



Stephan Ewen
@stephanewen

The Stream Processor as a Database

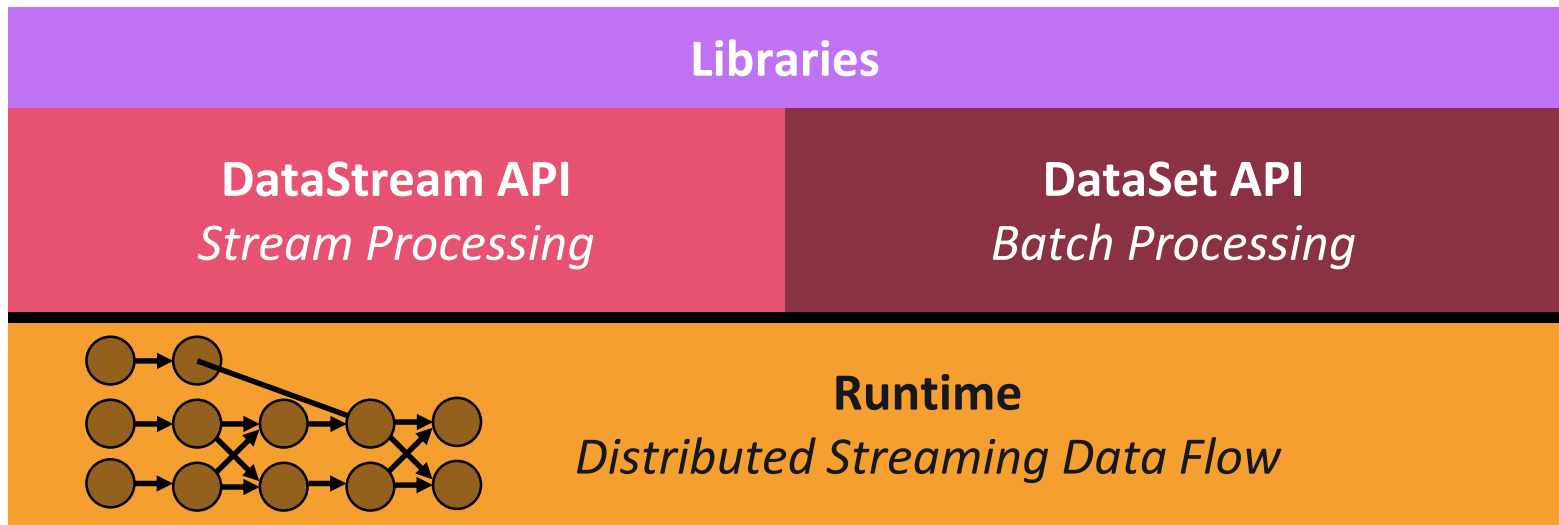
Apache Flink

dataArtisans



Streaming technology is enabling the obvious:
continuous processing on data that is
continuously produced

Apache Flink Stack



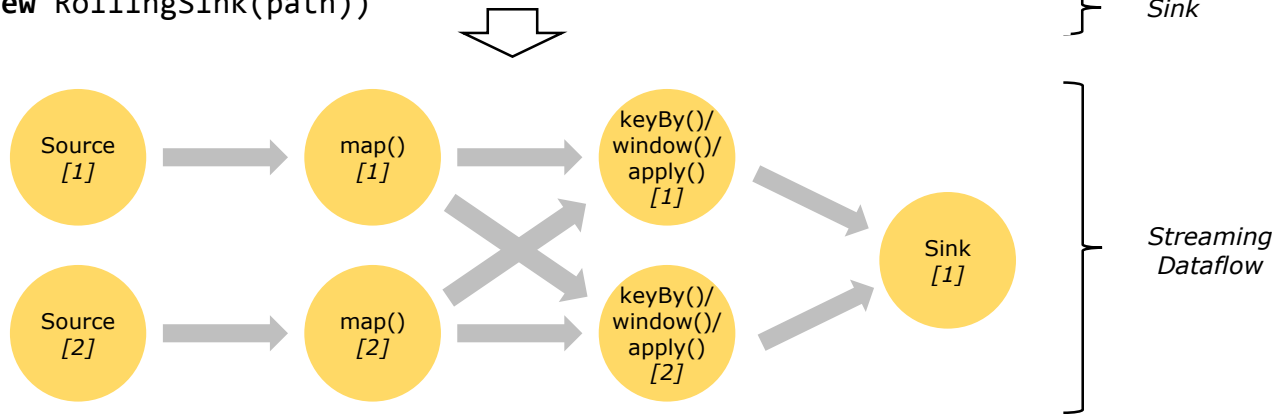
Streaming and batch as first class citizens.

Programs and Dataflows

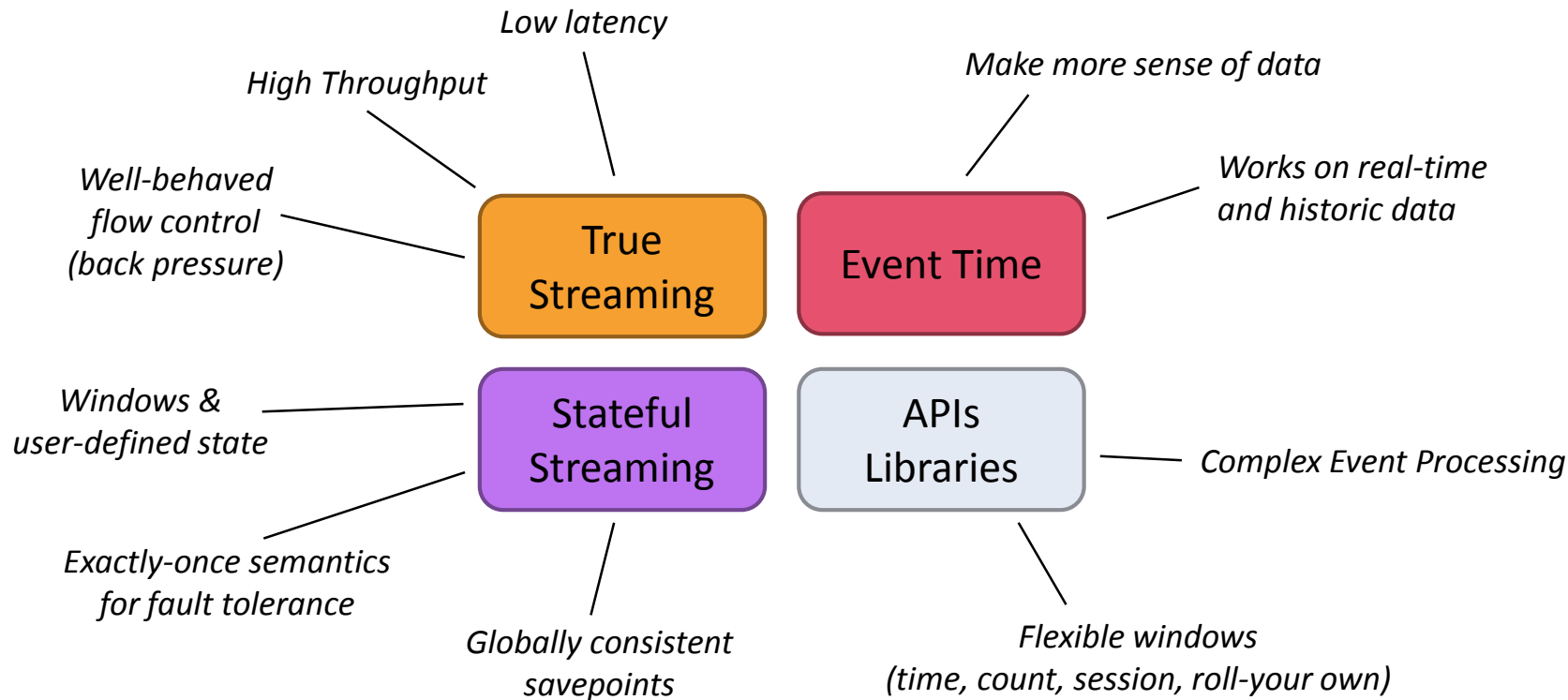


```
val lines: DataStream[String] = env.addSource(new FlinkKafkaConsumer09(...))  
  
val events: DataStream[Event] = lines.map((line) => parse(line))  
  
val stats: DataStream[Statistic] = stream  
  .keyBy("sensor")  
  .timeWindow(Time.seconds(5))  
  .apply(new MyAggregationFunction())  
  
stats.addSink(new RollingSink(path))
```

