



HAL
open science

Load-Aware Shedding in Stream Processing Systems

Nicolò Rivetti, Yann Busnel, Leonardo Querzoni

► **To cite this version:**

Nicolò Rivetti, Yann Busnel, Leonardo Querzoni. Load-Aware Shedding in Stream Processing Systems. Transactions on Large-Scale Data- and Knowledge-Centered Systems, 2020, pp.121-153. 10.1007/978-3-662-62386-2_5 . hal-03115253

HAL Id: hal-03115253

<https://imt-atlantique.hal.science/hal-03115253v1>

Submitted on 19 Jan 2021

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Load-Aware Shedding in Stream Processing Systems ^{*} ^{**}

Nicoló Rivetti¹, Yann Busnel², and Leonardo Querzoni^{3*} ^{**} ^[0000-0002-8711-4216]

¹ Independent researcher, nicolo.rivetti@gmail.com

² IMT Atlantique, IRISA (Rennes, France), yann.busnel@imt-atlantique.fr

³ DIAG, Sapienza University of Rome (Italy), querzoni@diag.uniroma1.it

Abstract. Distributed stream processing systems are today gaining momentum as a tool to perform analytics on continuous data streams. Load shedding is a technique used to handle unpredictable spikes in the input load whenever available computing resources are not adequately provisioned. In this paper, we propose Load-Aware Shedding (LAS), a novel load shedding solution that, unlike previous works, does not rely neither on a pre-defined cost model nor on any assumption on the tuple execution duration. Leveraging *sketches*, LAS efficiently estimates the execution duration of each tuple with small error bounds and uses this knowledge to proactively shed input streams at any operator to limiting queuing latencies while dropping as few tuples as possible. We provide a theoretical analysis proving that LAS is an (ϵ, δ) -approximation of the optimal online load shedder. Furthermore, through an extensive practical evaluation based on simulations and a prototype, we evaluate its impact on stream processing applications.

Keywords: Load-Shedding; Stream Processing; Data Streaming; Distributed systems

1 Introduction

Distributed stream processing systems (DSPS) and Complex Event Processing (CEP) are today considered as a mainstream technology to build architectures for the real-time analysis of big data. An application running in a DSPS, or a query executed by a CEP engine, is typically modeled as a directed acyclic graph (a topology) where data operators, represented by nodes, are interconnected by streams of tuples containing data to be analyzed, the directed edges. The success of such systems can be traced back to their ability to run complex applications at scale on clusters of commodity hardware or in the cloud.

^{*} This work has been partially funded by the MIUR SCN-00064 project RoMA and by Sapienza University of Rome through the project RM11916B75A3293D.

^{**} A preliminary short version of this work appeared in the *Proceedings of the 10th ACM International Conference on Distributed and Event-based Systems*.

^{***} *Corresponding author.*

Correctly provisioning computing resources for DSPS or CEP engines however is far from being a trivial task. System designers need to take into account several factors: the computational complexity of the operators, the overhead induced by the framework, and the characteristics of the input streams. This latter aspect is often the most critical, as input data streams may unpredictably change over time both in rate and in content. Over-provisioning is not economically sensible, thus system designers are today moving toward approaches based on elastic scalability [11], where an underlying infrastructure can tune at runtime the available resources in response to changes in the workload. This represents a desirable solution when coupled with on-demand provisioning offered by many cloud platforms, but still may be affected by transient overloads [3], caused for example by unexpected load spikes, that could temporarily degrade performance below the desired SLA.

Bursty input load represents a problem for both DSPS and CEP engines as it may create unpredictable bottlenecks within the system that lead to an increase in queuing latencies, pushing the system in a state where it cannot deliver the expected quality of service (typically expressed in terms of tuple completion latency). *Load shedding* is generally considered a practical approach to handle bursty traffic. It consists of dropping a subset of incoming tuples as soon as a bottleneck is detected in the system. As such, load shedding is a solution that is complementary [24] and must coexist with resource shaping techniques (like elastic scaling), rather than being an alternative.

Existing load shedding solutions either randomly drop tuples when bottlenecks are detected [1] or apply a pre-defined model of the application and its input that allows them to deterministically take the best shedding decision [25]. In any case, all the existing solutions assume that incoming tuples all impose the same computational load. However, such assumption does not hold for many practical use cases; tuple execution duration, in fact, may depend on the tuple content itself. This is often the case whenever the receiving operator implements a logic with branches where only a subset of the incoming tuples travels through every single branch. If the computation associated with each branch generates different loads, then the execution duration will change from tuple to tuple. A tuple with a large execution duration may delay the execution of subsequent tuples in the same stream, thus increasing queuing latencies. If further tuples are enqueued with large execution durations, this may bring to the emergence of a bottleneck.

As an example, consider the reach of a tweet, *i.e.*, the number of users that may receive the re-tweets of a given tweet. This computation entails counting the number of users that have a direct and un-direct follower relationship (until a given depth) with the tweet author. Then, depending on the size of the sub-graph rooted in the author node, the execution times vary. For instance, in our experiments the execution time belongs to the interval [0.01, 70] ms, the most frequent execution time was 65 ms, while the average per execution time was 20 ms. In the experimental evaluation, we provide a second use-case exhibiting the same phenomena.

Based on this simple observation, we introduce Load-Aware Shedding (LAS), a novel solution for load shedding in DSPS (or CEP engines) engines. LAS gets rid of the aforementioned assumptions and provides efficient shedding aimed at matching given queuing latency targets while dropping as few tuples as possible. To reach this goal LAS leverages a smart combination of *sketch* data structures to efficiently collect at runtime information on the time needed to compute tuples. This information is used to build and maintain, at runtime, a cost model that is then exploited to take decisions on when input tuples must be shed. LAS has been designed as a flexible solution that can be applied on a per-operator basis, thus allowing developers to target specific critical stream paths in their applications. The proposed solution provides predictable per operator queuing latencies, an extremely important feature in several application scenarios where the stream processing system is expected to deliver results to users in a quasi-real-time fashion. Furthermore, LAS implements an efficient load shedding solution that perfectly fits the characteristics of settings where scarce resources are available (e.g. fog-computing). Finally, LAS can be complemented by an output quality model that allows blocking it from dropping tuples that may significantly degrade the final output quality.

The contributions provided by this paper are:

- the introduction of LAS, the first solution for load shedding in DSPS (or CEP engines) that proactively drops tuples to avoid bottlenecks without requiring a predefined cost model and without any assumption on the distribution of tuples;
- a theoretical analysis of LAS that points out how it is an (ϵ, δ) -approximation of the optimal online shedding algorithm;
- an experimental evaluation that illustrates how LAS can provide predictable queuing latencies that approximate a given threshold while dropping a small fraction of the incoming tuples.

Below, the next section states the system model we consider. Afterward, Section 3 details LAS whose behavior is then theoretically analyzed in Section 4. Section 5 reports on our experimental evaluation and Section 6 analyzes the related works. Finally, Section 7 concludes the paper.

2 System Model and Problem Definition

We consider a distributed stream processing system (DSPS) or Complex Event Processing (CEP) engine deployed on a cluster where several computing nodes exchange data through messages sent over a network. The stream processing application (or query) executed by the DSPS (or CEP engine) can be represented by a *topology*: a directed acyclic graph interconnecting operators, represented by vertices, with data streams (DS), represented by edges. Each topology contains at least a *source*, *i.e.*, an operator connected only through outbound DSs, and a *sink*, *i.e.*, an operator connected only to inbound DSs.

Symbol	Description
t	Tuple
σ	Stream of tuples
$[n]$	Universe of possible tuples
f_t	Number of occurrences of t in σ
$w(t)$	Execution duration of tuple t on operator O
$q(i)$	Queuing latency of the i -th tuple of the stream
$\mathcal{D}(j)$	Set of dropped tuples
$d(j)$	Number of dropped tuples
$\bar{Q}(j)$	Average queuing latency
τ	Average queuing latency threshold
\hat{C}	Estimation of the total operator execution duration
\mathcal{F}	Count Min sketch that tracks tuple frequencies
\mathcal{W}	Count Min sketch that tracks tuple cumulated execution durations
N	Window size parameter
\mathcal{S}	Snapshot
η	Relative error between consecutive snapshots
μ	Error threshold

Table 1. Symbols used in the text.

Data injected by the source is encapsulated in units called tuples (or events) and each data stream is an unbounded sequence of tuples. Without loss of generality, here we assume that each tuple t is a finite set of key/value pairs that can be customized to represent complex data structures. To simplify the discussion, in the rest of this work, we deal with streams of unary tuples each representing a single non-negative integer value.

For the sake of clarity, and without loss of generality, here we restrict our model to a topology with an operator LS (*load shedder*) that decides which tuples of its outbound DS σ consumed by a downstream operator O shall be dropped. The actual positioning of LS within a real topology may be tuned, depending on where bottlenecks are expected to appear within the topology itself. Nevertheless, we assume that LS is never deployed as a source or sink in any topology. Tuples in σ are drawn from a large universe $[n] = \{1, \dots, n\}$ and are ordered, *i.e.*, $\sigma = \langle t_1, \dots, t_m \rangle$. Therefore $[m] = 1, \dots, m$ is the index sequence associated with the m tuples contained in the stream σ . Both m and n are unknown. We denote with f_t the unknown frequency of tuple t , *i.e.*, the number of occurrences⁴ of t in σ .

We assume that the execution duration of tuple t on operator O , denoted as $w(t)$, depends on the content of the tuple t . We simplify the model assuming that w depends on a single, fixed and known attribute value of tuple t . Cases in which this assumption does not hold, e.g. w depends on multiple attributes can be simply treated by concatenating their values and considering them as a sin-

⁴ In the data streaming literature, the frequency is the number of occurrences *not* divided by time, which differs from the classical (physics) definition [17].

gle multiplexed attribute [7,5,15]. The probability distribution of such attribute values, as well as the function w are unknown, may differ from operator to operator and may change over time. However, we assume that subsequent changes are interleaved by a large enough time frame such that an algorithm may have a reasonable amount of time to adapt. On the other hand, the input throughput of the stream may vary, even with a large magnitude, at any time.

Let $q(i)$ be the queuing latency of the i -th tuple of the stream, *i.e.*, the time spent by the i -th tuple in the inbound buffer of operator O before being processed. Let us denote as $\mathcal{D}(j) \subseteq [j], j \leq m$, the set of dropped tuples in a stream of length m , *i.e.*, dropped tuples are thus represented in $\mathcal{D}(j)$ by their indices in $[j] \subseteq [m]$. Moreover, let $d(j) \leq j \leq m$ be the number of dropped tuples in a stream prefix of length j , *i.e.*, $d(j) = |\mathcal{D}(j)|$. Then we can define the average queuing latency as: $\bar{Q}(j) = \sum_{i \in [j] \setminus \mathcal{D}(j)} q(i) / (j - d(j))$ for all $j \in [m]$.

The goal of the load shedder is to maintain at any point in the stream the average queuing latency smaller than a given threshold τ by dropping as few tuples as possible. The quality of the shedder can be evaluated both by comparing the resulting $\bar{Q}(j)$ against τ and by measuring the number of dropped tuples $d(j)$. More formally, the load shedding problem can be defined as follows⁵.

