

# Scaling Stream Processing with Transactional State Management on Multicores

Shuhao Zhang<sup>1</sup>, Yingjun Wu<sup>2</sup>, Feng Zhang<sup>3</sup>, Bingsheng He<sup>1</sup>

<sup>1</sup>National University of Singapore, <sup>2</sup>IBM Almaden Research Center, <sup>3</sup>Renmin University of China

## ABSTRACT

Data stream processing system (DSPS) with transactional state management relieves users from managing state consistency by themselves. This paper introduces **TStream**, a highly scalable DSPS with built-in transactional state management. **TStream** is specifically designed for modern shared-memory multicore architectures. **TStream**'s key contribution is a novel asynchronous state management paradigm. By detaching and postponing state accesses from the stream application logic, **TStream** minimizes stream computation stalling caused by state access, which is commonly found in existing approaches. The postponed state accesses naturally form a batch, and we further propose an operation-chain based execution model that allows **TStream** to aggressively extract parallelism opportunities within each batch of state accesses while guaranteeing transactional state consistency. To confirm the effectiveness of our proposal, we compared **TStream** against four alternative designs on a 40-core machine. Our extensive experiment study show that **TStream** yields much higher throughput with limited latency penalty comparing to existing approaches.

## PVLDB Reference Format:

Shuhao Zhang, Yingjun Wu, Feng Zhang, Bingsheng He. Scaling Stream Processing with Transactional State Management on Multicores. *PVLDB*, 12(xxx): xxxx-yyyy, 2019.  
DOI: <https://doi.org/10.14778/xxxxxxx.xxxxxxx>

## 1. INTRODUCTION

Data stream processing systems (DSPSs) are gaining their popularities in powering modern IoT (Internet-of-Things) and data streaming applications [15, 43, 21, 1, 3, 29]. Thanks to the popularity of modern commodity machines with massively parallel processors, researchers and practitioners nowadays can perform ultra-fast stream processing on a single multicore machine [47, 28, 33, 48, 45]. To support complex real-time analytics, several prior studies [39, 14, 32, 7] have pointed out that stream

applications often need to manage (read/update) large *shared mutable states* (e.g., user account) during stream processing. As a result, any uncoordinated accesses to the same state (e.g., read and update the same account at the same time) can cause state inconsistencies.

The usage of shared mutable states in stream processing are prevalent. To relieve application developers from managing state consistency by themselves, DSPS with transactional state management has recently received many attentions from both academia [32, 7] and industry community [4]. Scaling stream processing while providing transactional state management is challenging. To achieve both low latency and high throughput, DSPSs are designed to process multiple input events concurrently [47, 48, 33, 28]. However, parallel stream processing may lead to conflict accesses to the same application state, hence leading to higher chances of violating transactional state consistency. Simply relying on a third-party DBMS for managing stream application states not only involves high inter-process communication overhead degrading stream processing performance but also may still leads to consistency violation [16] (see Section 2).

Due to the complexity, today's DSPSs often force user to select from either strict state consistency guarantee with limited concurrency or high concurrency with limited state consistency guarantee (e.g., Storm [3], Flink [1]). Recent studies on transactional state management in DSPSs [32, 16, 14, 39, 7] typically choose the prior one. The limited scalability of prior works comes from mainly two folds (detailed in Section 3). Firstly, they are commonly based on a synchronized execution model, where stream operator (carried by a thread) continuously performs states read/update (triggered by processing of one input event) along with stream processing. Despite being a straightforward implementation, each operator has to wait for its current state accesses to finish before processing on more input events. This can cause extra long execution stalls and reduce stream processing concurrency, consequently degrade the overall system throughput. Secondly, in order to guarantee state consistency, they heavily rely on synchronization primitives (e.g., locks need to be strictly inserted in the order of event timestamp [39]) in performing shared state read/update. As a result, parallelism opportunities among state accesses execution are overlooked preventing prior systems from scaling well.

Witnessing limitations of prior works, we propose **TStream**, a new DSPS that can support highly scalable stream processing with transactional state consistency

This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>. For any use beyond those covered by this license, obtain permission by emailing [info@vldb.org](mailto:info@vldb.org). Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.

*Proceedings of the VLDB Endowment*, Vol. 12, No. xxx  
ISSN 2150-8097.

DOI: <https://doi.org/10.14778/xxxxxxx.xxxxxxx>

Table 1: Summary of Terminologies

Term	Definition
Event	Input stream event with a monotonically increasing timestamp
Watermark	A special event guaranteeing no future event has smaller timestamp than it
Operator	Basic unit of a stream application, continuously process event streams
Execution instance	Each operator can be carried by multiple execution instance (i.e., threads)
Shared states	Mutable states that are concurrently accessible by multiple execution instance of an operator
State transaction	A set of state access operations triggered by a single input event
ACID+O	A consistency property satisfying ACID properties and event timestamp ordering

guarantee. Targeting on shared-memory multicore architectures, **TStream** achieves high performance via two key designs (detailed in Section 4):

(1) **TStream** performs its internal state management *asynchronously* by detaching state access operations from the application computation logic. Through postponing state accesses, **TStream** minimizes unnecessary stalls caused by state management in stream processing. **TStream** adopts *watermark* as periodic signal to trigger the actual state accesses. Postponed state accesses between two subsequent watermarks hence naturally form a batch, and the batch length can be then tuned to tradeoff latency and throughput.

(2) To take advantage of modern multicore architectures, **TStream** adopts an *operation-chain* based execution model that aggressively extracts parallelism opportunities within each batch of state transactions. We further apply evaluation pushdown and NUMA-aware placement to improve its processing efficiency.

To confirm **TStream**'s effectiveness, we compared it against four alternative schemes on a high-performance 40-core machine. Our extensive experimental study shows that **TStream** achieves more than two times higher throughput on average over existing solutions with similar or even smaller processing latency.

We organize the paper as follows: Section 2 introduces the computational model of DSPS with transactional state management. Section 3 reviews limitations of existing solutions. Section 4 discusses design of **TStream**. Additional implementation details are described in Section 5. We report extensive experiment results in Section 6. Section 7 reviews related works and Section 8 concludes this paper.

## 2. STREAM MODEL AND BACKGROUND

In this section, we discuss the stream processing model and terminology through a running example. We summarize terminologies in Table 1.

We describe the execution model of stream processing with a general definition [33]. We briefly present here for completeness. Stream processing continuously processes on one or more streams<sup>1</sup> of *events*  $E$ . Each event  $e_{ts} \in E$  has a timestamp  $ts$  that indicates its temporal sequence. For simplicity, we assume events arrived at the system has a monotonically increasing timestamp. A

<sup>1</sup>In this paper, we only consider single-input stream operator, and defer multi-streams operator as future work.

*watermark* [33] is a special input event guaranteeing that all later input events must have larger timestamp than it. A streaming application contains a sequence of *operators*  $O = O_1, O_2, \dots, O_n$ . Each operator continuously processes events and optionally emits events. To sustain high input stream ingress rate, each operator can be further carried in multiple *execution instances* (e.g., Java threads), which handle multiple input events concurrently.

**Transactional state management.** We refer the management of *shared mutable state* while guaranteeing their consistency as transactional state management. The usage of shared mutable states in stream processing are prevalent and has been studied in several prior works [39, 14, 32, 7]. We use a simplified stream application as depicted in Figure 1 as a running example. The application represents a simplified bidding system [37], where users bid for items online. As multiple users may compete for the same item, the system needs to ensure the request made earlier be processed earlier.

The application can be implemented with the following four operators. **Parser** continuously emits events describing requests from either buyers or sellers. **Auth.** authenticates the requests such as validating the request's ip address and dispatches valid requests for further process. **Process** handles valid requests. There are three types of requests. (1) *bid request* for one item from buyers, if the bid hits the ask price for an item, which has sufficient quantity available, the request is processed (i.e., reduce the quantity of the item), otherwise rejected. (2) *alter request* modify prices of a list of items. (3) *top request* tops up quantity of a list of items. We use Sink to record the output stream from **Process** to monitor system performance. We discuss configurations such as operator selectivity later in Section 6.

During each request processing, **Process** may need to read/update item table, which contains item id, price and quantity information. The table has to be shared among all execution instances of **Process** as each request may touch any items – it is difficult to partition the table and store it as private state in each instance. As a result, any uncoordinated accesses (i.e., update and read) to the same record can cause state inconsistencies. This problem is exacerbated if more complex state storage and retrieval queries such as scan and range lookup is required, which is not supported in most existing DSPSs, such as Storm [3], Flink [1], and has recently motivated a number of systems [39, 14, 7].

**DEFINITION 1.** We define the set of state access operations triggered by a single input event  $e_{ts}$  in an operator as one *state transaction*  $txn_{ts}$ , due to its similarity to conventional database transaction. Timestamp  $ts$  of a state transaction is assigned to be the same of its triggering event.

As multiple state transactions may be concurrently generated from multiple execution instances of an operator, their execution must also satisfy ACID properties. Formally,

- *Atomicity*: state access operations of one state transaction shall be processed all at once or none.
- *Consistency*: the state changes as a result of evaluating state transaction shall not violate any state constraint. For example, user account can not become negative.
- *Isolation*: multiple conflicting state transactions shall run independently without interfering each other.

