# Scaling Ordered Stream Processing on Shared-Memory Multicores

Guna Prasaad
University of Washington
guna@cs.washington.edu

G. Ramalingam
Microsoft Research India
grama@microsoft.com

Kaushik Rajan
Microsoft Research India
krajan@microsoft.com

## ABSTRACT

Many modern applications require real-time processing of large volumes of high-speed data. Such data processing needs can be modeled as a streaming computation. A streaming computation is specified as a dataflow graph that exposes multiple opportunities for parallelizing its execution, in the form of data, pipeline and task parallelism. On the other hand, many important applications require that processing of the stream be ordered, where inputs are processed in the same order as they arrive. There is a fundamental conflict between ordered processing and parallelizing the streaming computation. This paper focuses on the problem of effectively parallelizing ordered streaming computations on a shared-memory multicore machine.

We first address the key challenges in exploiting data parallelism in the ordered setting. We present a low-latency, non-blocking concurrent data structure to order outputs produced by concurrent workers on an operator. We also propose a new approach to parallelizing partitioned stateful operators that can handle load imbalance across partitions effectively and mostly avoid delays due to ordering. We illustrate the trade-offs and effectiveness of our concurrent data-structures on micro-benchmarks and streaming queries from the TPCx-BB [16] benchmark. We then present an adaptive runtime that dynamically maps the exposed parallelism in the computation to that of the machine. We propose several intuitive scheduling heuristics and compare them empirically on the TPCx-BB queries. We find that for streaming computations, heuristics that exploit as much pipeline parallelism as possible perform better than those that seek to exploit data parallelism.

## KEYWORDS

stream processing systems, streaming dataflow graph, continuous queries, data parallelism, partitioned parallelism, pipeline parallelism, dynamic scheduling, ordered processing, runtime, concurrent data structures

## 1 INTRODUCTION

Stream processing as a computational model has a long history dating back to Petri Nets in the 1960s, Kahn Process Networks and Communicating Sequential Processes in the 1970s, and Synchronous Dataflow in the 1980s [39]. Practical applications of stream processing were, for a long time, limited to audio, video and digital signal processing that typically involves deterministic, high-performance computations. Several languages and compilers such as StreamIt [22], Continuous Query Language(CQL) [9] and Imagine [28] were designed to specify and optimize the execution of such programs on single and shared-memory architectures.

The emergence of sensors and similar small-scale computing devices that continuously produce large volumes of data led to the rise of many new applications such as surveillance, fraud detection, environment monitoring, etc. The scale and distributed nature of the problem spurred the interest of several research communities that further gave rise to large scale distributed stream processing systems such as Aurora [2], Borealis [1], STREAM [11] and TelegraphCQ [19].

In today's highly connected world, data is of utmost value as it arrives. The advent of Big Data has further increased the importance of realtime stream processing. Several modern use-cases like shopping cart abandonment analysis, ad serving, brand monitoring on social media, and leader board maintenance in online games require realtime processing of large volumes of high-speed data. In the past decade, there has been a tremendous increase in the number of products (e.g. IBM Streams [37], Millwheel [5], Spark Streaming [41], Apache Storm [8], S4 [35], Samza [7], Heron [29], Microsoft Stream Insight [32], TRILL [18]) that cater to such data processing needs and is evidence of its ever-growing importance.

In this paper, we focus on scaling stream processing on the shared-memory multicore architecture. We consider this an important problem for several reasons. Streaming pipelines generally have a low memory footprint as most of the operators are either stateless or have a small bounded state. With increasing main memory sizes and prevalence of multi-core architectures, the bandwidth and parallelism offered by a single machine today is often sufficient to deploy pipelines with large number of operators [37]. So, most streaming workloads can be efficiently handled in a single multicore machine without having to distribute it over a cluster of machines. In fact, systems like TRILL [18] run streaming computations entirely on a single multicore machine at scales sufficient for several important applications. This is unlike batch processing systems, where the input, intermediate results and output data are often large and hence the pipeline needs to be split into many stages.

Even in workloads where distribution across a cluster is important (e.g. for fault tolerance), typically the individual nodes in the cluster are shared-memory multicores themselves. Most distributed streaming systems [8, 41] today assume each core in a multicore node as an individual executor and fail to exploit the advantages of low overhead communication offered by shared-memory. We believe that considering a multicore machine as a single powerful node rather than as a set of independent nodes can help better exploit shared-memory parallelism. A streaming computation can be split into multiple stages and each stage can be deployed on a shared-memory node [37]. Most prior work in this area have not studied the shared-memory multicore setting in depth - they either focus on the single core [10, 17, 27] or distributed shared-nothing architectures [1, 5, 6, 8, 29, 41].

Further, we are interested in ordered stream processing. The stream of events/tuples usually have an associated notion of *temporal ordering* such as in user click-streams from an online shopping session or periodic sensor readings. In many scenarios the application logic depends on this temporal order. For example, clustering click-streams into sessions based on timeouts between two consecutive events and computing time-windowed aggregates over streams of data. Implementing such logic on systems that do not provide ordered semantics is complicated and often turns out to be a performance bottleneck.

Ordered processing also enables our parallelization framework to be deployed easily on individual multicore nodes in a distributed stream processing cluster. This guarantee is important especially when a large stream processing query is divided into sub-queries, each allotted to a multicore node and one of them contains a non-commutative operation. Moreover, fault tolerance techniques such as active replication depend on state of the pipeline on two replicas being the same and it cannot be guaranteed without ordered processing.

Many stream processing systems today provide mechanisms to support ordered stream processing. Most of them based on a micro-batch architecture [1, 5, 18, 41], in which the input stream is broken down into streams of smaller batches and each batch is processed like in a batch processing system such as Map Reduce [20] or Apache Spark [40]. They support order sensitive pipelines by periodically sending *watermarks* denoting that all events less than a specific timestamp have been received. However, These techniques are not suitable for latency critical applications mainly due to the batching delays. We show that it is possible to achieve this guarantee at much lower latencies without constraining execution of the pipeline excessively.

## 1.1 Background and Challenges

A *streaming computation* can be specified as a dataflow graph, where each vertex is associated with an operator and directed edges represent flow of input into and out of the operators. At runtime, every vertex receives a stream of values (which we refer to as *tuples* henceforth) along each of its incoming edges. These tuples are then processed by the operator to produce zero or more output tuples that are then sent along its outgoing edges.

A unary operator processes an input (of some type $T_{in}$) and produces a sequence of zero or more outputs (of some type $T_{out}$). Every vertex with a single incoming edge has an associated unary operator that specifies the computation to be performed at that vertex. map, filter and windowed-aggregate are examples of such unary operators. A vertex with $n$ edges abstractly represents an $n-$ary operator, with $n$ inputs of types $T_1, T_2, ..., T_n$. In the streaming setting, the semantics of an $n$-ary operator too can be specified as a function that maps a tuple on a (specified) incoming edge to a sequence of zero or more output tuples.

Some operators are pure functions that do not have any state associated with its computation and hence called *stateless operators*. Some operators have an internal state that is accessed and updated during the computation - for example, windowed-count maintains the count of tuples in the current window as internal state. Such operators are called *stateful operators* . In some cases, the operator accesses only a part of the state during the computation, which is pre-determined by a *key* associated with every input tuple. These are called *partitioned stateful operators*, as the state can be partitioned by the key. windowed-group-by-count is an operator of this type.

*1.1.1 Opportunities for Parallelization.* A streaming dataflow graph exposes various opportunities for parallelizing the computation efficiently. We elucidate this using an example: figure 1 represents an algorithm to detect high-mobility fraud using call data records as a streaming dataflow graph.

Call data records (CDR) are generated by every call between two mobile phones and it contains information such as time, duration of the call, location and phone number of the caller and the callee. In the detection algorithm, a CDR is first filtered (1, fig. 1) on the interested area code and the caller/callee's time and location information is projected (2) as a record. These location records are then grouped by phone number to compute (3) the speed at which a user must have traveled between locations. Phone numbers that have a speed greater than $T$, are then filtered (4), and the number of such cases in a given time window are counted (5).

An operator is said to be *data parallel*, if its inputs can be processed concurrently. Stateless operators such as (1, 2, 4) in the example are data parallel. On the other hand, inputs to a partitioned stateful operator can be processed in parallel only if they belong to different partitions. Hence, they are said to exhibit *partitioned parallelism*. In our example, computing the speed based on location records (3) for two different phone numbers can be done in parallel. Non-commutative stateful operators do not exhibit any data parallelism.

Further, when two operators are connected to each other such that the output of one forms the input to another, they are said to exhibit *pipeline parallelism*. In that case, these two operators can be processed concurrently. For example, one worker can compute the speed (3) for a particular phone number, while another filters (4) some phone numbers based on the speed already computed and sent to be filtered. Finally, a dataflow graph also exposes *task parallelism*. If two operators are not connected to each other via an input-output relationship, directly or indirectly, they can be processed concurrently. For example, operators on two sibling nodes in a DAG exhibit this kind of parallelism.

*1.1.2 Ordered Processing.* Ordered processing specifies that processing of inputs to a streaming pipeline must be semantically equivalent to executing them serially one at a time in the order of their arrival. We achieve this by ensuring that each individual operator implementation guarantees ordered processing and hence by extension any pipeline built by composing these implementations provide the ordering guarantee.

There is a fundamental conflict between data parallelism and ordered processing. Data parallelism seeks to improve the throughput of an operator by letting more than one worker operate on the inputs from the worklist concurrently. On the other hand, ordered processing requires to process them in the order of their arrival. The key observation here is that depending on the type of the operator, a concurrent execution might still be semantically equivalent to a serial single-threaded execution.
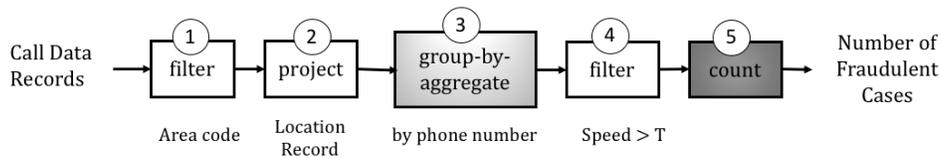
**Figure 1: Algorithm for high-mobility fraud detection expressed as a streaming dataflow graph.**

Ordered processing for a stateful operator is straightforward as its maximum allowed degree of parallelism is 1. In case of stateless and partitioned stateful operators, however, multiple workers process inputs concurrently and so they need special constructs to ensure that their concurrent execution is equivalent to a serial single-threaded execution.

