch of  $T_{\mathsf{Owner}(a)}$  (line 4.2) and in relevant delegation filters of other threads (line 4.3). Skew in the input data is beneficial to f(a)'s accuracy; frequent keys are often present in the filters where they are counted accurately, akin to ASketch. Based on the IVL definition and Obs. 2, along with the fact that entries in the partition-local sketch are increasing with subsequent updates, and that only the thread owning a key performs updates and queries on this sketch (hence cannot miss completed updates or double-count any), we have:

▶ Lemma 5. Delegation Sketch-based PQ is an IVL implementation of ASketch point query.

#### 4.1 Global Query Trade-offs

This approach of assigning query work to the owner thread will however not work for global queries which span all partitions. We investigate the trade-offs relating concurrency and accuracy for such bulk queries, as identified in § 3: the extremes of the spectrum serve as baselines, followed by our balancing approach.

STRICT: A baseline for strong consistency, using a Readers-Writer lock to allow shared access for updates and point queries, but exclusive access for global queries (acting as writers) to see a consistent global state.  $F_1$  is estimated by summing all counters on any row of each  $T_i$ . Sk and counts in all ASketch and delegation filters, yielding the number of updates completed before the query.  $F_2$  can be estimated using  $CM^+$  by merging all sketches and flushing filter contents (treating them as updates), thus preserving its error bounds. Note several drawbacks: updates incur locking overhead even when no global query is taking place; further, while thread-safe, the result can be stale due to "stopping the world" – input tuples keep arriving but are not processed or visible to the query.

<u>NOSYNC</u> (§ 4.2): This approach targets optimistic synchronization, for exploring freshness and concurrency maximization, at the cost of consistency, possibly risking mis-calculations.

<u>LAGOM</u> (§ 4.3): Our balanced method for concurrent queries with low latency, "just-enough" calculation, and lightweight synchronization. Query results have stronger semantics than NOSYNC, and are based on fresher state than STRICT.

#### 4.2 Nosync for Bulk Queries

We explore a baseline with no synchronization between updates and global queries, to explore freshness and concurrency maximization (at the cost of consistency).

 $F_1$ . To estimate  $F_1$ , improving efficiency and freshness compared to an approach as in STRICT, we introduce partial results in form of per-thread counters ( $T_i$ .F1partial) for the number of performed updates (line 2.12). The  $\hat{F}_1^{\text{Nosync}}$  query (Eq. 1) sums these counters, improving locality, efficiency, and NUMA-friendliness compared to STRICT.

$$\hat{F}_1^{\text{NOSYNC}} = \sum_{i}^{P} T_i.\text{F1partial} \tag{1}$$

▶ Lemma 6.  $\hat{F}_1^{NOSYNC}$  estimates  $F_1$  with IVL semantics.

 $F_2$ . A global  $\hat{F}_2$  equals the sum of per-partition  $\hat{F}_2$ , as partitions contain independent, non-overlapping parts of the input. However, an approach as for  $F_1$  with per-partition partial results poses obstacles: threads would update the partial result at the *owning partition*, reintroducing contention and serialization on shared data into the delegation design.

Instead, LMQ- $all^+$  (Eq. 2) uses  $CM^+$  to estimate  $F_2$  for each  $T_i$ .Sk: for each key a in  $T_i$ .AF,  $\mathcal{H}^{\text{Nosync}}$  (Eq. 3) performs a point query (Alg. 4, using the fast path on line 1.19) to obtain  $\hat{f}(a)^2$ , the  $F_2$  contribution of a. However, some occurrences of a may be stored in

#### Algorithm 6 $\hat{F}_2^{\text{LAGOM}}$ on own thread. **Algorithm 4** Point query on $T_{\mathsf{Owner}(a)}$ . 6.1 Function queryF2Lagom Function pointQuery(key a) 4.1 result $\leftarrow 0$ ; 6.2 result $\leftarrow T_{\mathsf{Owner}(a)}.\mathsf{ASketchQuery}(a);$ foreach $T_j$ do result $+=T_j.\mathsf{DF}_{\mathsf{Owner}(a)}[a];$ 4.2 foreach $T_i$ do 6.3 6.4 $T_i$ .beingScanned $\leftarrow$ true; return result: 6.5 repeat **Algorithm 5** $\hat{F}_2^{\text{NOSYNC}}$ on own thread. 6.6 $localResult \leftarrow min(T_i.Sk.F2partials);$ 6.7 foreach a in $T_i$ . AF do Function queryF2Nosync 5.1 6.9 $aFreq \leftarrow T_i.AF[a];$ 5.2 $\operatorname{arreq} \leftarrow I_i.\mathsf{AF}[a];$ $\operatorname{aFreq} += P \cdot T_i.\mathsf{AF}^{\operatorname{mavg}}[a] / 2;$ $\operatorname{localResult} += \operatorname{aFreq}^2 - (T_i.\mathsf{AF}^{\operatorname{old}}[a])^2;$ for each $T_i$ do 6.10 result $+=\hat{F_2}^{CM^+}(T_i.Sk);$ 6.11 result $+= T_2$ ( $T_i$ .SK), foreach a in $T_i$ .AF do aFreq $\leftarrow T_i$ .AF[a]; foreach $T_j$ do aFreq $+= T_j$ .DF<sub>i</sub>[a]; $v1 \leftarrow T_i.V1;$ 6.12 5.6 ⊳Retry until match 6.13 5.7 $T_i$ .beingScanned $\leftarrow$ false; result $+= aFreq^2 - (T_i.AF^{old}[a])^2;$ result += localResult; 5.8 return result; return result: 6.16 5.9

 $T_i$ .Sk and their contribution included in the  $CM^+$  estimate; this needs to be subtracted. We provide a geometric illustration, elaborating on the argument for this calculation, in § 5.3 and Fig. 2c. Alg. 5 shows the synchronization (or, rather, lack thereof) for concurrent LMQ- $all^+$ .

$$LMQ\text{-}all^{+} = \sum_{i}^{P} \left( \hat{F}_{2}^{CM^{+}}(T_{i}.\mathsf{Sk}) + \sum_{a}^{T_{i}.\mathsf{AF}} \mathcal{H}^{\mathsf{NOSYNC}}(a) \right)$$
 (2)

$$\mathcal{H}^{\text{Nosync}}(a) = \left(T_o.\mathsf{AF}[a] + \sum_{j}^{P} T_j.\mathsf{DF}_o[a]\right)^2 - (T_o.\mathsf{AF}^{\text{old}}[a])^2 \quad \text{where } o = \mathsf{Owner}(a) \quad (3)$$

▶ **Observation 7.** An  $\hat{F}_2^{NOSYNC}$  query Q can miss or double-count updates, due to data movement by overlapping processing of delegated updates (Alg. 3).

While NOSYNC avoids the overhead of STRICT, serving as a baseline for maximal concurrency and exploration of freshness, the lack of synchronization leaves weak consistency guarantees for  $F_2$  estimations<sup>2</sup>, implying arbitrary fluctuations from the accuracy bounds.

#### 4.3 Lagom for Bulk Queries – $F_2$

We now present our design to enable high throughput with clear concurrency semantics. We build upon some of the described components, i.e. support updates via Alg. 2, with highlighted enhancements, point queries via Alg. 4 (§ 3.11) and  $\hat{F}_1$  via Eq. 1 (§ 4.2). The crux to efficient, concurrent  $\hat{F}_2$  queries with consistency and accuracy guarantees lies in two key ideas: (1) efficient maintenance of partial results, which, based on a geometric observation, allows to argue about accuracy relative to the sequential bounds; and (2) lightweight synchronization implying IVL, over only few variables to scan, due to how partial results are maintained.

**Partial results for F\_2.** While a global  $\hat{F}_2$  can be obtained from summing per-partition  $\hat{F}_2$  values as in  $\hat{F}_2^{\text{Nosync}}$ , maintaining these per partition is, unlike  $F_1$ , not straightforward. Instead, we compute partial results to simplify the heaviest operations in Eq. 2, 3, which are: (1)  $CM^+$  on the underlying CMS, that requires reading all  $H \times K$  counters of the sketch, and (2)  $\mathcal{H}^{\text{Nosync}}$  (Eq. 3), that reads all delegation filters for each heavy key.