Source
Transformation
Transformation
Sink



What makes Flink flink?





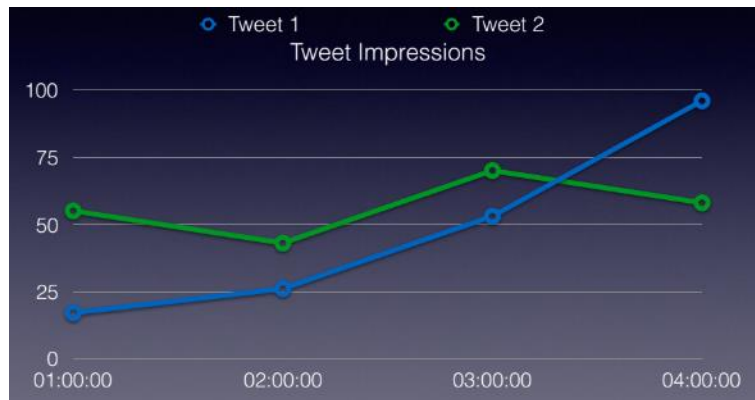
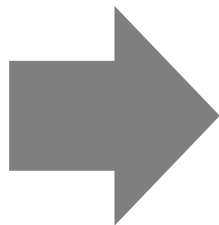
Realtime Counts and Aggregates

The (Classic) Use Case

(Real)Time Series Statistics



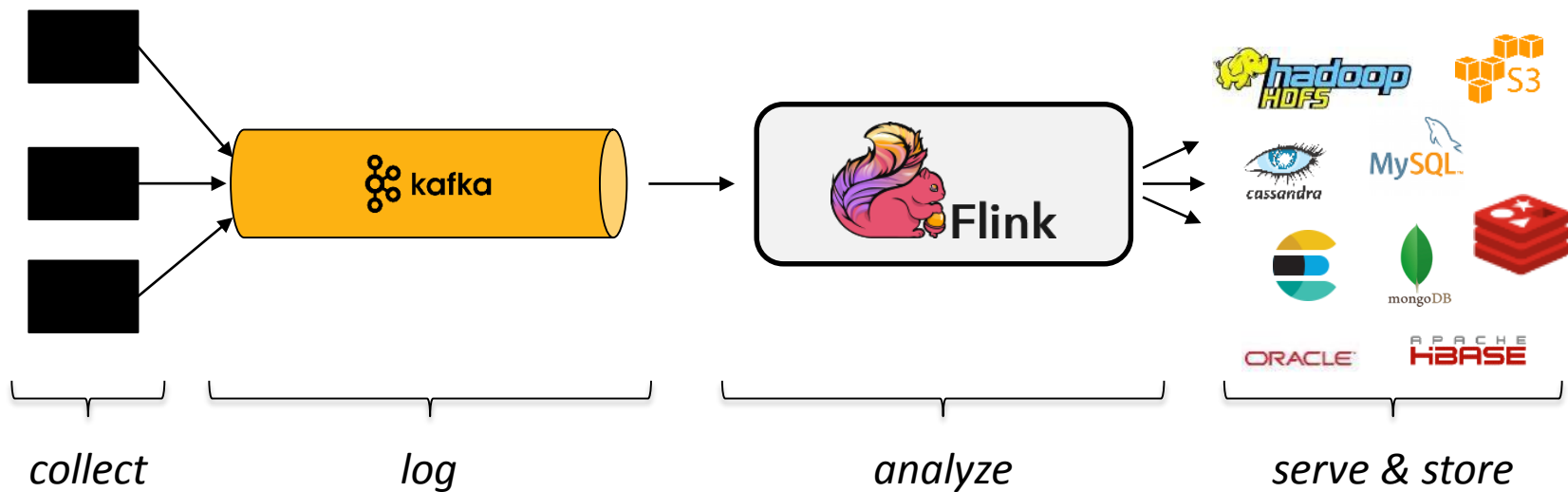
Stream
tweet-id: 1, event: url-click, time: 01:01:01
tweet-id: 2, event: url-click, time: 01:01:02
tweet-id: 1, event: impression, time: 01:01:03
tweet-id: 2, event: url-click, time: 02:01:01
tweet-id: 1, event: impression, time: 02:02:02