There are essentially two kinds of ordering requirements that must be handled correctly. The first kind is *processing order*: for some operators we need to ensure that the processing logic of the operator is executed on the inputs in the same order as they arrive. This is a key requirement for non-commutative stateful operators. On the other hand, there is no such constraint for stateless operators. Partitioned stateful operators present an interesting middle-ground where it is enough to guarantee that tuples with the same key are processed in their arrival order.

The second kind of requirement is *output ordering*, which specifies that the outputs of an operator are sent to the downstream operator in the same order as its inputs. In particular, even when inputs $i_1$ and $i_2$ can be processed concurrently (when they belong to different partitions or when the operator is stateless), we still need to ensure that the outputs $o_1$ and $o_2$ produced by these inputs respectively are sent out in the right order. We guarantee this property for both stateless and partitioned stateful operators using special concurrent data structures. We describe a low-overhead, non-blocking solution to this problem in Sec. 3.

Output ordering is innately a blocking constraint: even if $o_2$ is produced before $o_1$, it gets blocked until $o_1$ is produced and sent downstream. This manifests as an implicit advantage for parallelization schemes that processes inputs from the worklist *almost* in the order of their arrival even though the semantics does not impose a restriction on this order. For stateless operators, having a shared worklist directly enables this execution pattern. However, it is non-trivial to achieve this for partitioned stateful operators. We present an adaptive partitioning scheme that supports this notion of almost ordered processing in Sec. 4.

### 1.2 Contributions

In this paper, we make the following contributions:

(1) (Sec. 3) We present a low-overhead non-blocking reordering scheme to order outputs of an operator that are produced concurrently. We observe that it scales better than a standard lock-based scheme and overall provides better throughput for long pipeline queries.

(2) (Sec. 4) We propose a novel scheme for exploiting partitioned parallelism in the ordered setting. We observe that our scheme achieves better speedup than the predominantly used strategy for partitioned parallelism during partition-induced skews and

mostly avoids delay due to ordering constraints leading to much lower latencies.

(3) (Sec. 6) We propose several intuitive scheduling heuristics that can be used to dynamically schedule operators at runtime. We identify a single heuristic that produces the best throughput and near best latency.

(4) (Sec. 7) We evaluate our runtime on streaming queries from TPCx-BB[16] and demonstrate that we can provide a throughput of millions of tuples per second on some queries with latency in the order of few milliseconds.

## 2 SOLUTION OVERVIEW

A generalized solution model for executing a stream processing query comprises of two components: a *compiler* and a *runtime*. The compiler is responsible for static optimizations, while the runtime takes this compiled representation and executes it on the machine, potentially with dynamic optimizations.

The relative roles of the compiler and runtime are determined by the type of streaming computation. For example, the streaming computations in signal processing are deterministic, and operator characteristics (such as per-tuple processing cost, selectivity) are known a priori. Such workloads provide more opportunities for static compiler optimizations, and the runtime is a straightforward execution of the produced scheduling plan. This model is exemplified by systems like StreamIt [22] and Brook[15]. In other applications like monitoring, fraud detection or shopping cart analysis there is little to no information about the operator characteristics at compile-time and hence the scope of static optimizations are fewer. So, systems like Borealis[1] and STREAM [11] designed for these workloads rely heavily on dynamic optimizations. However, even in such dynamic workloads there is some scope for static optimizations like coarsening of operators, pushing up filters. Refer [26] for a detailed catalog of such optimizations.

In our system, we target dynamic workloads to support use-cases that have risen in many new Big Data applications. We assume that a stream processing query is initially compiled into an optimal pipeline using some of the known techniques. We then deploy this optimal version of the pipeline on a runtime that seeks to efficiently parallelize its execution with the ordered processing guarantee. Here we focus only on the design of runtime, as the compilation stage is quite well studied in earlier works. We limit our discussion to linear chain pipelines, which is the predominant structure present in most stream processing queries. We believe our ideas can be generalized to other DAG structures as well, but we do not specifically address them here.

## 2.1 Problem Definition

The system accepts a pipeline that consists of operators connected to each other as a linear chain. The operators are specified to be either of stateless, stateful or partitioned stateful type. In case of partitioned stateful operators, the user also specifies a key selector that can be used to associate tuples with keys and a partitioning strategy such as hash or range partitioning to further map keys to partitions. The goal of the runtime is to execute this linear pipeline efficiently on a shared-memory multicore machine by exploiting various forms of parallelism as described in Sec. 1.

There are two dimensions of performance for a stream processing system that we are interested in. First is the throughput, by which we refer to the number of tuples processed to completion every second. The second dimension is latency. There are two notions of latency prevalent in the literature: end-to-end latency, which is the time duration between entry at ingress and exit at egress, and processing latency, which is the time since the first operator begins processing a tuple until exit at egress. The difference between them is that end-to-end latency includes the time spent by the tuple in the input queue for the overall pipeline. In the rest of this paper, we refer to processing latency when we say latency. The objective here is to maximize the throughput to handle high-speed data while minimizing the latency to process them in realtime.

## 2.2 Runtime Design

Our runtime is based on an *asynchronous* model of execution. We first decouple the pipeline into individual operators and compile them to independently schedulable units, one for every operator. We do this by associating every operator with a worklist(s). Inputs to an operator are simply added to its worklist instead of executing the operator logic synchronously. When the operator is scheduled, it obtains inputs from the worklist, processes them and adds the outputs to the worklist of the downstream operator.

The goal of the runtime is to choose which operator to choose, at what time and on which core? The two essential components of our runtime are *worker threads* and the *scheduler* data structure. The worker threads are the work horses of our runtime and responsible for advancing the progress of operators. Worker threads periodically query the scheduler for work. A worker, when allotted to an operator, dequeues an input from the operator's worklist, performs the operation and adds the output(s) produced to worklist of the next operator in the pipeline. A worker is specified with the maximum number of tuples to process in an operator and when allotted it processes as many tuples before deciding which operator to work on next. The worker additionally collects runtime information about each operator such as number of inputs consumed, outputs produced, time taken to process them. This is then used to estimate operator characteristics like average per-tuple processing cost and average selectivity.

Scheduling decisions regarding which operator must be scheduled next are made by a central scheduler data structure. This decision is made using estimated operator characteristics, current worklist sizes, and possibly observed throughput and latency measurements. We achieve this using scheduling heuristics - we discuss several of them in Sec. 6. When a heuristic chooses to schedule two different operators on different cores, it seeks to exploit pipeline parallelism. When it schedules the same operator on different cores it exploits data parallelism ingrained in the operator. Overall the goal of the scheduler is use to *dynamically* determine an ideal combination of data and pipeline parallelism among operators to achieve optimal performance.

The scheduler in our runtime can dynamically schedule more than one worker on an operator. This is applicable only to data or partition parallel operators as the maximum degree of parallelism allowed by a stateful operator is 1. The implementation of these operators internally handle the required concurrency control to ensure correct and ordered processing (refer Sec. 5). This is unlike many other architectures [37], where a single logical operator is replicated into a statically determined number of physical operators that are then scheduled independently.

## 3 REORDERING SCHEME

In this section, we handle the problem of ordering outputs produced by concurrent workers before they are sent to downstream operator(s). Most prior solutions to this problem are restricted to the micro-batching architecture: the input tuple stream is considered as a stream of batches, where tuples in a batch are executed in parallel and their outputs are finally sorted before sending them downstream. The notion of batching has some advantages including amortizing the cost involved in sorting, admitting columnar-based and operator-specific batch optimizations. However, these solutions are predominantly known to trade off latency for throughput. Our approach seeks to perform this reordering incrementally using low overhead non-blocking concurrent data structures.

For stateless and partitioned stateful operators, multiple workers can consume inputs from their worklist producing outputs concurrently. Each input is associated with a unique serial number (starting from 1) denoting its arrival order into the worklist of the operator. This serial number is assigned using an atomic counter at the time of enqueueing them to the worklist(s) of the operator. In some cases, a single input can produce more than one outputs. However, they are considered together as one unit and is associated with a single serial number. The schemes we describe below are concerned only with ordering outputs based on this serial number.

Specifically, the ordering constraint requires that for all $t$, the output $o_t$ produced by a tuple $i_t$ be sent downstream (either to an operator or egress) only after $o_1, o_2, ..., o_{t-1}$ are sent downstream. Since, these outputs are produced by concurrent workers, they are produced in no predetermined order. So, $o_{t+1}$ might be produced before $o_t$ and in that case $o_{t+1}$ has to wait until $o_t$ is produced and sent. We first describe a lock-based solution that implements this waiting scheme. We show that such a straight-forward design could lead to sub-optimal performance. Then, we present our improved low-latency, non-blocking solution.

### 3.1 Lock-Based Solution.

A standard approach would be to use a waiting buffer and a counter. The counter keeps track of the serial number of next output to be sent. Whenever the corresponding output is available it is sent downstream immediately and the counter incremented. If an output is not the next one to be sent, we simply add it to the waiting buffer and return to process more inputs. So, when an output is sent, we

```
1    void send(o_t) {
2      lock();
3      if (t == next) {
4        send_downstream(o_t);
5        next++;
6        while(buffer has o_next) {
7          send_downstream(o_next);
8          next++;
9      } } else {
10         add o_t to buffer
11     }
12     unlock();
13   }
```

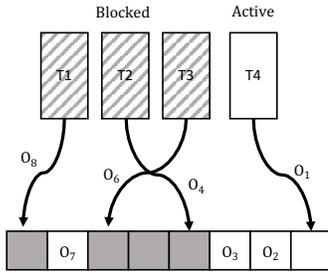**Figure 2: Lock-based Scheme: The global lock used here induces unnecessary blocking behavior**



**Figure 3: Unnecessary Blocking:** $T_1, T_2$ and $T_3$ **are blocked until** $T_4$ **sends outputs** $O_1, O_2$ **and** $O_3$ **downstream.**

must check the waiting buffer for the next output and if present, send that to the downstream operator and repeat. Further, we do not want multiple workers to send the output(s) downstream or increment the counter concurrently as that will violate our ordering guarantee. So, we protect the overall logic using a global lock to ensure correctness and progress. This scheme is listed in fig. 2

However, this scheme results in sub-optimal performance due to unnecessary blocking of workers. Consider the scenario shown in fig. 3: worker thread $T_4$ produces $O_1$, which is the next output to send downstream, obtains the lock and keeps sending outputs $O_2$ and $O_3$ as they are already available in the waiting buffer. Meanwhile, workers $T_1, T_2, T_3$ that produced outputs $O_4, O_6$ and $O_8$ respectively get blocked trying to acquire the lock. However, we know that outputs have pre-allotted serial numbers. So, adding them to the waiting buffer can be totally independent of sending them downstream. Ideally, workers $T_1, T_2$ must be able to add their outputs to the waiting buffer while another worker is sending outputs downstream and return back to do useful work.