Problem 1 (Load Shedding). Given a data stream $\sigma = \langle t_1, \dots, t_m \rangle$, find the smallest set $\mathcal{D}(j)$ such that

$$\forall j \in [m] \setminus \mathcal{D}(j), \quad \bar{Q}(j) \leq \tau.$$

3 Load Aware Shedding

This section introduces the Load-Aware Shedding algorithm by first providing an overview, then detailing some background knowledge, and finally describing the details of its functioning.

3.1 Overview

Load-Aware Shedding (LAS) is based on a simple, yet effective, idea: if we assume to know the execution duration $w(t)$ of each tuple t on the operator, then we can foresee the queuing time for each tuple of the operator input stream and then drop all tuples that will cause the queuing latency threshold τ to be violated. However, the value of $w(t)$ is generally unknown. A possible solution to this problem is to build a static cost model for tuple execution duration and then use it to proactively shed load. However, building an accurate cost model usually requires a large amount of *a priori* knowledge on the system. Furthermore, once a model has been built, it can be hard to handle changes in the system or input stream characteristics at runtime.

⁵ This is not the only possible definition of the load shedding problem. Other variants are briefly discussed in section 6.

LAS overcomes these issues by building and maintaining at run-time a cost model for tuple execution durations. It takes shedding decision based on the estimation $\widehat{\mathcal{C}}$ of the total execution duration of the operator: $\mathcal{C} = \sum_{i \in [m] \setminus \mathcal{D}(m)} w(t_i)$. To do so, LAS computes an estimation $\widehat{w}(t)$ of the execution duration $w(t)$ of each tuple t . Then, it computes the sum of the estimated execution durations of the tuples assigned to the operator, *i.e.*, $\widehat{\mathcal{C}} = \sum_{i \in [m] \setminus \mathcal{D}(m)} \widehat{w}(t)$. At the arrival of the i -th tuple, subtracting from $\widehat{\mathcal{C}}$ the (physical) time elapsed from the emission of the first tuple provides us with an estimation $\widehat{q}(i)$ of the queuing latency $q(i)$ for the current tuple.

To enable this approach, LAS builds a sketch on the operator (*i.e.*, a memory efficient data structure) that will track the execution duration of the tuples it processes. Using a sketch allows LAS to efficiently track this data independently from the amount of possibly different tuples handled by the operator. When a change in the stream or operator characteristics affects the tuples execution durations $w(t)$, *i.e.*, the sketch content changes, the operator will forward an updated version to the load shedder, which will then be able to (again) correctly estimate the tuples execution durations. This solution does not require any *a priori* knowledge on the stream or system and is designed to continuously adapt to changes in the input stream or on the operator characteristics.

Shedding tuples from an incoming stream has in general a negative impact on the stream processing output quality. LAS approach is focussed on discarding tuples whose contribution to operator overload is larger, independently from their content. This approach is meaningful only under the assumption that the contribution to the stream output is the same for each input tuple. Applications, where this assumption does not hold, can be managed in LAS by building up a model for output degradation caused by shedding and then using this model to check for any candidate tuple if its contribution to the output quality is compatible with a given constraint.

3.2 Background

2-Universal Hash Functions — Our algorithm uses hash functions randomly picked from a 2-universal hash functions family. A collection \mathcal{H} of hash functions $h : \{1, \dots, n\} \rightarrow \{0, \dots, c\}$ is said to be 2-universal if for every two different items $x, y \in [n]$, for any $h \in \mathcal{H}$, $\mathbb{P}\{h(x) = h(y)\} \leq \frac{1}{c}$, which is the probability of collision obtained if the hash function assigned truly random values to any $x \in [n]$. Carter and Wegman [4] provide an efficient method to build large families of hash functions approximating the 2-universality property.

Count Min sketch algorithm — Cormode and Muthukrishnan have introduced in [6] the **Count Min** sketch that provides, for each item t in the input stream an (ε, δ) -additive-approximation \widehat{f}_t of the frequency f_t . The **Count Min** sketch consists of a two-dimensional matrix \mathcal{F} of size $r \times c$, where $r = \lceil \log \frac{1}{\delta} \rceil$ and $c = \lceil \frac{n}{\varepsilon} \rceil$. Each row is associated with a different 2-universal hash function $h_i : [n] \rightarrow [c]$. When the **Count Min** algorithm reads sample t from the input stream,

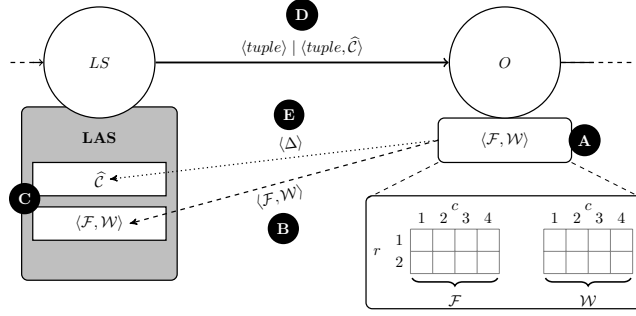


Fig. 1. Load-Aware Shedding design with $r = 2$ ($\delta = 0.25$), $c = 4$ ($\varepsilon = 0.70$).

it updates each row: $\forall i \in [r], \mathcal{F}[i, h_i(t)] \leftarrow \mathcal{F}[i, h_i(t)] + 1$. Thus, the cell value is the sum of the frequencies of all the items mapped to that cell. Upon request of f_t estimation, the algorithm returns the smallest cell value among the cells associated with t : $\hat{f}_t = \min_{i \in [r]} \{\mathcal{F}[i, h_i(t)]\}$.

Fed with a stream of m items, the space complexity of this algorithm is $O(\frac{1}{\varepsilon} \log \frac{1}{\delta} (\log m + \log n))$ bits, while update and query time complexities are $O(\log 1/\delta)$. The **Count Min** algorithm guarantees that the following bound holds on the estimation accuracy for each item read from the input stream: $\mathbb{P}\{|\hat{f}_t - f_t| \geq \varepsilon(m - f_t)\} \leq \delta$, while $f_t \leq \hat{f}_t$ is always true.

This algorithm can be easily generalized to provide (ε, δ) -additive-approximation of point queries on a stream of updates, *i.e.*, a stream where each item t carries a positive integer update value v_t . When the **Count Min** algorithm reads the pair $\langle t, v \rangle$ from the input stream, the update routine changes as follows: $\forall i \in [r], \mathcal{F}[i, h_i(t)] \leftarrow \mathcal{F}[i, h_i(t)] + v$.

3.3 LAS design

The operator stores two **Count Min** sketch matrices (Figure 1.A): the first one, denoted as \mathcal{F} , tracks the tuple frequencies f_t ; the second one, denoted as \mathcal{W} , tracks the tuple cumulated execution durations $W_t = w(t) \times f_t$. Both **Count Min** matrices share the same sizes, controlled by parameters ε and δ , and hash functions. The latter is the generalized version of the **Count Min** (Section 3.2) where the update value is the tuple execution duration when processed by the instance (*i.e.*, $v = w(t)$). The operator updates (Listing 3.1 lines 24-27) both matrices after each tuple execution.

The operator is modeled as a finite state machine (Figure 2) with two states: **START** and **STABILIZING**. The **START** state lasts as long as the operator has executed N tuples, where N is a user defined window size parameter. The transition to the **STABILIZING** state (Figure 2.A) triggers the creation of a new snapshot \mathcal{S} . A snapshot is a matrix of size $r \times c$ where $\forall i \in [r], j \in [c] : \mathcal{S}[i, j] = \mathcal{W}[i, j] / \mathcal{F}[i, j]$ (Listing 3.1 lines 15-16). We say that the \mathcal{F} and \mathcal{W}

Listing 3.1: Operator

```

1: init do
2:    $\mathcal{F} \leftarrow 0_{r,c}$  ▷ zero matrices of size  $r \times c$ 
3:    $\mathcal{W} \leftarrow 0_{r,c}$ 
4:    $\mathcal{S} \leftarrow 0_{r,c}$ 
5:    $r$  hash functions  $h_1, \dots, h_r : [n] \rightarrow [c]$  from a 2-universal family.
6:    $m \leftarrow 0$ 
7:    $state \leftarrow \text{START}$ 
8: end init
9: function UPDATE(tuple:  $t$ , execut. time:  $l$ , request:  $\widehat{\mathcal{C}}$ )
10:   $m \leftarrow m + 1$ 
11:  if  $\widehat{\mathcal{C}}$  not null then
12:     $\Delta \leftarrow \mathcal{C} - \widehat{\mathcal{C}}$ 
13:    send  $\langle \Delta \rangle$  to  $LS$ 
14:  if  $state = \text{START} \wedge m \bmod N = 0$  then ▷ Figure 2.A
15:    update  $\mathcal{S}$ 
16:     $state \leftarrow \text{STABILIZING}$ 
17:  else if  $state = \text{STABILIZING} \wedge m \bmod N = 0$  then ▷ Figure 2.C
18:    if  $\eta \leq \mu$  (Eq. 1) then
19:      send  $\langle \mathcal{F}, \mathcal{W} \rangle$  to  $LS$ 
20:       $state \leftarrow \text{START}$ 
21:      reset  $\mathcal{F}$  and  $\mathcal{W}$  to  $0_{r,c}$ 
22:    else ▷ Figure 2.B
23:      update  $\mathcal{S}$ 
24:    for  $i = 1$  to  $r$  do
25:       $\mathcal{F}[i, h_i(t)] \leftarrow \mathcal{F}[i, h_i(t)] + 1$ 
26:       $\mathcal{W}[i, h_i(t)] \leftarrow \mathcal{W}[i, h_i(t)] + l$ 
27:    end for
28:  end function
29: end function

```

matrices are stable when the relative error η between the previous snapshot and the current one is smaller than a parameter μ , *i.e.*,

$$\eta = \frac{\sum_{\forall i,j} |\mathcal{S}[i,j] - \frac{\mathcal{W}[i,j]}{\mathcal{F}[i,j]}|}{\sum_{\forall i,j} \mathcal{S}[i,j]} \leq \mu \quad (1)$$

is satisfied. Then, each time the operator has executed N tuples (Listing 3.1 lines 17-23), it checks whether Equation 1 is satisfied. **(i)** In the negative case \mathcal{S} is updated (Figure 2.B). **(ii)** In the positive case, the operator sends the \mathcal{F} and \mathcal{W} matrices to the load shedder (Figure 1.B), resets their content, and moves back to the START state (Figure 2.C). This approach allows to limit the amount of data sent from the operator to LS , and resembles what was proposed in [12].

There is a delay between any change in $w(t)$ and when LS receives the updated \mathcal{F} and \mathcal{W} matrices. This introduces a skew in the cumulated execution duration estimated by LS . To compensate this skew, we introduce a synchronization mechanism that kicks in whenever the LS receives a new pair of matrices from the operator.

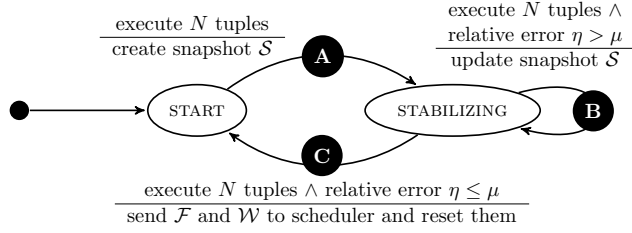


Fig. 2. Operator finite state machine.

