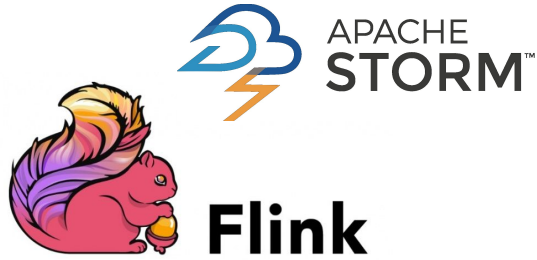


Distributed Stream Processing with Query Compilation

Ivan Ilic, Till Lehmann, Tobias Niedling, Youri Kaminsky

26.04.2021



+ Usability

+ Efficiency



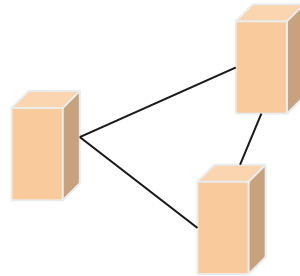
Outline

1. Technical Background
2. Workload Design
3. Implementation
4. Evaluation

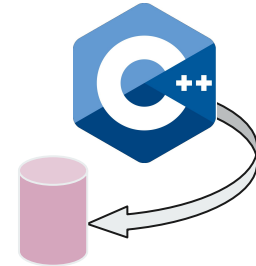
Technical Background



Stream Processing



Distributed Computing



Query Compilation



Assumptions

1. We focus on processing time windowing only.
2. We solely implement sliding windows.
3. We assume a fixed number of nodes to distribute across for the entire query runtime. Thus, we don't support re-scaling while executing a query.
4. We only support decomposable aggregation operations.
5. We do not ensure fault tolerance.

Workload Design



Sample Query

Stream 1:

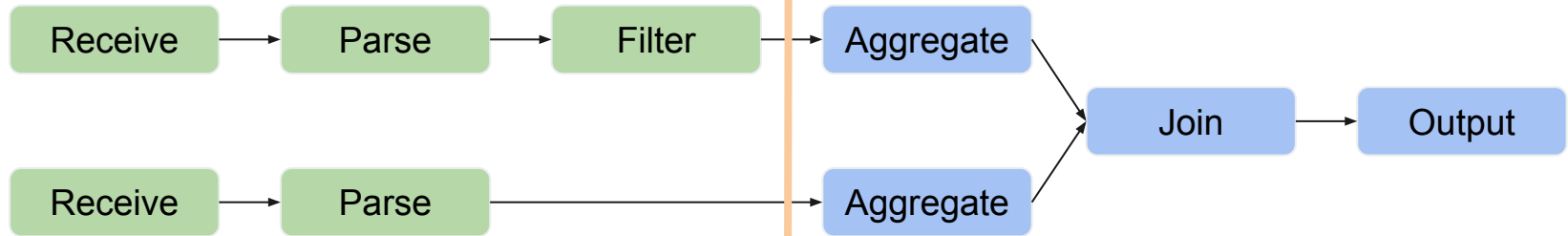
```
ADS (  
    ad_id: Int,  
    user_id: Int,  
    cost: Double  
)
```

Stream 2:

```
PURCHASES (  
    purchase_id: Int,  
    user_id: Int,  
    ad_id: Int,  
    value: Double  
)
```