## 3.2 Non-Blocking Solution

We improve this version by replacing the lock with an atomic flag, essentially to provide try_lock semantics. This scheme is listed elaborately in fig. 4. Any worker $w$ seeking to send an output

downstream, first tries to add it in a bounded circular buffer. The buffer is used to store available outputs that are not yet ready to be sent. This step can either fail or succeed based on the size of buffer and current value of the next counter. If it fails, the worker tries again with the same output, after it exits the send function.

Before exiting, irrespective of success or failure in the add step, $w$ tries to send pending outputs in the buffer to downstream operator(s). It can do so only when it can test_and_set a global atomic flag. If it cannot set the flag, it means that another worker $w'$ is performing this step. In that case, $w$ simply exits the function instead of getting blocked, unlike in the lock-based scheme.

If $w$ can set the flag, it has exclusive access to send the buffered outputs. First, it obtains the current value of next counter and the corresponding value from the buffer array. If this value is not EMPTY, it sends the output downstream, increments the counter and repeats this again for the new value of next. If the obtained value is EMPTY then, $w$ clears the flag and exits the loop. Further, to ensure that every output is sent downstream as soon as it is ready to be, $w$ checks the buffer array again and retries to send the previously unavailable output, if it is available now. This ensures that there is no ready-to-send output in the buffer, when there are no active workers inside send.

THEOREM 3.1 (CORRECTNESS OF NON-BLOCKING REORDERING SCHEME). *If all concurrent workers allotted to an operator send outputs to operators downstream by invoking the* send *procedure (fig. 4), then output $o_t$ (with serial number $t$) is sent downstream (by invocation of* send_downstream) *only after all outputs $o_1, o_2, ..., o_{t-1}$ are sent.*

PROOF. The outputs are sent downstream only inside the send-pending-outputs procedure, in which lines L27-36 (referred to as *exit section*) are protected from concurrent access by the atomic flag variable flag. Since this makes the exit section a critical section, at most one worker increments the next counter and sends the pending outputs in the buffer to the operator downstream. It is quite clear from the control flow in the exit section, that whatever non-EMPTY output is present in buffer[i], it is sent downstream as the output with serial number $n$, where $i = n \mod s$. Now, it suffices to prove that if the value of next is $n$ and $o$ is the value obtained by loading buffer[i] as in fig. 4, then the following two conditions hold:

(1) If $o_n$ has not been added to the buffer, then $o$ is EMPTY
(2) If $o$ is not EMPTY, then the value of $o$ is $o_n$

In order to prove this, we first define $T_k$ to be the time at which the next counter is atomically incremented from $k$ to $k + 1$. For simplicity of explanation, we assume $T_k$, for $k < 0$, to be some global initialization time when buffer array is initialized with EMPTY.

The condition at L16 (referred as *entry condition*) determines whether an output $o_t$ (with serial number $t$) can be added to the buffer at $i = (t \mod s)$ or not. This condition enforces that $o_t$ can be added only when next $\in (t - s, t]$, which in turn can happen only during the time interval $(T_{t-s}, T_t)$.

Since all updates to the global data fields are atomic, they are sequentially consistent. So, the value of buffer[i] (where $i = t \mod s$) is set to EMPTY in the exit section before $T_{t-s}$. We also know that this will definitely remain EMPTY until $T_{t-s}$. This is because

```
1    //data fields
2    atomic_long next;
3    atomic<output*> buffer[s];
4    atomic_flag flag;
5
6    //invoked by workers
7    bool send(o_t) {
8      bool success = try_add(o_t);
9      send_pending_outputs();
10     return success;
11   }
12
13   //helper functions
14   bool try_add(o_t) {
15     n = next.load();
16     if(t ≥ n and t < n + s) {
17       i = t mod s;
18       buffer[i].set(o_t);
19       return true;
20     } else {
21       return false;
22   } }
23
24   void send_pending_outputs() {
25     if (not flag.test_and_set()) {
26       //send as many outputs as possible
27       while(true) {
28         n = next.load();
29         i = n mod s;
30         o = buffer[i].load();
31         if (o is not EMPTY) {
32           send_downstream(o);
33           buffer[i].set(EMPTY);
34           next.fetch_add(1);
35         } else {
36           flag.clear();
37           break;
38       } }
39       //re-check if next output is available
40       o = buffer[i].load();
41       if (o is not EMPTY) {
42         send_pending_outputs();
43   } } }
```

Figure 4: Non-blocking Reordering Scheme

the only valid output that can be added at $i$ during that intermittent time is $o_{t-s}$ due to the entry condition. But, we know $o_{t-s}$ has already been added once and by uniqueness of serial numbers, we can assert it will not be added again. From just after $T_{t-s}$, this value will still remain EMPTY until some worker adds an output into that slot. Again, the entry condition now ensures that only $o_t$ can be added to buffer[$i$] during $(T_{t-s}, T_t)$. Hence, (1) holds.

Further, control flow in the exit section necessitates that $o_t$, if available, is read into $o$ before $T_t$. This together with the entry

condition, ensures that $o_t$ is not overwritten before being read back from buffer[$i$]. Hence, (2) holds. Both conditions (1) and (2), in addition with the guarantee that next cannot have a value $k + 1$ before $k$, we can assert that the outputs are indeed sent downstream in the serial order.                                                    □

*Progress.* In the above scheme, none of the concurrent workers get blocked due to another worker sending outputs. However, a worker can get blocked due to limited size of the waiting buffer: when it tries to send an output that corresponds to input with a serial number much higher than the current value of next, it can potentially get blocked trying and failing repeatedly to add the output. This is because the entry condition prevents this output to be added until some earlier outputs are sent and the buffer makes space for this output. Meanwhile, this worker repeatedly tries to send it and fails.

One simple way to handle this would be to use a non-blocking concurrent map instead of a bounded array. However, the overheads in a simple array are much lesser compared to the alternatives and hence we chose such a design. Even though, we can never eliminate this scenario with a bounded buffer, we can try to avoid its occurrence as much as possible. One could use an appropriately sized waiting buffer. Further, we could employ design strategies such that concurrent workers working on a data parallel operator would produce outputs almost in-order of their serial numbers. This is a key design strategy in exploiting partitioned parallelism in the ordered setting, which we present in the next section.

## 4 PARTITIONED PARALLELISM

The essence of partitioned parallelism is that every input to be processed has a *key*, and the state required to process inputs with different keys are disjoint. This allows us to process tuples with different keys in parallel, though those with the same key must be processed sequentially, in order.

The key space can be statically partitioned into many disjoint buckets based on a strategy such as range or hash partitioning. The system treats tuples belonging to the same bucket as potentially having the same key and processes them sequentially. If the number of buckets is $p$, it limits the degree of parallelism to $p$. Ideally, we would like to have as many buckets as the number of keys to exploit as much parallelism as possible even during load imbalance induced by the partitioning strategy. But, scheduling overheads and complexities in the key space force us to have a fewer, fixed number of buckets. However, a more fine-grained partitioning strategy is still preferable, given the overheads are admissible. Profiling data gathered from sample runs can be used to determine both $p$ and the partitioning of key space into $p$ buckets.

Further, we would like to design a flexible scheme where workers can be dynamically allotted to operators. This is necessary to support a dynamic scheduling based runtime that allots workers to operators based on current status of the pipeline. As we saw in Sec. 3, we would also like the processing order of inputs belonging to different buckets to be as close to arrival order as possible. This is because reordering of outputs will lead to unnecessary blocking if processed too much out-of-order. In the rest of this section, we describe the concurrent data-structure and strategy we employ to achieve ordered partition parallelism. We first describe two simpler
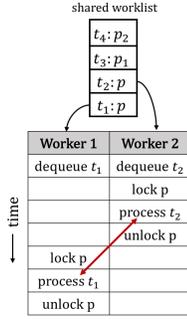
**Figure 5: Shared-Queue Approach: Each worker must dequeue the tuple and obtain a lock on the tuple's bucket atomically, otherwise a concurrent execution might violate the processing order constraint as shown above.**
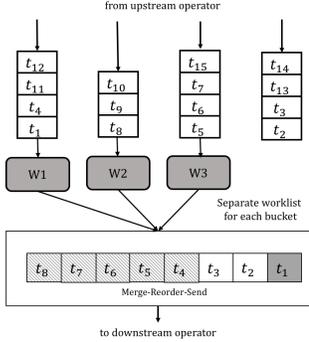


**Figure 6: Partitioned-Queue Approach: Since no worker is allotted to the last bucket which contains tuple $t_2, t_3$, the outputs of tuples $t_4$ to $t_{12}$ will get blocked from flowing to the downstream operator, limiting pipelined parallelism.**

strategies for implementing such an operator before presenting our approach.

## 4.1 Shared-Queue Approach.

In the first approach, which we refer to as the *shared-queue* approach, the producers (preceding operators in the dataflow graph) enqueue their outputs to a single queue (the worklist), and all concurrent workers extract tuples from the same queue and process them. This is a fairly straightforward strategy when the operator is stateless. We can use any linearizable concurrent queue to support multiple producers and consumers.