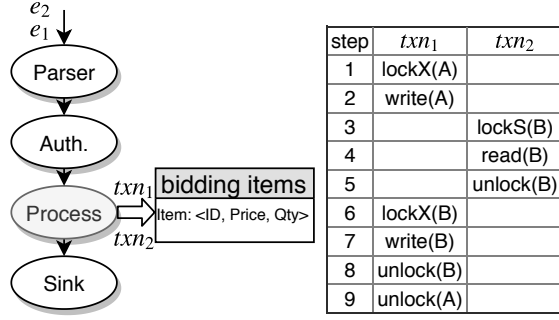


Figure 1: Running example Figure 2: Example application.

step	$txn_1$	$txn_2$
1	lockX(A)	
2	write(A)	
3		lockS(B)
4		read(B)
5		unlock(B)
6	lockX(B)	
7	write(B)	
8	unlock(B)	
9	unlock(A)	

- **Durability:** the changes made to the state shall be made persistent.<sup>2</sup>

Although this is a pretty common problem in traditional database system, we are targeting at a streaming context, where state accesses are triggered by input streaming events rather than user. Naively locking the state or relying on a third-party database to manage states not only degrades the performance but also leads to consistency violation [16] because of the additional requirement of guaranteeing event timestamp ordering, e.g., bidding request with smaller timestamp should be processed earlier.

**DEFINITION 2.** Following previous works [32, 7], we define a consistency property satisfying both *ACID* properties and event timestamp ordering constraint as *ACID+O* properties. Formally, a *DSPS* under *ACID+O* needs to ensure the state transaction schedule must be conflict equivalent to  $txn_1 \prec \dots \prec txn_n$ .

Unfortunately, such state access ordering guarantee is not well-supported in conventional database systems as also pointed out by a number of previous works [39, 14, 32, 7]. We use an example execution trace to illustrate the failure of simply using a conventional database to manage item table in our running example. Let us assume that two events  $e_1$  and  $e_2$  ( $e_1$  has a smaller timestamp) are concurrently processed by two instances of **Process** independently. Assume the processing of  $e_1$  triggers two updates of records of item  $A$  and  $B$ , i.e.,  $txn_1:W(A), W(B)$ , while  $e_2$  triggers a read of  $B$ , i.e.,  $txn_2:R(B)$ . System under *ACID+O* properties needs to ensure the processing of  $txn_2$  successfully obtain the correct updated value of  $B$  due to  $txn_1$ . However, conventional CC protocols serialize transactions in an order that is conflict-equivalent to *any* certain serial schedule. Taking 2-phase locking (2PL) [10] as an example, a possible schedule is shown in Figure 2. The resulting serial order is  $txn_2 \prec txn_1$ , which leads to a wrong reading of  $B$  in processing  $e_2$ .

### 3. EXISTING SOLUTIONS AND LIMITATIONS

Scaling stream processing with transactional state management is challenging and has motivated a number of

related systems [39, 7, 14, 4, 32]. The fundamental challenge is that the state access operations triggered by one event is unaware of the access pattern of operations concurrently triggered by other events.

Processing each input event serially following their timestamp order (e.g., SEI [39]) eliminates the challenge of preserving state consistency. For example, one can restrict the **Process** operator to have only one execution instance (e.g., single thread). However, it severely restricts system concurrency. To achieve both low latency and high throughput, DSPSs need to exploit parallelism aggressively, and operators are typically carried by multiple instances. However, there is a *potential conflict* between those concurrently running instances as they may access the same shared state at the same time. As a result, higher concurrency also intensifies the state access collision, and consistency needs to be guaranteed.

Prior solutions [39, 14, 32] can still confront severe scaling bottlenecks, primarily due to the requirement of tracking order of each state transaction. Figure 3 compares how different algorithms process three concurrently generated state transactions. In addition to  $txn_1$  and  $txn_2$  of our running example shown previously, let  $txn_3$  contain a read access to record of item  $C$ , that is  $txn_3:R(C)$ .

**Lock Ahead 2PL (LAL).** An earlier work from Wang et al. [39] described a strict two-phase locking (S2PL) based algorithm that allows multiple state transactions run concurrently while maintaining state consistency. Extended from the original S2PL [10], it relies on a *lock-ahead* mechanism to make sure that locks on each state are granted in the correct order (see the **Lock(order)** in Figure 3(a)) to preserve the aforementioned order-preserving property. As the lock-ahead operations from different transactions have to be performed serially, LAL severely restrict system concurrency.

**Low-Water-Mark (LWM).** Wang et al. [39] propose an improved version of LAL by adding multi-versioning, called *low water mark (LWM)*. LWM leverages a status variable (i.e., **lwm**) of each state to guard the processing sequence: write must be performed monotonically (by sequentially increasing **lwm** during transaction commits), and will not be blocked by read locks; a read may not be blocked by write locks as long as it is able to read the correct version of shared states. That is, the read operation has a timestamp larger than **lwm** of the corresponding state so it can read an *readily* version of state. As shown in Figure 3(b), **LockS(B)** will not be blocked by  $txn_1$ 's lock on  $B$ , which is its major improvement compared to LAL. However, **Read(B)** has to wait for **lwm** to be updated by  $txn_1$  during its commit. Furthermore, LWM still requires the **Lock(order)** mechanism to ensure locks are inserted at a correct order following event sequence.

**Partition-based approach (PAT).** S-Store [32] fuses OLTP and streaming into one engine, which is built based on H-Store – a shared-nothing deterministic database system [27]. Similar to how H-Store splits database, S-Store splits the shared states into multiple disjoint *partitions*, and hence only needs to guard accessing order in each partition. Unfortunately, the original implementation of S-Store uses a single-thread [32] for all shared state accesses, which severely limits system concurrency. We hence implement a mechanism called *PAT* to represent the similar idea of S-Store without such restriction. Despite being partitioned,

<sup>2</sup>In this paper, we assume shared states can be kept in memory and we defer the durability in future work.

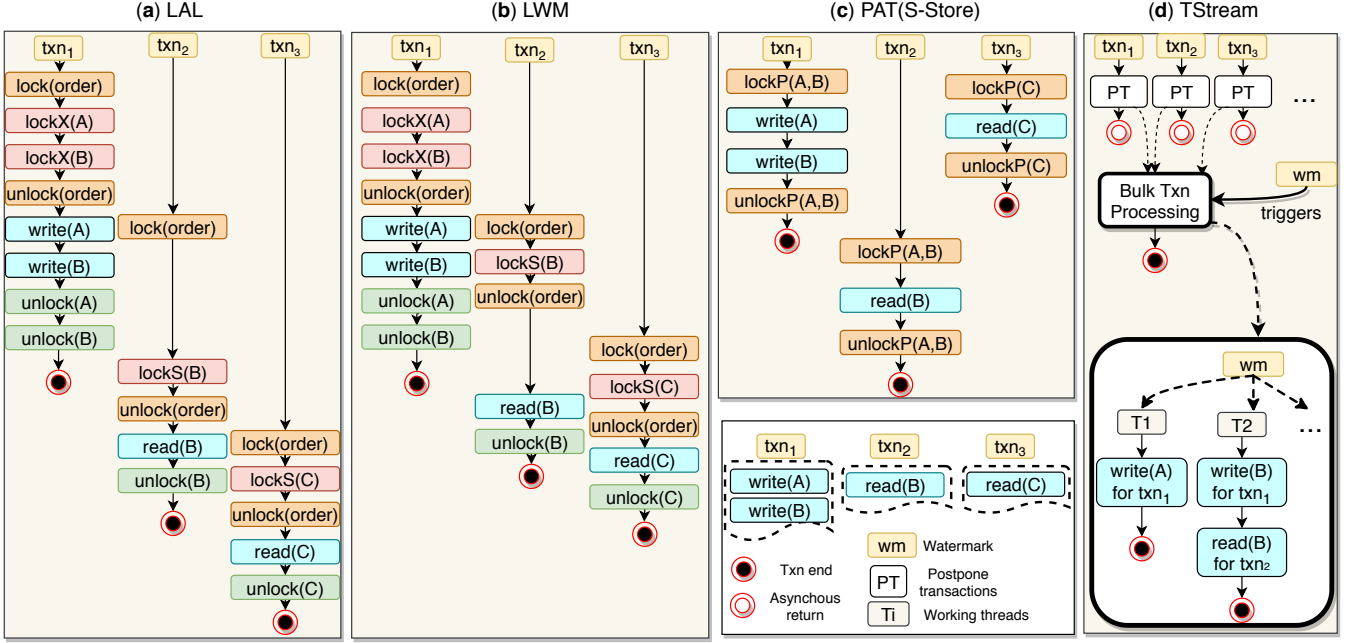


Figure 3: Examples of how different algorithms process three concurrent state transactions.

all execution instances of an operator may access any partitions of shared states during stream processing. To guarantee state consistency, PAT requires partition-level locks (*LockP*) to synchronize the accesses to each partition, which act much like an array of *Lock(order)* of multiple partitions.

Let us assume record *A* and *B* are grouped into the same partition and *C* is grouped into different partitions. As shown in Figure 3(c), the execution of *txn<sub>3</sub>* (only touches *C*) will never compete with the other two transactions achieving good system concurrency. In contrast, as *txn<sub>1</sub>* and *txn<sub>2</sub>* target at the same partition, their execution must be performed serially similar to LAL. To achieve good performance, PAT requires the system to perfectly partition shard states beforehand and each state transaction only needs to touch on a single partition with minimum conflicts among each other. However, multi-partition transactions can significantly degrade the overall performance – a common problem for all partition-based approaches [34]. For example, consider *txn<sub>3</sub>* also needs to access *B*, or *C* is grouped in the same partition of *A* and *B*. In either case, all three transactions will be executed serially under PAT.