To be precise,  $\hat{F}_2^{\text{Nosync}}$  is quiescence-consistent [25]: in absence of concurrent updates, the "bad things" in Obs. 7 cannot happen. But this property is to little purpose in high-rate scenarios.

For (1), we maintain the per-row sums in Def. 1 for each sketch row, labeled F2partials in Fig. 1. This is done within the enhanced CMS update in Alg. 1, with incremental (associative) calculations.  $\hat{F}_2^{CM^+}$  can then return the min of these H values (highlighted in Eq. 4).

For (2), we aim to avoid scanning delegation filters when estimating the frequency of heavy keys. However, as filters can buffer a significant number of updates, particularly for skewed streams, their contents should not be ignored to avoid excessive underestimation (analyzed in § 5.3). To compensate, we devise  $\mathcal{H}^{\text{Lagom}}$  (Eq. 5) as a lightweight estimation of the number of heavy-key occurrences buffered in delegation filters. For each slot, ASketch filters maintain an exponentially weighted moving average of the number of occurrences received during filter flushes (line 1.4, invoked from line 3.7). We replace the scan over delegation filters in  $\mathcal{H}^{\text{Nosync}}$  with the following idea (highlighted in Eq. 5), requiring only one read from the ASketch filter: there are P delegation filters for a partition, projected to contain  $AF^{\text{mavg}}[a]$  counts of a once full and getting flushed. At any point when the query scans a partition, filters are on average half full. This yields LMQ-proj<sup>+</sup> (Eq. 5). In § 5 we analyze the accuracy implications and the association with the sequential bounds.

$$LMQ\text{-}proj^{+} = \sum_{i}^{P} \left( \min(T_{i}.\mathsf{Sk}.\mathsf{F2partials}) + \sum_{a}^{T_{i}.\mathsf{AF}} \mathcal{H}^{\mathsf{LAGOM}}(a) \right)$$

$$\mathcal{H}^{\mathsf{LAGOM}}(a) = \left( T_{o}.\mathsf{AF}[a] + P \cdot \frac{T_{o}.\mathsf{AF}^{\mathsf{mavg}}[a]}{2} \right)^{2} - \left( T_{o}.\mathsf{AF}^{\mathsf{old}}[a] \right)^{2} \quad \text{where } o = \mathsf{Owner}(a) \tag{5}$$

$$\mathcal{H}^{\text{Lagom}}(a) = \left(T_o.\mathsf{AF}[a] + P \cdot \frac{T_o.\mathsf{AF}^{\text{mavg}}[a]}{2}\right)^2 - \left(T_o.\mathsf{AF}^{\text{old}}[a]\right)^2 \quad \text{where } o = \mathsf{Owner}(a) \tag{5}$$

**Lightweight Synchronization.** The  $\hat{F}_2^{\text{LAGOM}}$  query (Alg. 6) uses lightweight synchronization to safely calculate LMQ-proj<sup>+</sup> concurrently with updates. To achieve the IVL-semantics goal (G1) set in § 3 and avoid the safety problems of  $\hat{F}_2^{\text{NOSYNC}}$  seen in Obs. 7,  $\hat{F}_2^{\text{LAGOM}}$  synchronizes with each partition i being scanned (i.e., with thread  $T_i$ ) via a handshake, involving a pair of initially equal version numbers (as in [30]) for detecting a concurrent filter flush, and an atomic flag set by the query, to signal when the partition is being scanned.

In the common case, when  $T_i$  processes a tuple (a, v), Alg. 2 with enhancements finds available space in  $T_i.\mathsf{DF}_{\mathsf{Owner}(a)}$  (condition on line 2.3);  $T_i.\mathsf{DF}_{\mathsf{Owner}(a)}[a]$  is then incremented by v (line 2.6 or 2.9), but a handover is not triggered immediately (line 2.13 is false).  $\hat{F}_2^{\text{LAGOM}}$ queries cannot directly account for this update before  $T_i.\mathsf{DF}_{\mathsf{Owner}(a)}$  is flushed to  $T_{\mathsf{Owner}(a)}$ .  $\hat{F}_2^{\text{Lagom}}$  targets a consistent view of each partition's relevant information independently. Concurrent filter flushes (Alg. 3) are detected if version numbers mismatch (line 6.13), causing the query to retry, though only for at most one concurrent flush per partition, as subsequent flushes will be stalled by the flag (line 3.3). The sum of the collected per-partition partial results is then returned as  $\hat{F}_2$ . The update-query interaction implies:

▶ Lemma 8.  $\hat{F}_2^{LAGOM}$  gets an atomic snapshot per partition, and cannot deadlock with updates.

Optimizations on Filters: Bounds and Self-Delegation. To bound the interval between filter flushes, a parameter B limits the buffering capacity of filters (line 2.7 and 2.11). Lower values for B trade performance for accuracy – both for updates, as filters are unavailable more often, and for queries, where more overlapping flushes interfere with them. The impact of B is studied in  $\S$  5.3, according to this observation, motivating tracking of updates in  $\hat{F}_2^{\text{LAGOM}}$ :

**Description 9.** At most r = PB single increment updates can be buffered in delegation filters and may thus not be explicitly observed by  $\hat{F}_2^{\text{LAGOM}}$  queries.

Further, in Delegation Sketch, a thread processing an update for a key it owns updates its ASketch directly. However, in LAGOM, ASketch updates require synchronizing with global queries, causing delays particularly in partitions owning heavy keys. To improve efficiency, each partition is extended with a *self-delegation filter* to buffer local updates and process them in bulk, also interfering less with queries.

#### 4.4 Multiple Query Types - Relative Consistency

In § 3 we identify monotonicity of scans [17] (G2) as intuitive and desirable behavior – a query should observe the updates observed by a previously performed query. Having a single data structure implies certain monotonicity compared to a disjoint per-sketch data structure. As LMQ-Sketch consists of multiple components, an update can become visible to queries of different types at slightly different points. Here we bound how much/little "the world can differ" between query types. To this end, we identify as sub-operations of an update U the atomic operations, denoted here  $U|_{op}$ , after which U becomes visible to a query of the respective type:  $U|_{PQ}$  at line 2.6 or 2.9,  $U|_{F_1}$  at line 2.12, and  $U|_{F_2}$  at line 3.11. We have  $U|_{PQ} \longrightarrow U|_{F_1} \longrightarrow U|_{F_2}$  from program order. The cross-query consistency properties are:

▶ **Lemma 10.** For queries  $Q_1 \longrightarrow Q_2$  the monotonicity of scans for each combination follows

	$Q_2$ :	PQ	$F_1$	$F_{2}$
	PQ	Monotonic	$P ext{-}relaxed\ monotonic$	$r ext{-}relaxed\ monotonic$
$Q_1$ :	$F_1$	Monotonic	Monotonic	$rP ext{-}relaxed\ monotonic$
	$F_{2}$	Monotonic	Monotonic	Monotonic

Note that: (1) The bounds are pessimistic, e.g. for  $PQ \longrightarrow F_I$  to deviate by P updates, all must occur in the partition queried by  $Q_1$ , while for all other updates except these, i.e., where  $U|_{F_1} \longrightarrow Q_2$ , full monotonicity of scans applies. (2) P is small relative to update rates. (3) Due to the compensation scheme of  $\hat{F}_2^{\text{LAGOM}}$ , even in the very unlikely case that the maximum deviation in operations observed occurs, actual results differ significantly less.

# 5 Accuracy of $F_2$ Estimation

We analyze the consistency and accuracy of LMQ-Sketch for  $\hat{F}_2$ , as done for PQ and  $F_1$  queries in § 3.11 and § 4.2. We develop a sequence of auxiliary designs (Table 2) and argue for their properties, ultimately arriving at  $\hat{F}_2^{\text{LAGOM}}$ . We develop a geometric interpretation of  $\hat{F}_2$  in terms of areas of squares, illustrating accuracy properties. For each step, we describe the organization of the data structure and the  $F_2$  query on it in a sequential setting along with its accuracy, followed by a parallelized construction and its concurrency semantics.

#### 5.1 Count-Min Sketch & Partitioning

We begin by comparing a "wide" CMS with a partitioned one, on fixed memory budget.

Wide: A Count-Min Sketch with $H$ row	and $P \times K$ columns.
A toy example is shown to the side.	