A partitioned operator, however, introduces a key challenge: we need to ensure that the items with same key are processed sequentially and in order. A naive approach would be as follows: Each worker first dequeues an item $t$ and then acquires a lock (or use any equivalent mechanism to ensure isolation) on the item's key so that two items with the same key are not processed concurrently. However, these two actions must be performed *atomically*: otherwise, two workers could concurrently dequeue items $t_1$ and $t_2$ with the same key $k$, but end up acquiring the lock on $k$ out-of-order and thus process them out-of-order as shown in fig. 5. This necessitates

quite complex and expensive concurrency control. Furthermore, this also introduces potentially blocking behavior when one worker waits for another, which is processing an input tuple with the same key. A naive implementation could aggravate this, causing all workers to be blocked, if a global lock is used to ensure the atomicity of the sequence of these two actions.

## 4.2 Partitioned-Queue Approach.

The second approach, which we refer to as the *partitioned-queue* approach avoids this problem. We use separate queues (worklist) for each bucket and the producers enqueue each tuple into the queue corresponding to the tuple's key. Different workers process different queues and hence there is no need for explicit concurrency control. However, this approach has its own set of drawbacks: Consider the scenario shown in figure 6, where the number of workers assigned to an operator is less than $p$ (number of buckets). In this case, the workers may make progress processing a subset of the $p$ queues. However, the outputs produced by these workers will be blocked by the reordering scheme that merges the outputs produced from the $p$ queues in order. This can cause further sub-optimal performance downstream as this behavior limits available pipelined parallelism between this and the downstream operator.

## 4.3 Hybrid-Queue Approach

We propose a hybrid approach that combines techniques from both these strategies. We use separate queues, one for each bucket as described above. In addition, we utilize a master queue which is analogous to the single queue of the former approach. Actual tuples are stored in individual bucket queues while, the master queue stores the *key* of each tuple. We list the execution model in fig. 7.

Every worker $w$ dequeues a key $k$ from the master queue, and then tries to gain exclusive access to the queue $Q_k$ that corresponds to $k$. If some other worker $w'$ already has exclusive access to queue $Q_k$, then worker $w$ delegates the responsibility of processing the corresponding tuple to $w'$, by incrementing a concurrent counter $count_k$ associated with the key $k$. The counter $count_k$ denotes the number of tuples from $Q_k$ to be processed before the active worker of key $k$ ($w'$ in this case) tries to dequeue the next key from master queue. The same counter is used to provide exclusive access to the queue $Q_k$. Having delegated the responsibility of processing the dequeued tuple to $w'$, worker $w$ can return to process the next key from the master queue.

If, on the other hand, worker $w$ gains exclusive access to queue $Q_k$, it dequeues the next tuple from $Q_k$ and processes it. However, after processing it, the worker needs to check if there are any delegated tuples that it needs to process from the same queue $Q_k$. As long as the concurrent counter $count_k$ indicates there are delegated items, the worker continues to dequeue tuples from $Q_k$ and processes them. When the counter becomes zero, the worker returns to processing the master queue. We prove the correctness of this scheme in the theorem below.

THEOREM 4.1 (CORRECTNESS OF HYBRID-QUEUE ALGORITHM). *If inputs to a partitioned stateful operator **o** are added using the* addInput *(fig. 7) procedure and workers allotted to **o**, consume inputs by invoking the* consumeInputs *procedure (fig. 7), then the following properties hold:*

```
1   //invoked by producers
2   void addInput(tuple) {
3     p = getPartition(tuple);
4     partitionQueues[p].enqueue(msg);
5     masterQueue.enqueue(p);
6   }
7   //invoked by workers
8   void consumeInputs() {
9     while(masterQueue.tryDequeue(p)) {
10      if(count[p].fetch_add(1) == 0) {
11        do {
12          partitionQueues[p].tryDequeue(tuple);
13          operate(tuple);
14        } while(count[p].fetch_sub(1) > 1) ;
15  } } }
```

**Figure 7: Hybrid Queue Approach: The `addInput` procedure is invoked by upstream operators and `consumeInputs` procedure is invoked by workers allotted to the partitioned stateful operator**

(1) No two workers can `operate` on tuples having the same key $k$ concurrently

(2) All tuples that have the same key $k$ are processed exactly once and in the order of their arrival

Proof. Any worker allotted to the operator first dequeues a partition $p$ from the master queue. The condition at L10 in fig. 7 ensures that a worker can obtain a tuple from the partition queue for $p$ only when value of count[$p$] (counter for $p$) is zero before the atomic increment. Now, to prove (1), it is enough to assert that the value of count[$p$] is never zero when a tuple belonging to $p$ is being actively processed by a worker. In the do-while loop (L11-14), the counter is decremented only after the dequeued tuple is processed completely. Note that the control flow in the addInput procedure ensures that tryDequeue at L12 always succeeds. Since the counter is decremented only at L14, it is clear that only the active worker of $p$ can reduce the value to zero, after which any other worker can enter L11-14. The atomic decrement and the condition at L14 ensures that the current active worker does not process any more tuples when count[$p$] becomes zero. Hence, at most one worker operates on tuples belonging to the same key.

Further, the FIFO guarantee of the linearizable concurrent queues in partitionQueues and the constraint that at most one worker can enter L11-14 for a particular $p$ (proved above) ensure that tuples belonging to the same key are processed exactly once and in order of their arrival into partitionQueues[$p$]. □

*Progress.* No worker can get blocked in the hybrid-queue approach. Adding inputs happen only by a single worker (of the operator upstream) due to the execution model employed in the reordering scheme. When consuming inputs, a worker that dequeues a tuple with same key as one being concurrently processed by another worker will simply delegate it to the active worker. So, this worker does not get blocked and moves on to the next key in the master queue. In this approach, outputs are also produced almost

in their arrival order, which avoids blocking of outputs (sometimes workers themselves) by the reordering scheme.

## 5 CORRECTNESS OF IMPLEMENTATION

In this section, we describe how we use the concurrent data structures described in Sec. 3 and 4 to implement the operators and prove that our implementation in combination with the runtime always guarantees ordered processing.

We start by defining the notion of correctness on the concurrent execution resulting from any implementation of the streaming computation.

*Definition 5.1 (Ordered execution).* A concurrent execution $E$ of a streaming computation on any input sequence $i_1, i_2, i_3, \ldots$ is *ordered*, if and only if, the output sequence $o_1, o_2, o_3, \ldots$ produced by the execution is the same output sequence produced by a sequential execution of the pipeline on $i_1, i_2, i_3, \ldots$.

We would like to ensure our implementation of the streaming computation is correct with respect to the above definition of ordering. An implementation of a streaming computation is said to be *ordered*, if and only if, any concurrent execution of the implementation is ordered.

There are three types of operators supported in our system: stateful, stateless and partitioned stateful operators. The implementation of a stateful operator is straight-forward. A worker of the upstream operator adds an input tuple to its worklist(a single-producer single-consumer concurrent queue). Only a single worker is allotted to a stateful operator at any time and this worker consumes these inputs serially and adds the corresponding outputs to the worklist of the downstream operator. A stateless operator is built using a shared-worklist (a multi-producer multi-consumer concurrent queue) and our non-blocking reordering buffer (Sec. 3). Input tuples are added to the shared-worklist and every tuple is allotted a unique serial number using an atomic counter. Worker(s) allotted to this stateless operator dequeue an input from this worklist and process it to produce output, which are then sent to the downstream operator by invoking the send method of the reordering buffer(fig. 4). If it fails, the worker tries again until it successfully adds the output to this buffer.

We implement the partitioned stateful operator by composing the hybrid partitioning scheme we described in Sec. 4 with our non-blocking reordering buffer. Inputs are allotted a unique increasing serial number in the order of their arrival and added by invoking the addInput method (fig. 7). Workers alloted to this operator consumes inputs using the consumeInputs method and invokes the send method of our reordering buffer (fig. 4) to send outputs downstream.

Theorem 5.2 (Correctness of pipeline implementation). *Any pipeline built by composing the above operator implementations and executed using our dynamic runtime only allows ordered executions.*

Proof. It is easy to see that the above theorem holds for a pipeline composed only of a single stateful operator. For a pipeline composed only of a stateless operator, even though $i_1, i_2, i_3 \ldots$ may be processed in any order, the corresponding output produced for $i_t$ is $o_t$ since the operator is stateless. The reordering buffer (sec. 3) guarantees that reordered sequence sent out is $o_1, o_2, o_3 \ldots$.

Now, let us consider a pipeline composed only of a partitioned stateful operator. For any two inputs $i_k$ and $i_l$ ($k < l$), if they belong to different keys, then irrespective of the order in which they are processed the corresponding outputs produced will be $o_k$ and $o_l$. When they have the same keys, the hybrid scheme guarantees that $i_k$ will be processed before $i_l$ and hence the outputs produced will be $o_k$ and $o_l$. So, the output produced for $i_t$ is $o_t$. Similar to a stateless operator pipeline, reordering scheme ensures that the output sequence produced is $o_1, o_2, o_3, \ldots$.

Since each of these single operator pipelines lead to correct executions, it is straightforward to see any linear composition of these operators will always lead to correct executions in our runtime. □

## 6 DYNAMIC SCHEDULING

Our system consists of many workers that consume inputs from the worklist(s) of an operator to produce outputs using the user-specified operator logic. The number of such workers is the same as number of cores available on the multicore machine. Each worker queries a central scheduler data structure to obtain some work and returns back for more, after it finishes the work allotted previously. The scheduler is responsible for answering two questions: (1) which operator to work on and (2) how many tuples to process from its worklist(s) before returning back. In this section, we propose some scheduling heuristics to perform this dynamic work allotment to worker threads.

We say an operator is *schedulable*, if the currently allotted number of workers is less than its maximum allowed degree of parallelism and its worklist is not empty. The theoretical maximum degree of parallelism of a stateful operator is 1, of a partitioned stateful operator is the number of partitions $p$, and that of a stateless operator is $\infty$ (essentially the number of available cores $n$). In all the heuristics we discuss below, we consider only those operators that are schedulable at the time we make the scheduling decision.