stream of events

realtime statistics

The Architecture



The Flink Job



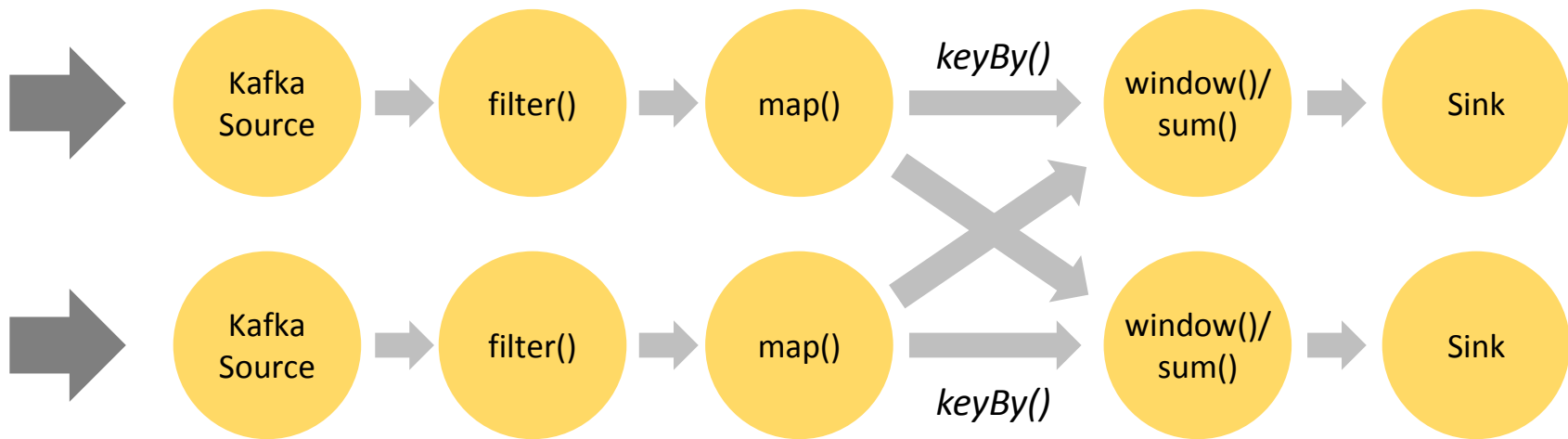
```
case class Impressions(id: String, impressions: Long)

val events: DataStream[Event] =
    env.addSource(new FlinkKafkaConsumer09(...))

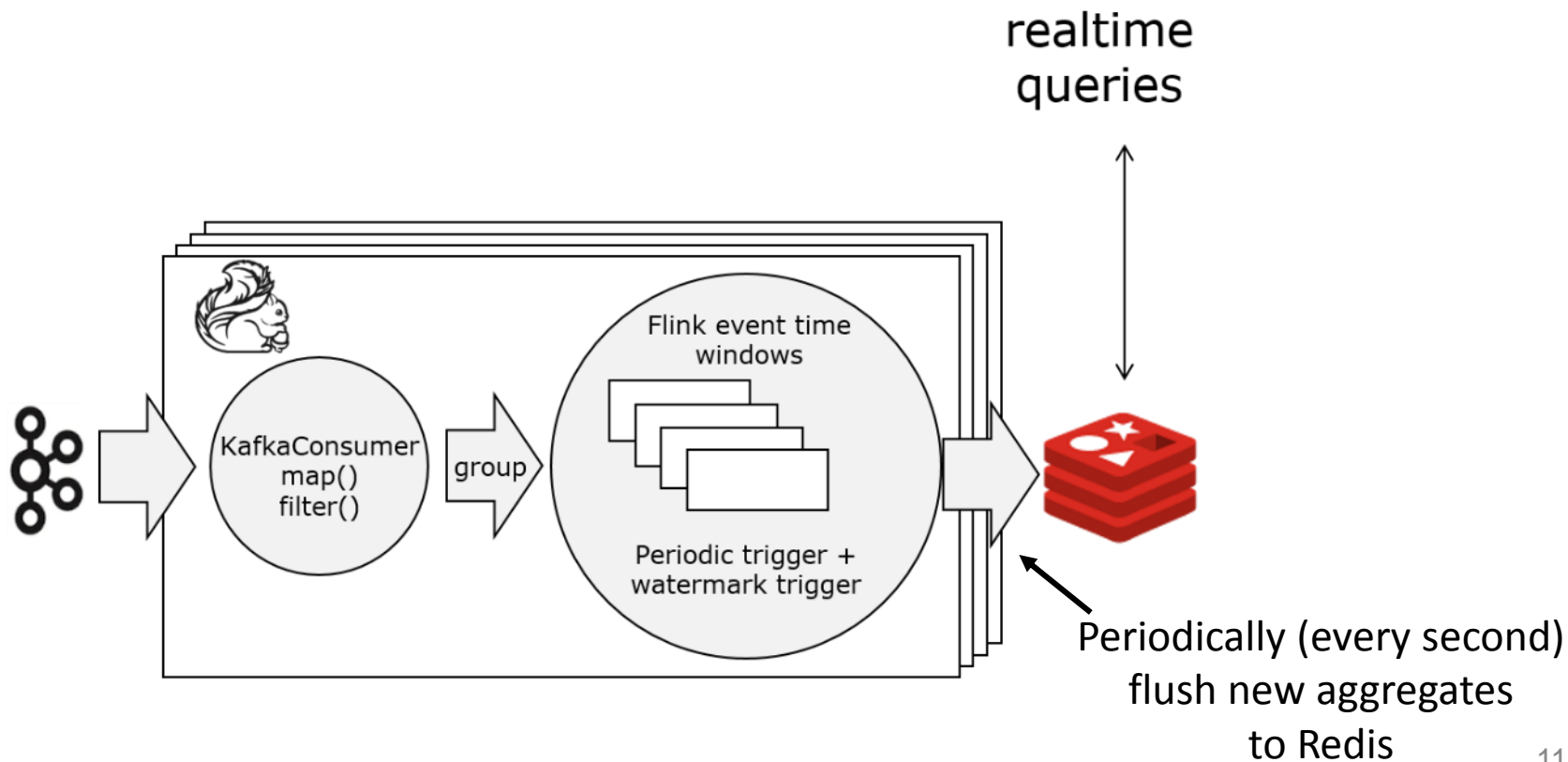
val impressions: DataStream[Impressions] = events
    .filter(evt => evt.isImpression)
    .map(evt => Impressions(evt.id, evt.numImpressions))

val counts: DataStream[Impressions]= stream
    .keyBy("id")
    .timeWindow(Time.hours(1))
    .sum("impressions")
```

The Flink Job



Putting it all together



How does it perform?