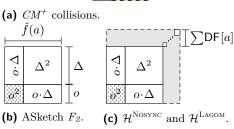
1		2, 3	4
4	1, 2		3

**Sequential.**  $CM^+$  (Def. 1) estimates  $F_2$  from WIDE. Consider the example CMS; each number represents a unique key hashed to that counter, and is simultaneously the frequency of that key (i.e., f(1) = 1, f(2) = 2, etc.). Colliding keys are separated by commas; CMS will only store their sum. The true  $F_2$  is  $1^2 + 2^2 + 3^2 + 4^2 = 30$ .  $CM^+$  computes

**Table 2** Sequential (for accuracy reasoning) & concurrency-aware  $\hat{F}_2$  (for IVL reasoning).

	$2^2$
$1^2$	

Data structure	Sequential	Concurrent
Wide (CMS)	$CM^+$	$CM_{conc}^{+}$
PARTCMS	$PARTCMS^+$	$PARTCMS_{conc}^{+}$
PARTAS	$PARTAS^+$	$PARTAS^{+}_{conc}$
LMQ-Sketch	LMQ- $all$ <sup>+</sup>	$\hat{F_2}^{ m Nosync}$
LMQ-Sketch	$PARTAS^+(\mathrm{LMQ})$	
LMQ-Sketch	$LMQ ext{-}proj^+$	$\hat{F_2}^{ ext{LAGOM}}$



**Figure 2** Geometric interpretation of  $\hat{F}_2$ .

 $1^2 + (2+3)^2 + 4^2 = 42$  from the first row and  $4^2 + (1+2)^2 + 3^2 = 34$  from the second, returning the smaller estimate. Overestimation arises from colliding keys. Fig. 2a illustrates the  $CM^+$  calculation for the counter containing (1,2). The true  $F_2$  contribution of the contained keys is 5, the sum of the plain areas.  $CM^+$  calculates  $(1+2)^2 = 9$ , the entire square, thus overestimating by an amount equal to the striped areas. The impossibility of underestimating is clear. Heavy keys induce particularly large extra areas, suggesting a need to treat them specially.

**Concurrent.** CMS updates and queries can be parallelized by, e.g., atomic fetch-and-add. In [39], Rinberg and Keidar show a similar construction to be an IVL implementation of CMS. On top of this parallel CMS, consider  $CM_{conc}^+$  as a concurrency-aware adaptation of  $CM^+$  using atomic reads; a similar argument as [39, Lemma 5.3] implies:

▶ **Lemma 11.**  $CM_{conc}^+$  is an IVL implementation of  $CM^+$ , preserving  $CM^+$ 's  $(\epsilon, \delta)$  bounds.

For reduced contention relative to WIDE, the space is now partitioned into P sketches:

	PARTCMS: $P$ partitions, each with a $H \times K$ CMS. Each partition
ı	sketches a subdomain of the input. (Toy example on the right).

1	1 2		3, 4			
		1, 2	4	3		

**Sequential.** Partitioning the memory of the example into P = 2 partitions (note each key is only present in one partition), to estimate  $F_2$  from PARTCMS,  $PARTCMS^+$  sums a  $CM^+$  estimate for each partition. In the example, 5 is calculated for the first partition and 25 for the second. Their sum is an  $\hat{F}_2$  for the complete PARTCMS.

▶ Observation 12. PARTCMS<sup>+</sup> can only improve accuracy of  $\hat{F}_2$  compared to CM<sup>+</sup> on WIDE, as it selects sketch rows with minimum overestimation independently for each partition.

**Concurrent.** Each partition is assigned a thread to perform updates. For now, assume that each thread only receives tuples for keys of its partition. Consider  $PARTCMS^+_{conc}$  as a parallelization of  $PARTCMS^+$  using atomic reads (similar to  $CM^+_{conc}$ ) to perform  $PARTCMS^+$  concurrently with updates. Similar reasoning as in Lemma 11 implies:

▶ Lemma 13. PARTCMS<sup>+</sup><sub>conc</sub> is an IVL implementation of PARTCMS<sup>+</sup>.

#### 5.2 Augmented Sketch & Partitioning

ASketch [41], § 2, introduced filters to track heavy keys more accurately. In a similar fashion to PARTCMS, partitioning can be applied to ASketch while keeping total memory constant: PARTAS: P ASketches, each of size  $H \times K'$  – reducing width to fit filters in same total memory.

- **Sequential.** Consider  $PARTAS^+$ : for each partition,  $CM^+$  is applied to its CMS; the  $F_2$ contribution of each key a in the filters is  $\hat{f}(a)^2$ , shown in Fig. 2b, where  $\hat{f}(a) = \mathsf{AF}[a]$ (line 1.19). Note,  $\hat{f}(a)$  contains  $o = \mathsf{AF}^{\mathrm{old}}[a]$  occurrences which are also in the CMS from times when a did not reside in the filter due to not being heavy enough. Hence,  $CM^+$  already included  $o^2$  (dotted area in Fig. 2b). To avoid double-counting this quantity,  $o^2$  is subtracted  $(\mathcal{H}^{\text{NOSYNC}}$  (Eq. 3) and  $\mathcal{H}^{\text{LAGOM}}$  (Eq. 5)). Finally, the contributions are summed to get an  $\hat{F}_2$ .
- $\blacktriangleright$  **Observation 14.** Higher P is beneficial to PARTAS<sup>+</sup>'s accuracy, reducing overestimation through P filters accurately counting heavy keys, while maintaining one-sided error as CM<sup>+</sup>.
- **Concurrent.** As in  $PARTCMS_{conc}^+$ , each partition is updated by a dedicated thread, (still) with the simplifying assumption that input tuples are distributed to the owning partition. Instead of unsafe scanning, consider  $PARTAS_{conc}^+$  which obtains an atomic snapshot of each partition-local ASketch (using any synchronization that can guarantee that) and computes an estimate as  $PARTAS^+$ , implying:
- ▶ **Lemma 15.** PartAS<sup>+</sup><sub>conc</sub> is an IVL implementation of PartAS<sup>+</sup>.

#### 5.3 **Concurrency Awareness – Delegation**

We now waive the simplification of input being distributed to the owning partition; threads delegate updates to the owner. We describe and compare three approaches to estimating  $F_2$ . LMQ-Sketch: Each partition of PARTAS uses P delegation filters to buffer updates, periodically flushed to the owning partition. K is adjusted to maintain the memory budget.

- All DFs. LMQ-all<sup>+</sup> (Eq. 2) extends  $PARTAS^+$  to include all delegation filters when estimating  $\hat{f}(a)$ , shown as  $\sum \mathsf{DF}[a]$  in Fig. 2c. Buffered occurrences of light keys are not considered, as they by definition do not significantly contribute to  $F_2$ , particularly for skewed streams which we target.  $\hat{F}_2^{\text{NOSYNC}}$  (§ 4.2) calculates this as a concurrent query, but lack of synchronization makes it unsafe (Obs. 7). Further,  $LMQ-all^+$  scales quadratically in P; all P delegation filters are read from other threads for each of C heavy keys in P ASketch filters.
- **No DFs.** To query more efficiently, we consider outright ignoring delegation filters, which will underestimate  $F_2$ , by applying  $PARTAS^+$  to the LMQ-Sketch data structure. However, the number of buffered, hence ignored, updates per partition is bounded (Obs. 9).
- $\blacktriangleright$  Lemma 16. PartAS $_{\mathrm{conc}}^+$  on LMQ-Sketch (i.e. its ASketch part) is an r-relaxed IVL implementation of LMQ-all<sup>+</sup> per partition, where r = PB.
- ▶ Corollary 17. Since IVL is a local property [39], PARTAS<sup>+</sup><sub>conc</sub> on LMQ-Sketch is an rP-relaxed IVL implementation of LMQ-all<sup>+</sup>.

Note the worst-case is extremely unlikely to happen, as all r updates would need to be for the same key (i.e. extreme skew), and filters do not become full simultaneously.

**Lagom.** To compensate for this relaxation yet maintain its high efficiency, LMQ- $proj^+$  (§ 4.3) projects the number of buffered occurrences of heavy keys to replace the exact calculation of  $\sum DF[a]$  in Fig. 2c and  $\mathcal{H}^{NOSYNC}$  (Eq. 3).

▶ Observation 18.  $\mathcal{H}^{LAGOM}$  (Eq. 5) projection compensates for up to r/2 ignored occurrences per heavy key, bringing the query result closer to the actual  $F_2$ , based on the reasoning of the geometric argument and observations 12 and 14.