Sample Query

Purchases



Ads

Processing Time Sliding Window



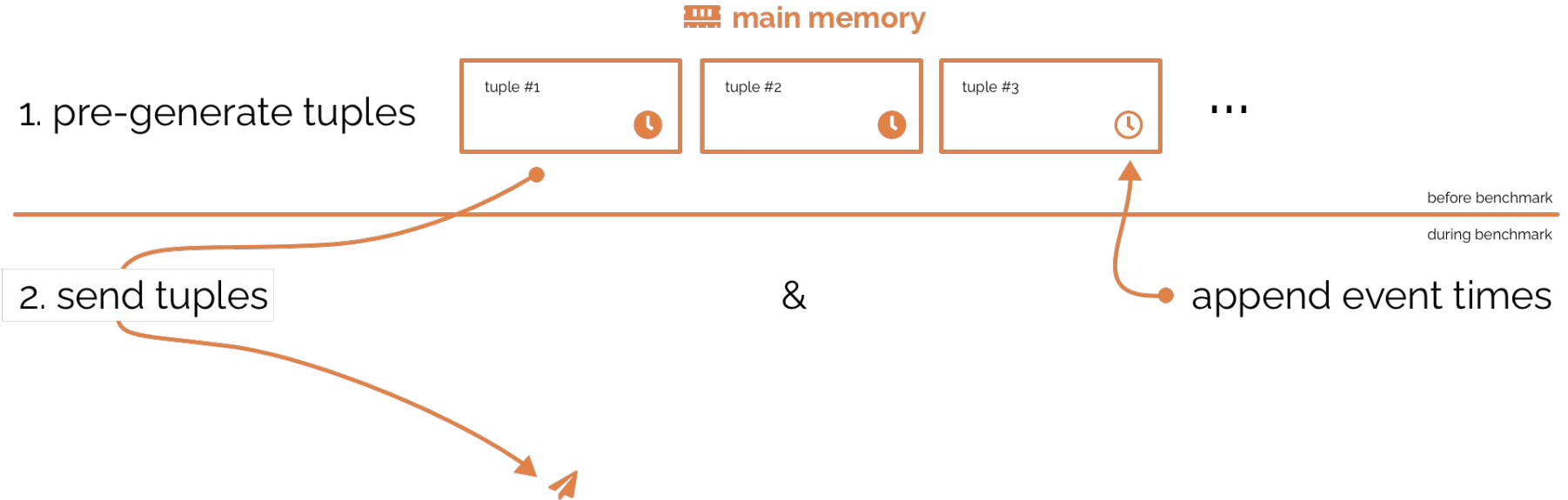
Sample Query

```
SELECT
    a.ad_id, p.sum_purchases - p.sum_costs
FROM
    (SELECT ad_id, SUM(value) as sum_purchases
    FROM PURCHASES GROUP BY ad_id) as p,

    (SELECT ad_id, SUM(cost) as sum_costs
    FROM ADS GROUP BY ad_id) as a
WHERE
    p.ad_id == a.ad_id AND p.ad_id != 0;
```

Implementation

Data Generator

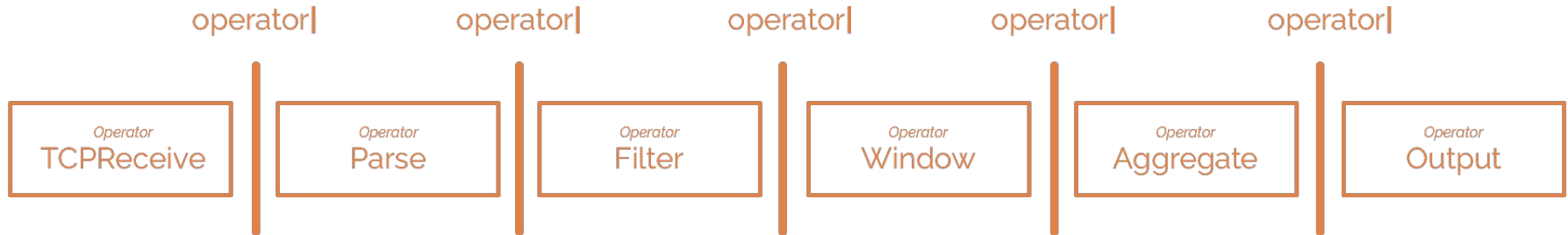


Streaming Engines

#	engine	purpose	expectation
1	query implemented into Apache Flink (JVM)	baseline	lower bound
2	query implemented into an iterator style C++ engine	baseline	lower bound
3	highly optimized hardcoded C++ query implementation	baseline	upper bound
4	C++ engine generating a distributed, compiled C++ query	<i>evaluation</i>	(best engine ever!)