A worker when allotted to an operator, operates on it for a constant time slice $s$. The maximum number of tuples that must be processed by the worker can be computed using this constant $s$ and $c_i$, the cost of processing a single input tuple by $o_i$. If the worklist of an allotted operator becomes empty before processing the specified number of tuples, the worker does not get blocked; instead returns back to query the scheduler for more work. There are several alternatives for choosing the time for which an operator should be scheduled. However, we focus on constant time slices in order to study characteristics of the heuristics we propose without interference from these changes. Nevertheless, one has to be careful in choosing $s$. Higher the value of $s$, lower the contention for querying the scheduler and better amortization of scheduling overheads. On the other hand, a larger value of $s$ impedes the responsiveness of the system to dynamic changes, as it can get stuck on a previous scheduling decision for a long time.

We first propose some intuitive heuristics based on the idea of orchestrating the flow of tuples through a pipeline. There are two simple ways to enable this flow: one is to provide a thrust from ingress towards egress or use a suction pressure from egress to pull items from ingress. The following two heuristics are based on this key idea.

### 6.1 Queue-size-throttling (QST)

In this heuristic, we push tuples from the entry point towards the exit point and try to focus on one operator at a time. We schedule an operator until it generates enough inputs for the downstream operators and then go on to schedule the next one in the pipeline. We implement this scheme using queue throttling: each operator has an upper bound on its output queue (worklist of the downstream operator) size and is not scheduled if current size is higher than this threshold. In short, the heuristic always picks the earliest operator in the pipeline that has current output queue size less than its threshold.

Further, each operator $o_i$ has a selectivity, denoted by $s_i$, which is the average number of outputs produced by $o_i$ on processing a single input tuple. For example, selectivity is 1 for a `map` operator that maps each input tuple to a single output tuple, while it is less than 1 for a `filter` and more than 1 for `flat-map`, which maps a single input tuple to more than one output tuples. Due to difference in selectivities, having a uniform threshold for all operators could potentially create a slack in the pipeline. So, we set the output queue size threshold $T_i$ for an operator $o_i$ as follows, where $cs_i$ is the cumulative selectivity of operator $o_i$ since ingress ($cs_i = \prod_{k=1}^{i} s_k$) and $C$ is a constant that can be imagined as capacity of the system.

$$T_i = \frac{C * cs_i}{\sum_{i=1}^{n} cs_i} \tag{1}$$

Note that $T_i$ is proportional to the expected number of tuples produced by $o_i$ as input to $o_{i+1}$, when $\frac{C}{\sum_{i=1}^{n} cs_i}$ tuples are processed in the overall pipeline.

### 6.2 Last-in-pipeline (LP)

This heuristic is based on the complementary idea of pulling tuples from the exit point. In contrast to QST, this heuristic seeks to schedule operators later in the pipeline. Whenever an operator is not schedulable, this heuristic moves to its upstream operator and schedules that. This scheme depends entirely on the imminent dataflow between the operators and not on any of the operator characteristics. So, LP chooses the latest operator in the pipeline that has a non-empty input queue. An alternative could be to have a minimum worklist size, in which case only operators with worklist at least as big as this threshold would be considered for scheduling. But, in our empirical evaluation we consider only the simpler case where this threshold is 1.

The next set of heuristics take a slightly different approach to scheduling by *prioritizing* operators based on a certain measure of priority. This priority is computed using operator characteristics and current status of the pipeline. Essentially, these heuristics answer the question: which operator in the pipeline currently needs the most worker time to reach our performance goals? We discuss two heuristics designed using this strategy below.

### 6.3 Estimated-time (ET)

In this heuristic, we prioritize operators based on the estimated time it would take to process its current worklist, if we allot a new worker to it. We compute priority $p_i$ of an operator $o_i$, as follows, where $I_i$ denotes the current size of its worklist, $c_i$ denotes the cost of processing a single tuple by $o_i$, $w_i$ denotes the number of

workers currently assigned to $o_i$, and $M_i$ is its maximum allowed degree of parallelism:

$$p_i = \begin{cases} \frac{I_i * c_i}{w_i + 1} & \text{if } w_i < M_i \\ 0 & \text{otherwise} \end{cases} \qquad (2)$$

This strategy is based on the intuition that an operator that needs more worker time will lag behind and have a worklist that will take longer to complete.

## 6.4 Current-throughput (CT)

The key idea here is to choose the operator with the lowest throughput, as it is likely to be the *bottleneck* in the pipeline. We have to normalize the throughput to account for non-unit selectivities. We divide the time dimension into windows of size $w$ and compute the effective number of tuples processed by an operator in that time window as a measure of its throughput. The effective number of tuples $n_i^w$ that would be processed in the current window under current allocation of workers can be computed approximately as follows:

$$n_i^w = \frac{T_i^w + w_i * s}{c_i * cs_i} \qquad (3)$$

where, $T_i^w$ is the total worker time spent on $o_i$ in the current window $w$, $w_i$ is the number of workers alloted to $o_i$ currently and $s$ is the time slice for which each of the $w_i$ workers are allotted to $o_i$. CT chooses the operator with the lowest $n_i^w$ value. Another critical issue in the above heuristic is deciding on the window size $w$. It is possible for the scheduler to make sub-optimal decisions if the window size $w$ is too low. Ideally, we would like to use a window size that would have same $n_i^w$ for all the operators at the end of the window. This is similar to the *period* of a static schedule.

We evaluate these heuristics on real-world streaming queries from the TPCx-BB benchmark and discuss the pros and cons of choosing one over another in the next section.

## 7 EVALUATION

In this section, we present results of evaluation of the different scheduling heuristics and highlight benefits of our design of the parallelization framework for ordered stream processing empirically.

*Experimental Setup.* We perform all our experiments on Intel Xeon E5 family 2698B v3 series which runs the Windows Server 2012 R2 Datacenter operating system. It has 16 physical cores, with L1, L2 and L3 cache of size 32 KB, 256 KB and 40 MB. We implemented our research prototype in C++ on Windows using standard library implementations of concurrent queues and other atomic primitives. We measure throughput and latency by sending marker wrappers over tuples at equal tuple intervals, which carry information about entry and exit times. We ran all experiments for 2-10 mins and report the mean over 3 runs. For measurements, we consider only markers in the 20th to 80th percentile range, to eliminate starting up and shutting down interferences. Average throughput is computed by obtaining the ratio of number of tuples to the total time taken to process them and latency by averaging the processing latency of each marker in the range.

*Benchmark.* We use queries from the TPCx-BB benchmark[16], which is a modern Big Data benchmark that covers various categories of data analytics. We use all queries (Q1, Q2, Q3, Q4 and Q15) that correspond to stream processing workloads from TCxBB to compare our heuristics and evaluate various aspects of the runtime design. These queries and their implementation details are summarized in table 1. Web clickstreams are generated by every click made by a user on the online shopping portal and every item purchase in a retail store generates a store sales tuple.

## 7.1 Comparison of scheduling heuristics

We discussed several heuristics for dynamically scheduling operators in a stream processing pipeline in Sec. 6, namely normalized-current-throughput (CT), estimate-completion-time (ET), last-in-pipeline (LP) and queue-size-throttling (QST). We present results of experiments comparing their performance in terms of throughput and latency for the above queries in figures 8a and 8b respectively. We increase the number of cores until peak throughput of the best heuristic, beyond which the performance drops when the overheads of parallelization outweigh its benefits. We can use existing techniques in the literature [37] to identify this break-even point automatically. So, we do not focus on that aspect here.
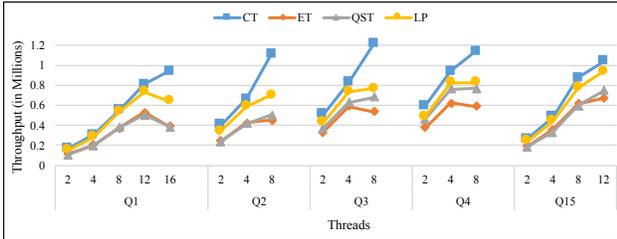
**Throughput.** We observe that heuristic CT scales almost linearly up to 16 cores for Q1, 12 cores for Q15 and 8 cores for Q2, Q3 and Q4. It achieves a peak throughput of approximately millions of web clickstreams and store sales tuples per second. We observe that this is the best possible throughput based on the per-tuple processing costs and selectivities of operators in the pipeline for corresponding degrees of parallelism. Among other heuristics, LP performs as well as CT for queries Q1 and Q15, but achieves sub-optimal performance for the others. Both ET and QST are observed to follow a similar trend in speedup achieved, however, they do not perform as well as CT or LP in terms of absolute throughput.

**Latency.** LP is the best heuristic for low-latency processing, followed closely by CT. It achieves latencies as low as a few milliseconds, which is the best known for stream processing systems. CT, which yields the best throughput, also processes tuples with such low latencies in many cases while it shoots up to 100s of milliseconds in some cases. Note that this is still quite low compared to other stream processing systems, which are based on batched stream processing [18, 41]. On the other hand, ET and QST have quite high latencies. This increase in latency for QST maybe due to a higher value of $M$ (refer Sec. 6), while ET is heavily influenced by the throughput of input stream to the overall pipeline.
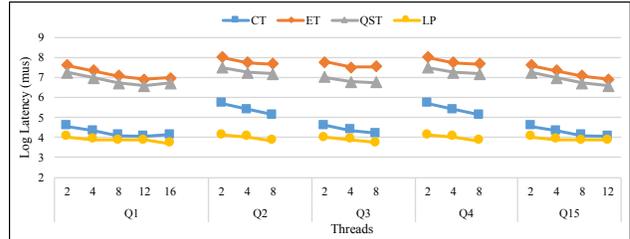
**Analysis.** From our analysis of the experimental results, we observe that there is a difference in performance among the heuristics even when their worker time distribution (ratio of total worker time spent on each operator in the pipeline) is almost similar. Heuristics that distribute workers across operators in the pipeline simultaneously tend to establish a continuous pipelined flow and are seen to yield much better throughput and latency. Those that focus on a single operator by exploiting maximum data parallelism at a time lead to increased per-tuple processing cost due to overheads at higher degrees of parallelism.

| | Pipeline | Brief Description |
|---|---|---|
| 1 | SS → SL → PS → PS → SF | Find top 100 pairs of items that are sold together frequently in the retail stores every hour |
| 2 | WC → SL → PS → SL → PS → SF | Find top 30 products that are viewed together online. Viewed together relates to a click-session of a user with session time-out of 60 mins |
| 3 | WC → SL → PS → PS | Find top 30 list of items (sorted by number of views) which are the last 5 products (in the past 10 days) that are mostly viewed before an item was purchased online |
| 4 | WC → SL → PS → SL → SF | Shopping cart abandonment analysis: For users who added products in their shopping cart but did not check out, find average number of pages they visited during their session |
| 15 | SS → SL → SL → PS | Find item categories with flat or declining sales for in-store purchases |