From the above and Lemma 8, we have:

▶ **Lemma 19.**  $\hat{F}_2^{LAGOM}$  is an IVL implementation of LMQ-proj<sup>+</sup>.

The behavior of  $\hat{F}_2^{\text{LAGOM}}$  within the bounds is data- and execution-dependent and is evaluated in the next section. Summarizing the properties of LMQ-Sketch, we have:

▶ Corollary 20. LMQ-Sketch is a composite concurrent sketch data structure for multiple queries: (1) PQ (§ 3.11), an IVL implementation of the ASketch point query; (Lemma 5); (2)  $\hat{F}_1^{\text{NOSYNC}}$ , an IVL implementation of exact  $F_1$  (Lemma 6); (3)  $\hat{F}_2^{\text{LAGOM}}$ , an IVL implementation of LMQ-proj<sup>+</sup> (Obs. 18, Lemma 19). Queries follow the monotonicity properties in Lemma 10.

#### 6 Evaluation

**Baselines.** We study LMQ-Sketch relative to: (1) SW-SKT [9], the software implementation of SKT which targets FPGA-acceleration for high-rate sketching. SKT uses separate sketches for multiple queries ( $F_0$ , point queries, and  $F_2$  via HLL, CMS, and Fast-AGMS) updated in parallel, but no concurrent queries (first merging all thread-local sketches before querying). This comparison is for insights about answering mixed queries in a single sketch, from a memory, accuracy and scalability point of view. (2) Delegation Sketch [42], which supports concurrent updates but with point queries only. By comparing, we aim to understand the synchronization overhead for concurrent global queries. (3) The elementary STRICT and NOSYNC designs are used to evaluate the balance achieved by LAGOM. Note that simply adding concurrent queries to SW-SKT would imply effects similar to NOSYNC and STRICT and hence are not studied separately.

**Datasets and Hardware Platform.** From the CAIDA [6] network packet traces we extract  $18.5\,\mathrm{M}$  tuples of headers. Besides, synthetic datasets with skew  $z=1,\,1.5,\,2$  and 3 (each having  $100\,\mathrm{M}$  tuples sampled from a domain of size  $1\,\mathrm{M}$ ) are used to explore performance characteristics up to extreme skew levels where a majority of input tuples are owned by a single partition. All experiments are conducted on a 128-core AMD EPYC 7954 running at  $2.25\,\mathrm{GHz}$ , without using SMT. The server runs openSUSE Tumbleweed 20250211. LMQ-Sketch is implemented in C++23 [26], compiled with GCC  $14.2.1\,\mathrm{and}$  -03 -march=native.

**Metrics.** To study the scalability and efficiency of LMQ-Sketch, we measure throughput (updates per unit-time) at different thread counts, without and with concurrent queries. To study accuracy in presence of concurrency, we focus on latency of global queries, which has implications on the IVL-permitted interval; we propose a methodology to evaluate accuracy of IVL-associated queries (cf. § 5), measuring the difference of returned values relative to the estimated IVL-bounds, giving insight into the synergy of IVL semantics, query latency and freshness. Finally, we compare the impact of memory budget on accuracy for  $\hat{F}_2^{\text{LAGOM}}$ , state-of-the-art sketching techniques, and reference methods in § 5.

**Parameters.** Threads, partitions & memory: we stress-test with high thread counts (P=1 to 128). Memory budget per partition is constant at 32 KiB:  $H \times K = 8 \times 1024$ , while reducing K when having filters (cf. evaluations in [41, 42]), and keeping sketch rows (hash function computations) fixed. Filter size: C is kept at 16 as in [42], delegation filter buffer capacity is bounded at B=1000. Synchronization: we compare LAGOM with the aforementioned baselines, including two Delegation Sketch configurations: the plain one (unbounded B) and one with B=1000. Skewness of input data: z>1 for the synthetic data, and  $z\approx 1.4$  for the CAIDA data. Frequency of concurrent queries: stress testing with global queries at rates of 1000/s and point queries for 0.1% of update tuples, following the benchmark in [42].

**Experiments.** We begin by comparing the impact of synchronization on update operation's throughput when processing a complete input dataset from system memory (§ 6.1). Next, we investigate the accuracy and freshness of global query results, in relation to memory budget and concurrency. Higher weight lies on the more complicated  $F_2$  queries,  $\hat{F_1}$ 's behavior being studied mainly in terms of query latency and PQ's behavior evaluated in [42] (§ 6.2, 6.3, 6.4). For brevity, the results are shown in figures with summary descriptions in the caption, and main takeaways in associated paragraphs. More detailed discussions appear in Appendix B.

#### 6.1 Concurrency and Synchronization

**Design and Parameters.** Each dataset is processed with varying numbers of partitions and threads, P. First, we measure update throughput with no concurrent queries. We then repeat the measurement with a constant high load of concurrent queries:  $F_1$  and  $F_2$  at 1000/s each, and point queries at 0.1% of updates, for the designs that support them.

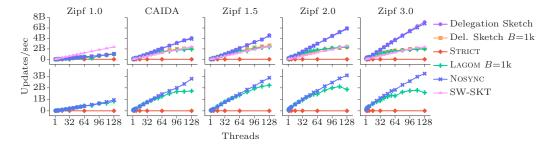


Figure 3 Mean update throughput, without (upper row) and with concurrent queries (lower row, for designs that support them). Fluctuations are small and omitted for clarity. Note the minimal overhead of LAGOM for global queries as it matches Delegation Sketch with *B*=1k. STRICT does not scale with increasing threads or skew. With concurrent queries, LAGOM stays close to NOSYNC, affirming the lightweight-ness of the synchronization design; increasing skew improves performance, as the latency hiding of delegation counteracts and postpones saturation even at extreme skew.

**Takeaways.** The results, shown in Fig. 3, imply that LMQ-Sketch with LAGOM is able to maintain state in a way that allows independence by different threads, enabling high processing throughput, yet still imparting the necessary consistency to serve the purpose of sketching for continuously and concurrently answering multiples queries in a streaming setting.

#### 6.2 Memory and Accuracy

We consider a sequential setting, to explore the impact of sketch memory budget on  $F_2$  estimation accuracy for several baselines and our methods. Note: There are two main parameters for memory budget, but with differing impact on accuracy: per-partition memory (given by H and K) and the total number of partitions P (also number of threads when parallelizing). In the partitioned design, an increase in any of these improves accuracy (§ 5).

**Design and Parameters.** The methods in the analysis (§ 5) form points of reference; additionally, we compare with Fast-AGMS [10], used in SW-SKT. For each method and skew  $z=1,\,1.5$  and 2, five synthetic datasets (length 100M, cardinality 1M) are processed under various memory budgets, for  $\hat{F}_2$  and the mean absolute percent error (MAPE) relative to the true  $F_2$ , at the end of the execution (i.e., based on the global state for the filterenabled designs). For each method in § 5, we use constant per-partition memory 32 KiB  $(H \times K = 8 \times 1024 \text{ with adjustments in presence of filters, as before); total memory increases with <math>P$ .

Fig. 4 shows the MAPE in log-scale, along the outcome for a single-partition Fast-AGMS at several budgets:  $6 \times 2^{13}$  (as in the original evaluation of [9]),  $8 \times 1024$  (same P=1 in the partitioned approaches), and  $128 \times 8 \times 1024$  (identical to the total budget of LMQ-Sketch at P=128). Increasing parallelism for thread-local designs such as SW-SKT will require P times this memory at runtime, but the accuracy bounds of queries remain the same as for a single sketch of 1/P of the total memory (since the local sketches simply get merged). Fig. 5 shows memory budget scaling for SW-SKT and LMQ-Sketch when the number of partitions/local-sketches follow the number of threads when run in parallel.

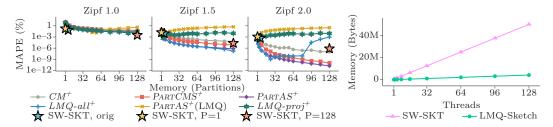


Figure 4 Non-concurrent MAPE for  $F_2$  estimations at Figure 5 Total memory scalvarious z and memory budgets, in multiples of 1 partition ing linearly with P. LMQ-Sketch occupying  $H \times K = 8 \times 1024$  counters (32 KiB). Thread-local uses 32 KiB per thread/partition. approaches (SW-SKT) must merge local sketches before query- SW-SKT, with  $6 \times 2^{13}$  counters for its ing, and achieve the accuracy of this single sketch, various sizes CMS and Fast-AGMS as in the eval-of which are shown here. For such designs, increasing memory uation of [9] (but without the HLL budget with additional threads does not improve accuracy. for  $F_0$ ) uses 393 kB per thread.