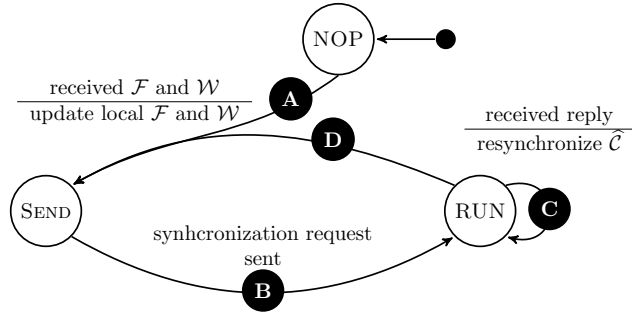


Fig. 3. Load shedder LS finite state machine.

The LS (Figure 1.C) maintains the estimated cumulated execution duration of the operator \hat{C} and a pair of initially empty matrices $\langle \mathcal{F}, \mathcal{W} \rangle$. LS is modeled as a finite state machine (Figure 3) with three states: NOP, SEND, and RUN. The LS executes the code reported in Listing 3.2. In particular, every time a new tuple t arrives at the LS , the function `SHED` is executed. The LS starts in the NOP state where no action is performed (Listing 3.2 lines 15-17). Here we assume that in this initial phase, *i.e.*, when the topology has just been deployed, no load shedding is required. When LS receives the first pair $\langle \mathcal{F}, \mathcal{W} \rangle$ of matrices (Figure 3.A), it moves into the SEND state and updates its local pair of matrices (Listing 3.2 lines 7-9). While being in the SEND states, LS sends to O the current cumulated execution duration estimation \hat{C} (Figure 1.D) piggybacking it with the first tuple t that is not dropped (Listing 3.2 lines 22-24) and moves in the RUN state (Figure 3.B). This information is used to synchronize the LS with O and remove the skew between O 's cumulated execution duration C and the estimation \hat{C} at LS . O replies to this request (Figure 1.E) with the difference $\Delta = C - \hat{C}$ (Listing 3.1 lines 11-13). When the load shedder receives the synchronization reply (Figure 3.C) it updates its estimation $\hat{C} + \Delta$ (Listing 3.2 lines 11-13).

In the RUN state, the load shedder computes, for each tuple t , the estimated queuing latency $\hat{q}(i)$ as the difference between the operator estimated execution duration \hat{C} and the time elapsed from the emission of the first tuple

Listing 3.2: Load shedder

```

1: init do
2:    $\hat{C} \leftarrow 0$ 
3:    $\langle \mathcal{F}, \mathcal{W} \rangle \leftarrow \langle 0_{r,c}, 0_{r,c} \rangle$  ▷ zero matrices pair of size  $r \times c$ 
4:   Same hash functions  $h_1 \dots h_r$  of the operator
5:    $state \leftarrow \text{NOP}$ 
6: end init
7: upon  $\langle \mathcal{F}', \mathcal{W}' \rangle$  do ▷ Figure 3.A and 3.D
8:    $state \leftarrow \text{SEND}$ 
9:    $\langle \mathcal{F}, \mathcal{W} \rangle \leftarrow \langle \mathcal{F}', \mathcal{W}' \rangle$ 
10: end upon
11: upon  $\langle \Delta \rangle$  do ▷ Figure 3.C
12:    $\hat{C} \leftarrow \hat{C} + \Delta$ 
13: end upon
14: function SHED(tuple:  $t$ )
15:   if  $state = \text{NOP}$  then
16:     return false
17:    $\hat{q} \leftarrow \hat{C}$  - elapsed time from first tuple
18:   if CHECKLATENCY( $\hat{q}$ )  $\wedge$  CHECKUTILITY( $t$ ) then
19:     return true
20:    $i \leftarrow \arg \min_{i \in [r]} \{\mathcal{F}[i, h_i(t)]\}$ 
21:    $\hat{C} \leftarrow \hat{C} + (\mathcal{W}[i, h_i(t)] / \mathcal{F}[i, h_i(t)]) \times (1 + \varepsilon)$ 
22:   if  $state = \text{SEND}$  then ▷ Figure 3.B
23:     piggyback  $\hat{C}$  to operator on  $t$ 
24:      $state \leftarrow \text{RUN}$ 
25:   return false
26: end function
27: function CHECKLATENCY( $q$ )
28:   if  $(Q + q) / \ell > \tau$  then
29:     return true
30:    $Q \leftarrow Q + q$ 
31:    $\ell \leftarrow \ell + 1$ 
32:   return false
33: end function

```

(Listing 3.2 line 17). It then checks if the estimated queuing latency for t satisfies the CHECKLATENCY method (Listing 3.2 line 18).

This method encapsulates the logic for checking if a desired condition on queuing latencies is violated or not. In this paper, as stated in Section 2, we aim at maintaining the average queuing latency below a threshold τ . Then, CHECKLATENCY tries to add \hat{q} to the current average queuing latency (Listing 3.2 lines 28). If the result is larger than τ (i), it simply returns *true*; otherwise (ii), it updates its local value for the average queuing latency and returns *false* (Listing 3.2 lines 30-32). Note that different goals, based on the queuing latency, can be defined and encapsulated within CHECKLATENCY, *e.g.*, maintain the absolute per-tuple queuing latency below τ , or maintain the average queuing latency calculated on a sliding window below τ [21].

Function CHECKUTILITY evaluates the impact the output quality would incur by dropping t . This function encapsulates optional requirements on the maximum acceptable quality drop as defined by the semantics of the application. Considering that the quality definition is application dependent, we don't provide here a specific implementation. However, we assume that, independently of the implementation, it will return *true* if the t can be dropped with an acceptable quality loss.

If both CHECKLATENCY(\hat{q}) and CHECKUTILITY(t) return *true* (i) the load shedder returns *true* as well, *i.e.*, tuple t must be dropped. Otherwise (ii), the operator estimated execution duration \hat{C} is updated with the estimated tuple execution duration $\hat{w}(t)$, increased by a factor $1 + \varepsilon$ to mitigate potential underestimations⁶, and the load shedder returns *false* (Listing 3.2 line 25), *i.e.*, the tuple must not be dropped. Finally, if the load shedder receives a new pair $\langle \mathcal{F}, \mathcal{W} \rangle$ of matrices (Figure 3.D), it will update its local pair of matrices and move to the SEND state (Listing 3.2 lines 7-9).

Now we will discuss the complexity of LAS.⁷

Theorem 1 (Time complexity of LAS).

For each tuple read from the input stream, the time complexity of LAS for the operator and the load shedder is $\mathcal{O}(\log 1/\delta)$.

Theorem 2 (Space Complexity of LAS).

The space complexity of LAS for the operator and load shedder is

$$\mathcal{O}\left(\frac{1}{\varepsilon} \log \frac{1}{\delta} (\log m + \log n)\right) \text{ bits.}$$

Theorem 3 (Communication complexity of LAS).

The communication complexity of LAS is of $\mathcal{O}\left(\frac{m}{N}\right)$ messages and

$$\mathcal{O}\left(\frac{m}{N} \left(\frac{1}{\varepsilon} \log \frac{1}{\delta} (\log m + \log n) + \log m\right)\right) \text{ bits.}$$

Note that the communication cost is low with respect to the stream size since the window size N should be chosen such that $N \gg 1$ (*e.g.*, in our tests we have $N = 1024$).

4 Theoretical Analysis

This section provides an analysis of the quality of the shedding performed by LAS in two steps. First, we study the correctness and optimality of the shedding algorithm, under *full knowledge* assumption (*i.e.*, the shedding strategy is aware of the exact execution duration w_t for each tuple t). Then, in Section 4.2, we

⁶ This correction factor derives from the fact that $\hat{w}(t)$ is a (ε, δ) -approximation of $w(t)$ as shown in Section 4.

⁷ For readability reasons, proofs of these theorems are available in A.

provide a probabilistic analysis of the mechanism that LAS uses to estimate the tuple execution durations. For the sake of simplicity, in both sections, we assume CHECKUTILITY always returns *true*. The proofs of the theorem are available in Appendix A.

4.1 Correctness of LAS

We suppose that tuples cannot be preempted, that is they must be processed uninterruptedly on the available operator instance. As mentioned before, in this analysis we assume that the execution duration $w(t)$ is known for each tuple t . Finally, given our system model, we consider the problem of minimizing d , the number of dropped tuples, while guaranteeing that the average queuing latency $\overline{Q}(t)$ will be upper-bounded by τ , $\forall t \in \sigma$. The solution must work online, thus the decision of enqueueing or dropping a tuple has to be made only resorting to knowledge about tuples received so far in the stream.

Let OPT be the online algorithm that provides the optimal solution to Problem 1. We denote with \mathcal{D}_{OPT}^σ (resp. d_{OPT}^σ) the set of dropped tuple indices (resp. the number of dropped tuples) produced by the OPT algorithm fed by stream σ (cf., Section 2). We also denote with d_{LAS}^σ the number of dropped tuples produced by LAS introduced in Section 3.3 fed with the same stream σ .

Theorem 4 (Correctness and Optimality of LAS). *For any σ , we have $d_{LAS}^\sigma = d_{OPT}^\sigma$ and $\forall t \in \sigma, \overline{Q}_{LAS}^\sigma(t) \leq \tau$.*

This theorem establishes that LAS is optimal, given that its execution time is the same as that of the optimal OPT algorithm. Moreover, it is correct in the sense of the Definition 1 proposed in Section 2, namely that its average queuing latency will not exceed the predetermined threshold τ .

4.2 Execution Time Estimation

In this section, we analyze the approximation made on execution duration $w(t)$ for each tuple t when the assumption of full knowledge is removed. LAS uses two matrices, \mathcal{F} and \mathcal{W} , to estimate the execution time $w(t)$ of each tuple submitted to the operator. By the Count Min sketch algorithm (cf., Section 3.2) and Listing 3.1, we have that for any $t \in [n]$ and each row $i \in [r]$,

$$\mathcal{F}[i][h_i(t)](m) = f_t + \sum_{u=1, u \neq t}^n f_u \mathbf{1}_{\{h_i(u)=h_i(t)\}},$$

and

$$\mathcal{W}[i][h_i(t)](m) = f_t w_t + \sum_{u=1, u \neq t}^n f_u w_u \mathbf{1}_{\{h_i(u)=h_i(t)\}}.$$

Let us denote respectively by w_{\min} and w_{\max} the minimum and the maximum execution time of the items. For sake of clarity in the following equations, we denote the ratio

$$\mathcal{V}_{i,t} = \mathcal{W}[i][h_i(t)]/\mathcal{F}[i][h_i(t)].$$

We have trivially

$$w_{\min} \leq \mathcal{V}_{i,t} \leq w_{\max}.$$

We define $S = \sum_{\ell=1}^n w_{\ell}$. We then have

Theorem 5.

$$\mathbb{E}\{\mathcal{V}_{i,t}\} = \frac{S - w_t}{n - 1} - \frac{k(S - nw_t)}{n(n - 1)} \left(1 - \left(1 - \frac{1}{k}\right)^n\right).$$

The proof of this theorem is available in appendix. First, it important to note that this result does not depend on m . Moreover, we easily understand that the formula proposed in this last theorem may seem rather uninformative. Thus, we propose to present a numeric application of it to take the measure of the potential use of it for an end-user.