**Table 1: Summary of streaming queries in TPCx-BB. In the above table, WC = web clickstreams, SS = store sales, SL = stateless, PS = partitioned stateful and SF = stateful operator**



(a) Average Throughput



(b) Average latency

**Figure 8: Performance of the runtime when using different scheduling heuristics over TPCx-BB queries (We increase number of cores until peak throughput of the best heuristic - performance drops after that due to overheads of parallelization.)**

CT and LP seem to be exploiting this dichotomy quite efficiently. Choosing the operator with lowest estimated normalized through-put in the current window easily establishes this pipelined flow and hence uses an ideal combination of data, partitioned and pipeline parallelism. LP, that aims to always schedule operators later in the pipeline also establishes this continuous flow as follows: Initially, it is forced to schedule earlier operators in the pipeline as later ones are not schedulable; as they are scheduled it generates inputs for later operators and any worker that exits this operator is scheduled immediately on the next while some others are still processing the earlier operator. However, LP over-allots workers to operators later in the pipeline when they are schedulable which leads to sub-optimal performance in some queries above. The QST heuristic focuses on one operator at a time by design, similar to batched stream processing, thereby scheduling operators one-by-one along the pipeline. ET seems to be highly influenced by the input stream throughput as priority of the first operator depends on this. Hence, for a value of throughput higher than current system throughput, ET focuses mainly on the earliest operator and leads to sub-optimal performance as is evident from the results.

In the next two sub-sections, we discuss certain aspects of our parallelization framework that handles concurrent workers allotted to the same data or partitioned parallel operator. We designed *parametric operators* that can be used to create stateless and partitioned stateful operators with different computation profiles to help analyze their scalability in our framework. These operators are based on matrix computations on the input tuple. The per-tuple
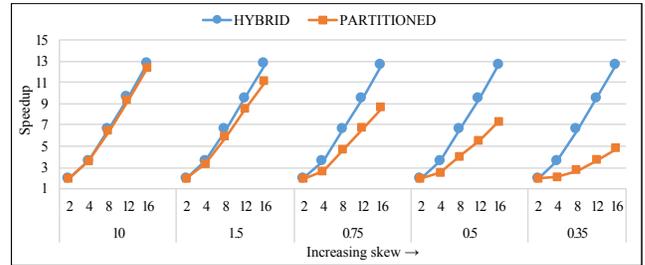


**Figure 9: Handling load imbalance across partitions**

processing cost, input tuple size, state size (for partitioned stateful) and selectivity can be varied by initializing these parametric operators with appropriate parameters.

## 7.2 Comparison of Partitioning Schemes

Now, we compare the two partitioning schemes we described in Sec. 4, PARTITIONED-QUEUE and HYBRID-QUEUE, that help achieve partitioned parallelism. Specifically, we compare their performance during load imbalance and in terms of latency with the constraint of ordered processing. Both schemes behave similarly in terms of per-operator throughput under uniform distribution, but hybrid scheme performs better in longer pipeline queries as it is more amenable to pipeline parallelism.

*7.2.1 Load Balancing.* Skewed distribution is known to highly limit partitioned parallelism. It is especially important to be able to balance load across workers in the stream processing setting
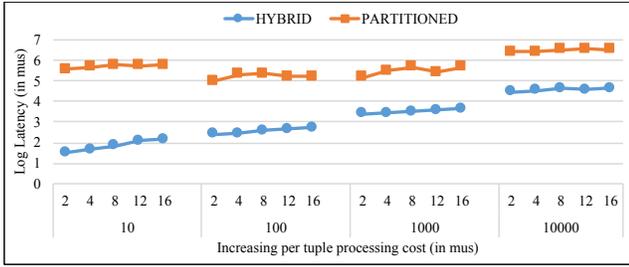
**Figure 10: Latency vs. Per-Tuple Processing Cost**



**Figure 11: Peak throughput (left) and latency (right) of the two partitioning schemes on TPCx-BB queries**

as they are expected to be long running continuous queries. In this experiment, we provide empirical evidence that the HYBRID-QUEUE approach can handle load imbalance much better than the PARTITIONED-QUEUE approach. In order to systematically induce skew in the distribution, we do range partitioning on keys sampled from a Gaussian distribution. We scale values in $[-1, 1]$ generated by $\mathcal{N}(0, \sigma)$ appropriately to fit the key space. We vary the value of $\sigma$ to vary the skew across partitions - higher the value of $\sigma$, closer the distribution is to a uniform distribution. The maximum number of partitions for PARTITIONED-QUEUE is limited to the number of workers, while the number of partitions in the HYBRID-QUEUE allows finer partitions and so is set to 100. The results of this experiment are presented in figure 9. We observe that both schemes perform similarly when the distribution is almost uniform. However, as we increase skew in the distribution, HYBRID-QUEUE performs consistently while scalability of the PARTITIONED-QUEUE approach drops heavily. This is because HYBRID-QUEUE admits finer partitions and hence leads to better load balancing.

*7.2.2 Latency.* We now compare the average processing latency in either schemes - the time between start of processing an input to the time at which its outputs exit the operator through the reordering scheme. We observe this for operators with various per-tuple processing costs (10, 100, 1000 and 10000 micro seconds) and a uniform distribution of tuples across partitions - the results are presented in figure 10. We can see that the average processing latency is much higher for PARTITIONED-QUEUE, while for HYBRID-QUEUE it is close to the corresponding operator's per-tuple processing cost. This is because the outputs produced through the PARTITIONED-QUEUE approach has to wait longer in the reordering buffer for outputs with earlier serial numbers. We do not report throughput comparisons between the two schemes here as both yield similar throughputs due to a uniformly random distribution of keys. However, this difference in individual operator processing latency leads to throughput differences in larger pipeline queries as we will see in the next experiment.

*7.2.3 Pipeline queries.* We compare performance of the two approaches on the TPCx-BB queries we described above. We use the CT scheduling heuristic, which yields the best performance among all the heuristics, and change only the partitioning scheme keeping the rest of the framework same. Peak throughput and latencies are reported in figure 11 as we vary the number of workers from 2 to 16. HYBRID-QUEUE is able to achieve much higher throughput than
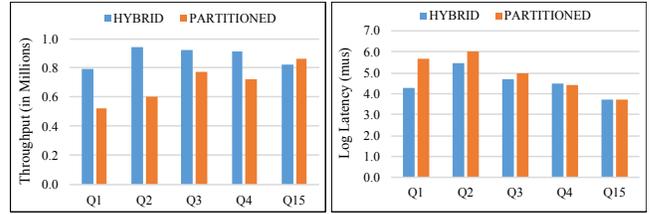
PARTITIONED-QUEUE in all queries. As expected, the difference is higher for queries that have more partitioned stateful operators. Query 15 contains only one such operator and is partitioned on item category id, where the total number of item categories in TPCx-BB is 10. It does not support higher degrees of parallelism and so the difference is unclear. HYBRID-QUEUE performs better than PARTITIONED-QUEUE also in terms of latency in 3 out of 5 queries, which have more partitioned stateful operators and almost similar for the rest.

## 7.3 Comparison of Reordering Schemes

We report the results of our empirical evaluation comparing the NON-BLOCKING scheme (fig. 4) and the LOCK-BASED scheme (fig. 2) in this subsection. We specifically highlight scenarios which are seen to be important in real-world queries from TPCx-BB using micro-benchmark experiments and also support it by evaluating them on the pipeline queries themselves.

*7.3.1 Light-weight Operators.* When the per-tuple processing cost of a stateless or partitioned-stateful operator is large and its computation profile is amenable to parallelization, the overhead of reordering outputs is relatively smaller and hence does not impede scalability of the operator. However, when this quantity is small, reordering could potentially become a huge bottleneck. We demonstrate that our NON-BLOCKING strategy minimizes this overhead leading to better scalability of such operators. We designed a stateless parametric operator with a per-tuple processing cost in the order of 10s of microseconds on a single core serial execution. Now, we varied the degree of parallelism of this operator and observed the increase in average per-tuple processing cost and the corresponding speedup achieved (fig. 12). Higher the reordering overhead, higher the average per-tuple processing cost and lower the speedup achieved. The results show that NON-BLOCKING reordering scheme scales better than the LOCK-BASED scheme. As expected, the average per-tuple processing cost of the operator, which includes the time for which a worker is blocked, increases more steeply for the LOCK-BASED strategy due to unnecessary blocking of workers when another worker is sending outputs downstream. This is avoided in our improved non-blocking design.

*7.3.2 High Selectivity Operators.* Similarly, when these operators have a huge selectivity (number of outputs per input), the amount of serial overhead involved in reordering is higher. In such cases, NON-BLOCKING strategy performs better in comparison to LOCK-BASED, even for operators with larger computation sizes . To illustrate this, we construct a pipeline that consists of two
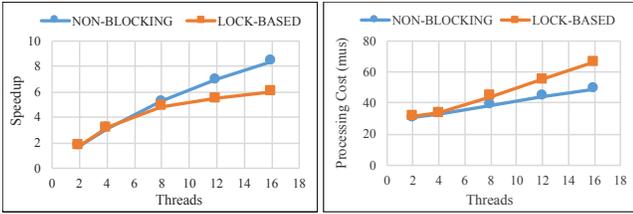
**Figure 12: (a) Speedup and (b) average processing cost for a light-weight operator**
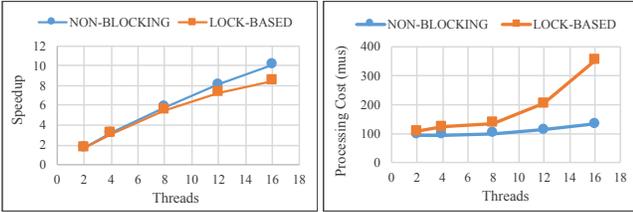


**Figure 13: (a) Speedup of the pipeline and (b) average processing cost of the first operator with LOCK-BASED and NON-BLOCKING reordering schemes**
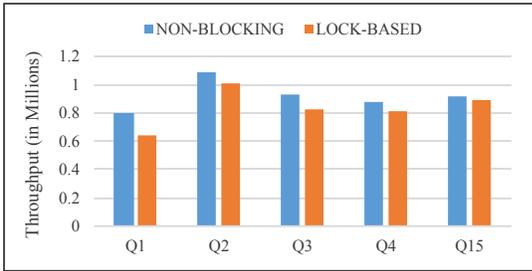


**Figure 14: Peak throughput comparison of the reordering schemes on TPCx-BB queries**

operators, a parametric stateless operator that is followed by a partitioned stateful operator. We use operators with a processing cost of approximately $100\mu s$ and the stateless operator has a selectivity of 50. Such high selectivity is not uncommon in real workloads. For example in query Q2, all clickstreams in a session are analyzed to produce a large number of item pairs viewed together. We report the average per-tuple processing cost and speedup achieved for this pipeline query in figure 13. Every tuple in the batch of outputs generated by the stateless operator has to be added into the appropriate queue of the partitioned stateful operator. To ensure ordering constraints, this operation is performed serially, which leads to blocking of workers in LOCK-BASED strategy, while in our scheme this is avoided.