**Takeaways.**  $PARTCMS^+$  and  $PARTAS^+$ , our stepwise enhancements for  $CM^+$  in § 5, improve estimation accuracy. Unlike for thread-local designs, the partitioned approach allows LMQ-Sketch to utilize the increased memory budget for a two-fold benefit: increasing per-partition memory improves (sequential) accuracy in isolation; increasing the number of partitions benefits both accuracy and parallelism. Although Fast-AGMS is one of the most accurate  $F_2$  sketches, its thread-local-oriented design does not see improved accuracy for the same per-partition memory budget, while LMQ-Sketch achieves the same or higher accuracy by effectively navigating the memory and concurrency trade-offs.

#### 6.3 **Query Latency and Accuracy**

Next, we consider concurrent queries. As described in (G1) (§ 3), IVL, while preserving bounds for sketches, cannot fully characterize the final accuracy of the result; query latency plays a key role. LMQ-Sketch addresses this by optimizing query latency using partial results at insertion. We evaluate the latency improvement due to this algorithmic design.

**Design and Parameters.** We measure the latency of global queries with different values of P, concurrent with updates. After  $100 \,\mathrm{ms}$  of warmup to populate delegation filters, queries are performed at a rate of 1000/s. For  $F_1$  queries, we focus on the  $\hat{F}_1^{\text{NOSYNC}}$  IVL design (§ 4.2). For  $F_2$ , the more complex query, we benchmark the synchronization designs.

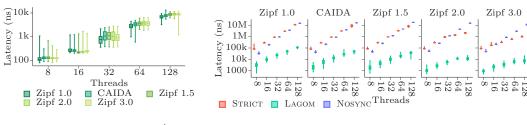


Figure 6 Latency of  $\hat{F}_1^{\text{NOSYNC}}$ 

**Figure 7** Latencies of 100  $F_2$  queries concurrent with queries concurrent with updates scales updates. LAGOM is 2-3 orders of magnitude faster than linearly with P. No impact from skew. other designs, which struggle with high thread counts.

**Takeaways.**  $\hat{F}_1^{\text{NOSYNC}}$  query duration ranges from 100 ns to 10 µs at high thread counts. Our optimizations in  $\hat{F}_2^{\text{LAGOM}}$  yield significant improvements in query latency, which stays below 100 µs. Such short latencies suggest the suitability of our queries for accurate estimation under IVL, due to reduced number of overlapping updates, evaluated next.

#### 6.4 **Concurrency and Accuracy**

Building on insights from measuring concurrent query latencies, we now evaluate the impact on accuracy. Query duration in conjunction with update throughput determine the IVL interval size. Although IVL does not exactly target accuracy guarantees, to complement, we explore the admissible freshness of this interval.

**Design and Parameters.** We compare the return value of concurrent  $\hat{F}_2^{\text{LAGOM}}$  with the IVL-permitted interval. We adopt the following methodology to determine the relative error induced by concurrent executions, without interfering with execution patterns, e.g., as an approach based on stopping and starting updater threads to record measurements would do.

Methodology to determine bounds of return interval of concurrent query Q: We reconstruct ideal query return values at  $Q^{start}$  and  $Q^{end}$ . A threshold  $F_1$ -value T is selected (beyond the warmup period), as a trigger point for performing Q which overlaps an arbitrary number of updates n. To reconstruct  $Q^{start}$ , updates are stopped when  $F_1 = T$  and a sequential query is performed. In a new execution, reaching  $F_1 = T$  instead triggers Q concurrently. Upon the completion of Q, updates are stopped and  $Q^{end}$  is recorded in a sequential setting.

We set T = 10M tuples and perform 50 repetitions for each dataset (CAIDA and synthetic with z=1, 1.5 and 2) and increasing P. We study  $\hat{F}_2^{\text{Nosync}}$  and  $\hat{F}_2^{\text{Lagom}}$  using the same methodology to determine the effect of latency and weaker semantics on accuracy.

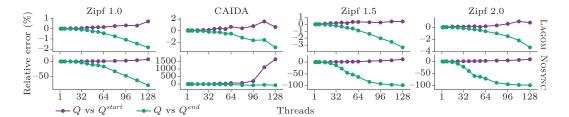


Figure 8 Comparison of query return value against bounds of IVL-permitted interval. For both designs, Q observes more updates than  $Q^{start}$  (which is what  $\hat{F}_2^{STRICT}$  would return) but misses some updates observed by  $Q^{end}$ . LAGOM consistently yields (orders of magnitude) narrower return value intervals compared to NOSYNC, which may miss many updates due to its long operation latency.

**Takeaways.** The concurrent  $F_2$  queries can observe overlapping updates (towards improving STRICT's freshness). However, the interval of return values for NOSYNC is arbitrarily large, while LAGOM (with agile calculation and IVL guarantees) imparts a seriously smaller uncertainty, particularly at low–intermediate P, demonstrating that LAGOM's properties preserve, and are useful for, accuracy.

#### 7 Other Related Work

Rinberg et al. in [40] present a snapshot-based methodology for sketch querying concurrently with updates, building on a thread-local design with limited local buffering and a global propagator which merges local sketches into a global one when full. Concurrent queries require a snapshot of this global sketch, which gets more costly with the size of the state, also leading to a tradeoff with respect to accuracy. Using a partitioning design as in LMQ-Sketch allows efficient queries taking snapshots per partition, while supporting distributed state and fine-tuning of freshness. Concurrent queries and updates are also targeted in [18] for estimating quantiles, although not using the generic framework of [40] due to risk of sequential bottlenecks.

Several recent orthogonal works such as Hydra [32] and OmniSketch [38] explore using sketches for *multiple querying* of multidimensional data streams. These works are not targeting queries concurrent with updates, though, instead using an online sketching phase followed by an offline querying phase, during which the sketch is no longer updated.

Universal Sketching appears in [31, 32] as a useful technique to estimate sums of monotonic functions applied to stream frequency vectors, such as  $F_1$  and  $F_2$ , from a single sketch in small (polylogarithmic) memory. As its core functions associate with the ones in LMQ-Sketch, the latter can be a possible candidate for supporting concurrency in it.

#### 8 Conclusions

LMQ-Sketch is a concurrent, multi-query sketch in a single, low memory-footprint object, with explicit concurrency semantics that lead to predictably high accuracy even with very demanding stream properties. The analytical and detailed empirical insights using real-world and synthetic data showed that (1) having a single data structure balances multiple targets: accuracy (Fig. 4), timeliness (Fig. 6, 7), memory footprint (Fig. 5), freshness (Fig. 8), concurrency, and throughput (Fig. 3); (2) besides the memory budget, the way that the memory is used is catalytic for the achievable concurrency and accuracy, both for local and

global queries; (3) employing IVL, in conjunction with low query latency and concurrency-aware compensation, enables accuracy aligning with that of the sequential sketches, even under very high-rate, skewed streams.

LMQ-Sketch can be a useful component in systems such as Redis, Apache Druid/Spark/DataSketches and more. The method facilitates continued research towards supporting additional queries concurrently, which all need to deal with the same challenges regarding efficiency, consistency, and accuracy; e.g. top-k elements, possibly supporting wavelets, quantiles, and associating with queries required for universal sketching.