We take for instance $k = 55$, $n = 4096$ and the distinct values of w_u equal to $1, 2, 3, \dots, 64$, each item being present 64 times in the input stream, we get for $t = 1, \dots, 64$, $\mathbb{E}\{\mathcal{V}_{i,t}\} \in [32.08, 32.92]$. Note also from above that we have $1 \leq \mathcal{V}_{i,t} = \mathcal{W}[i][h_i(t)]/\mathcal{F}[i][h_i(t)] \leq 64$.

From the Markov inequality, we have, for every $x > 0$,

$$\mathbb{P}\{\mathcal{V}_{i,t} \geq x\} \leq \frac{\mathbb{E}\{\mathcal{V}_{i,t}\}}{x}.$$

By taking $x = 64a$, with $a \in [0.6, 1)$, we obtain

$$\mathbb{P}\{\mathcal{V}_{i,t} \geq 64a\} \leq \frac{\mathbb{E}\{\mathcal{V}_{i,t}\}}{64a} \leq \frac{33}{64a}.$$

Recall that r denotes the number of rows of the system; we then have by the independence of the h functions,

$$\begin{aligned} \mathbb{P}\{\min_{i=1,\dots,r} (\mathcal{V}_{i,t}) \geq 64a\} \\ = (\mathbb{P}\{\mathcal{V}_{i,t} \geq 64a\})^r \leq \left(\frac{33}{64a}\right)^r. \end{aligned}$$

By taking for instance $a = 3/4$ and $r = 10$, we get

$$\mathbb{P}\{\min_{i=1,\dots,r} (\mathcal{V}_{i,t}) \geq 48\} \leq \left(\frac{11}{16}\right)^{10} \leq 0.024.$$

5 Experimental Evaluation

In this section, we evaluate the performance obtained by using LAS to perform load shedding. We first describe the general setting used to run the tests and then discuss the results obtained through simulations (Section 5.2) and with a prototype of LAS integrated within Apache Storm (Section 5.3).

5.1 Setup

Datasets — In our tests we consider both synthetic and real datasets. Synthetic datasets are built as streams of integer values (items) representing the values of the tuple attribute driving the execution duration when processed on the operator. We consider streams of $m = 32,768$ tuples, each containing a value chosen among $n = 4,096$ distinct items. Streams have been generated using the Uniform and Zipfian distributions with different values of $\alpha \in \{0.5, 1.0, 1.5, 2.0, 2.5, 3.0\}$, denoted respectively as Zipf-0.5, Zipf-1.0, Zipf-1.5, Zipf-2.0, Zipf-2.5, and Zipf-3.0. We define w_n as the number of distinct execution duration values that the tuples can have. These w_n values are selected at a *constant* distance in the interval $[w_{min}, w_{max}]$. We ran experiments with $w_n \{1, 2, \dots, 64\}$, however, due to space constraints, we only report results for $w_n = 64$, and with $w_{max} \in \{0.1, 0.2 \dots, 51.2\}$ milliseconds. Tests performed with different values for w_n did not show unexpected deviations from what is reported in this section. Unless otherwise specified, the frequency distribution is Zipf-1.0 and the stream parameters are set to $w_n = 64$, $w_{min} = 0.1$ ms and $w_{max} = 6.4$ ms; this means that the $w_n = 64$ execution durations are picked in the set $\{0.1, 0.2, \dots, 6.4\}$ ms.

Let \bar{W} be the average execution duration of the stream tuples, then the stream maximum theoretical input throughput sustainable by the setup is equal to $1/\bar{W}$. When fed with an input throughput smaller than $1/\bar{W}$ the system will be over-provisioned (*i.e.*, possible underutilization of computing resources). Conversely, an input throughput larger than $1/\bar{W}$ will result in an underprovisioned system. We refer to the ratio between the maximum theoretical input throughput and the actual input throughput as the percentage of underprovisioning that, unless otherwise stated, was set to 25%.

To generate 100 different streams, we randomize the association between the w_n execution duration values and the n distinct items: for each of the w_n execution duration values, we pick *uniformly* at random n/w_n different values in $[n]$ that will be associated to that execution duration value. This means that the 100 different streams we use in our tests do not share the same association between execution duration and item as well as the association between frequency and execution duration (thus each stream has also a different average execution duration \bar{W}). Each of these permutations has been run with 50 different seeds to randomize the stream ordering and the generation of the hash functions used by LAS. This means that each single experiment reports the mean outcome of 5,000 independent runs.

We considered two types of constraints defined on the queuing latency:

ABS(τ): requires that the queuing latency per tuple does not exceed τ milliseconds: $\forall i \in [m] \setminus D, q(i) \leq \tau$.

AVG(τ): requires that the total average queuing latency does not exceed τ milliseconds: $\forall i \in [m] \setminus D, \overline{Q}(i) \leq \tau$.

While not being a realistic requirement, the straightforwardness of the **ABS(τ)** constraint allowed us to grasp a better insight of the mechanisms of the algorithm. However, in this section, we only show results for the **AVG(6.4)** constraint as it is a much more sensible requirement with respect to a real setting.

The LAS operator window size parameter N , the tolerance parameter μ and the number of rows of the \mathcal{F} and \mathcal{W} matrices δ were set to $N = 1024$, $\mu = 0.05$ and $\delta = 0.1$ (*i.e.*, $r = 4$ rows) respectively. By default, the LAS precision parameter (*i.e.*, the number of columns of the \mathcal{F} and \mathcal{W} matrices) was set to $\varepsilon = 0.05$ (*i.e.*, $c = 54$ columns), however in one of the test we evaluated LAS performance using several values: $\varepsilon \in [0.001, 1.0]$. To evaluate LAS performance without other external factors, in all our experiments we set **CHECKUTILITY** to always return *true*.

For the real data, we used a dataset containing a stream of preprocessed tweets related to the 2014 European elections. Among other information, the tweets are enriched with a field *mention* containing the *entities* mentioned in the tweet. These entities can be easily classified into *politicians*, *media*, and *others*. We consider the first 500,000 tweets, mentioning roughly $n = 35,000$ distinct entities and where the most frequent entity has an empirical probability of occurrence equal to 0.065.

Tested Algorithms —We compare LAS performance against three other algorithms:

Base Line The Base Line algorithm takes as input the percentage of under-provisioning and drops at random an equivalent fraction of the tuples.

Straw-Man The Straw-Man algorithm uses the same shedding strategy of LAS, however, it uses the average execution duration \overline{W} as the estimated execution duration $\hat{w}(t)$ for each tuple t .

Full Knowledge The Full Knowledge algorithm uses the same shedding strategy of LAS, however, it feeds it with the exact execution duration w_t for each tuple t as they were provided by an omniscient oracle.

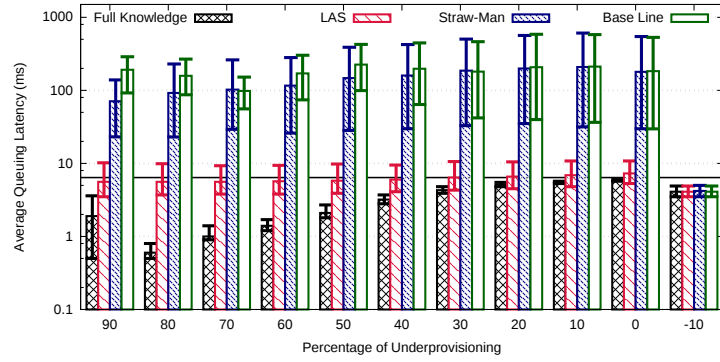
Evaluation Metrics —The evaluation metrics we used are:

- the dropped ratio: $\alpha = d/m$.
- the ratio of tuples dropped by algorithm *alg* with respect to Base Line: $\lambda = (d^{alg} - d^{\text{Base Line}})/d^{\text{Base Line}}$. In the following, we refer to this metric as shedding ratio.
- the average queuing latency: $\overline{Q} = \sum_{i \in [m] \setminus \mathcal{D}} q(i)/(m - d)$.
- the average completion latency, *i.e.*, the average time it takes for a tuple from the moment it is injected by the source in the topology, till the moment operator O concludes its processing.

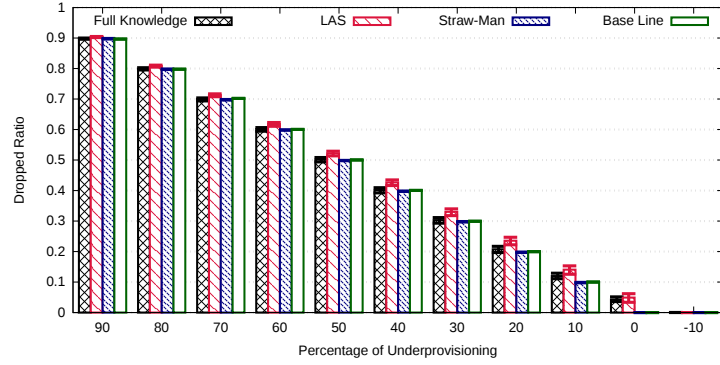
Whenever applicable we provide the maximum, mean, and minimum figures over the 5,000 runs.

5.2 Simulation Results

In this section, we analyze, through a simulator built ad-hoc for this study, the sensitivity of LAS while varying several characteristics of the input load. The simulator faithfully simulates the execution of LAS and the other algorithms and simulates the execution of each tuple t on O doing busy waiting for $w(t)$ milliseconds.



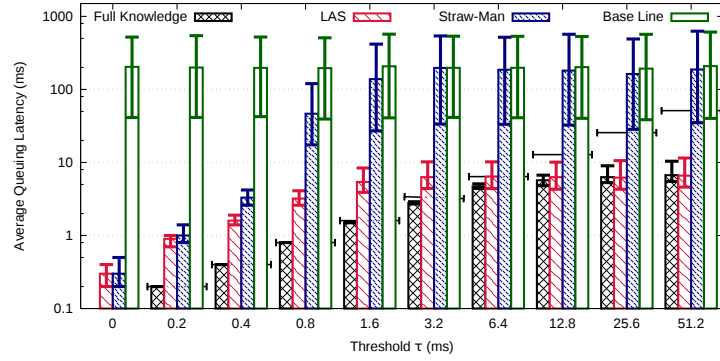
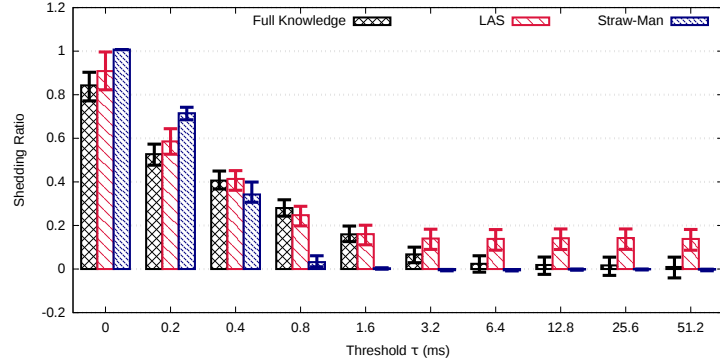
(a) Average queuing latency \bar{Q}



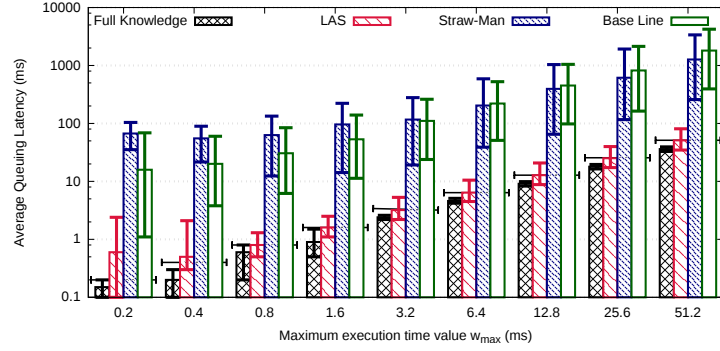
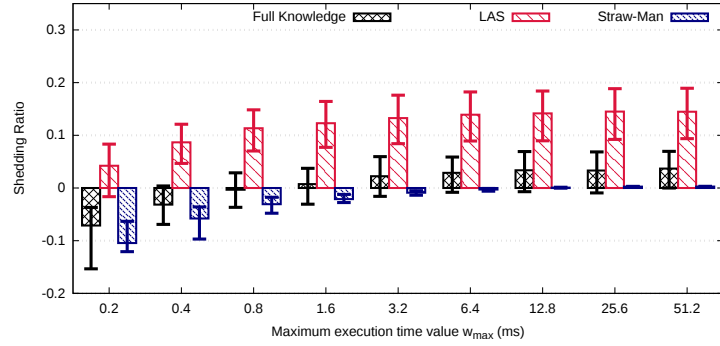
(b) Dropped ratio α

Fig. 4. LAS performance varying the amount of underprovisioning.