*7.3.3 Pipeline Queries.* To further validate the benefits of our non-blocking reordering scheme, we compare it against the LOCK-BASED scheme on TPCx-BB queries. We report the peak throughput of the runtime for each of the queries using the best heuristic (CT) and by varying just the reordering scheme in fig. 14. We can clearly see that NON-BLOCKING scheme consistently yields a better throughput than the LOCK-BASED strategy. They do not differ much in processing latency and hence we do not present them here.

# 8 RELATED WORK

In this section, we review prior work related to concurrent data structures that we designed for ordered processing and scheduling of streaming computations.

## 8.1 Concurrent Data Structures

Our system uses several non-blocking concurrent data structures that have been proposed in the literature [25] such as single-producer single-consumer FIFO queues and multi-producer, multi-consumer queues. The reordering scheme we presented has very specific requirements (non-blocking, low-latency buffering), which are not directly met by any other data structure. The pre-allotted monotonically increasing serial numbers enabled further optimizations that would be inaccessible to a generic data structure such as concurrent priority queues.

Most partitioned parallelism implementations are based on the partitioned-queue approach we presented in Sec. 4, initially proposed in the Volcano [23] model of query evaluation for databases. In such a design, the degree of parallelism associated with the operator is determined statically and cannot be controlled by a dynamic scheduler. In case of shared-nothing architectures, some techniques [38] exist that adaptively repartitions the query during runtime. The trade-offs with respect to communication and repartitioning overheads are very different in a shared-memory architecture, so those techniques do not apply here directly. In addition, we address partition parallelism in the presence of ordering constraints, which to the best of our knowledge, none of the existing concurrent data structures address.

## 8.2 Static Scheduling

Static schedulers assume that the per-tuple processing cost and selectivity of the operators are known at compile time. Early streaming systems designed for applications from the digital signal processing domain focused on compiling down synchronous dataflow graphs (SDF), to single and multicores [12, 14]. For their application domain a purely static solution is not unreasonable as operator characteristics are largely fixed. There is a huge body of literature on scheduling SDF graphs to optimize various metrics such as throughput, memory and cache locality [4, 34, 36]. StreamIt [22], Brook [15] and Imagine [28] are some of the early systems designed based on this model of execution. However, none of these works address the case when operator characteristics change during runtime.

## 8.3 Dynamic Solutions

Aurora [2], its distributed counterpart, Borealis [1] and STREAM [11] are some of the early prototypes of stream processing engines that make dynamic scheduling decisions. Many recent stream processing engines (NaiagraST [31], Nile [24], Naiad [33],Spark Streaming[41], Storm [8], S4 [35]) also scheduling decisions during

runtime. All these systems either focus on single core or shared-nothing architectures. Even distributed solutions composed of individual shared-memory multicores consider each core as a separate executor and hence fail to exploit the advantages of a fast shared-memory.

IBM Streams [37] is one of the systems that target shared-memory architecture. Their runtime focuses on two issues: First, they design a mechanism to dynamically determine the maximum number of cores needed by a pipeline. This work is orthogonal to our work and can be easily adapted to our system. Second, they design a scalable scheduler that can schedule a large pipeline on a multicore; the focus is on scalability of the scheduler (number of scheduling decisions made) and not necessarily overall performance of the pipeline as they assume manual fine-tuning. Another key difference is that their scheduler works on an expanded pipeline, where each logical operator is duplicated a number of times specified through user annotations. This limits the flexibility of scheduler while also increasing the scheduling overhead. Their system also does not natively support totally ordered processing making a direct comparison infeasible.

Other systems such as TRILL [18] and Spark Streaming [41] are based on the micro-batch architecture. The idea is to execute a batch of inputs on an operator to completion before starting the next operator, thus relying primarily on the (data) parallelism within an operator. At any given time, a bulk of the workers are involved in executing instances of a single operator. Batching of streams is known to increase latency. We believe that such systems can be built on top of our parallelization and scheduling framework without much effort. We also note that several architectural proposals [3, 13, 21, 30] exist in the literature for a shared-memory streaming parallelization framework, but none of them address dynamic scheduling in the ordered setting or compare different scheduling heuristics empirically, which is a key contribution of this paper.

The approach we present in this paper is based on dynamic scheduling. Process/thread scheduling in operating systems is an example of this type of scheduling. We seek to develop a customized solution for the streaming setting taking advantage of the extra information available in the form of a dataflow graph. Further in a typical task graph, total amount of work to be done is fixed and the scheduler just needs to pick the right order once whereas in a streaming setting the scheduler has to continuously choose based on the status of pipeline. So, classical notions like work stealing do not apply to our setting [37].

## 9 CONCLUSIONS AND FUTURE WORK

We presented a new design for a dynamic runtime that executes streaming computations on a shared-memory multicore machine, along with the guarantee of ordered processing of tuples. We empirically demonstrated that our runtime is able to achieve good throughput (in millions of tuples per second) without compromising on the latency (a few milliseconds) on some TPCx-BB queries. We presented a couple of concurrent data structures that help achieve data and partitioned parallelism in the ordered setting, proved their correctness and showed their usefulness empirically using micro-benchmarks and on TPCx-BB queries.

In our current scheme, we assume all worker threads are uniform. However, in reality a worker is closer to some workers than others due to the hierarchical cache architecture and more so in modern non-uniform memory access (NUMA) architectures. An important extension to our work is to design scheduling heuristics that discriminate workers based on their spatial distribution.

## REFERENCES

[1] Daniel J Abadi, Yanif Ahmad, Magdalena Balazinska, Ugur Cetintemel, Mitch Cherniack, Jeong-Hyon Hwang, Wolfgang Lindner, Anurag S Maskey, Alexander Rasin, Esther Ryvkina, Nesime Tatbul, Ying Xing, and Stan Zdonik. 2005. The Design of the Borealis Stream Processing Engine. In *Second Biennial Conference on Innovative Data Systems Research (CIDR 2005)*. Asilomar, CA.
[2] Daniel J. Abadi, Donald Carney, Ugur Çetintemel, Mitch Cherniack, Christian Convey, Sangdon Lee, Michael Stonebraker, Nesime Tatbul, and Stanley B. Zdonik. 2003. Aurora: a new model and architecture for data stream management. *VLDB J.* 12, 2 (2003), 120–139. https://doi.org/10.1007/s00778-003-0095-z
[3] Kunal Agrawal, Jeremy T. Fineman, Jordan Krage, Charles E. Leiserson, and Sivan Toledo. 2012. Cache-conscious Scheduling of Streaming Applications. In *Proceedings of the Twenty-fourth Annual ACM Symposium on Parallelism in Algorithms and Architectures (SPAA '12)*. ACM, New York, NY, USA, 236–245. https://doi.org/10.1145/2312005.2312049
[4] Waheed Ahmad, Robert de Groote, Philip K. F. Hölzenspies, Mariëlle Stoelinga, and Jaco van de Pol. 2014. Resource-Constrained Optimal Scheduling of Synchronous Dataflow Graphs via Timed Automata. In *Proceedings of the 2014 14th International Conference on Application of Concurrency to System Design (ACSD '14)*. IEEE Computer Society, Washington, DC, USA, 72–81. https://doi.org/10.1109/ACSD.2014.13
[5] Tyler Akidau, Alex Balikov, Kaya Bekiroglu, Slava Chernyak, Josh Haberman, Reuven Lax, Sam McVeety, Daniel Mills, Paul Nordstrom, and Sam Whittle. 2013. MillWheel: Fault-Tolerant Stream Processing at Internet Scale. In *Very Large Data Bases*. 734–746.
[6] Tyler Akidau, Robert Bradshaw, Craig Chambers, Slava Chernyak, Rafael J. FernÃ¡ndez-Moctezuma, Reuven Lax, Sam McVeety, Daniel Mills, Frances Perry, Eric Schmidt, and Sam Whittle. 2015. The Dataflow Model: A Practical Approach to Balancing Correctness, Latency, and Cost in Massive-Scale, Unbounded, Out-of-Order Data Processing. *Proceedings of the VLDB Endowment* 8 (2015), 1792–1803.
[7] Apache. 2016. Apache Samza. http://samza.apache.org/. (2016). Accessed: 2017-01-16.
[8] Apache. 2016. Apache Storm. http://storm.apache.org/. (2016). Accessed: 2017-01-16.
[9] Arvind Arasu, Shivnath Babu, and Jennifer Widom. 2006. The CQL Continuous Query Language: Semantic Foundations and Query Execution. *The VLDB Journal* 15, 2 (June 2006), 121–142. https://doi.org/10.1007/s00778-004-0147-z
[10] Brian Babcock, Shivnath Babu, Rajeev Motwani, and Mayur Datar. 2003. Chain: Operator Scheduling for Memory Minimization in Data Stream Systems. In *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data (SIGMOD '03)*. ACM, New York, NY, USA, 253–264. https://doi.org/10.1145/872757.872789
[11] Shivnath Babu and Jennifer Widom. 2001. Continuous Queries over Data Streams. *SIGMOD Rec.* 30, 3 (Sept. 2001), 109–120. https://doi.org/10.1145/603867.603884
[12] Shuvra S. Battacharyya, Edward A. Lee, and Praveen K. Murthy. 1996. *Software Synthesis from Dataflow Graphs*. Kluwer Academic Publishers, Norwell, MA, USA.
[13] Jonathan C. Beard, Peng Li, and Roger D. Chamberlain. 2015. RaftLib: A C++ Template Library for High Performance Stream Parallel Processing. In *Proceedings of the Sixth International Workshop on Programming Models and Applications for Multicores and Manycores (PMAM '15)*. ACM, New York, NY, USA, 96–105. https://doi.org/10.1145/2712386.2712400
[14] Shuvra S. Bhattacharyya, Praveen K. Murthy, and Edward A. Lee. 1999. Synthesis of Embedded Software from Synchronous Dataflow Specifications. *J. VLSI Signal Process.* 21, 2 (June 1999), 151–166. https://doi.org/10.1023/A:1008052046456
[15] Ian Buck, Tim Foley, Daniel Horn, Jeremy Sugerman, Kayvon Fatahalian, Mike Houston, and Pat Hanrahan. 2004. Brook for GPUs: Stream Computing on Graphics Hardware. In *ACM SIGGRAPH 2004 Papers (SIGGRAPH '04)*. ACM, New York, NY, USA, 777–786. https://doi.org/10.1145/1186562.1015800
[16] Paul Cao, Bhaskar Gowda, Seetha Lakshmi, Chinmayi Narasimhadevara, Patrick Nguyen, John Poelman, Meikel Poess, and Tilmann Rabl. 2017. *From BigBench to TPCx-BB: Standardization of a Big Data Benchmark*. Springer International Publishing, Cham, 24–44. https://doi.org/10.1007/978-3-319-54334-5_3
[17] Don Carney, Uğur Çetintemel, Alex Rasin, Stan Zdonik, Mitch Cherniack, and Mike Stonebraker. 2003. Operator Scheduling in a Data Stream Manager. In *Proceedings of the 29th International Conference on Very Large Data Bases - Volume 29 (VLDB '03)*. VLDB Endowment, 838–849. http://dl.acm.org/citation.cfm?id=1315451.1315523