#### - References -

- 1 Noga Alon, Phillip B. Gibbons, Yossi Matias, and Mario Szegedy. Tracking join and self-join sizes in limited storage. In *Proceedings of the Eighteenth ACM SIGMOD-SIGACT-SIGART Symposium on Principles of Database Systems*, PODS '99, pages 10–20, New York, NY, USA, May 1999. Association for Computing Machinery. doi:10.1145/303976.303978.
- 2 DataSketches | Apache® DataSketches. URL: https://datasketches.apache.org/.
- 3 Apache Druid | Apache® Druid. URL: https://druid.apache.org/.
- Vladimir Braverman and Rafail Ostrovsky. Zero-one frequency laws. In Proceedings of the Forty-Second ACM Symposium on Theory of Computing, STOC '10, pages 281–290, New York, NY, USA, June 2010. Association for Computing Machinery. doi:10.1145/1806689.1806729.
- 5 Chiranjeeb Buragohain and Subhash Suri. Quantiles on Streams. In *Encyclopedia of Database Systems*, pages 2235–2240. Springer US, Boston, MA, 2009. doi:10.1007/978-0-387-39940-9\_290.
- 6 The CAIDA UCSD Anonymized Internet Traces 2018. URL: https://www.caida.org/catalog/datasets/passive\_dataset.
- 7 Lidia Ceriani and Paolo Verme. The origins of the Gini index: Extracts from Variabilità e Mutabilità (1912) by Corrado Gini. *The Journal of Economic Inequality*, 10(3):421–443, September 2012. doi:10.1007/s10888-011-9188-x.
- 8 Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding Frequent Items in Data Streams. In *Automata, Languages and Programming*, Lecture Notes in Computer Science, pages 693–703, Berlin, Heidelberg, 2002. Springer. doi:10.1007/3-540-45465-9\_59.
- 9 Monica Chiosa, Thomas B. Preußer, and Gustavo Alonso. SKT: A one-pass multi-sketch data analytics accelerator. Proceedings of the VLDB Endowment, 14(11):2369–2382, July 2021. doi:10.14778/3476249.3476287.
- Graham Cormode and Minos Garofalakis. Sketching streams through the net: Distributed approximate query tracking. In *Proceedings of the 31st International Conference on Very Large Data Bases*, VLDB '05, pages 13-24, Trondheim, Norway, August 2005. VLDB Endowment. URL: http://www.vldb.org/archives/website/2005/program/paper/tue/p13-cormode.pdf.
- Graham Cormode and Minos Garofalakis. Join sizes, frequency moments, and applications. In *Data Stream Management: Processing High-Speed Data Streams*, pages 87–102. Springer, 2016. doi:10.1007/978-3-540-28608-0\_4.
- Graham Cormode, Minos Garofalakis, and Dimitris Sacharidis. Fast Approximate Wavelet Tracking on Streams. In *Advances in Database Technology EDBT 2006*, pages 4–22, Berlin, Heidelberg, 2006. Springer. doi:10.1007/11687238\_4.
- Graham Cormode and S. Muthukrishnan. An improved data stream summary: The countmin sketch and its applications. *Journal of Algorithms*, 55(1):58–75, April 2005. doi: 10.1016/j.jalgor.2003.12.001.
- Graham Cormode and S. Muthukrishnan. Summarizing and Mining Skewed Data Streams. In Proceedings of the 2005 SIAM International Conference on Data Mining, pages 44–55. Society for Industrial and Applied Mathematics, April 2005. doi:10.1137/1.9781611972757.5.

- Graham Cormode and Ke Yi. Small Summaries for Big Data. Cambridge University Press, Cambridge, 2020. doi:10.1017/9781108769938.
- Alin Dobra, Minos Garofalakis, Johannes Gehrke, and Rajeev Rastogi. Sketch-Based Multi-Query Processing over Data Streams. In *Data Stream Management: Processing High-Speed Data Streams*, Data-Centric Systems and Applications, pages 241–261. Springer, Berlin, Heidelberg, 2016. doi:10.1007/978-3-540-28608-0\_12.
- 17 Cynthia Dwork, Maurice Herlihy, Serge Plotkin, and Orli Waarts. Time-Lapse Snapshots. SIAM Journal on Computing, 28(5):1848–1874, January 1999. doi:10.1137/S0097539793243685.
- 18 Shaked Elias Zada, Arik Rinberg, and Idit Keidar. Quancurrent: A Concurrent Quantiles Sketch. In *Proceedings of the 35th ACM Symposium on Parallelism in Algorithms and Architectures*, SPAA '23, pages 15–25, New York, NY, USA, June 2023. Association for Computing Machinery. doi:10.1145/3558481.3591074.
- Minos Garofalakis. Discrete Wavelet Transform and Wavelet Synopses. In Encyclopedia of Database Systems, pages 857–863. Springer US, Boston, MA, 2009. doi:10.1007/ 978-0-387-39940-9\_539.
- 20 Anna C. Gilbert, Yannis Kotidis, S. Muthukrishnan, and Martin J. Strauss. How to summarize the universe: Dynamic maintenance of quantiles. In *Proceedings of the 28th International Conference on Very Large Data Bases*, VLDB '02, pages 454–465, Hong Kong, China, August 2002. VLDB Endowment. doi:10.1016/B978-155860869-6/50047-0.
- Anna C. Gilbert, Yannis Kotidis, S. Muthukrishnan, and Martin J. Strauss. One-pass wavelet decompositions of data streams. *IEEE Transactions on Knowledge and Data Engineering*, 15(3):541–554, May 2003. doi:10.1109/TKDE.2003.1198389.
- 22 Corrado Gini. Variabilità e Mutabilità. Reprinted in Memorie Di Metodologia Statistica (Ed. E. Pizetti and T. Salvemini.) 1955. Libreria Eredi Virgilio Veschi, Rome, 1912.
- 23 Itamar Haber. Count-Min Sketch: The Art and Science of Estimating Stuff, March 2022. URL: https://redis.io/blog/count-min-sketch-the-art-and-science-of-estimating-stuff/.
- 24 Thomas A. Henzinger, Christoph M. Kirsch, Hannes Payer, Ali Sezgin, and Ana Sokolova. Quantitative relaxation of concurrent data structures. In Proceedings of the 40th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL '13, pages 317–328, New York, NY, USA, January 2013. Association for Computing Machinery. doi:10.1145/2429069.2429109.
- 25 Maurice Herlihy, Nir Shavit, Victor Luchangco, and Michael Spear. The Art of Multiprocessor Programming. Elsevier, Morgan Kaufmann Publishers, Cambridge, MA, United States, second edition, 2021.
- 26 Martin Hilgendorf. LMQ-Sketch. URL: https://gitlab.com/marhilg/lmq-sketch.
- Victor Jarlow, Charalampos Stylianopoulos, and Marina Papatriantafilou. QPOPSS: Query and Parallelism Optimized Space-Saving for finding frequent stream elements. *Journal of Parallel and Distributed Computing*, 204:105134, October 2025. doi:10.1016/j.jpdc.2025.105134.
- 28 Lior Kogan. T-digest: A New Probabilistic Data Structure in Redis Stack, March 2023. URL: https://redis.io/blog/t-digest-in-redis-stack/.
- 29 Balachander Krishnamurthy, Subhabrata Sen, Yin Zhang, and Yan Chen. Sketch-based change detection: Methods, evaluation, and applications. In *Proceedings of the 3rd ACM SIGCOMM Conference on Internet Measurement*, IMC '03, pages 234–247, New York, NY, USA, October 2003. Association for Computing Machinery. doi:10.1145/948205.948236.
- 30 Leslie Lamport. Concurrent reading and writing. Communications of the ACM, 20(11):806–811, November 1977. doi:10.1145/359863.359878.
- Zaoxing Liu, Antonis Manousis, Gregory Vorsanger, Vyas Sekar, and Vladimir Braverman. One Sketch to Rule Them All: Rethinking Network Flow Monitoring with UnivMon. In Proceedings of the 2016 ACM SIGCOMM Conference, pages 101–114, Florianopolis Brazil, August 2016. ACM. doi:10.1145/2934872.2934906.