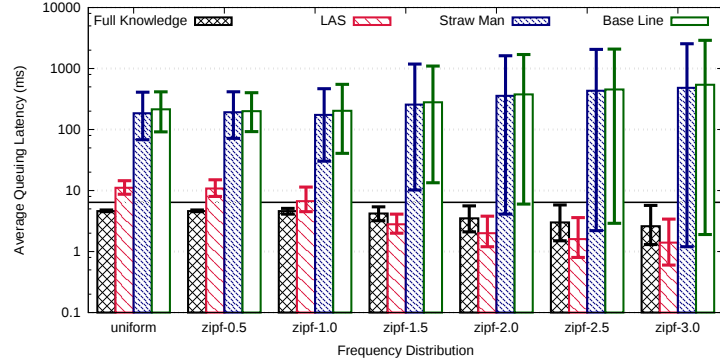
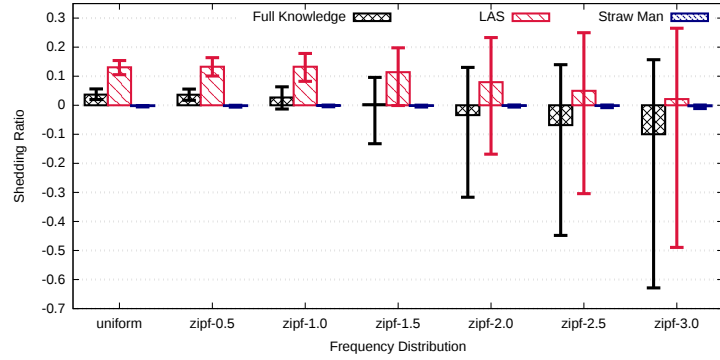
Input Throughput — Figure 4 shows the average queuing latency \bar{Q} (top) and dropped ratio α (bottom) as a function of the percentage of under-provisioning ranging from 90% to -10% (*i.e.*, the system is 10% overprovisioned with respect

(a) Average queuing latency \bar{Q} (b) Shedding ratio λ **Fig. 5.** LAS performance varying the threshold τ .

to the average input throughput). As expected, in this latter case all algorithms perform at the same level as load shedding is superfluous. In all the other cases both Base Line and Straw-Man do not shed enough load and induce a huge amount of exceeding queuing latency. On the other hand, LAS average queuing latency is quite close to the required value of $\tau = 6.4$ milliseconds, even if this threshold is violated in some of the tests. Finally, Full Knowledge always abide by the constraint and is even able to produce a much lower average queuing latency while dropping no more tuples that the competing solutions. Comparing the two plots we can see that the resulting average queuing latency is strongly linked to which tuples are dropped. In particular, Base Line and Straw-Man shed the same amount of tuples, LAS slightly more and Full Knowledge is in the middle. This result corroborates our initial claim that dropping tuples based on the load they impose allows designing more effective load shedding strategies.

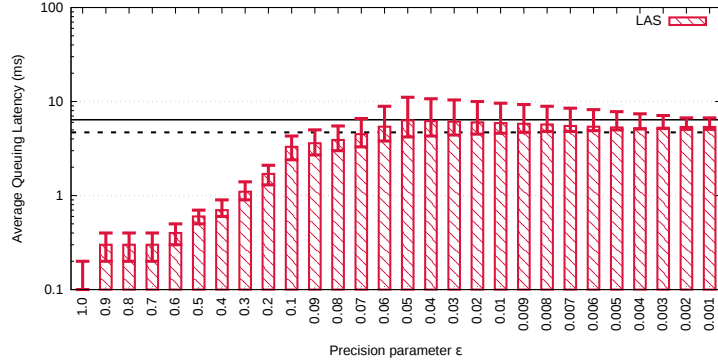
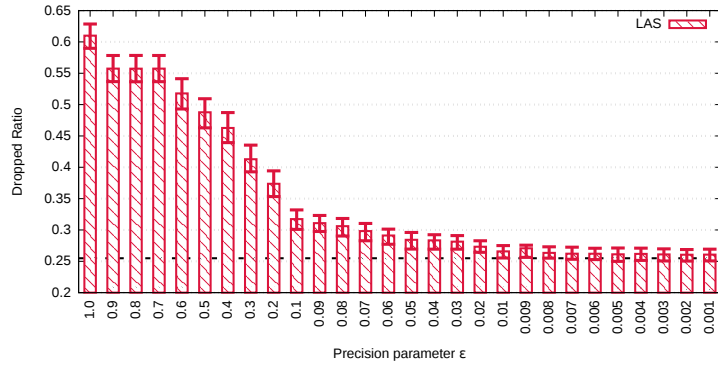
(a) Average queuing latency \bar{Q} (b) Shedding ratio λ **Fig. 6.** LAS performance varying the maximum execution duration value w_{max} .

Threshold τ — Figure 5 shows the average queuing latency \bar{Q} (top) and shedding ratio λ (bottom) as a function of the τ threshold. Notice that with $\tau = 0$ we do not allow any queuing, while with $\tau = 6.4$ we allow at least a queuing latency equal to the maximum execution duration w_{max} . In other words, we believe that with $\tau < 6.4$ the constraint is strongly conservative, thus representing a difficult scenario for any load shedding solution. Since Base Line does not take into account the latency constraint τ it always drops the same amount of tuples and achieves a constant average queuing latency. For this reason, Figure 5b reports the shedding ratio λ achieved by Full Knowledge, LAS, and Straw-Man against Base Line. The horizontal segments in Figure 5a represent the distinct values for τ . As the graph shows Full Knowledge always perfectly approaches the latency threshold, but for $\tau \geq 12.8$ where it is slightly smaller. Straw-Man performs reasonably well when the threshold is very small, but this is a consequence of the fact that it drops a large number of tuples when compared with Base Line as can be seen by Figure 5b. However, as τ becomes larger (*i.e.*, $\tau \geq 0.8$) Straw-Man

(a) Average queuing latency \bar{Q} (b) Shedding ratio λ **Fig. 7.** LAS performance varying the frequency probability distributions.

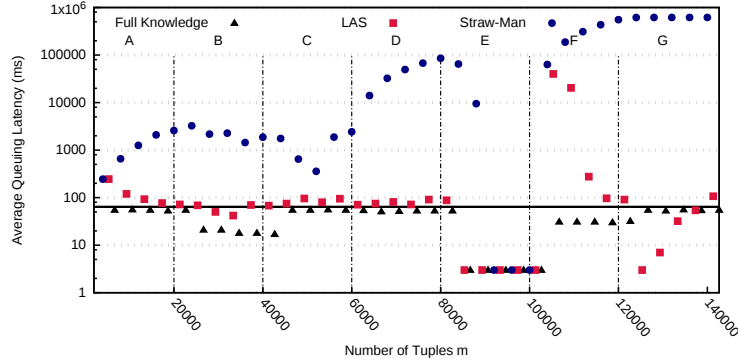
average queuing latency quickly grows and approaches the one from Base Line as it starts to drop the same amount of tuples. LAS, in the same setting, performs largely better, with the average queuing latency that for large values of τ approaches the one provided by Full Knowledge. While delivering this performance LAS drops a slightly larger amount of tuples compared to Full Knowledge, to account for the approximation in calculating tuple execution durations.

Maximum execution duration value w_{max} — Figure 6 shows the average queuing latency \bar{Q} (top) and dropped ratio λ (bottom) as a function of the maximum execution duration value w_{max} . Notice that in this test we varied the value for τ setting it equal to w_{max} . Accordingly, Figure 6a shows horizontal lines that mark the different thresholds τ . As the two graphs show, the behavior for LAS is rather consistent while varying w_{max} ; this means that LAS can be employed in widely different settings where the load imposed by tuples in the operator is not easily predictable. The price paid for this flexibility is in the shedding ratio

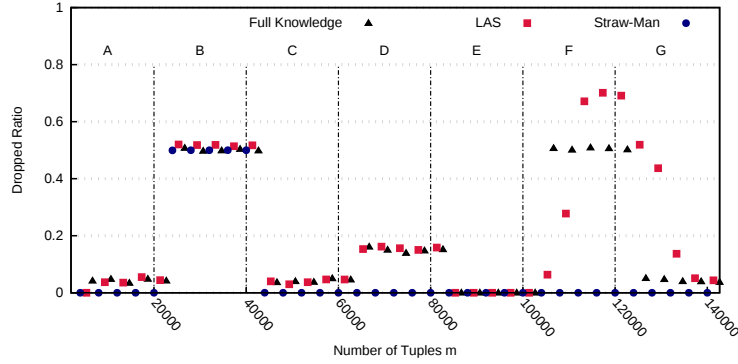
(a) Average queuing latency \bar{Q} (b) Dropped ratio α **Fig. 8.** LAS performance varying the precision parameter ϵ .

that, as shown in Figure 6b, is always positive.

Frequency Probability Distributions — Figure 7 shows the average queuing latency \bar{Q} (top) and dropped ratio λ (bottom) as a function of the input frequency distribution. As Figure 7a shows Straw-Man and Base Line perform invariably bad with any distribution. The span between the best and worst performance per run increases as we move from a uniform distribution to more skewed distributions as the latter may present extreme cases where tuple latencies match their frequencies in a way that is particularly favorable or unfavorable for these two solutions. Conversely, LAS performance improves the more the frequency distribution is skewed. This result stems from the fact that the sketch data structures tracing tuple execution durations perform at their best on strongly skewed distribution, rather than on uniform ones. This result is confirmed by the shedding ratio (Figure 7b) that decreases, on average, as α for the



(a) Average queuing latency \bar{Q}



(b) Dropped ratio α

Fig. 9. Simulator time-series.

distribution increases.