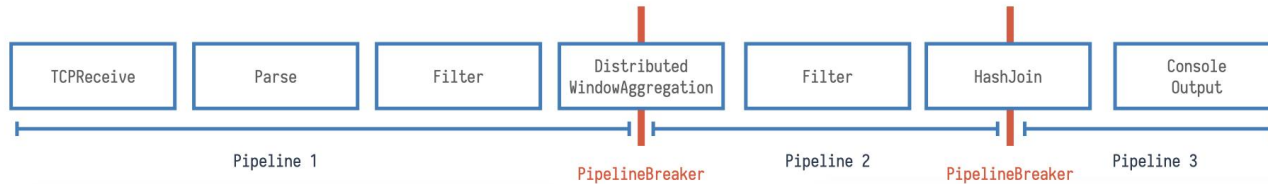
Query compilation



Query compilation

```
1 struct CheckoutTuple {
2     std::shared_ptr<Datatype> purchaseId,
3     userId,
4     adId,
5     value,
6     eventTime,
7     processingTime;
8 };
9
10 CheckoutTuple checkoutTuple{
11     POD("int", "purchaseId"),
12     POD("int", "userId"),
13     POD("int", "adId"),
14     POD("double", "value"),
15     POD("uint64_t", "eventTime"),
16     POD("uint64_t", "processingTime")
17 };
18
19 auto checkoutPipeline =
20     TCPReceiver(12346)
21     | Parse(Schema(std::vector{
22         checkoutTuple.purchaseId,
23         checkoutTuple.userId,
24         checkoutTuple.adId,
25         checkoutTuple.value,
26         checkoutTuple.eventTime}))
27     | Filter(checkoutTuple.adId, " == 0")
28     | AppendTimestamp(checkoutTuple.processingTime)
29     | DistributedWindowAggregation(1000, 5, aggregatedCheckoutTuple);
```

Task parallelism



```

1 thread pipeline1([]{
2   while (true) {
3     TCPReceive();
4     Parse();
5     Filter();
6     if (DistributedWindowAggregation::preprocess())
7       semaphore1.signal();
8   }
9 });

```

```

1 thread pipeline2([]{
2   while (true) {
3     semaphore1.wait();
4     DistributedWindowAggregation();
5     Filter();
6     if (HashJoin::preprocess())
7       semaphore2.signal();
8   }
9 });

```

```

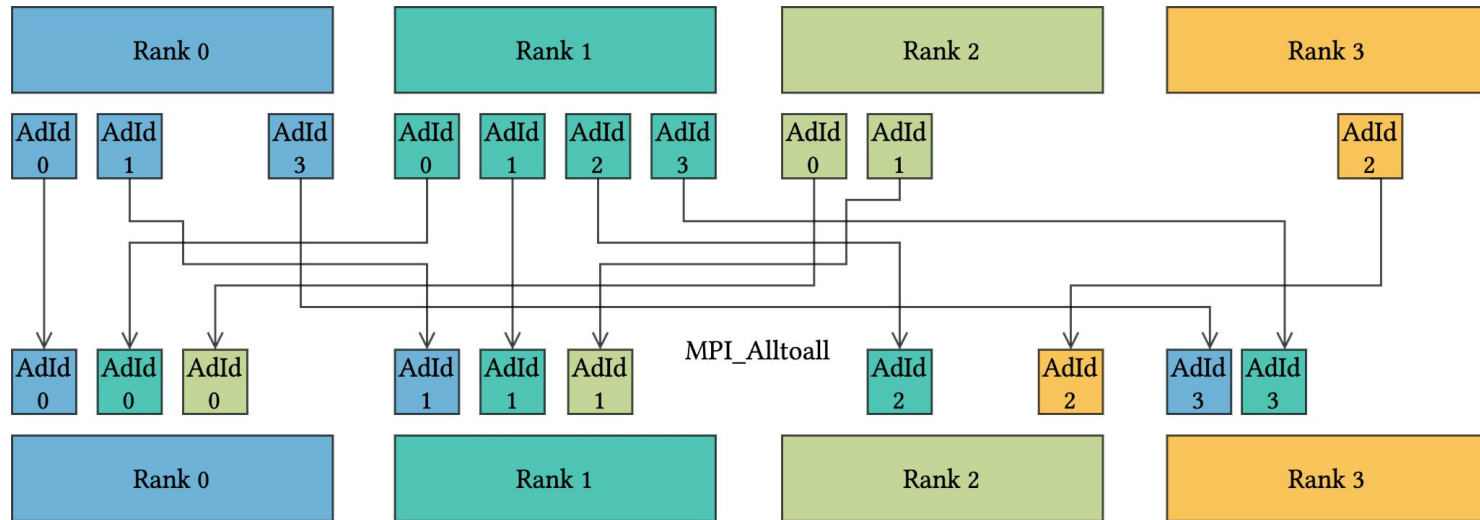
1 thread pipeline3([]{
2   while (true) {
3     semaphore2.wait();
4     semaphore3.wait(); // from other Stream
5     HashJoin();
6     ConsoleOutput();
7   }
8 });