- 32 Antonis Manousis, Zhuo Cheng, Ran Ben Basat, Zaoxing Liu, and Vyas Sekar. Enabling efficient and general subpopulation analytics in multidimensional data streams. *Proceedings of the VLDB Endowment*, 15(11):3249–3262, July 2022. doi:10.14778/3551793.3551867.
- Gonçalo Matos, Salvatore Signorello, and Fernando M. V. Ramos. Generic change detection (almost entirely) in the dataplane. In *Proceedings of the Symposium on Architectures for Networking and Communications Systems*, ANCS '21, pages 113–120, New York, NY, USA, January 2022. Association for Computing Machinery. doi:10.1145/3493425.3502767.
- Vinh Quang Ngo and Marina Papatriantafilou. Cuckoo heavy keeper and the balancing act of maintaining heavy hitters in stream processing. *Proceedings of the VLDB Endowment*, 18(9):3149–3161, August 2025. doi:10.14778/3746405.3746434.
- Yiannis Nikolakopoulos, Anders Gidenstam, Marina Papatriantafilou, and Philippas Tsigas. A Consistency Framework for Iteration Operations in Concurrent Data Structures. In 2015 IEEE International Parallel and Distributed Processing Symposium, pages 239–248, May 2015. doi:10.1109/IPDPS.2015.84.
- 36 Savannah Norem. Probabilistic Data Structures in Redis, August 2022. URL: https://redis.io/blog/streaming-analytics-with-probabilistic-data-structures/.
- 37 Erez Petrank and Shahar Timnat. Lock-Free Data-Structure Iterators. In *Proceedings of the* 27th International Symposium on Distributed Computing Volume 8205, DISC 2013, pages 224–238, Berlin, Heidelberg, October 2013. Springer-Verlag. doi:10.1007/978-3-642-41527-2\_16.
- Wieger R. Punter, Odysseas Papapetrou, and Minos Garofalakis. OmniSketch: Efficient Multi-Dimensional High-Velocity Stream Analytics with Arbitrary Predicates. Proc. VLDB Endow., 17(3):319–331, November 2023. doi:10.14778/3632093.3632098.
- Arik Rinberg and Idit Keidar. Intermediate Value Linearizability: A Quantitative Correctness Criterion. *Journal of the ACM*, 70(2):17:1–17:21, April 2023. doi:10.1145/3584699.
- 40 Arik Rinberg, Alexander Spiegelman, Edward Bortnikov, Eshcar Hillel, Idit Keidar, Lee Rhodes, and Hadar Serviansky. Fast Concurrent Data Sketches. *ACM Transactions on Parallel Computing*, 9(2):6:1–6:35, April 2022. doi:10.1145/3512758.
- Pratanu Roy, Arijit Khan, and Gustavo Alonso. Augmented Sketch: Faster and More Accurate Stream Processing. In *Proceedings of the 2016 International Conference on Management* of Data, SIGMOD '16, pages 1449–1463, New York, NY, USA, June 2016. Association for Computing Machinery. doi:10.1145/2882903.2882948.
- 42 Charalampos Stylianopoulos, Ivan Walulya, Magnus Almgren, Olaf Landsiedel, and Marina Papatriantafilou. Delegation sketch: A parallel design with support for fast and accurate concurrent operations. In *Proceedings of the Fifteenth European Conference on Computer Systems*, EuroSys '20, pages 1–16, New York, NY, USA, April 2020. Association for Computing Machinery. doi:10.1145/3342195.3387542.
- Da Tong and Viktor K. Prasanna. Sketch Acceleration on FPGA and its Applications in Network Anomaly Detection. *IEEE Transactions on Parallel and Distributed Systems*, 29(4):929–942, April 2018. doi:10.1109/TPDS.2017.2766633.
- R. Kent Treiber. Systems Programming: Coping with Parallelism. Number 5118 in Research Report. International Business Machines Incorporated Research Division, Thomas J. Watson Research Center, 1986. URL: https://dominoweb.draco.res.ibm.com/reports/rj5118.pdf.
- 45 David P. Woodruff. Sketching as a Tool for Numerical Linear Algebra. Foundations and Trends® in Theoretical Computer Science, 10(1–2):1–157, October 2014. doi:10.1561/040000060.
- 46 David P. Woodruff. Frequency Moments. In Encyclopedia of Database Systems, pages 1518–1519. Springer, New York, NY, 2018. doi:10.1007/978-1-4614-8265-9\_167.

# A Proof arguments

Due to space limitations, arguments that support the claims in the main part of the paper are presented in more detail in this appendix.

#### A.1 From § 4

▶ Lemma 5. Delegation Sketch-based PQ is an IVL implementation of ASketch point query.

**Proof sketch.** While  $T_{\mathsf{Owner}(a)}$  executes PQ, other threads may concurrently update relevant delegation filters. However, only updates overlapping PQ can be missed; updates completed before PQ are reflected in the returned  $\hat{f}(a)$  [42, Claim 2]. Double-counting of updates is not possible [42, Claim 3], as this would mean PQ observed an update both while it is buffered in a delegation filter, as well as when it has been flushed to the partition-local ASketch. This would require a flush of the delegation filter, which can only be performed by the same thread as PQ itself, and hence cannot overlap PQ. Therefore, the maximal return value is the one for a linearization which observes all overlapping updates.

▶ **Lemma 6.**  $\hat{F}_1^{NOSYNC}$  estimates  $F_1$  with IVL semantics.

**Proof sketch.** The partial results form jointly a shared counter, where each thread maintains a local counter of completed updates. The query, performing one atomic read per partition, can neither double-count updates not omit completed ones.

▶ **Observation 7.** An  $\hat{F}_2^{Nosync}$  query Q can miss or double-count updates, due to data movement by overlapping processing of delegated updates (Alg. 3).

This can occur when Q overlaps with processing delegated updates (Alg. 3), which non-atomically moves occurrences of keys from delegation filters to the ASketch of the owning partition. Updates can be missed when Q scans a sketch  $T_i$ .Sk (line 5.4), whereupon a filter  $T_j$ .DF<sub>i</sub> is handed over to  $T_i$  and flushed (Alg. 3), followed by Q reading the now-empty filter (line 5.7), thus missing all the concurrently flushed updates, regardless of whether they were completed before Q began. Or, if Q overlaps a flush of  $T_j$ .DF<sub>i</sub>, updates may be seen both in the ASketch of  $T_i$  and in  $T_j$ .DF<sub>i</sub> before the latter is cleared.<sup>3</sup>

▶ **Lemma 8.**  $\hat{F}_2^{LAGOM}$  gets an atomic snapshot per partition, and cannot deadlock with updates.

**Proof sketch.** Atomic per-partition snapshot – When a full delegation filter is flushed, buffered updates move to the owning partition's ASketch and associated partial results are updated (Alg. 1). As this operation is not atomic, the query synchronizes with it to avoid mis-calculations, such as double-counting updates or omitting completed ones:

- A flush begins by incrementing V1 and increments V2 upon completion (line 3.4 and 3.11). A query reads these in reverse order (V2 then V1) to detect concurrent modification of the data it read, and the scan is retried is they do not match (condition on line 6.13).
- To prevent unbounded retrying in the (unlikely) event that repeated flushes occur in a query's duration, the query flags the partition it is scanning (line 6.4). If the flush finds the flag set (line 3.3), it waits until the query has cleared the flag (line 6.14).
- If the flag is not set, the flush can proceed to increment V1 (line 3.4) and begin updating the ASketch contents. If a query now sets the flag, it will either find the version numbers match and knows that a consistent view was obtained, or detect a mismatch, indicating the observed state may be inconsistent, and retry the scan (line 6.13). The query can be blocked by at most one concurrent filter flush, as subsequent ones will be stalled by the flag.

<sup>&</sup>lt;sup>3</sup> Similarly, there can there be omissions and double-counting if we swap the relative order of scanning of sketches and delegation filters.

The synchronization design performs a double-collect of the version numbers (line 6.13) which indicate concurrent modification of the ASketch of the partition being scanned. As a scan gets an image of a partition's partial results without interferences from filter flushes, it sees a linearizable view of that information in the partition.

No deadlock – Deadlock entails one or more updater threads and the  $F_2$  query thread being unable to proceed because they are waiting for one another. The query algorithm interacts with only one partition at a time; all remaining updater threads for other partitions are not involved and continue executing independently. If the query has set the flag for updater thread  $T_i$  (line 6.4),  $T_i$  will not proceed past line 3.3. Within bounded steps, the query will complete its scan of partition i, as  $T_i$  is currently blocked and cannot update the version numbers to mismatch, and ultimately reset the flag (line 6.14), thus unblocking  $T_i$ . Similarly, if the query is attempting obtain a consistent scan while  $T_j$  is performing a filter flush operation and has incremented V1 – causing the query to retry –, from our system properties,  $T_i$  is guaranteed to make progress and complete the flush within bounded time, finally incrementing V2 to match. Again, the query will be blocked by at most one concurrent filter flush, as subsequent ones will be stalled by the flag.