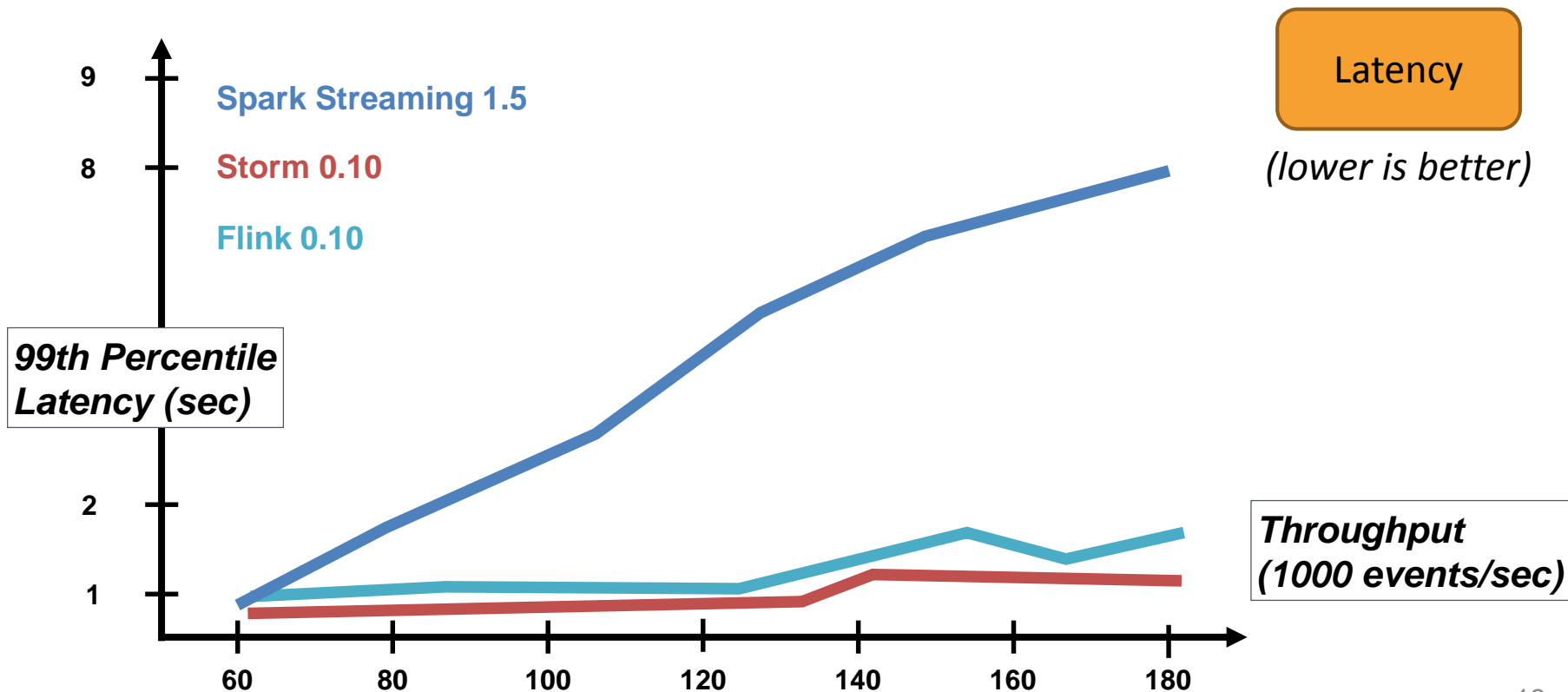


Latency

Throughput

Number of
Keys

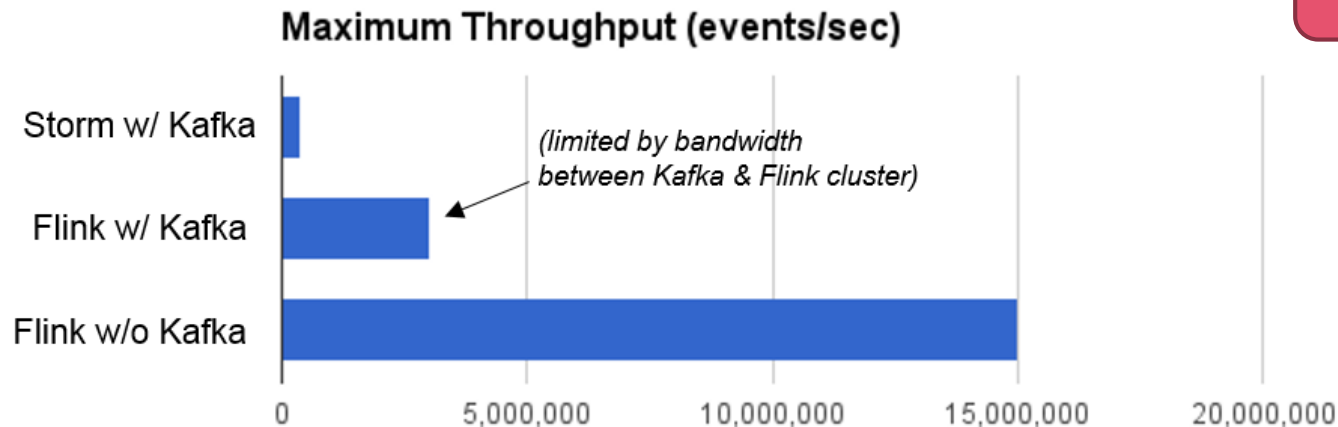
Yahoo! Streaming Benchmark



Extended Benchmark: Throughput



Throughput



- 10 Kafka brokers with 2 partitions each
- 10 compute machines (Flink / Storm)
 - Xeon E3-1230-V2@3.30GHz CPU (4 cores HT)
 - 32 GB RAM (only 8GB allocated to JVMs)
- 10 GigE Ethernet between compute nodes
- 1 GigE Ethernet between Kafka cluster and Flink nodes