```


Distribution



Distribution



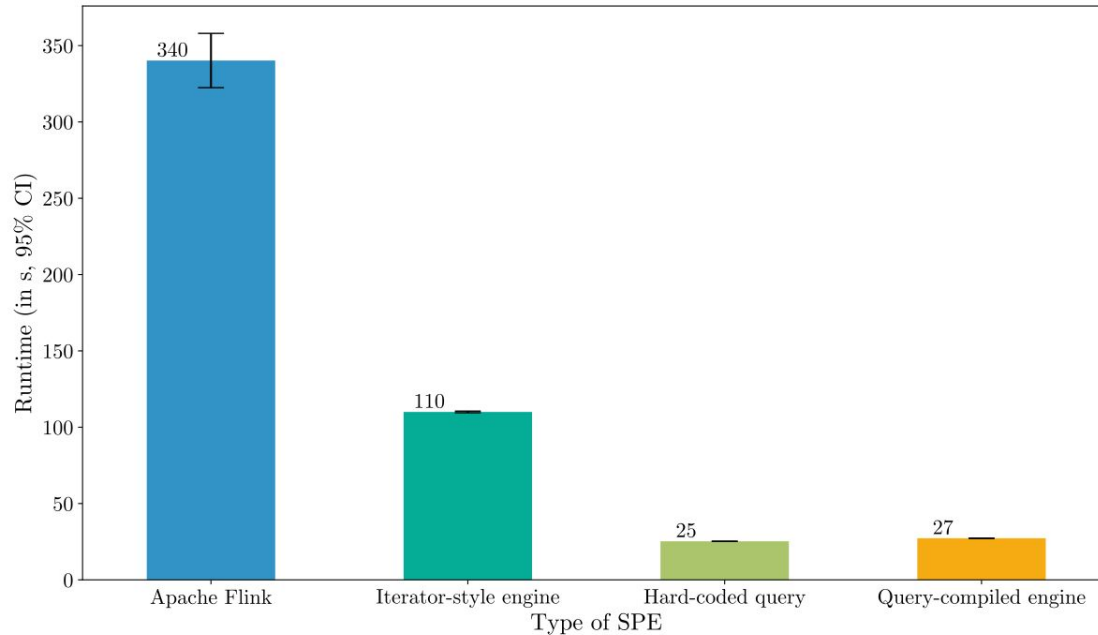
Evaluation



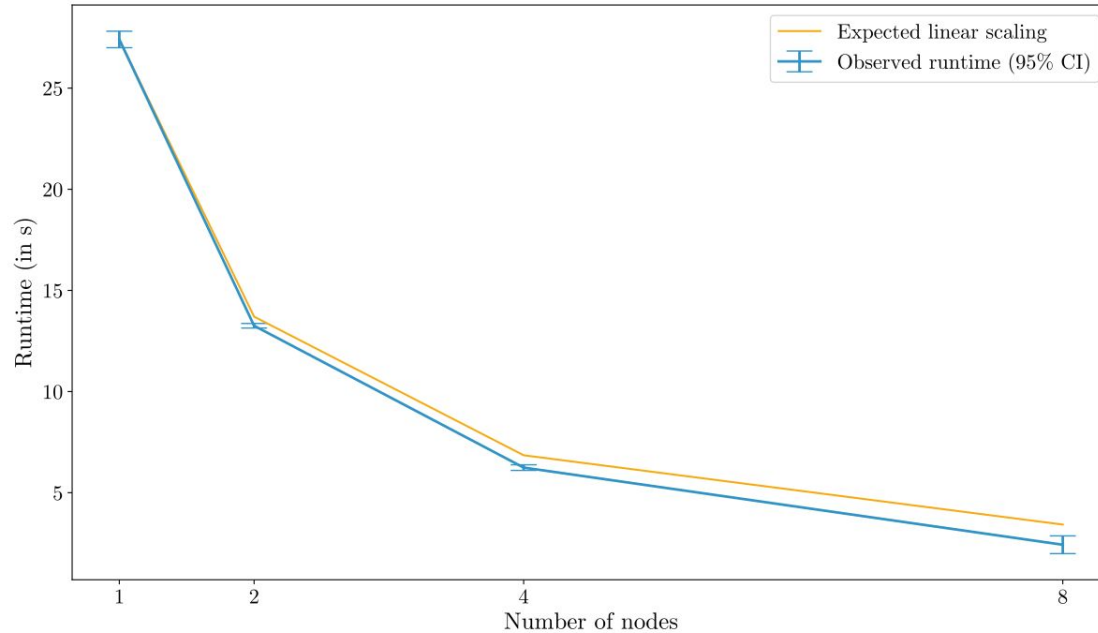
Experimental Setup

- 16 nodes cluster of Score Lab
 - 2x Intel Xeon Gold 5220S CPU
 - 95 GB RAM
 - 25 Gbit/s Ethernet networking
- numactl to bind process and memory allocation
 - Different NUMA nodes for generator and SPE
- Each experiment conducted 5 times