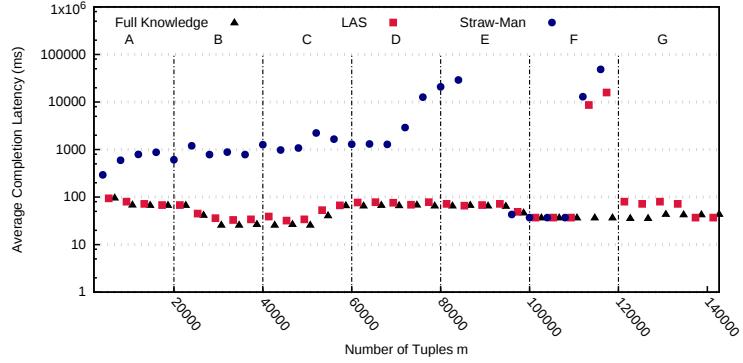
Precision parameter ε — Figure 8 shows the average queuing latency \bar{Q} (top) and dropped ratio α (bottom) as a function of the precision parameter ε . This parameter controls the trade-off between the precision and the space complexity of the sketches maintained by LAS. As a consequence, it has an impact on LAS performance. In particular, for large values of ε (left side of the graph), the sketch data structures are extremely small, thus the estimation $\hat{w}(t)$ is extremely unreliable. The corrective factor $1 + \varepsilon$ (see Listing 3.2 line 21) in this case is so large that it pushes LAS to largely overestimate the execution duration of each tuple. As a consequence LAS drops a large number of tuples while delivering average queuing latencies that are close to 0. By decreasing the value of ε (*i.e.*, $\varepsilon \leq 0.1$), sketches become larger and their estimation more reliable. In this configuration LAS performs at its best delivering average queuing latencies that are always below or equal to the threshold $\tau = 6.4$ while dropping a smaller

number of tuples. The dotted lines in both graphs represent the performance of Full Knowledge and are provided as a reference.

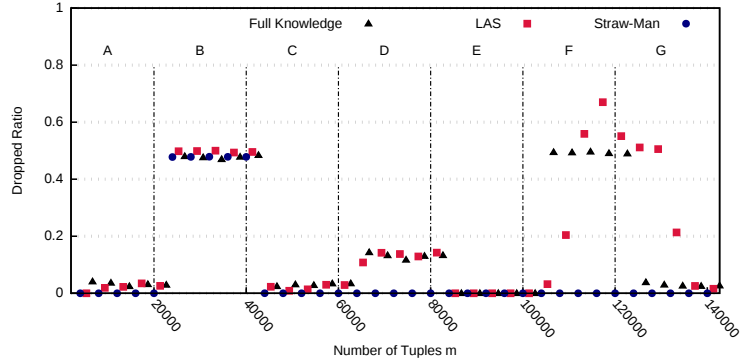
Time Series — Figure 9 shows the average queuing latency \bar{Q} (top) and dropped ratio α (bottom) as the stream unfolds (x -axis). Both metrics are computed on a jumping window of 4.000 tuples, *i.e.*, each dot represents the mean queuing latency \bar{Q} or the dropped ratio α computed on the previous 4.000 tuples. Notice that the points for Straw-Man, LAS and Full Knowledge related to the same value of the x -axis are artificially shifted to improve readability. In this test, we set $\tau = 64$ milliseconds. The input stream is made of 140,000 tuples and is divided into phases, from a A through G, each lasting 20,000 tuples. At the beginning of each phase we inject an abrupt change in the input stream throughput and distribution, as well as in $w(t)$ as follows:

- phase A** : the input throughput is set according to the provisioning (*i.e.*, 0% underprovisioning);
- phase B** : the input throughput is increased to induce 50% of underprovisioning;
- phase C** : same as phase A;
- phase D** : we swap the most frequent tuple t with a less frequent tuple t' such that $w(t') = w_{max}$, inducing an abrupt change in the tuple values frequency distribution and in the average execution duration \bar{W} ;
- phase E** : the input throughput is reduced to induce 50% of overprovisioning;
- phase F** : the input throughput is increased back to 0% underprovisioning and we also double the execution duration $w(t)$ for each tuple, simulating a change in the operator resource availability;
- phase G** : same as phase A.

As the graphs show, during phase A the queuing latencies of LAS and Straw-Man diverge: while LAS quickly approaches the performance provided by Full Knowledge, Straw-Man average queuing latencies quickly grow. In the same timespan, both Full Knowledge and LAS drop slightly more tuples than Straw-Man. All the three solutions correctly manage phase B: their average queuing latencies see slight changes, while, correctly, they start to drop larger amounts of tuples to compensate for the increased input throughput. The transition to phase C brings the system back in the initial configuration, while in phase D the change in the tuple frequency distribution is managed very differently by each solution: both Full Knowledge and LAS compensate this change by starting to drop more tuples, but still maintaining the average queuing latency close to the desired threshold τ . Conversely, Straw-Man cannot handle such change, and its performance incurs a strong deterioration as it drops still the same amount of tuples. In phase E the system is strongly overprovisioned, and, as it was expected, all three solutions perform equally well as no tuple needs to be dropped. The transition to phase F is extremely abrupt as the input throughput is brought back to the equivalent of 0% of underprovisioning, but the cost to handle each tuple on the operator is doubled. At the beginning of this phase, both Straw-Man and LAS perform badly, with queuing latencies that are largely above τ . However,



(a) Average completion latency



(b) Dropped tuples d

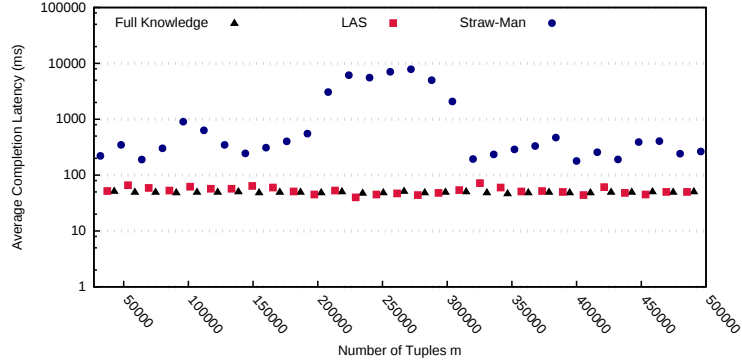
Fig. 10. Prototype time-series

while the phase unfolds LAS quickly updates its data structures and converges toward the given threshold, while Straw-Man diverges as tuples continue to be enqueued on the operator worsening the bottleneck effect. Bringing back the tuple execution durations to the initial values in phase G has little effect on LAS, while the bottleneck created by Straw-Man cannot be recovered as it continues to drop an insufficient number of tuples.

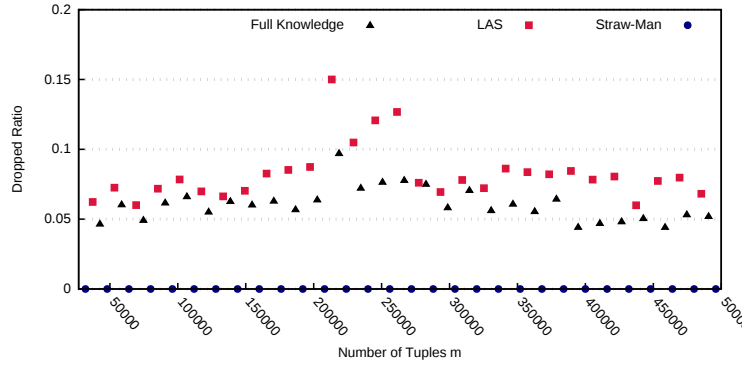
5.3 Prototype

To evaluate the impact of LAS on real applications we implemented it as a bolt within the Apache Storm [27] framework. We have deployed our cluster on Microsoft Azure cloud service, using a Standard Tier A4 VM (4 cores and 7 GB of RAM) for each worker node, each with a single available slot.

The test topology is made of a source (*spout*) and two operators (*bolts*) LS and O . The source generates (reads) the synthetic (real) input stream and emits the tuples consumed by bolt LS . Bolt LS uses either Straw-Man, LAS or Full



(a) Average completion latency

(b) Dropped tuples d **Fig. 11.** Prototype use case

Knowledge to perform the load shedding on its outbound data stream consumed by bolt O . Finally operator O implements the logic.

Time Series — In this test we ran the simulator using the same synthetic load used for the time series discussed in the previous section. The goal of this test is to show how our simulated tests capture the main characteristic of a real run. Notice, however, that plots in Figure 10 report the average completion latency per tuple instead of the queuing latency. This is due to the difficulties in correctly measuring queuing latencies in Storm. Furthermore, the completion latency is, from a practical point of view, a more significant metric as it can be directly perceived on the output. From this standpoint, the results, depicted in Figure 10, report the same qualitative behavior already discussed with Figure 9. Two main differences are worth to be discussed: firstly, the behaviors exposed by the shedding solution in response to phase transitions in the input load are in general shifted in time (with respect to the same effects reported in Figure 9) as

a consequence of the general overhead induced by the software stack. Secondly, several data points for Straw-Man are missing in phases E and G. This is a consequence of failed tuples that start to appear as soon as the number of enqueued tuples is too large to be managed by Storm. While this may appear as a sort of “implicit” load shedding imposed by Storm, we decided not to consider these tuples in the metric calculation as they have not been dropped as a consequence of a decision taken by the Straw-Man load shedder.

Simple Application with Real Dataset — In this test we pretended to run a simple application on a real dataset: for each tweet of the twitter dataset mentioned in Section 5.1 we want to gather some statistics and decorate the outgoing tuples with some additional information. However, the statistics and additional information differ depending on which class the entities mentioned in each tweet belong. We assumed that this leads to a long execution duration for *media* (e.g., possibly caused by access to an external DB to gather historical data), an average execution duration for *politicians* and a fast execution duration for *others* (e.g., possibly because these tweets are not decorated). We modeled execution durations with 25 milliseconds, 5 milliseconds, and 1 millisecond of busy waiting respectively. Each of the 500,000 tweets may contain more than one mention, leading to $w_n = 110$ different execution duration values from $w_{min} = 1$ millisecond to $w_{max} = 152$ milliseconds, among which the most frequent (36% of the stream) execution duration is 1 millisecond. The average execution time \bar{W} is equal to 9.7 millisecond, the threshold τ is set to 32 milliseconds and the under-provisioning is set to 0%.

Figure 11 reports the average completion latency (top) and dropped ratio λ (bottom) as the stream unfolds. As the plots show, LAS provides completion latencies that are extremely close to Full Knowledge, dropping a similar amount of tuples. Conversely, Straw-Man completion latencies are at least one order of magnitude larger. This is a consequence of the fact that in the given setting Straw-Man does not drop tuples, while Full Knowledge and LAS drop on average a steady amount of tuples ranging from 5% to 10% of the stream. These results confirm the effectiveness of LAS in keeping close control on queuing latencies (and thus provide more predictable performance) at the cost of dropping a fraction of the input load.

6 Related Work

Aurora [1] is the first stream processing system where shedding has been proposed as a technique to deal with bursty input traffic. Aurora employs two different kinds of shedding, the first and better detailed being random tuple dropping at strategic places in the application topology to satisfy QoS constraints.

A large number of works proposed solutions aimed at reducing the impact of load shedding on the quality of the system output. These solutions fall under the name of *semantic* load shedding, as drop policies are linked to the significance of each tuple with respect to the computation results. Tatbul et al. first introduced

in [26] the idea of semantic load shedding. Babcock et al. in [2] provided an approach tailored to aggregation queries. Tatbul et al. in [25] ported the concept of semantic load shedding in the realm of DSPS. GrubJoin [8] is a solution tailored for shedding load in multiway windowed stream joins while minimizing output degradation. Finally, Kalyvianaki et al. in [14] contextualized the problem to the realm of federated DSPS, and provided a solution for shedding fairness. Several solutions assume that the utility of an event depends on the event type and its frequency in the input event stream [26], *i.e.* they assume a static model for quality degradation; other works propose solutions to build and maintain at runtime a model for event utility [18,16]. All the previous works are based on the same goal, *i.e.*, to reduce the impact of load shedding on the semantics of the queries deployed in the stream processing system, while avoiding overloads. We believe that avoiding excessive degradation in the performance of the DSPS and in the semantics of the deployed query output are two orthogonal facets of the load shedding problem. In our work, we did not consider the latter and focused on the former while including in our solution the possibility to limit output quality degradation.