[18] Badrish Chandramouli, Jonathan Goldstein, Mike Barnett, Robert DeLine, John C. Platt, James F. Terwilliger, and John Wernsing. 2014. Trill: A High-Performance Incremental Query Processor for Diverse Analytics. *PVLDB* 8, 4 (2014), 401–412. http://www.vldb.org/pvldb/vol8/p401-chandramouli.pdf

[19] Sirish Chandrasekaran, Owen Cooper, Amol Deshpande, Michael J. Franklin, Joseph M. Hellerstein, Wei Hong, Sailesh Krishnamurthy, Samuel R. Madden, Fred Reiss, and Mehul A. Shah. 2003. TelegraphCQ: Continuous Dataflow Processing. In *Proceedings of the 2003 ACM SIGMOD International Conference on Management of Data (SIGMOD '03)*. ACM, New York, NY, USA, 668–668. https://doi.org/10.1145/872757.872857

[20] Jeffrey Dean and Sanjay Ghemawat. 2008. MapReduce: Simplified Data Processing on Large Clusters. *Commun. ACM* 51, 1 (Jan. 2008), 107–113. https://doi.org/10.1145/1327452.1327492

[21] Mark A. Franklin, Eric J. Tyson, James Buckley, Patrick Crowley, and John Maschmeyer. 2006. Auto-pipe and the X Language: A Pipeline Design Tool and Description Language. In *Proceedings of the 20th International Conference on Parallel and Distributed Processing (IPDPS'06)*. IEEE Computer Society, Washington, DC, USA, 117–117. http://dl.acm.org/citation.cfm?id=1898953.1899049

[22] Michael I. Gordon, William Thies, Michal Karczmarek, Jasper Lin, Ali S. Meli, Andrew A. Lamb, Chris Leger, Jeremy Wong, Henry Hoffmann, David Maze, and Saman P. Amarasinghe. 2002. A stream compiler for communication-exposed architectures. In *Proceedings of the 10th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS-X), San Jose, California, USA, October 5-9, 2002*. 291–303. https://doi.org/10.1145/605397.605428

[23] G. Graefe. 1994. Volcano&#151 An Extensible and Parallel Query Evaluation System. *IEEE Trans. on Knowl. and Data Eng.* 6, 1 (Feb. 1994), 120–135. https://doi.org/10.1109/69.273032

[24] M. A. Hammad, M. F. Mokbel, M. H. Ali, W. G. Aref, A. C. Catlin, A. K. Elmagarmid, M. Eltabakh, M. G. Elfeky, T. M. Ghanem, R. Gwadera, I. F. Ilyas, M. Marzouk, and X. Xiong. 2004. Nile: a query processing engine for data streams. In *Proceedings. 20th International Conference on Data Engineering*. 851–. https://doi.org/10.1109/ICDE.2004.1320080

[25] Maurice Herlihy and Nir Shavit. 2008. *The art of multiprocessor programming*. Morgan Kaufmann.

[26] Martin Hirzel, Robert Soulé, Scott Schneider, BuÄ§ra Gedik, and Robert Grimm. 2014. A catalog of stream processing optimizations. *Comput. Surveys* 46, 4 (mar 2014), 1–34. https://doi.org/10.1145/2528412

[27] Qingchun Jiang and Sharma Chakravarthy. 2004. *Scheduling Strategies for Processing Continuous Queries over Streams*. Springer Berlin Heidelberg, Berlin, Heidelberg, 16–30. https://doi.org/10.1007/978-3-540-27811-5_3

[28] Ujval Kapasi, William J. Dally, Scott Rixner, John D. Owens, and Brucek Khailany. 2002. The Imagine Stream Processor. In *Proceedings 2002 IEEE International Conference on Computer Design*. 282–288.

[29] Sanjeev Kulkarni, Nikunj Bhagat, Maosong Fu, Vikas Kedigehalli, Christopher Kellogg, Sailesh Mittal, Jignesh M. Patel, Karthik Ramasamy, and Siddarth Taneja. 2015. Twitter Heron: Stream Processing at Scale. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data (SIGMOD '15)*. ACM, New York, NY, USA, 239–250. https://doi.org/10.1145/2723372.2742788

[30] I-Ting Angelina Lee, Charles E. Leiserson, Tao B. Schardl, Zhunping Zhang, and Jim Sukha. 2015. On-the-Fly Pipeline Parallelism. *ACM Transactions on Parallel Computing* 2, 3, Article 17 (Sept. 2015), 42 pages. https://doi.org/10.1145/2809808

[31] David Maier, Jin Li, Peter Tucker, Kristin Tufte, and Vassilis Papadimos. 2005. Semantics of Data Streams and Operators. In *Proceedings of the 10th International Conference on Database Theory (ICDT'05)*. Springer-Verlag, Berlin, Heidelberg, 37–52. https://doi.org/10.1007/978-3-540-30570-5_3

[32] Microsoft. 2016. Microsoft StreamInsight. https://msdn.microsoft.com/en-us/library/ee362541(v=sql.111).aspx. (2016). Accessed: 2017-01-16.

[33] Derek G. Murray, Frank McSherry, Rebecca Isaacs, Michael Isard, Paul Barham, and Martín Abadi. 2013. Naiad: A Timely Dataflow System. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles (SOSP '13)*. ACM, New York, NY, USA, 439–455. https://doi.org/10.1145/2517349.2522738

[34] P. K. Murthy, S. S. Bhattacharyya, and E. A. Lee. 1994. Minimizing memory requirements for chain-structured synchronous dataflow programs. In *Acoustics, Speech, and Signal Processing, 1994. ICASSP-94., 1994 IEEE International Conference on*, Vol. ii. II/453–II/456 vol.2. https://doi.org/10.1109/ICASSP.1994.389625

[35] Leonardo Neumeyer, Bruce Robbins, Anish Nair, and Anand Kesari. 2010. S4: Distributed Stream Computing Platform. In *ICDMW 2010, The 10th IEEE International Conference on Data Mining Workshops, Sydney, Australia, 13 December 2010*. 170–177. https://doi.org/10.1109/ICDMW.2010.172

[36] Jose L Pino, Shuvra S. Bhattacharyya, and Edward A. Lee. 1999. *A Hierarchical Multiprocessor Scheduling Framework for Synchronous Dataflow Graphs*. Technical Report. Berkeley, CA, USA.

[37] Scott Schneider and Kun-Lung Wu. 2017. Low-synchronization, Mostly Lock-free, Elastic Scheduling for Streaming Runtimes. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI 2017)*. ACM, New York, NY, USA, 648–661. https://doi.org/10.1145/3062341.3062366

[38] Mehul A. Shah, Joseph M. Hellerstein, Sirish Chandrasekaran, and Michael J. Franklin. 2003. Flux: An Adaptive Partitioning Operator for Continuous Query Systems. In *Proceedings of the 19th International Conference on Data Engineering, March 5-8, 2003, Bangalore, India*. 25–36. https://doi.org/10.1109/ICDE.2003.1260779

[39] William Thies. 2009. *[StreamIt] Language and Compiler Support for Stream Programs*. Ph.D. Dissertation.

[40] Matei Zaharia, Mosharaf Chowdhury, Michael J. Franklin, Scott Shenker, and Ion Stoica. 2010. Spark: Cluster Computing with Working Sets. In *Proceedings of the 2Nd USENIX Conference on Hot Topics in Cloud Computing (HotCloud'10)*. USENIX Association, Berkeley, CA, USA, 10–10. http://dl.acm.org/citation.cfm?id=1863103.1863113

[41] Matei Zaharia, Tathagata Das, Haoyuan Li, Scott Shenker, and Ion Stoica. 2012. Discretized Streams: An Efficient and Fault-tolerant Model for Stream Processing on Large Clusters. In *Proceedings of the 4th USENIX Conference on Hot Topics in Cloud Ccomputing (HotCloud'12)*. USENIX Association, Berkeley, CA, USA, 10–10. http://dl.acm.org/citation.cfm?id=2342763.2342773