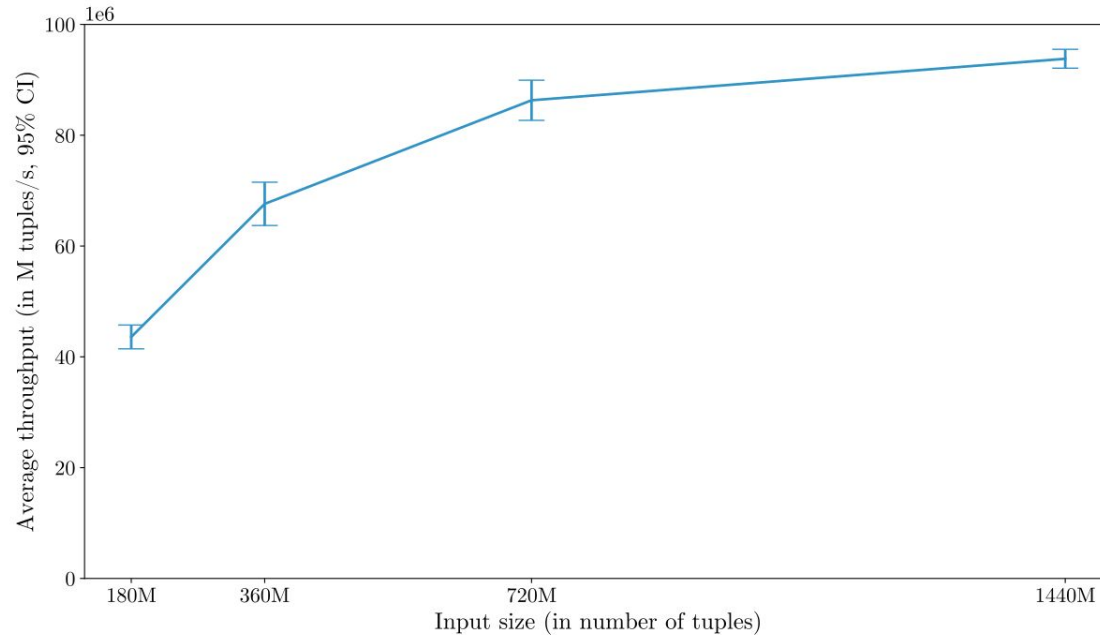
Experiment: Comparing SPEs



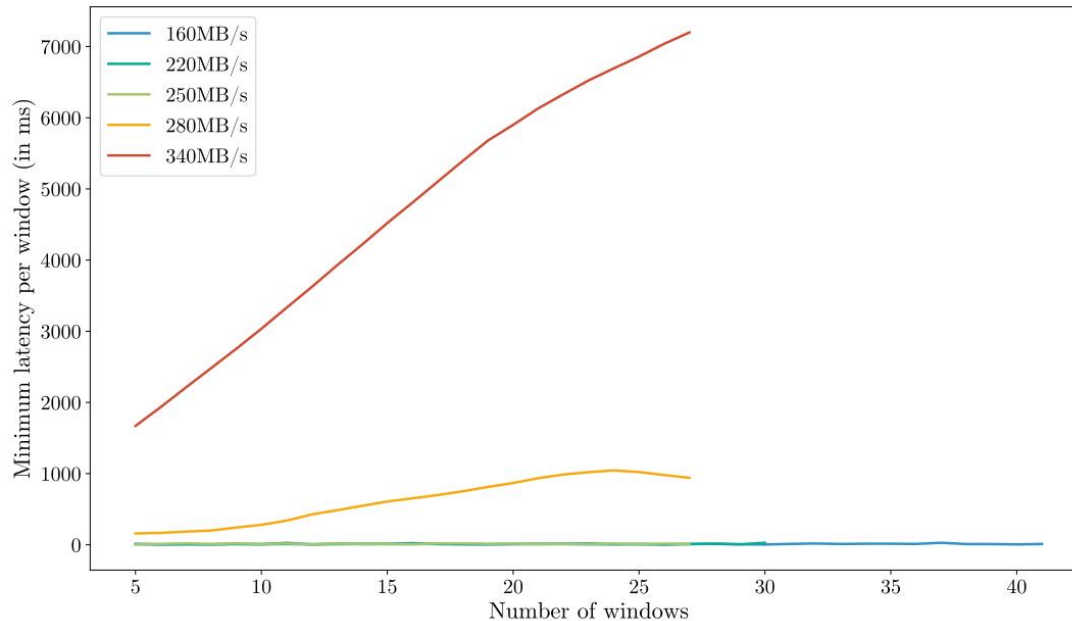
Experiment: Scaling number of nodes



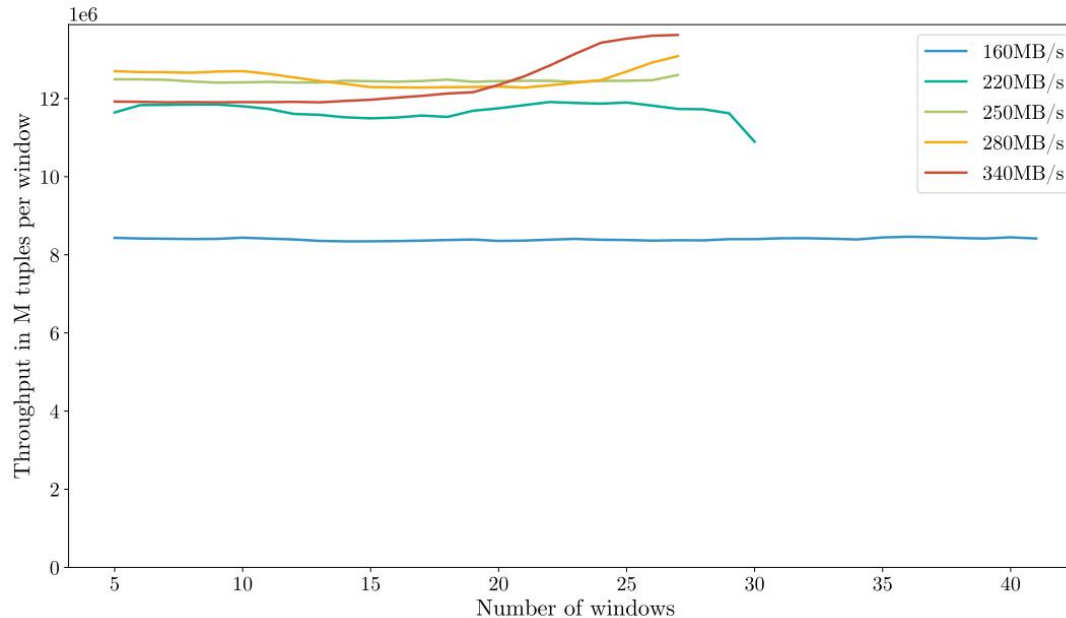
Experiment: Scaling data sizes



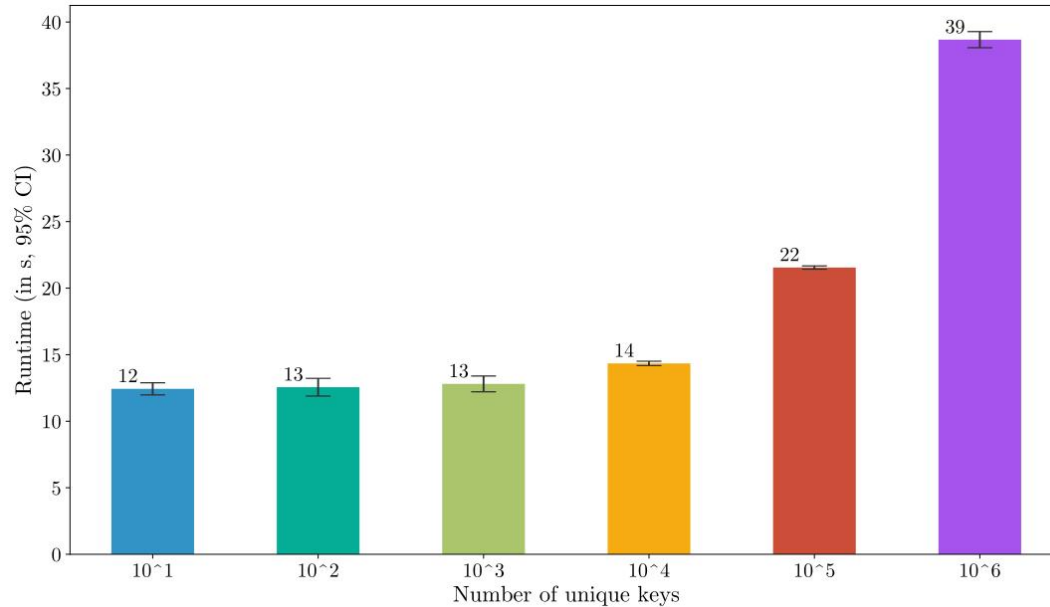
Experiment: Sustainable Throughput (1)



Experiment: Sustainable Throughput (2)



Experiment: Varying Key Ranges





Discussion

- ✓ Included: Performance comparison with Flink and baseline approaches on a single node
 - Scale-Out experiments show real-world behaviour (multi-node cluster and unbounded stream)
 - Performance impact of data rate and key range

- ✗ To Do: Distribute Apache Flink and compare to our engine prototype on n nodes
 - Evaluate other workloads (i.e., different user-defined queries)

Conclusion



Conclusion

- Combining **query compilation** and **distribution** ...
 - ... is practically feasible
 - ... shows significant performance improvements
- Our prototype achieves 12.6× higher throughput than Flink, and scales well when distributing

Future Work:

- Further extend our evaluation as discussed previously
- Move away from a prototype towards a more complete streaming engine