A different approach has been proposed in [20], with a system that builds summaries of dropped tuples to later produce approximate evaluations of queries. The idea is that such approximate results may provide users with useful information about the contribution of dropped tuples. A similar approach is adopted in StreamApprox [19] where the authors designed an online stratified reservoir sampling algorithm to produce approximate output with rigorous error bounds. A similar approach was also adopted in [28].

A classical control theory approach based on a closed control loop with feedback has been considered in [13,29,30]. In all these works the focus is on the design of the loop controller, while data is shed using a simple random selection strategy. In all these cases the goal is to reactively feed the stream processing engine system with a bounded tuple rate, without proactively considering how much load these tuples will generate.

Finally, a few works have recently appeared that address the problem of shedding load in Complex Event Processing (CEP) applications [9,10,23,22,31]. While these solution leverage techniques similar to those discussed in the previous paragraphs, they provide specific adaptations to the CEP context where input load can be shed both in the form of events and partial pattern matches.

7 Conclusions

In this paper, we introduced Load-Aware Shedding (LAS), a novel solution for load shedding in DSPS. LAS exploits a characteristic of many stream-based applications, *i.e.*, the fact that load on operators depends both on the input rate and on the content of tuples, to smartly drop tuples and avoid the appearance of performance bottlenecks. In particular, LAS leverages sketch data structures to efficiently collect at runtime information on the operator load characteristics and then use this information to implement a load shedding policy aimed at

maintaining the average queuing latencies close to a given threshold. Through a theoretical analysis, we proved that LAS is an (ϵ, δ) -approximation of the optimal algorithm. Furthermore, we extensively tested LAS both in a simulated setting, studying its sensitivity to changes of several characteristics of the input load, and with a prototype implementation integrated within the Apache Storm DSPS. Our tests confirm that by taking into account the specific load imposed by each tuple, LAS can provide performance that closely approaches a given target, while dropping a limited number of tuples.

References

1. D. J. Abadi, D. Carney, U. Çetintemel, M. Cherniack, C. Convey, S. Lee, M. Stonebraker, N. Tatbul, and S. Zdonik. Aurora: a new model and architecture for data stream management. *The International Journal on Very Large Data Bases (VLDB Journal)*, 12(2):120–139, 2003.
2. B. Babcock, M. Datar, and R. Motwani. Load shedding for aggregation queries over data streams. In *Proceedings of the 20th International Conference on Data Engineering (ICDE '04)*, pages 350–361. IEEE, 2004.
3. M. Borkowski, C. Hochreiner, and S. Schulte. Minimizing cost by reducing scaling operations in distributed stream processing. *Proc. VLDB Endow.*, 12(7):724–737, Mar. 2019.
4. J. L. Carter and M. N. Wegman. Universal classes of hash functions. *Journal of Computer and System Sciences*, 18, 1979.
5. G. Cormode. Sketch techniques for approximate query processing. In *Synopses for Approximate Query Processing: Samples, Histograms, Wavelets and Sketches, Foundations and Trends in Databases*. NOW publishers, 2011.
6. G. Cormode and S. Muthukrishnan. An improved data stream summary: The count-min sketch and its applications. *Journal of Algorithms*, 55, 2005.
7. A. Dobra, M. Garofalakis, J. Gehrke, and R. Rastogi. Sketch-based multi-query processing over data streams. In E. Bertino, S. Christodoulakis, D. Plexousakis, V. Christophides, M. Koubarakis, K. Böhm, and E. Ferrari, editors, *Advances in Database Technology - EDBT 2004*, pages 551–568, Berlin, Heidelberg, 2004. Springer Berlin Heidelberg.
8. B. Gedik, K. Wu, P. S. Yu, and L. Liu. Grubjoin: An adaptive, multi-way, windowed stream join with time correlation-aware cpu load shedding. *IEEE Transactions on Knowledge and Data Engineering*, 19(10):1363–1380, 2007.
9. Y. He, S. Barman, and J. F. Naughton. On load shedding in complex event processing. *arXiv preprint arXiv:1312.4283*, 2013.
10. Y. He, S. Barman, and J. F. Naughton. On load shedding in complex event processing. In *Proceedings of the 17th International Conference on Database Theory (ICDT '14)*, pages 213–224. OpenProceedings.org, 2014.
11. T. Heinze, L. Aniello, L. Querzoni, and Z. Jerzak. Cloud-based data stream processing. In *Proceedings of the 8th ACM International Conference on Distributed Event-Based Systems (DEBS '14)*, pages 238–245. ACM, 2014.
12. S. Ilarri, O. Wolfson, E. Mena, A. Illarramendi, and P. Sistla. A query processor for prediction-based monitoring of data streams. In *Proceedings of the 12th International Conference on Extending Database Technology: Advances in Database Technology*, EDBT '09, pages 415–426, New York, NY, USA, 2009. Association for Computing Machinery.

13. E. Kalyvianaki, T. Charalambous, M. Fiscato, and P. Pietzuch. Overload management in data stream processing systems with latency guarantees. In *7th IEEE International Workshop on Feedback Computing (Feedback Computing'12)*, 2012.
14. E. Kalyvianaki, M. Fiscato, T. Salonidis, and P. Pietzuch. Themis: Fairness in federated stream processing under overload. In *Proceedings of the 2016 International Conference on Management of Data*, pages 541–553. ACM, 2016.
15. A. Kammoun. *Enhancing Stream Processing and Complex Event Processing Systems*. PhD thesis, Université Jean Monnet, Saint-Etienne, 2019.
16. N. R. Katsipoulakis, A. Labrinidis, and P. K. Chrysanthis. Concept-driven load shedding: Reducing size and error of voluminous and variable data streams. In *2018 IEEE International Conference on Big Data (Big Data)*, pages 418–427, 2018.
17. Muthukrishnan. *Data Streams: Algorithms and Applications*. Now Publishers Inc., 2005.
18. C. Olston, J. Jiang, and J. Widom. Adaptive filters for continuous queries over distributed data streams. In *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data*, SIGMOD '03, pages 563–574, New York, NY, USA, 2003. Association for Computing Machinery.
19. D. L. Quoc, R. Chen, P. Bhatotia, C. Fetzer, V. Hilt, and T. Strufe. Streamapprox: Approximate computing for stream analytics. In *Proceedings of the 18th ACM/I-FIP/USENIX Middleware Conference*, Middleware '17, pages 185–197, New York, NY, USA, 2017. Association for Computing Machinery.
20. F. Reiss and J. M. Hellerstein. Data triage: An adaptive architecture for load shedding in TelegraphCQ. In *Proceedings of the 21st International Conference on Data Engineering (ICDE '05)*, pages 155–156. IEEE, 2005.
21. N. Rivetti, Y. Busnel, and A. Mostefaoui. Efficiently summarizing data streams over sliding windows. In *Proc. of the 14th IEEE International Symposium on Network Computing and Applications (NCA'15)*, Boston, USA, *Best Student Paper Award*, Sept. 2015.
22. A. Slo, S. Bhowmik, A. Flaig, and K. Rothermel. pSPICE: Partial match shedding for complex event processing. In *2019 IEEE International Conference on Big Data (Big Data)*, pages 372–382. IEEE, 2019.
23. A. Slo, S. Bhowmik, and K. Rothermel. eSPICE: Probabilistic load shedding from input event streams in complex event processing. In *Proceedings of the 20th International Middleware Conference*, pages 215–227, 2019.
24. I. Stanoi, G. Mihaila, T. Palpanas, and C. Lang. Whitewater: Distributed processing of fast streams. *IEEE Transactions on Knowledge and Data Engineering*, 19(9):1214–1226, 2007.
25. N. Tatbul, U. Çetintemel, and S. Zdonik. Staying fit: Efficient load shedding techniques for distributed stream processing. In *Proceedings of the 33rd international conference on Very large data bases*, pages 159–170. VLDB Endowment, 2007.
26. N. Tatbul, U. Çetintemel, S. Zdonik, M. Cherniack, and M. Stonebraker. Load shedding in a data stream manager. In *Proceedings of the 29th international conference on Very large data bases (VLDB '03)*, pages 309–320. VLDB Endowment, 2003.
27. The Apache Software Foundation. Apache Storm. <http://storm.apache.org>.
28. W. H. Tok, S. Bressan, and M.-L. Lee. A stratified approach to progressive approximate joins. In *Proceedings of the 11th International Conference on Extending Database Technology: Advances in Database Technology*, EDBT '08, pages 582–593, New York, NY, USA, 2008. Association for Computing Machinery.

29. Y.-C. Tu, S. Liu, S. Prabhakar, and B. Yao. Load shedding in stream databases: a control-based approach. In *Proceedings of the 32nd international conference on Very large data bases (VLDB '06)*, pages 787–798. VLDB Endowment, 2006.
30. Y. Zhang, C. Huang, and C. Huang. A novel adaptive load shedding scheme for data stream processing. In *Future Generation Communication and Networking (FGCN '07)*, pages 378–384. IEEE, 2007.
31. B. Zhao, N. Q. V. Hung, and M. Weidlich. Load shedding for complex event processing: Input-based and state-based techniques. In *Proceedings of ICDE (to appear)*, 2020.

A Theoretical Analysis

Data streaming algorithms strongly rely on pseudo-random functions that map elements of the stream to uniformly distributed image values to keep the essential information of the input stream, regardless of the stream elements frequency distribution.

This appendix extends with the proofs the theoretical analysis of the quality of the shedding performed by LAS in two steps provided in Section 4 as well as the complexities presented in Section 3.

First we study the correctness and optimality of the shedding algorithm, under *full knowledge* assumption (*i.e.*, the shedding strategy is aware of the exact execution duration w_t for each tuple t). Then, in A.3, we provide a probabilistic analysis of the mechanism that LAS uses to estimate the tuple execution durations.

A.1 Time, Space and Communication Complexities

In this section we provide the proofs of the time, space and communication complexities presented in Section 3.

Theorem 1 [Time complexity of LAS] For each tuple read from the input stream, the time complexity of LAS for the operator and the load shedder is $\mathcal{O}(\log 1/\delta)$.

Proof. By Listing 3.1, for each tuple read from the input stream, the algorithm increments an entry per row of both the \mathcal{F} and \mathcal{W} matrices. Since each has $\log 1/\delta$ rows, the resulting update time complexity is $\mathcal{O}(\log 1/\delta)$. By Listing 3.2, for each submitted tuple, the scheduler has to retrieve the estimated execution duration for the submitted tuple. This operation requires to read entry per row of both the \mathcal{F} and \mathcal{W} matrices. Since each has $\log 1/\delta$ rows, the resulting query time complexity is $\mathcal{O}(\log 1/\delta)$. \square

Theorem 2 [Space Complexity of LAS] The space complexity of LAS for the operator and load shedder is $\mathcal{O}\left(\frac{1}{\varepsilon} \log \frac{1}{\delta} (\log m + \log n)\right)$ bits.