▶ **Lemma 10.** For queries  $Q_1 \longrightarrow Q_2$  the monotonicity of scans for each combination follows

	$Q_2$ :	PQ	$F_1$	$F_2$
	PQ	Monotonic	$P ext{-}relaxed\ monotonic$	$r ext{-}relaxed\ monotonic$
$Q_1$ :	$F_1$	Monotonic	Monotonic	$rP ext{-}relaxed\ monotonic$
	$F_2$	Monotonic	Monotonic	Monotonic

**Proof sketch.** If  $Q_1$  and  $Q_2$  are of the same type or if  $U\big|_{Q_2} \longrightarrow U\big|_{Q_1}$  according to program order, then monotonicity of scans is immediate. There are three remaining cases where it is possible that  $U\big|_{Q_1} \longrightarrow Q_1 \longrightarrow Q_2 \longrightarrow U\big|_{Q_2}$  when U overlaps both queries:

- $PQ \longrightarrow F_1$  At most one update operation per updater thread can overlap both  $Q_1$  and  $Q_2$ . Therefore, a PQ is limited to observing at most P updates that are not yet visible to a subsequent  $F_1$  query.
- $PQ \longrightarrow F_2$  Since  $Q_1$  is a local query, hence only observes one partition, at most r = PB updates buffered in delegation filters for this partition may be observed by  $Q_1$  but not by  $Q_2$  (Obs. 9).
- $F_1 \longrightarrow F_2$  Similar to the previous case; however,  $Q_1$  is now a global query and may in the worst case observe up to r = PB buffered updates per partition not yet visible to  $Q_2$ . Hence, the queries may deviate from each other by up to rP updates, but no more.

#### A.2 From § 5

▶ **Lemma 11.**  $CM_{conc}^+$  is an IVL implementation of  $CM^+$ , preserving  $CM^+$ 's  $(\epsilon, \delta)$  bounds.

**Proof sketch.** Following a similar argument as [39, Lemma 5.3]; as updates only increment counters, each counter read by a  $CM_{conc}^+$  query Q will return a value at least as large as the value at the start of Q, and no larger than the value at the end of Q. Transitively, this holds also for the sum of squared values of said counters as computed by  $CM_{conc}^+$ . The returned  $\hat{F}_2$  will be the minimum of these per-row results, and it cannot be lower than the least per-row result at the start of the query, or larger than the least per-row result at the end of the query, as required for IVL of monotonically increasing quantities (Obs. 2).

▶ Lemma 13. PARTCMS<sup>+</sup> is an IVL implementation of PARTCMS<sup>+</sup>.

**Proof sketch.** Similar reasoning as for Lemma 11.

▶ **Lemma 15.** PartAS<sup>+</sup><sub>conc</sub> is an IVL implementation of PartAS<sup>+</sup>.

**Proof sketch.** Similar to earlier lemmas. By virtue of being atomic, the snapshots for each partition will linearize with updates by the corresponding updater thread. Since counter values read by the query are non-decreasing, the return value is bounded between an ideal return value at the start of the query (observing all preceding updates) and an ideal return value at the end of the query (additionally observing all concurrent updates).

▶ **Lemma 16.** PartAS<sup>+</sup><sub>conc</sub> on LMQ-Sketch (i.e. its ASketch part) is an r-relaxed IVL implementation of LMQ-all<sup>+</sup> per partition, where r = PB.

**Proof sketch.** For a partition, up to r = PB updates buffered in delegation filters (Obs. 9) may be missing from the atomic snapshot taken by  $PARTAS_{conc}^+$ . Hence, the return value is bounded between the value at query start excluding the r buffered updates, and the value at the end, including all r buffered updates.

#### B Evaluation & Discussion

Due to space limitations, the discussion supporting the main takeaways of the empirical study is presented in more detail in this appendix.

#### B.1 Concurrency and Synchronization (§ 6.1) – Results

Fig. 3 shows the mean rate of update operations for sketching the complete input data sets. Lagom exhibits very similar scaling as Delegation Sketch with B=1000, processing approximately 1.5 billion updates/second on real-world data, demonstrating the very low overhead of its synchronization design for consistent global queries. Strict cannot scale and performance worsens with more threads as contention around the global RW lock increases. Nosync scales similarly to the plain Delegation Sketch (neither of them support consistent global queries). As expected, SW-SKT throughput grows linearly with the number of threads, as threads process updates entirely independently. However, concurrent queries cannot be supported. With increasing skewness, a clear upwards trend in throughput of the delegation design is seen, as delegation filters are able to buffer more updates locally, reducing inter-thread communication.

On the lower row, Fig. 3 shows the mean rate of update operations with concurrent queries for designs which support them. STRICT shows similar scaling to before. LAGOM achieves very similar performance in presence of concurrent queries. Decreases in update throughput are explained by the fact that updater threads are responsible for serving point queries alongside their update workload. NOSYNC clearly shows this effect as global queries have no impact in this design, but the overhead of point query work leads to a reduction in peak throughput not seen in LAGOM by around 50 % compared to Fig. 3 at 0.1 % point queries (exactly reproducing the result in [42]).

#### B.2 Memory and Accuracy (§ 6.2) – Results

Fig. 4 shows the accuracy of various  $F_2$  methods, without concurrent updates during query execution. As discussed in § 5,  $PARTCMS^+$  and  $PARTAS^+$ , our enhanced versions of  $CM^+$  for more efficient data structures, improve estimation accuracy. Further, for partitioning-based approaches, increasing the memory budget by increasing the number of partitions leads to improved accuracy, as more memory is available for (1) avoiding hash collisions (Obs. 12) and

(2) accurate tracking of heavy keys in ASketch filters (Obs. 14). Although Fast-AGMS, used in SW-SKT, is one of the most accurate techniques for  $F_2$  estimation for arbitrary skewness, a thread-local design does not see improved accuracy for the same memory budget, while LMQ-Sketch's partition-targeting methods, navigating memory and concurrency trade-off challenges, achieve higher accuracy. Of course, to be fair, one should observe that SW-SKT was not designed with concurrent queries as a target.

#### B.3 Latency and Accuracy (§ 6.3) – Results

Fig. 6 shows  $\hat{F}_1^{\text{NOSYNC}}$  query latency to be small – implying a low number of overlapping updates – growing linearly with the number of threads and partitions and not impacted by skew.

Similarly, Fig. 7 shows the distribution of  $F_2$  query durations for various synchronization designs. The operational complexity of  $\hat{F}_2^{\text{NOSYNC}}$  scales with the number of delegation filters and ASketch filters, as for each slot in ASketch filters, all delegation filters of that partition are read ( $\mathcal{H}^{\text{NOSYNC}}$ , Eq. 3).  $\hat{F}_2^{\text{LAGOM}}$  is significantly faster, both absolutely and comparatively, at less than 100 µs, 2 to 3 orders of magnitude faster than  $\hat{F}_2^{\text{STRICT}}$  and  $\hat{F}_2^{\text{NOSYNC}}$ , regardless of skew, demonstrating the impact of our projection-based latency optimization.

#### B.4 Concurrency and Accuracy (§ 6.4) – Results

Fig. 8 shows how query return values relate to IVL interval boundaries for  $\hat{F}_2^{\text{LAGOM}}$  and  $\hat{F}_2^{\text{NOSYNC}}$ . In all cases, the concurrent query Q observes more updates than  $Q^{start}$  – which is what  $\hat{F}_2^{\text{STRICT}}$  would return – but ignores some updates that  $Q^{end}$  observes, as expected. Thus, the truly concurrent  $\hat{F}_2$  query results are more fresh than STRICT synchronization. However, the size of the interval for Nosync is large compared to LAGOM, due to its significantly longer query latency (seen in Fig. 7) and the weakness of its semantics (described in Obs. 7), while LAGOM imparts a much smaller IVL-interval on the return value of Q, demonstrating that our lightweight synchronization preserves accuracy of results close to the sequential expectation.

#### **B.5** Overall Takeaways

While IVL, as a useful correctness criterion, allows reasoning about semantics of concurrent queries and can preserve  $(\epsilon, \delta)$  bounds of sketches, it cannot alone fully characterize the accuracy of results. Our study shows that  $\hat{F}_2^{\text{LAGOM}}$ 's compensation scheme is very close to accurate estimations in a sequential setting, utilizing available memory for improved parallelism and accuracy in a concurrent one, where efficient synchronization permits low query latency, implying freshness of returned results. These observations illustrate the challenges tackled in this work.