Scaling Number of Users

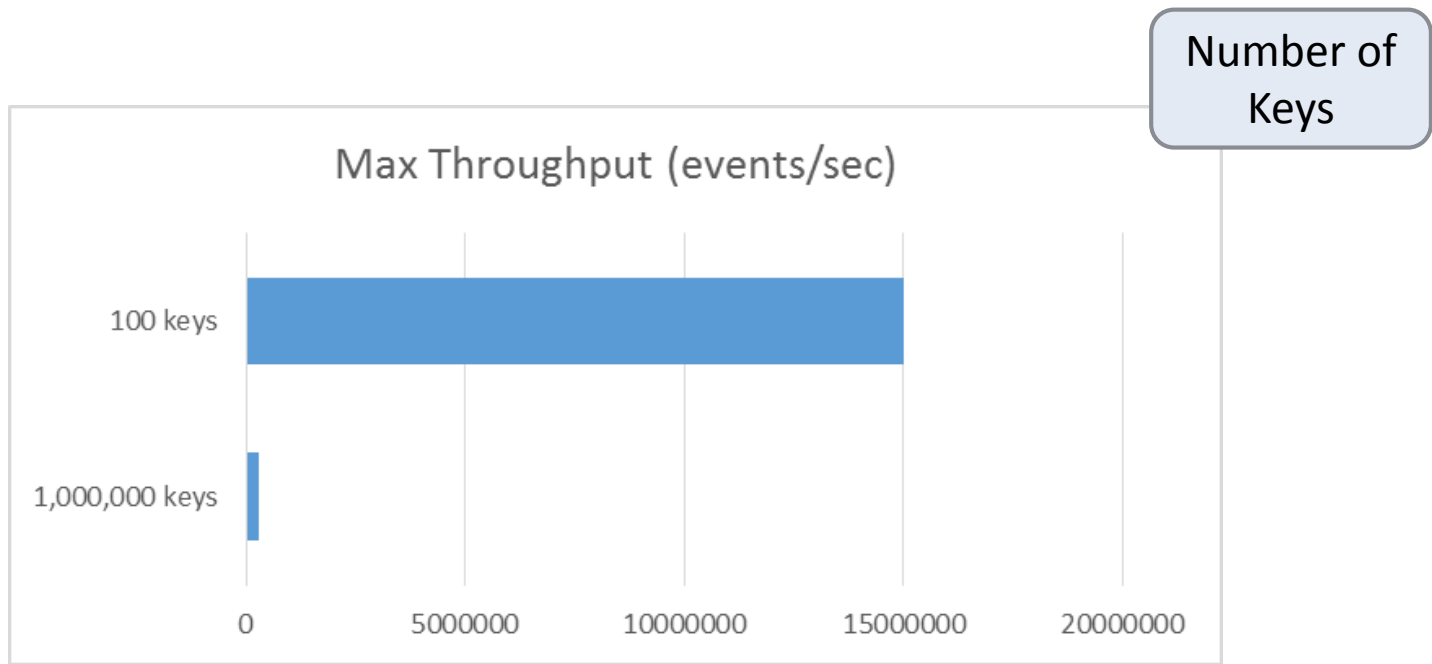


- Yahoo! Streaming Benchmark has 100 keys only
 - Every second, only 100 keys are written to key/value store
 - Quite few, compared to many real world use cases

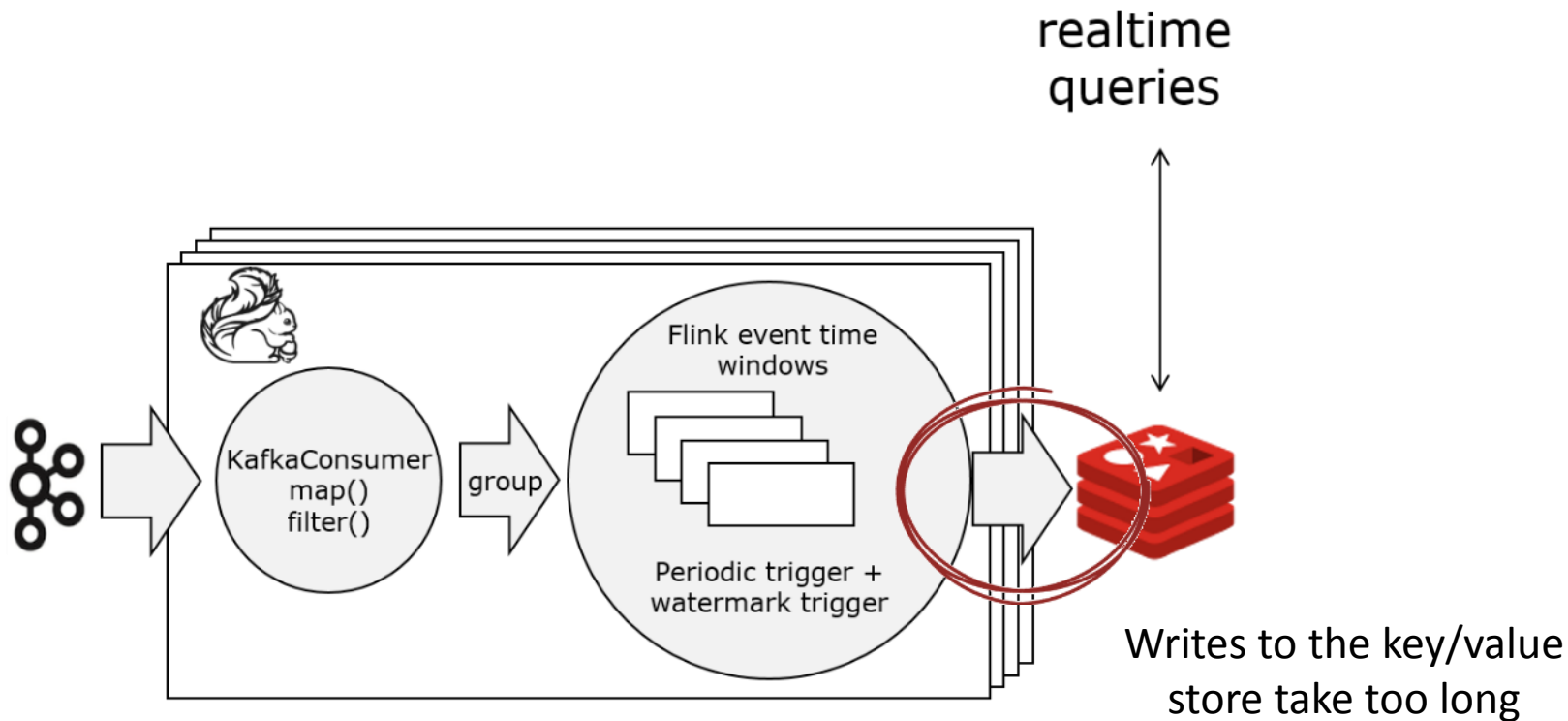
- Tweet impressions: millions keys/hour
 - Up to millions of keys updated per second