There are a few more related studies [14, 7], which are omitted here as they are similarly based on locks but missing further details on how they execute state transactions. In summary, prior solutions [32, 16, 14, 39] commonly adopt a synchronized execution model (shown in Figure 4a), which severely restrict system concurrency and scalability. *Firstly*, executing the state access operations strictly following the dependencies encoded in the streaming computation logic can cause extra long execution stalls in stream processing, consequently degrading the overall system throughput. *Secondly*, each transaction is synchronized with locks to guard correct locking sequence, which can cause high performance overhead under highly contended workloads. Although, PAT (i.e., S-Store) can potentially reduce such

overhead by careful partition shared states beforehand, it quickly falls back to LAL with more multi-partition transactions. Therefore, we need a new solution for scaling the transactional data management in DSPSs on modern multicore architectures.

## 4. SYSTEM DESIGN

In this section, we discuss the design of **TStream**, and highlight how it addresses the scalability limitation in existing solutions.

### 4.1 Overview

Targeting at shared-memory multicore architectures, we propose **TStream** with two key novel designs.

**Asynchronous state management.** Synchronized state accesses are not only *harmful* leaving many parallelism opportunities unexplored but also often *unnecessary*. To allow streaming operators to continue process more input events without waiting for its currently issued state transaction to finish, **TStream** adopts an asynchronous execution model that detaches state transaction processing from application computation logic. Through postponing state accesses by recording necessary information (e.g., read/write sets) into a conceptual data structure – placeholders, **TStream** minimizes unnecessary stalls caused by state management in stream processing. **TStream** adopts *watermark* as periodic signal to trigger the actual process of state transactions. Postponed state transactions between two subsequent watermarks hence naturally form a batch, and the batch length can be then tuned to tradeoff system latency and throughput.

**Operation-chains parallel execution.** To take advantage of modern multicore architectures, we propose an *operation-chain* based execution model, that greedily extract parallelism opportunists among each batch transactions. As shown in Figure 3(d), multiple operation

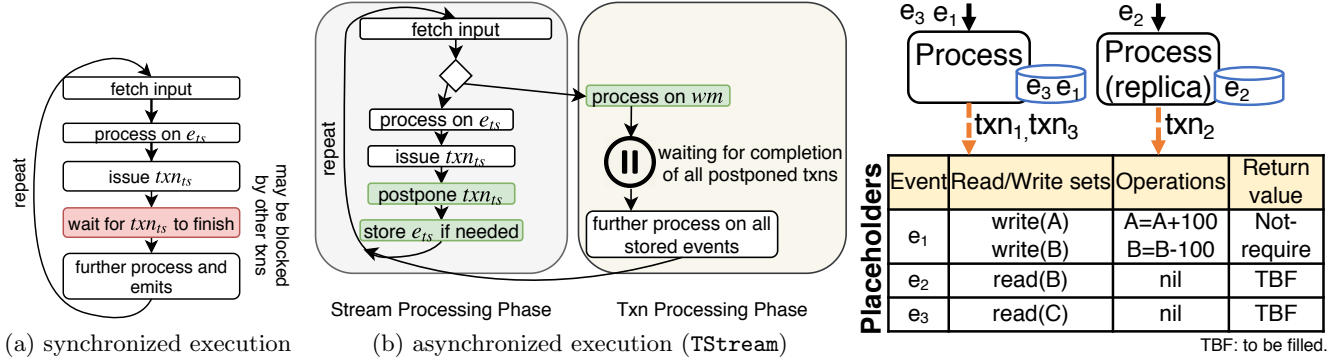


Figure 4: Workflow of an execution instance under different execution models. Figure 5: Example updates of placeholders.

chains are built and can be (ideally) concurrently processed with no locks while still guaranteeing the aforementioned consistency properties. We further apply dynamic work-scheduling and NUMA-aware placement techniques to improve its processing efficiency.

## 4.2 Asynchronous State Management

Figure 4b shows the workflow of an execution instance (of one operator) in **TStream**, which contains two phases. (1) During stream processing, after operator issues state transaction, it does not wait for transaction to be executed. Instead, it only *postpones* the transaction and immediately works on more input events. This minimize unnecessary computation stalls caused by state transactions in stream processing. (2) When an operator receives a *watermark* (a special event to trigger transaction processing), it enters *transaction processing phase* and waits for all postponed transactions to be executed.

**Stream processing phase.** In order to postpone state transaction, necessary input information including read/write sets of each state transaction are recorded in a conceptual data structure called *placeholders*. Figure 5 shows an example process of postponing two state transactions. After  $txn_1$  is recorded in placeholders, **Process** can continue work on other input events (e.g.,  $e_3$ ). This eliminates expensive synchronization on each individual transaction.

The placeholder is implemented as an auxiliary data structure of each input event in **TStream**, hence no centralized contention. Events may need to be stored temporarily as they may require further process depending on evaluation results of corresponding state transaction. For example,  $e_2$  has to be marked as *in-complete* since it requires the value of state *B* for further computation. A reference (i.e., pointer) to the input event (and its placeholder) will be embedded with transaction to be evaluated later. For example, during the processing of  $txn_2$ ,  $e_2$ 's placeholder will be updated with the value of state *B* through the reference. When stream operator is resumed after all transactions are processed,  $e_2$ 's placeholder will then contain the desired value to support further stream computation.

**Transaction processing phase.** Instead of eagerly evaluate each state transaction, **TStream** postpones and accumulates them into a batch. Subsequently, **TStream** adopts *watermark* as periodic synchronization signal to

trigger the process of a batched transactions. Following previous work, we rely on stream sources and operators (e.g., **Parser**) to create watermarks based on their knowledge of the stream data [33].

Figure 6 shows an example workflow of receiving watermarks. We assume a watermark with timestamp of 5 is arrived, and it triggers **Process** operator to start transaction processing. Note that, the same watermark must be broadcast to all instances of **Process** before transaction processing can start. This is implemented by a *Countdown Latch*<sup>3</sup> in **TStream**. During transaction processing, all instances are paused and no further input events (e.g.,  $e_6$ ,  $e_7$ ) are allowed to enter the system. When all postponed transactions are evaluated, operator will resume the process of future input events, and the current watermark is forward to all instances of downstream operator.

The interval size of two subsequent watermark plays an important role in tuning system throughput and processing latency. Having a large interval, the system will wait for longer to start transaction processing (i.e., synchronization overhead in **TStream**), which potentially worse processing latency. Conversely, having a small interval size, the system throughput may drops with insufficient parallelism (too few transactions in the batch), which may end up with more input events queued up, and eventually leads to high processing latency. We evaluate the effect of watermark interval in our experiments.

## 4.3 Operation-Chains Parallel Execution

Enabling more concurrency has become particularly important with the proliferation of multicore and large-scale systems. As a watermark guarantees that no subsequent input events has a timestamp smaller than any events a prior of it, **TStream** has the complete knowledge of all state transactions to be handled *without worrying future* transactions. In other words, **TStream** only needs to concentrate on improving processing throughput of every batch of state transactions arrived between two consecutive watermarks.

Observing that the order-preserving property essentially determines the execution sequences of conflicting operations of each state. For example, a write operation must precede a read operation of the same state if the write

<sup>3</sup><https://www.baeldung.com/java-countdown-latch>

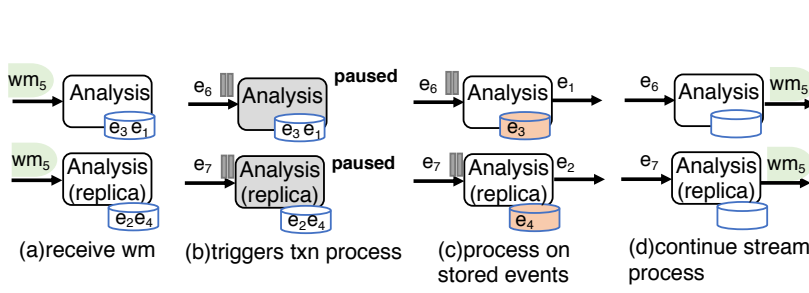


Figure 6: Example workflow of receiving watermarks.

operation is triggered by an earlier input event. We propose an *operation-chain* based execution model to process every batch of state transactions. Specifically, all state transactions are decomposed into atomic state access operations. Each operation, either read/write targets at one state. Operations targeting the same state form an **operation chain** sorted in the timestamp order, which naturally preserves ACID+O properties. These operation chains can be then (ideally) executed concurrently and independently, and hence allows TStream to achieve a maximum concurrency of the number of operation chains constructed for every batch of state transactions.

**Construct operation chains.** Figure 7 shows an example of decomposing three transactions. Two operation chains are formed as a result. First,  $txn_{1,2,3}$  are decomposed into atomic operations, where each operation is annotated with timestamp of its original transaction, targeting state, access type, and optional parameters (e.g., value to update state). Then, those decomposed operations are hash partitioned based on the targeting state. Finally, operations of the same partition group form an operation chain that is sorted by timestamp. Each operation chain is implemented by a concurrent ordered data structure (e.g., concurrent skip-list), allowing concurrent insertion from multiple execution instances while still guaranteeing that items (i.e., operations) are ordered by the triggering event timestamp.