Proof. The operator stores two matrices of size $\log(\frac{1}{\delta}) \times \frac{\varepsilon}{\delta}$ of counters of size $\log m$. In addition, it also stores a hash function with a domain of size n . Then the space complexity of LAS on the operator is $\mathcal{O}\left(\frac{1}{\varepsilon} \log \frac{1}{\delta} (\log m + \log n)\right)$ bits. The load shedder stores the same matrices, as well as a scalar. Then the space complexity of LAS on the load shedder is also $\mathcal{O}\left(\frac{1}{\varepsilon} \log \frac{1}{\delta} (\log m + \log n)\right)$ bits. \square

Theorem 3 [Communication complexity of LAS] The communication complexity of LAS is of $\mathcal{O}\left(\frac{m}{N}\right)$ messages and $\mathcal{O}\left(\frac{m}{N} \left(\frac{1}{\varepsilon} \log \frac{1}{\delta} (\log m + \log n) + \log m\right)\right)$ bits.

Proof. After executing N tuples, the operator may send the \mathcal{F} and \mathcal{W} matrices to the load shedder.

This generates a communication cost of $\mathcal{O}\left(\frac{m}{N} \frac{1}{\varepsilon} \log \frac{1}{\delta} (\log m + \log n)\right)$ bits via $\mathcal{O}\left(\frac{m}{N}\right)$ messages. When the load shedder receives these matrices, the synchronization mechanism kicks in and triggers a round trip communication (half of which is piggybacked by the tuples) with the operator. The communication cost of the synchronization mechanism is $\mathcal{O}\left(\frac{m}{N}\right)$ messages and $\mathcal{O}\left(\frac{m}{N} \log m\right)$ bits. \square

Note that the communication cost is low with respect to the stream size since the window size N should be chosen such that $N \gg 1$ (e.g., in our tests we have $N = 1024$).

A.2 Correctness of LAS

We suppose that tuples cannot be preempted, that is they must be processed uninterruptedly on the available operator instance. As mentioned before, in this analysis we assume that the execution duration $w(t)$ is known for each tuple t . Finally, given our system model, we consider the problem of minimizing d , the number of dropped tuples, while guaranteeing that the average queuing latency $\bar{Q}(t)$ will be upper-bounded by τ , $\forall t \in \sigma$. The solution must work online, thus the decision of enqueueing or dropping a tuple has to be made only resorting to knowledge about tuples received so far in the stream.

Let OPT be the online algorithm that provides the optimal solution to Problem 1. We denote with \mathcal{D}_{OPT}^σ (resp. d_{OPT}^σ) the set of dropped tuple indices (resp. the number of dropped tuples) produced by the OPT algorithm fed by stream σ (cf., Section 2). We also denote with d_{LAS}^σ the number of dropped tuples produced by LAS introduced in Section 3.3 fed with the same stream σ .

Theorem 4 [Correctness and Optimality of LAS] For any σ , we have $d_{LAS}^\sigma = d_{OPT}^\sigma$ and $\forall t \in \sigma, \bar{Q}_{LAS}^\sigma(t) \leq \tau$.

Proof. Given a stream σ , consider the sets of indices of tuples dropped by respectively OPT and LAS, namely \mathcal{D}_{OPT}^σ and \mathcal{D}_{LAS}^σ . Below, we prove by contradiction that $d_{LAS}^\sigma = d_{OPT}^\sigma$.

Assume that $d_{LAS}^\sigma > d_{OPT}^\sigma$. Without loss of generality, we denote $i_1, \dots, i_{d_{LAS}^\sigma}$ the ordered indices in \mathcal{D}_{LAS}^σ , and $j_1, \dots, j_{d_{OPT}^\sigma}$ the ordered indices in \mathcal{D}_{OPT}^σ . Let us define a as the largest natural integer such that $\forall \ell \leq a, i_\ell = j_\ell$ (i.e., $i_1 = j_1, \dots, i_a = j_a$). Thus, we have $i_{a+1} \neq j_{a+1}$.

- Assume that $i_{a+1} < j_{a+1}$. Then, according to Section 3.3, the i_{a+1} -th tuple of σ has been dropped by LAS as the method CHECK returned *true*. Thus, as $i_{a+1} \notin \mathcal{D}_{OPT}^\sigma$, the OPT run has enqueued this tuple violating the constraint τ . But this is in contradiction with the definition of OPT.
- Assume now that $i_{a+1} > j_{a+1}$. The fact that LAS does not drop the j_{a+1} tuple means that CHECK returns *false*, thus that tuple does not violate the constraint on τ . However, as OPT is optimal, it may drop some tuples for which CHECK is *false*, just because this allows it to drop an overall lower

number of tuples. Therefore, if it drops this j_{a+1} tuple, it means that OPT knows the future evolution of the stream and takes a decision on this knowledge. But, by assumption, OPT is an online algorithm, and the contradiction follows.

Then, we have that $i_{a+1} = j_{a+1}$. By induction, we iterate this reasoning for all the remaining indices from $a+1$ to d_{OPT}^σ . We then obtain that $\mathcal{D}_{OPT}^\sigma \subseteq \mathcal{D}_{LAS}^\sigma$.

As by assumption $d_{OPT}^\sigma < d_{LAS}^\sigma$, we have that $\exists \ell \in \mathcal{D}_{LAS}^\sigma \setminus \mathcal{D}_{OPT}^\sigma$ such that ℓ has been dropped by LAS. This means that, with the same tuple index prefix shared by OPT and LAS, the method CHECK returned *true* when evaluated on ℓ , and OPT would violate the condition on τ by enqueueing it. That leads to a contradiction. Then, $\mathcal{D}_{LAS}^\sigma \setminus \mathcal{D}_{OPT}^\sigma = \emptyset$, and $d_{OPT}^\sigma = d_{LAS}^\sigma$.

Furthermore, by construction, LAS never enqueues a tuple that violates the condition on τ because CHECK would return *true*.

Consequently, $\forall t \in \sigma, \overline{Q}_{LAS}^\sigma(t) \leq \tau$, which concludes the proof. \square

A.3 Execution Time Estimation

In this section, we analyze the approximation made on execution duration $w(t)$ for each tuple t when the assumption of full knowledge is removed. LAS uses two matrices, \mathcal{F} and \mathcal{W} , to estimate the execution time $w(t)$ of each tuple submitted to the operator. By the Count Min sketch algorithm (*cf.*, Section 3.2) and Listing 3.1, we have that for any $t \in [n]$ and for each row $i \in [r]$,

$$\begin{aligned} \mathcal{F}[i][h_i(t)](m) &= \sum_{u=1}^n f_u \mathbf{1}_{\{h_i(u)=h_i(t)\}} \\ &= f_t + \sum_{u=1, u \neq t}^n f_u \mathbf{1}_{\{h_i(u)=h_i(t)\}}. \end{aligned}$$

and

$$\mathcal{W}[i][h_i(t)](m) = f_t w_t + \sum_{u=1, u \neq t}^n f_u w_u \mathbf{1}_{\{h_i(u)=h_i(t)\}},$$

Let us denote respectively by w_{\min} and w_{\max} the minimum and the maximum execution time of the items. We have trivially

$$w_{\min} \leq \frac{\mathcal{W}[i][h_i(t)]}{\mathcal{F}[i][h_i(t)]} \leq w_{\max}.$$

We define $S = \sum_{\ell=1}^n w_\ell$. We then have

Theorem 5

$$\begin{aligned} &\mathbb{E}\{\mathcal{W}[i][h_i(t)]/\mathcal{F}[i][h_i(t)]\} \\ &= \frac{S - w_t}{n - 1} - \frac{k(S - nw_t)}{n(n - 1)} \left(1 - \left(1 - \frac{1}{k}\right)^n\right). \end{aligned}$$

It important to note that this result does not depend on m .

Proof.

For any $t = 1, \dots, n$, $\ell = 0, \dots, n-1$ and $A \in U_\ell(t)$, we introduce the event $B(t, \ell, A)$ defined by

$$B(t, \ell, A) = \{h_i(u) = h_i(t), \forall u \in A \text{ and} \\ h_i(u) \neq h_i(t), \forall u \in \{1, \dots, n\} \setminus (A \cup \{t\})\}.$$

From the independence of the hash function h_i , we have

$$\mathbb{P}\{B(t, \ell, A)\} = \left(\frac{1}{k}\right)^\ell \left(1 - \frac{1}{k}\right)^{n-1-\ell}.$$

Let us consider the ratio

$$\mathcal{V}_{i,t} = \mathcal{W}[i][h_i(t)] / \mathcal{F}[i][h_i(t)].$$

For any $i = 0, \dots, n$, we define

$$R_\ell(t) = \left\{ \frac{f_t w_t + \sum_{u \in A} f_u w_u}{f_t + \sum_{u \in A} f_u}, A \in U_\ell(t) \right\}.$$

We have $R_0(t) = \{w_t\}$. We introduce the set $R(t)$ defined by

$$R(t) = \bigcup_{\ell=0}^{n-1} R_\ell(t).$$

Thus with probability 1,

$$\mathcal{W}[i][h_i(t)] / \mathcal{F}[i][h_i(t)] \in R(t).$$

Let $x \in R(t)$. We have

$$\begin{aligned} & \mathbb{P}\{\mathcal{V}_{i,t} = x\} \\ &= \sum_{\ell=0}^{n-1} \sum_{A \in U_\ell(t)} \mathbb{P}\{\mathcal{V}_{i,t} = x \mid B(t, \ell, A)\} \mathbb{P}\{B(t, \ell, A)\} \\ &= \sum_{\ell=0}^{n-1} \left(\frac{1}{k}\right)^\ell \left(1 - \frac{1}{k}\right)^{n-1-\ell} \sum_{A \in U_\ell(t)} 1_{\{x=X(t,A)\}}. \end{aligned}$$

where $X(t, A)$ is the fraction:

$$X(t, A) = \frac{f_t w_t + \sum_{u \in A} f_u w_u}{f_t + \sum_{u \in A} f_u}$$

Thus

$$\begin{aligned}
\mathbb{E}\{\mathcal{V}_{i,t}\} &= \sum_{\ell=0}^{n-1} \left(\frac{1}{k}\right)^\ell \left(1 - \frac{1}{k}\right)^{n-1-\ell} \sum_{A \in U_\ell(t)} \sum_{x \in R(t)} x 1_{\{x=X(t,A)\}} \\
&= \sum_{\ell=0}^{n-1} \left(\frac{1}{k}\right)^\ell \left(1 - \frac{1}{k}\right)^{n-1-\ell} \sum_{A \in U_\ell(t)} X(t, A).
\end{aligned}$$

Let us assume that all the f_u are equal, that is for each u , we have $f_u = m/n$. The experimental evaluation tends to show that the worst case scenario of input streams is exhibited when all the items show the same number of occurrences in the input stream. We get

$$\begin{aligned}
\mathbb{P}\{\mathcal{V}_{i,t} = x\} &= \sum_{\ell=0}^{n-1} \left(\frac{1}{k}\right)^\ell \left(1 - \frac{1}{k}\right)^{n-1-\ell} \sum_{A \in U_\ell(t)} 1_{\{x = \frac{w_t + \sum_{u \in A} w_u}{\ell+1}\}}
\end{aligned}$$

that concludes the proof. \square