Number of
Keys

Performance



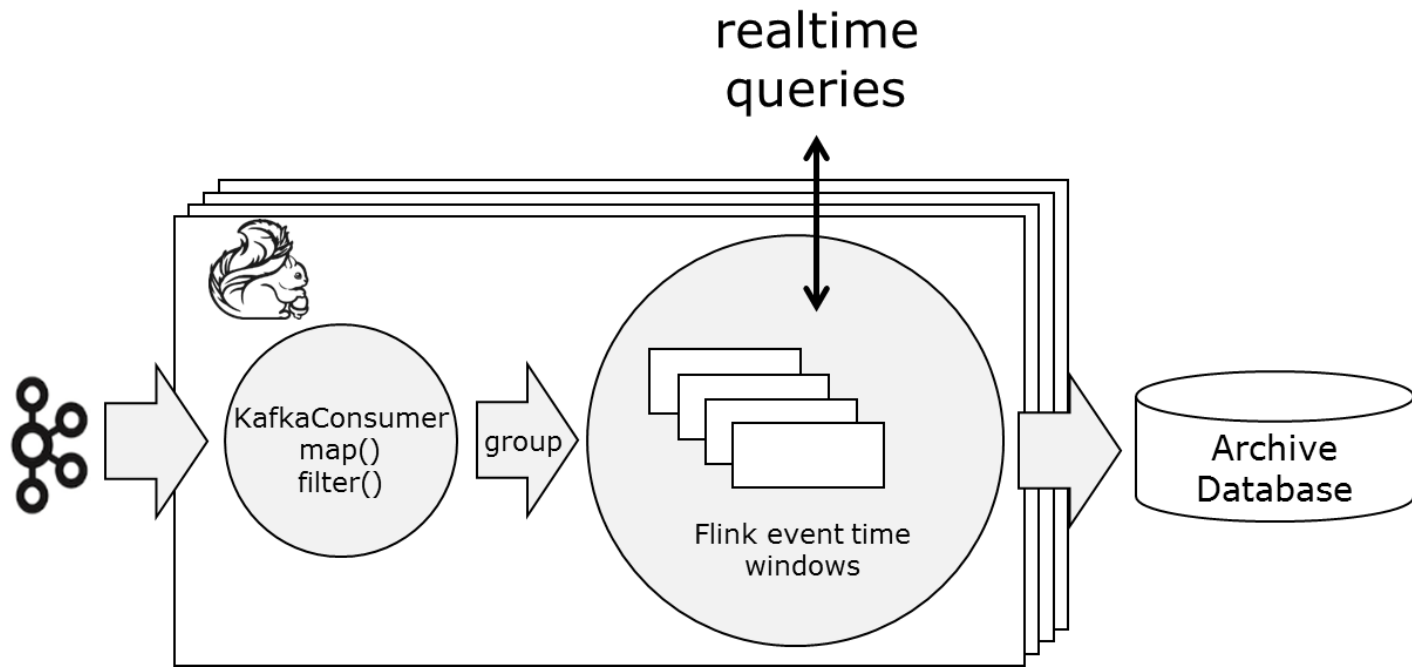
The Bottleneck



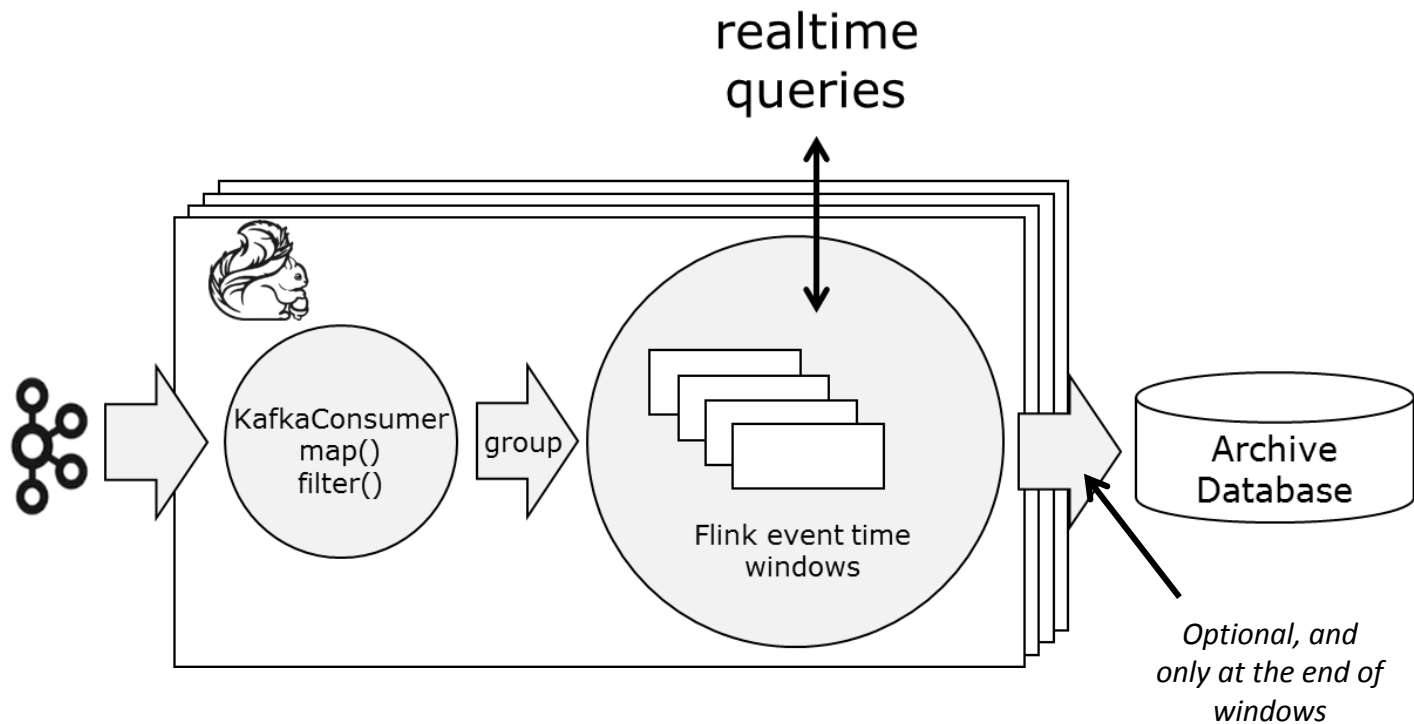


Queryable State

Queryable State



Queryable State



Queryable State Enablers



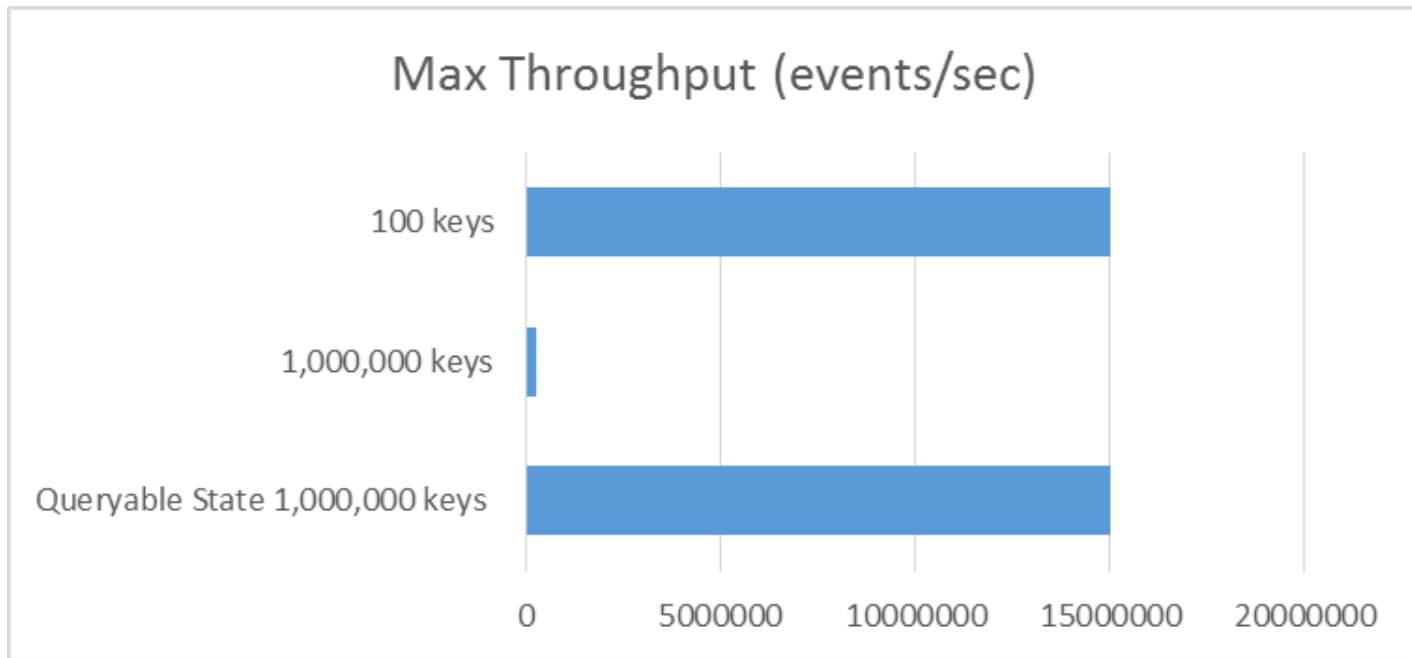
- Flink has state as a **first class citizen**
- State is **fault tolerant** (exactly once semantics)
- State is **partitioned** (sharded) together with the operators that create/update it
- State is **continuous** (not mini batched)
- State is **scalable** (e.g., embedded RocksDB state backend)