**Evaluate operation chains.** Constructed operation chains can be then (ideally) processed independently by multiple threads concurrently. Recall that all instances of the operator entering transaction processing phase are paused (i.e., waiting to be notified), hardware resources (e.g., CPU cores) can be utilized to process transactions. In TStream, we use a dedicated thread pool with predefined number of threads to evaluate operation chains.

A simple strategy is to statically assign an equal number of operation chains (as tasks) to every working thread. However, such static approach may not always achieve good load balancing. Runtime overhead of more complicated load-balancing scheduling algorithm may offsets the gain as every batch of transactions need to be rescheduled. TStream adopts dynamic work-scheduling [12] to achieve better load balancing. Specifically, TStream maintains a pool of operation chains, where multiple threads (from the thread pool) can fetch one operation chain (as a task) from the pool to work with. Once a thread finishes its task (e.g., a short operation chain with few operations), it fetches another until there is no task left in the pool. Furthermore, once an operation chain has only read-only

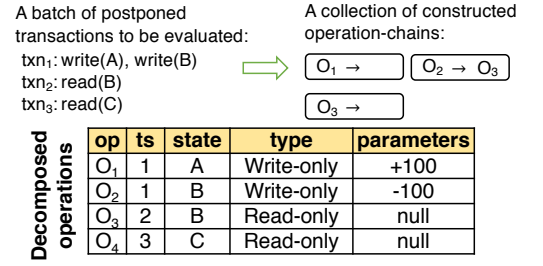


Figure 7: Example operation chains construction.

operations unprocessed, it can be cooperatively processed by multiple threads.

**Handling cross-states dependency.** The aforementioned evaluation approach unfortunately does not work correctly when operations on one state may depend on another state. Assuming  $txn_7$  is to update  $C$  with the value of  $B$  plus one (i.e.,  $txn_7: C=B+1$ ). Denote the corresponding decomposed operation as  $O_7$ . To support such cross-chain dependency is tricky as we need to make sure all (at least those a-priori of  $O_7$ ) updates on  $B$  are applied before execute  $O_7$ . Otherwise,  $O_7$  may read a wrong version of  $B$  violating ACID+O properties.

TStream handles such issue from two aspects. An example of handling  $txn_7$  is illustrated in Figure 8.

(1) *Encoding.* TStream encodes dependency information and update operations on the depended state inside each operation. For example,  $O_7$ 's dependency on  $B$  is encoded as its parameters. It also needs to record all update operations of depended state with a smaller timestamp. This can be done by traversing operation chain of depended state. In this example,  $O_6$  is recorded and  $O_9$  is not.

Such encoding step has to be performed only when watermark is received by all instances and before actual transaction execution starts. This ensures the complete knowledge of all operations to handle.

(2) *Querying.* During transaction execution, one thread will eventually encounter operation with dependency information encoded. It then enquires the correct value of depended state. Figure 8 shows that  $O_7$  fails to obtain correct version of state  $B$  as  $O_6$  has not been executed. The thread can switch out and evaluate other operation chains. After  $O_6$  is executed and the value of  $B$  is updated,  $O_7$  can then continue its process.

Note that, TStream keeps *multiple versions* of each state updated at different timestamp to allow write without waiting for read. Therefore, the process of  $O_7$  will not block the process of  $O_9$ . After the current batch of transactions are processed, all versions except the latest version are expired and can be garbage collected due to the ordering property of watermarks.

## 4.4 Further Optimizations

**Evaluation pushdown.** Consider a write operation, which requires to *double* the value of a state. Naive implementation would require the system return the state value to stream operator, who will multiply the value by two, and write the value back to state. To reduce such round-trip communication, TStream encodes the multiplication with operation itself, similar to how TStream handles cross-states



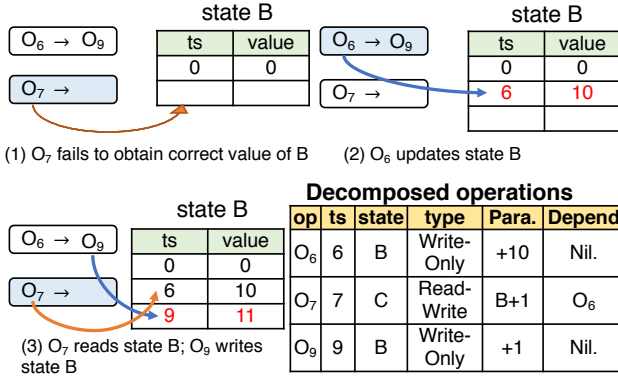


Figure 8: Handling cross-state operation.

dependency. During the processing of the write operation, the multiplication and write will be executed consecutively after read. Conversely, if an operator requires the state value for further computation (e.g., in subsequent stream computation), then **TStream** has to write back the state value to placeholders during transaction processing.

**NUMA-awareness.** Following previous work [35], we configure the data placement of the shared state in an island-aware manner for NUMA-awareness. Specifically, we first distribute the shared states evenly to each NUMA node, so that the workload assigned to each node is (ideally) balanced. Then, the working threads for transaction processing is divided into several groups, where each group is responsible for workloads of part of the CPU sockets (i.e., island). We evaluate the effect of different grouping and select the best one for all algorithms in our experiments.

## 4.5 Discussion

**Limitations.** Compared to existing synchronized execution model, the main disadvantage of **TStream** is the potentially higher stream processing latency. This is caused by two reasons. First, earlier arrived events, which triggers transaction to *read* shared states, are stored and not fully processed as their triggering transactions are postponed until subsequent watermark arrived. Second, to ensure correctness, *all* execution instances of an operator are paused until *all* postponed transactions are evaluated.<sup>4</sup> Nevertheless, as discussed before, it brings many opportunities in extracting higher concurrency while preserving **ACID+0** properties. Our experimental results also show that the significant improved performance minimizes the latency penalty. Another limitation is that, the read/write-set of a state transaction must be deducible before the transaction begins. Fortunately, all the tested workloads (including both of our two use cases), that we are going to discuss shortly later, provides read/write sets before execution and **TStream**'s limitation does not apply here. Specifically, stream operator already knows which state(s) to access according to the input event it processes. Note that, use case 1 is a real-world application suggested by a commercial transactional stream processing system, namely Ledger [4]. We admit that this is not a must, and we plan to investigate the situation when read/write sets are not

<sup>4</sup>Earlier/speculatively resume processing is deferred as a future work to explore.

Table 2: Comparisons in related systems

	ACID	ACID +O	Large State	Multicore/ Stream Manycore	Processing
Storm[3]/Flink[1]	×	×	×	×	✓
BriskStream[47]/StreamBox[32]	×	×	×	✓	✓
Cavalier[41]	✓	×	✓	✓	×
SDG[20]	×	×	✓	✓	✓
LWM[38]/PAT[31]	✓	✓	✓	limited	✓
TStream	✓	✓	✓	✓	✓

available. A potential solution is to rely on pre-analysis [36], and we defer it as future work to explore.

**System Comparison.** We list the comparison to some related systems<sup>5</sup> in Table 2. DSPSs like Storm [3] and Flink [1] provide simple API to express streaming application but scale poorly on modern multicore processors [47]. StreamBox [33], BriskStream [48], and Saber [28] are able to achieve very high throughput and low latency of stream processing, but they are lack of a built-in support of transactional state management. State-of-the-art OLTP database such as Cavalier [42] does not provide order-preserving consistency nor stream processing capability. SDG [20] manages large mutable state distributedly but does not provide transactional guarantees nor order-preserving consistency. The closest works to **TStream** are DSPSs supporting **ACID+0** properties [39, 32], which are however limited at their scalability on modern multicores. as we have discussed in Section 3. In the next section, we compare **TStream** with those prior solutions guaranteeing **ACID+0** properties, and thus omit the comparison with others (e.g., Storm/Flink and SDG).

## 5. IMPLEMENTATION DETAILS

**TStream** is built on BriskStream [48] – a multicore optimized stream processing system. We extend BriskStream with state transaction APIs (to be discussed shortly later) and a transaction processing engine, which is developed based on Cavalier [42] database. For comparison, we implement multiple related algorithms including LAL, LWM and PAT (S-Store) into **TStream**. This allows us to abstract away the implementation details of each approach and concentrate on the algorithm itself.

**No Order Preserving CC (NOCC).** At its minimum, the system needs to ensure that accesses to shared states are processed with **ACID** guarantees. In addition to the four **ACID+0** guaranteed algorithms (i.e., LAL, LWM, PAT and **TStream**), we additionally implement NOCC (w/o order-preserving consistency) as a baseline based on the two-phase-locking with no-wait deadlock prevention strategy [18]. Other conventional protocols (such as MVCC [41]) are also applicable but similarly *do not* guarantee order-preserving consistency as we have discussed before, we hence omit them for brevity.