Queryable State Status



- [FLINK-3779] / Pull Request #2051 :
Queryable State Prototype
- Design and implementation under evolution
- Some experiments were using earlier versions of the implementation
- Exact numbers may differ in final implementation, but order of magnitude is comparable

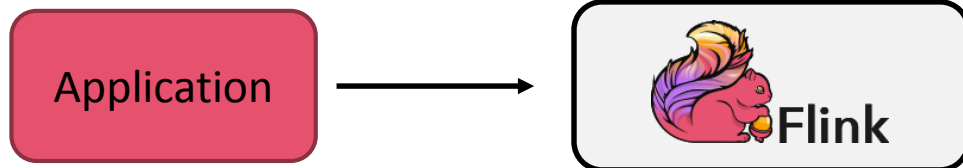
Queryable State Performance



Queryable State: Application View



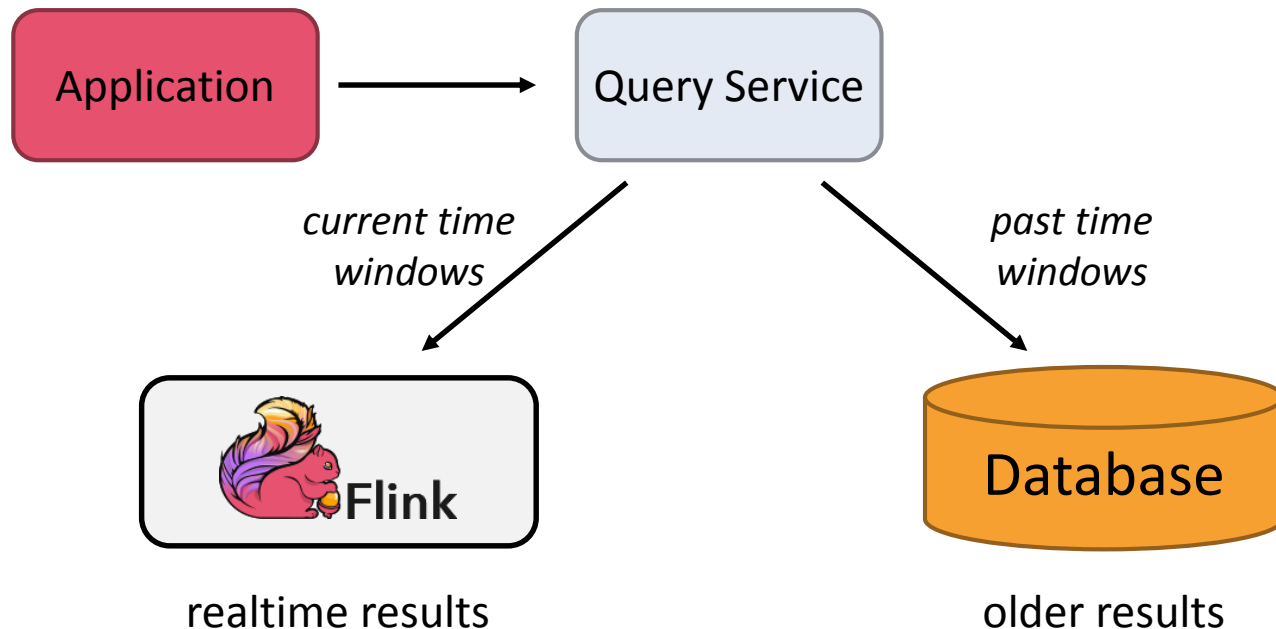
Application only interested in latest realtime results