**Transactional APIs.** In **TStream**, developers can declare transaction manager in each operator that allows

<sup>5</sup>Ledger [4] recently introduces transactional state management to Flink, but is unfortunately close-sourced.

multiple instances of the operator to access *shared mutable states* without worrying state consistency violation. This is exemplified by the following code that defines how **Process** manages item table as shared states discussed in Section 2.

```

1 private TxnManager m; //Interface of shared states.
2 public void pre_process(Event e); //user defined
   function (UDF)
3 public void post_process(Event e); //UDF depended on
   shared states
4 public void emit_output(Event e); //UDF generates
   output to downstream operator
5 public void initialize(/*config info*/) {
6     m = new TxnManagerLAL(); //creates a TxnManager of
   specific algorithm
7     m.initialize("ItemTable", config); //initialize
   shared states
8 }
9 public void execute(Event e) { //repeatedly invoked for
   each input event
   pre_process(e); //e.g., count, filter, etc.
10
11     //state accesses as one transaction
12     if(e.ReadRequest()){
13         m.Read("ItemTable", e.keys, e.holders); //e.
14         holders are going to be filled with the
           value of targetting state.
15     }else{
16         m.Write("ItemTable", e.keys, e.values);
17     }
18
19     post_process(e); //e.g., performing aggregation
           on state values in e.holders
20     emit_output(e);
21 }

```

Listing 1: Code template of **Process** operator of our running example

As shown in line 1, application developer needs to claim a **TxnManager** as the interface to access shared states. Different transaction execution algorithms including NOCC, LAL [39], LWM [39] and PAT [32] are supported in **TStream**, which are implemented as custom **TxnManagers** (e.g., **TxnManagerLAL**). The **execute** method in line 9 is repeatedly involved for every input event. At line 14 and 16, the processing of input event *e* triggers shared states read or write. The read/write sets are extract from input event directly (i.e., *e.keys* and *e.values*). *e.holders* is the placeholder data structure (Section 4.2) to be filled during transaction evaluation. Other parameters of state transaction can be similarly pushed to **TxnManagers**.

## 6. EVALUATION

In this section, we evaluate the effectiveness of **TStream**. We first study (Section 6.3) the case when shared states cannot be partitioned beforehand, and PAT cannot be applied. Then, we study the effect of states partitioning in Section 6.4, where we assume the shared states are pre-partitioned into disjoint subsets. We study two use case workloads in Section 6.5. Finally, we perform sensitivity study to understand the design trade-off in **TStream** with different workloads in Section 6.6.

### 6.1 Experimental Setup

We conduct the experiment on a 4-socket Intel Xeon E7-4820 server with 128 GB DRAM. Each socket contains ten 1.9GHz cores and 25MB of L3 cache. The operating system is Linux 4.4.0-62-generic. The number of cores devoted to the system and the size of watermark interval are system parameters, which can be varied by the

system administrator. We vary both parameters in our experiments. In our evaluation, there is a one-to-one binding between threads and cores. We reserve two cores in the system: one for running data stream producer (i.e., **Parser**), and the other for output stream receiver (i.e., **Sink**). Therefore, we devote 2 to 38 cores for running other operators of the application to evaluate the system scalability. To minimize cross-operator communication, we fuse [26] operators (e.g., **Access** and **Compute**) into a single joint operator (**AccessCompute**). Further fine-grained query optimization (e.g., [26, 48]) is orthogonal to this work and omitted for simplicity. Under **TStream**, we additionally configure a dedicated thread pool of 2 to 38 cores to perform the actual transaction execution (see Section 4.3). The thread pool never compete with stream processing as it is paused while stream operators are running and only involved when operators are paused.

Following the previous work [44], we also report how much time each transaction spends in different components of the system. Specifically, we divide the time spend during a transaction execution into the following components. *Useful* is the time that the transaction is really operating on records. *Abort* is the time spent due to transaction abort (e.g., failure for acquiring write locks). *Lock* stands for the total amount of time that a transaction spends due to lock acquisition. *Sync* is the time that a transaction spends due to synchronization to enforce order-preserving property. A transaction may need to synchronize for either 1) a global synchronization primitives, e.g., **Lock(order)** in LAL and LWM, **LockP** in PAT or 2) watermarks in **TStream**. *Others* stands for all other overheads including indexing, timestamp allocation, etc.

### 6.2 Benchmark Workloads

We use both microbenchmark and two use case workloads for performance evaluation.

**MicroBenchmark.** Figure 9(a) shows the topology of our microbenchmark containing four operators. **Parser** continuously feeds synthetic input events to **Access** operator. **Access** issues an state access transaction (read-only or write-only) to a table, which is shared among all execution instances. Information including both access type (read/write), targeting states (length is set to ten) and optional parameters (the value to write) are extracted from each input event. If an event triggers a read-only transaction, **Access** emits the returned state values as one event to **Compute**; otherwise, it simply forwards its input event (which triggers write-only transaction) to **Sink**. **Compute** performs a summation calculation of the returned state values from **Access**, and emits the results as one event to **Sink**. We use **Sink** to record the output stream from **Compute** and **Access** to monitor system performance. Despite its simplicity, the microbenchmark shall be applicable to cover a wide range of different workloads by varying different workload parameters such as the skew of keys (*theta*), read-write ratio and state partition. We vary all those factors in our experiments.

**Use case 1: Online Book Shopping.** Use case 1 processes events describing wiring money and assets between different user accounts. This application is suggested by a commercial DSPS, namely Streaming Ledger [4]. Detailed descriptions can be found in their white paper [4] and are omit here for brevity. We implement it with three



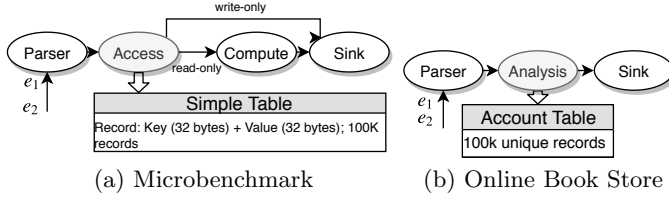


Figure 9: Application topology of two more workloads.

operators. **Parser** parses received input events and forward to **Analysis**. **Analysis** performs some analytics on the received events, and optionally need to read or update some user accounts. The process results are further passed to **Sink**. **Sink** reports the processing results to end users.

The user accounts must be shared among all execution instances of **Analysis**. We set a balanced ratio of transfer and deposit requests (i.e., 50% for each) in **Parser**. Transaction length is four for transfer request (i.e., each request touches four records) and is two for deposit request. We vary the skew of keys ( $\theta$ ) from 0.6 to 0.8. The table contains 100k unique records. The ratio of multi-partition transaction is 0.5 under PAT.

**Use case 2: Online Bidding System.** We have described use case 2 as our running example in Section 2, and we highlight here its configurations used in our experiments. **Process** has a selectivity of one in this work, that is, it always generates one output for each input event. We use a table with 10k records to represent the bidding items. The ratio of bid, alter, and top requests is configured as 6:1:1. Such configuration represents a rather competitive seller’s market scenario and transactions conflict severely. The ratio of multi-partition transaction is set to be 0.75 under PAT, and  $\theta$  is set to be 0.6.

### 6.3 Single Partition Analysis

We first use microbenchmark to evaluate different algorithms described in Section 3. We use a single table consisting of 100k records to represent a non-partitionable shared states. Each record has a size of 64 bytes, and the watermark interval is set to 100 ms in **TStream**. Records are stored in a single partition and all worker threads can access any states. Hence, the partition-based scheme (PAT) is excluded in this section.

**Read-Only Workload.** In the first scalability analysis, we configure input events to trigger only read requests. We set the key skew factor to be 0, and hence the state is accessed with uniform frequency. The read-only workload helps to stress the handover between stream computation and state management of **TStream** as much as possible. Specifically, all input events have to be marked as *incomplete* and the pending transactions need to be processed (all at once) once the watermark is arrived. Compute can only start its computation after **Access** successfully obtains state value.

Figure 10(a) shows the results of our experiment, and we have three main observations. First, the throughput of the most relaxed scheme *without order-preserving guarantee* – NOCC, increases almost linearly with the number of cores. This is expected as there is no lock contention in this workload. Second, lock-ahead based algorithms (LAL and LWM) perform well when the number of cores is small, but

they stop scaling when more than 8 cores are used. The time breakdown in Figure 11(a) indicates that *sync time* dominates the runtime of both LAL and LWM with a large core count. This is primarily caused by their lock-ahead mechanism. Third, **TStream** performs much better than LAL and LWM while guaranteeing the same consistency. However, Figure 11(a) indicates that **TStream** still spends a large portion of time in synchronization, mainly due to the waiting for watermarks. Furthermore, operation chain construction process also contributes to the high overhead in **TStream**. There is still large room for further improvement.

**Write-Intensive Workload.** We next study a write-intensive workload, where input events only trigger write requests to the shared states. To represent a more realistic scenario, we model the accessing distribution as Zipfian skew, where certain states are more likely to be accessed than others. The amount of skew in the workload is determined by the parameter,  $\theta$ . We use the medium and high contention levels for the transactions’ access patterns.

The medium contention results in Figure 10(b) shows a similar trend as in read-only case except that NOCC performs poorly with larger core counts. Figure 11(b) shows that there is a significant increase in management overhead. Under NOCC, transaction needs to hold a copy of records so that it can roll back the changes during aborting. Other order-preserving schemes are permissive in nature and does not involve any transaction aborts in this workload. This is because their lock accusations are performed sequentially due to the lock-ahead mechanism. Therefore, they do not bear such memory-copy overhead. The results of high contention case shown in Figure 10(c) show that only **TStream** is able to further improve its performance beyond 32 cores. The reasons are two folds. First, **TStream** is lock-free and no abort in acquiring locks, higher skew level only potentially brings higher workload unbalancing (i.e., some operation chains may be significant longer than others) to **TStream**. Second, **TStream**’s dynamic work-scheduling scheme relieves the increased workload unbalance issue. In contrast, the breakdown in Figure 11(c) indicates that NOCC is inhibited by excessive abort and redo. Again, LAL and LWM spend most of their execution time in synchronization for guaranteeing correct lock insertion sequence. LWM performs significantly worse than LAL, due to the additional overhead of accessing the `lwm` counter. Note that, `lwm` is read (at least once) during process and updated once during commit for each write-only transaction.

**Read/Write Mixture Workload.** An application may issue both read and write operations to the internal states. In our experiment, we vary the percentage of read operations executed by each transaction. Figure 12(a) shows that NOCC performs better with more read operations. The reason is that there is an increasing concurrency with lesser contention as it does not need to preserve ordering property. LAL and LWM are not affected by the read-write ratio because they are dominated by sync time for enforcing ordering. **TStream** also performs similar regardless of read-write ratio, which is mainly due to the its unique asynchronous execution model. On the one hand, the average processing time per transaction per core of **TStream** increases with more write requests, which is shown in Figure 12(b). On the other hand, each read request

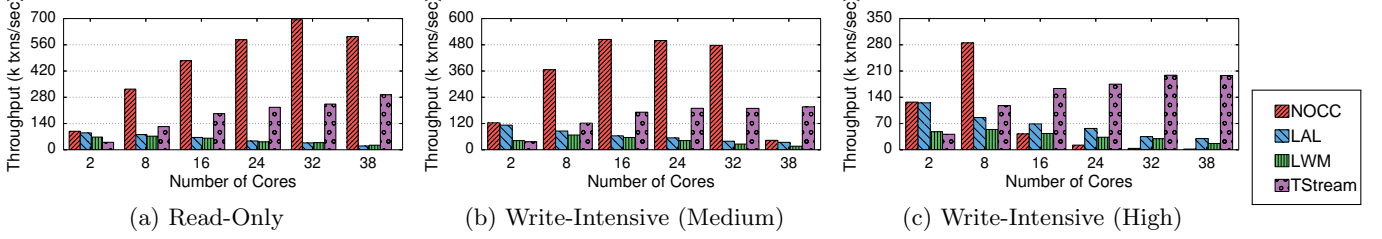


Figure 10: Throughput comparison for varying types of workload.

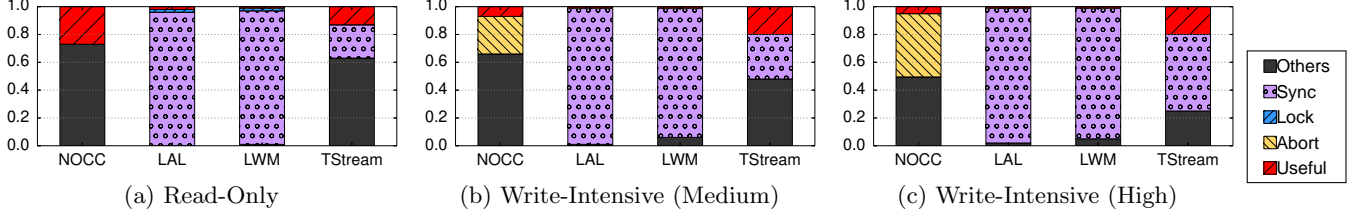


Figure 11: Runtime breakdown (38 cores) for varying types of workload.

requires **TStream** to fill up the *return value* of placeholders (see Figure 5 in Section 4.2), which also involves write operations.

**Small State Size.** The internal state of stream processing (e.g., IoT) can be relatively small compared to traditional relational database. Under such configuration, the chances of transaction contention is very high. We conduct such a study by reducing state size to be only 10k. We use read/write mixture workload (50% each) in this study. Figure 13 shows that NOCC cannot progress when more than 24 cores are used due to the excessive abort and redo process. Conversely, **TStream** shows similar and sometimes even better performance. The robust of **TStream** under highly contended workloads comes from two aspects. First, under such configuration, many transaction requests are targeting at the same state. The overhead of constructing operation chains are effectively amortized. Second, **TStream**'s dynamic work-scheduling scheme successfully reduces the potential workload unbalancing issue among working threads yielding better overall performance.

## 6.4 Multi-Partition Analysis

We now study the effect of shared states partitioning. We use a simple hashing strategy to assign the records to partitions based on their primary keys so that each partition stores approximately the same number of records.

As a common issue of all partition based algorithms [34], the performance of PAT is heavily depended on the length and ratio of multi-partition transactions. Hence, we vary both parameters in our experiments. We first set the length of multi-partition to be 6, that is each multi-partition transaction needs to touch 6 different partitions of the shared states. We vary the percentage of multi-partition transactions in the workload. The results are shown in Figure 14(a). As expected, the system's throughput degrades with more multi-partition transactions as they reduce the amount of parallelism. When there are more

than 25% of the transactions access multi-partition, PAT scheme performs worse than **TStream**. We next execute workload with 50% multi-partition transactions, and Figure 14(b) illustrates the results of varying the number of partitions that they access. The system performance drops significantly with multi-partition transactions accessing two or more partitions. In the worst case, PAT falls back to LAL when every transaction needs to access all partitions.

## 6.5 Use Case Studies

We now study the system performance for two use case workloads assuming each table is pre-partitioned. For this study, we exclude NOCC as it cannot ensure processing correctness. This applies to all other conventional CC protocols without order-preserving guarantee. The results shown in Figure 15 confirm that **TStream** outperforms all other schemes significantly at large core counts at both use cases, and performs especially better at highly contended workload. Comparing to LAL and LWM [39], **TStream** achieves higher concurrency while preserving the same consistency properties. The key reason is that **TStream** is able to explore the hidden parallelism opportunities among transactions at runtime even if they have conflicting operations. Comparing to PAT (i.e., S-Store), **TStream** does not need to statically partition shared states beforehand. Instead, it dynamically partitions operations inside each transaction at a fine-granularity of operation chain. This allows **TStream** to be able to explore more parallelism. However, such fine-grained partition also comes with high constant overhead of managing each operation chain. This is also the reason why **TStream** performs well, sometimes even better under highly contended workload, where management overhead of each operation chain is effectively amortized among more operations.

## 6.6 Sensitivity Study

**The Effect of Watermark Interval.** We now study the effect of varying size of watermark interval in **TStream**,

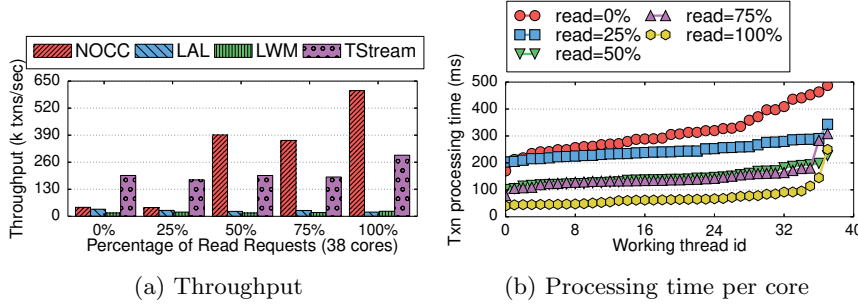


Figure 12: Read/write mixture evaluation.

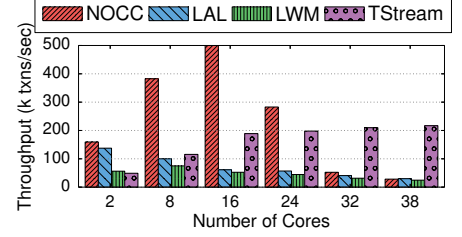


Figure 13: Small state size study (state size=10,000).

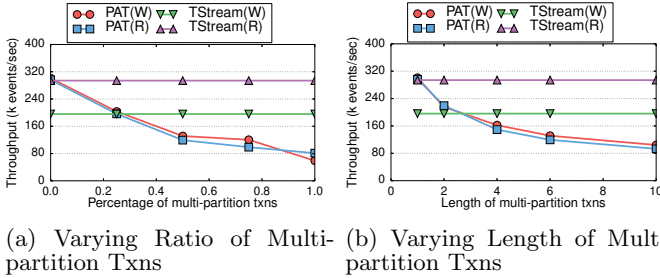


Figure 14: Sensitivity analysis of the PAT scheme for workloads with multi-partition transactions.

on both latency and throughput. Figure 16 shows that the relationship between system throughput and watermark interval is non-linear. Deciding a suitable batch size can be hence complex. Potential techniques have been studied in previous work [17], and we leave it as future work to enhance our system. Following previous work [17], we define the end-to-end latency of a streaming workload as the duration between the time when an input event enters the system and the time when the results corresponding to that event is generated. The results in Figure 17 show that the processing latency of **TStream** is comparable to other schemes and can be tuned to be even smaller. This is a highly desired feature as different applications can have different latency requirements, and **TStream** provides users the opportunity to tune the system flexibly.

**Factor Study.** Figure 18 shows the relative effectiveness of several optimization techniques of **TStream**. Changes are added left to right and are cumulative. *Simple* refers to **TStream** with an synchronous execution mode – it falls back to LAL. *+Asy-execute* enables asynchronous execution model with evaluation push down. *+Dynamic-scheduling* further enables dynamic task scheduling technique to resolve load unbalancing issue. *+Numa-aware* adds the consideration of different NUMA-aware thread and workload placement configurations.

There are three major observations in Figure 18. First, asynchronous state management brings remarkable performance improvement as it avoids unnecessary execution stalls between stream processing and transaction processing. Second, dynamic work scheduling plays a critical role in reducing workload unbalancing issue

under contented workloads. Third, NUMA-awareness brings minor performance improvement in all testing workloads. In **TStream**, each state is always processed by the same thread (as long as rescheduling is not involved for that state), which greedily minimizes NUMA traffic. Nevertheless, we plan to evaluate the impact of NUMA-awareness with more extensive applications in future work.

## 7. RELATED WORK

**Transactional DSPSs.** Most modern DSPSs either do not support transactional state management or have to rely on third-party data storage systems (e.g., [2, 5]), which not only degrades the system performance but also violates the state consistency [16]. Similar motivations have led to a very recent launch of a commercial system, called Ledger [4], developed by *Data Artisans*. It is close-sourced, and we cannot compare our system with it. Wang et al. [39] conducted an early study on the importance of supporting transactional state management in stream processing. They propose a new computing paradigm, called *active complex event processing* (ACEP), to enable complex interactive real-time analytics. Botan et al. [14] presented an *unified transactional model* for streaming applications, called UTM. Recently, MeeHan et al. [32] attempted to fuse OLTP and streaming processing together and developed the S-Store system. Different from previous implementation, **TStream**'s novel asynchronous execution design has shown to achieve much higher throughput and scalability at different types of workloads with limited penalty of higher processing latency.

**Concurrency Control.** Concurrency control (CC) protocols have been investigated widely in decades [11, 19, 44, 42]. Beyond guaranteeing ACID properties, DSPSs must also provide *order-preserving consistency* – a property that conventional CC protocols are not well-prepared for. Several prior works on transaction chopping and lazy evaluation [40, 19] inspired our design of the operation chain based processing model. The key contribution of our work is the application of transaction decomposition to scale transactional state management for stream processing with order-preserving consistency. In particular, we show that the overhead of maintaining sorted data structure (i.e., the operation chain) can be overcome by the performance gains from improved stream processing concurrency. Several prior works [44, 41, 23] studied the scalability bottlenecks in various aspects of concurrency control algorithms. Different from these works, our work studied the concurrency

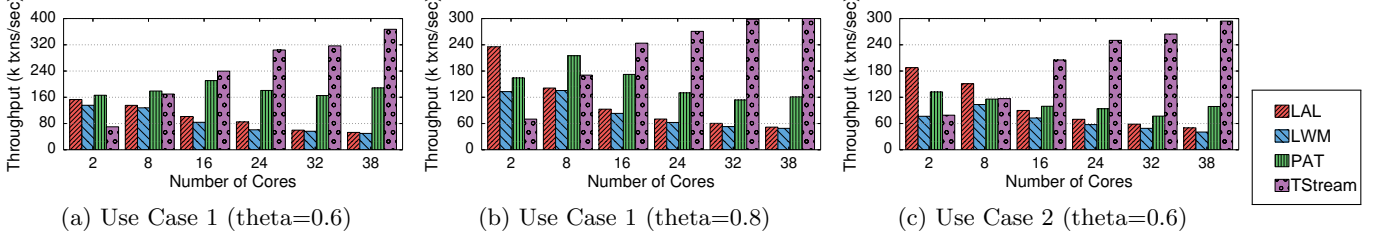


Figure 15: Performance evaluation of two use case workloads.

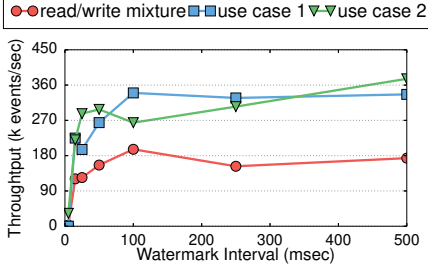


Figure 16: Throughput on varying watermark interval.

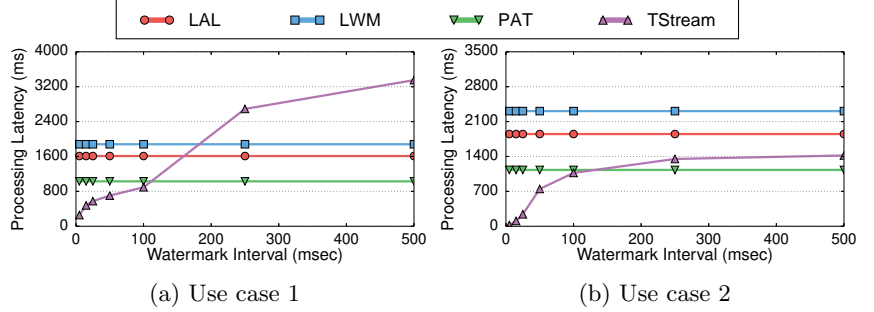


Figure 17: 90<sup>th</sup> Percentile end to end processing latency.

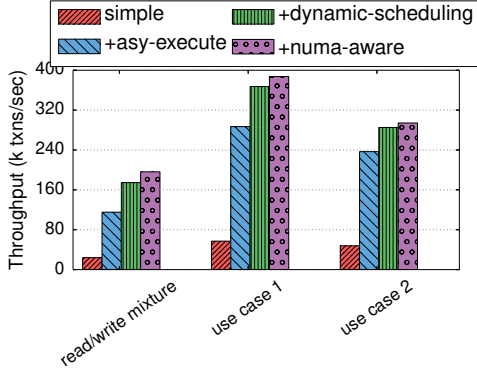


Figure 18: A factor analysis for **TStream** ( $\theta=0.6$ ).

control algorithms in the context of DSPSs with different consistency requirements (i.e., order-preserving). **TStream** is conceptually similar to some deterministic databases [6], which generates a dependency graph that deterministically orders transactions' conflicting record. However, the determinism of **TStream** is given by input events rather than the system itself. Furthermore, deterministic databases are designed for distributed environment, and the dependency analysis is performed by one thread (or one host) before the actual execution can start. Such dry-run phase is too heavy to be applied in **TStream**.

**Data Management Systems on Multicores.** Multicore architectures have brought many research challenges and opportunities for in-memory data management, as outlined in recent surveys [38, 46]. There have been studies on optimizing the instruction cache performance [49, 24], the memory and cache performance [8, 25, 13, 9] and NUMA [30, 31, 22]. The importance of scale

up stream processing is getting more and more attentions recently [47, 28, 33, 48, 45], and we believe it is an equally important optimization direction compared to scale out. **TStream** is built to improve multicore utilization standing on the shoulders of many valuable existing works. Particularly, its asynchronous state management and operation chain parallel evaluation improve system concurrency are inspired by past works such as [49, 40]. However, none of the previous work addresses the scalability bottlenecks that **TStream** faces. That is how to scale stream processing guaranteeing transactional semantics of shared states. To address the needs for a NUMA-aware OLTP system, Porobic et al. [35] proposed "hardware islands", in which NUMA nodes are grouped into logical partitions as islands and communicate through message passing among different islands. We applied the similar idea in **TStream** but it brings only minor improvement in the testing workloads.

## 8. CONCLUSION

This paper introduces **TStream** aiming at scaling stream processing with transactional state management on shared-memory multicore architectures. **TStream** achieves high scalability via two key designs including 1) asynchronous state management, which minimizes execution stalls caused by state access operations and 2) operation-chains parallel execution model, which maximize transaction execution concurrency while guaranteeing ACID+O properties. To the best of our knowledge, we also provide the first comprehensive study of different algorithms for supporting transactional state management in DSPS and results have confirmed the superiority of **TStream**'s designs.

## 9. REFERENCES

- [1] Apache flink, <https://flink.apache.org/>, 2018.

- [2] Apache smaza, <https://samza.apache.org/learn/documentation/0.7.0/container/state-management.html>, 2018.
- [3] Apache storm, <http://storm.apache.org/>, 2018.
- [4] Data Artisans Streaming Ledger Serializable ACID Transactions on Streaming Data, <https://www.da-platform.com/streaming-ledger>. 2018.
- [5] Stateful stream processing in flink, <https://cwiki.apache.org/confluence/display/FLINK/Stateful+Stream+Processing>, 2018.
- [6] D. J. Abadi, Y. Ahmad, M. Balazinska, U. Cetintemel, M. Cherniack, J.-H. Hwang, W. Lindner, A. Maskey, A. Rasin, E. Ryvkina, et al. The design of the borealis stream processing engine. In *CIDR'05*.
- [7] L. Affetti, A. Margara, and G. Cugola. Flowdb: Integrating stream processing and consistent state management. In *Proceedings of the 11th ACM International Conference on Distributed and Event-based Systems, DEBS '17*, pages 134–145, New York, NY, USA, 2017. ACM.
- [8] A. Ailamaki, D. J. DeWitt, M. D. Hill, and D. A. Wood. Dbmss on a modern processor: Where does time go? In *Proceedings of the 25th International Conference on Very Large Data Bases, VLDB '99*, pages 266–277, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc.
- [9] C. Balkesen, J. Teubner, G. Alonso, and M. T. Özsu. Main-memory hash joins on multi-core cpus: Tuning to the underlying hardware. In *2013 IEEE 29th International Conference on Data Engineering (ICDE)*, pages 362–373, April 2013.
- [10] P. Bernstein and E. Newcomer. *Principles of Transaction Processing: For the Systems Professional*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1997.
- [11] P. A. Bernstein and N. Goodman. Concurrency control in distributed database systems. *ACM Comput. Surv.* 1981.
- [12] R. D. Blumofe and C. E. Leiserson. Scheduling multithreaded computations by work stealing. *J. ACM*, 46(5):720–748, Sept. 1999.
- [13] P. A. Boncz, S. Manegold, and M. L. Kersten. Database architecture optimized for the new bottleneck: Memory access. In *Proceedings of the 25th International Conference on Very Large Data Bases, VLDB '99*, pages 54–65, San Francisco, CA, USA, 1999. Morgan Kaufmann Publishers Inc.
- [14] I. Botan, P. M. Fischer, D. Kossmann, and N. Tatbul. Transactional stream processing. In *Proceedings of the 15th International Conference on Extending Database Technology, EDBT '12*, pages 204–215, New York, NY, USA, 2012. ACM.
- [15] R. Castro Fernandez, M. Migliavacca, E. Kalyvianaki, and P. Pietzuch. Integrating scale out and fault tolerance in stream processing using operator state management. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data, SIGMOD '13*, pages 725–736, New York, NY, USA, 2013. ACM.
- [16] U. Cetintemel, J. Du, T. Kraska, S. Madden, D. Maier, J. Meehan, A. Pavlo, M. Stonebraker, E. Sutherland, N. Tatbul, K. Tufte, H. Wang, and S. Zdonik. S-store: A streaming newsql system for big velocity applications. *Proc. VLDB Endow.*, 7(13):1633–1636, Aug. 2014.
- [17] T. Das, Y. Zhong, I. Stoica, and S. Shenker. Adaptive stream processing using dynamic batch sizing. In *Proceedings of the ACM Symposium on Cloud Computing, SOCC '14*, pages 16:1–16:13, New York, NY, USA, 2014. ACM.
- [18] K. P. Eswaran, J. N. Gray, R. A. Lorie, and I. L. Traiger. The notions of consistency and predicate locks in a database system. *Commun. ACM*, 19(11):624–633, Nov. 1976.
- [19] J. M. Faleiro, A. Thomson, and D. J. Abadi. Lazy evaluation of transactions in database systems. In *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data, SIGMOD '14*, pages 15–26, New York, NY, USA, 2014. ACM.
- [20] R. C. Fernandez, M. Migliavacca, E. Kalyvianaki, and P. Pietzuch. Making state explicit for imperative big data processing. In *2014 USENIX Annual Technical Conference (USENIX ATC 14)*, pages 49–60, Philadelphia, PA, 2014. USENIX Association.
- [21] J. Ghaderi, S. Shakkottai, and R. Srikant. Scheduling storms and streams in the cloud. In *Proceedings of the 2015 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems, SIGMETRICS '15*, pages 439–440, New York, NY, USA, 2015. ACM.
- [22] J. Giceva, G. Alonso, T. Roscoe, and T. Harris. Deployment of query plans on multicores. *Proc. VLDB Endow.*, 8(3):233–244, Nov. 2014.
- [23] R. Harding, D. Van Aken, A. Pavlo, and M. Stonebraker. An evaluation of distributed concurrency control. *Proc. VLDB Endow.*, 10(5), Jan. 2017.
- [24] S. Harizopoulos and A. Ailamaki. Improving Instruction Cache Performance in OLTP. *ACM Trans. Database Syst.*, 2006.
- [25] B. He, Q. Luo, and B. Choi. Cache-conscious automata for xml filtering. In *21st International Conference on Data Engineering (ICDE'05)*, pages 878–889, April 2005.
- [26] M. Hirzel, R. Soulé, S. Schneider, B. Gedik, and R. Grimm. A catalog of stream processing optimizations. *ACM Comput. Surv.*, 46(4):46:1–46:34, Mar. 2014.
- [27] R. Kallman, H. Kimura, J. Natkins, A. Pavlo, A. Rasin, S. Zdonik, E. P. C. Jones, S. Madden, M. Stonebraker, Y. Zhang, J. Hugg, and D. J. Abadi. H-store: A high-performance, distributed main memory transaction processing system. *Proc. VLDB Endow.*, 1(2):1496–1499, Aug. 2008.
- [28] A. Koliosis, M. Weidlich, R. Castro Fernandez, A. L. Wolf, P. Costa, and P. Pietzuch. Saber: Window-based hybrid stream processing for heterogeneous architectures. In *Proceedings of the 2016 International Conference on Management of Data, SIGMOD '16*, pages 555–569, New York, NY, USA, 2016. ACM.
- [29] S. Kulkarni, N. Bhagat, M. Fu, V. Kedigehalli, C. Kellogg, S. Mittal, J. M. Patel, K. Ramasamy, and S. Taneja. Twitter heron: Stream processing at scale.



- In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data*, SIGMOD '15, pages 239–250, New York, NY, USA, 2015. ACM.
- [30] V. Leis, P. Boncz, A. Kemper, and T. Neumann. Morsel-driven parallelism: A numa-aware query evaluation framework for the many-core age. In *Proceedings of the 2014 ACM SIGMOD International Conference on Management of Data*, SIGMOD '14, pages 743–754, New York, NY, USA, 2014. ACM.
  - [31] Y. Li, I. Pandis, R. Mueller, V. Raman, and G. M. Lohman. Numa-aware algorithms: the case of data shuffling. In *CIDR*, 2013.
  - [32] J. Meehan, N. Tatbul, S. Zdonik, C. Aslantas, U. Cetintemel, J. Du, T. Kraska, S. Madden, D. Maier, A. Pavlo, M. Stonebraker, K. Tufte, and H. Wang. S-store: Streaming meets transaction processing. *Proc. VLDB Endow.*, 8(13):2134–2145, Sept. 2015.
  - [33] H. Miao, H. Park, M. Jeon, G. Pekhimenko, K. S. McKinley, and F. X. Lin. Streambox: Modern stream processing on a multicore machine. In *2017 USENIX Annual Technical Conference (USENIX ATC 17)*, pages 617–629, Santa Clara, CA, 2017. USENIX Association.
  - [34] A. Pavlo, C. Curino, and S. Zdonik. Skew-aware automatic database partitioning in shared-nothing, parallel oltp systems. In *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*, SIGMOD '12, pages 61–72, New York, NY, USA, 2012. ACM.
  - [35] D. Porobic, I. Pandis, M. Branco, P. Tözün, and A. Ailamaki. Oltp on hardware islands. *Proceedings of the VLDB Endowment*, 5(11):1447–1458, 2012.
  - [36] K. Ren, A. Thomson, and D. J. Abadi. An evaluation of the advantages and disadvantages of deterministic database systems. *Proc. VLDB Endow.*, 7(10):821–832, June 2014.
  - [37] J. Tan and M. Zhong. An online bidding system (obs) under price match mechanism for commercial procurement. *Applied Mechanics and Materials*, 556-562:6540–6543, 05 2014.
  - [38] K.-L. Tan, Q. Cai, B. C. Ooi, W.-F. Wong, C. Yao, and H. Zhang. In-memory databases: Challenges and opportunities from software and hardware perspectives. *SIGMOD Rec.*, 44(2):35–40, Aug. 2015.
  - [39] D. Wang, E. A. Rundensteiner, and R. T. Ellison, III. Active complex event processing over event streams. *Proc. VLDB Endow.*, 4(10):634–645, July 2011.
  - [40] Z. Wang, S. Mu, Y. Cui, H. Yi, H. Chen, and J. Li. Scaling multicore databases via constrained parallel execution. In *Proceedings of the 2016 International Conference on Management of Data*, SIGMOD '16, pages 1643–1658, New York, NY, USA, 2016. ACM.
  - [41] Y. Wu, J. Arulraj, J. Lin, R. Xian, and A. Pavlo. An empirical evaluation of in-memory multi-version concurrency control. *Proc. VLDB Endow.*, 10(7):781–792, Mar. 2017.
  - [42] Y. Wu, C.-Y. Chan, and K.-L. Tan. Transaction healing: Scaling optimistic concurrency control on multicores. In *Proceedings of the 2016 International Conference on Management of Data*, SIGMOD '16, pages 1689–1704, New York, NY, USA, 2016. ACM.
  - [43] Y. Wu and K. Tan. Chronostream: Elastic stateful stream computation in the cloud. In *2015 IEEE 31st International Conference on Data Engineering*, pages 723–734, April 2015.
  - [44] X. Yu, G. Bezerra, A. Pavlo, S. Devadas, and M. Stonebraker. Staring into the abyss: An evaluation of concurrency control with one thousand cores. *Proc. VLDB Endow.*, 8(3):209–220, Nov. 2014.
  - [45] S. Zeuch, B. D. Monte, J. Karimov, C. Lutz, M. Renz, J. Traub, S. Breß, T. Rabl, and V. Markl. Analyzing efficient stream processing on modern hardware. *Proc. VLDB Endow.*, 12(5):516–530, Jan. 2019.
  - [46] H. Zhang, G. Chen, B. C. Ooi, K. Tan, and M. Zhang. In-memory big data management and processing: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 27(7):1920–1948, July 2015.
  - [47] S. Zhang, B. He, D. Dahlmeier, A. C. Zhou, and T. Heinze. Revisiting the design of data stream processing systems on multi-core processors. In *2017 IEEE 33rd International Conference on Data Engineering (ICDE)*, pages 659–670, April 2017.
  - [48] S. Zhang, J. He, A. C. Zhou, and B. He. Briskstream: Scaling Data Stream Processing on Multicore Architectures. To appear in SIGMOD, 2019.
  - [49] J. Zhou and K. A. Ross. Buffering database operations for enhanced instruction cache performance. In *Proceedings of the 2004 ACM SIGMOD international conference on Management of data*, pages 191–202. ACM, 2004.