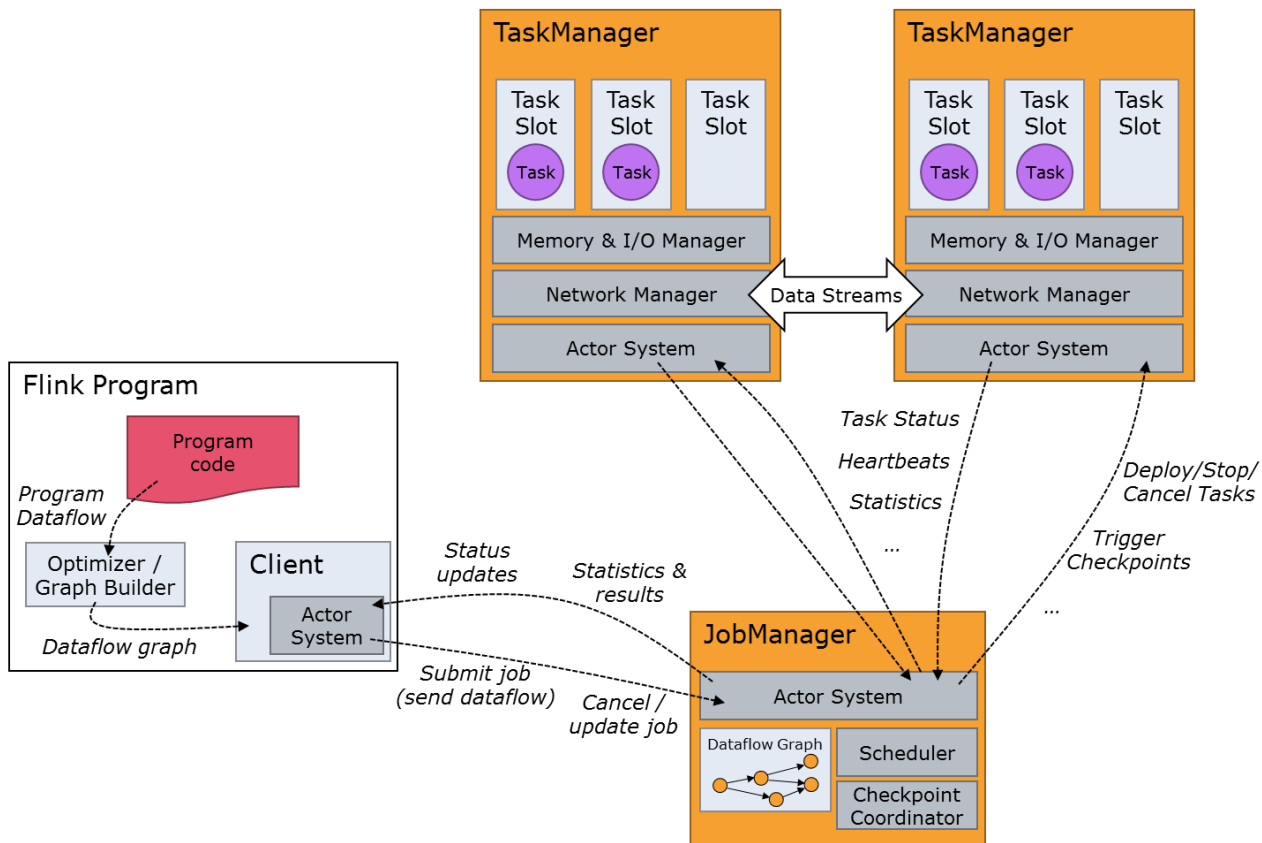
Queryable State: Application View



Application requires both latest realtime- and older results



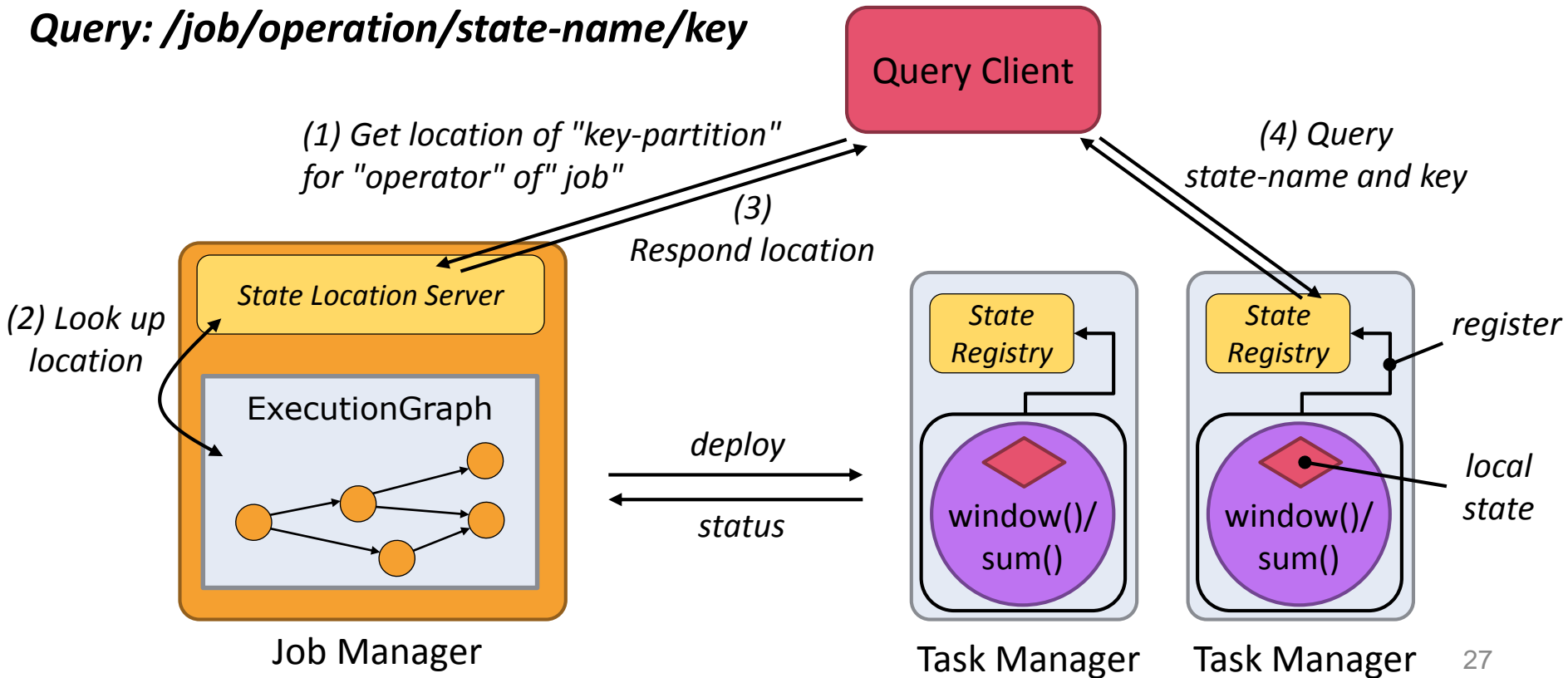
Apache Flink Architecture Review



Queryable State: Implementation



Query: /job/operation/state-name/key





Contrasting with key/value stores

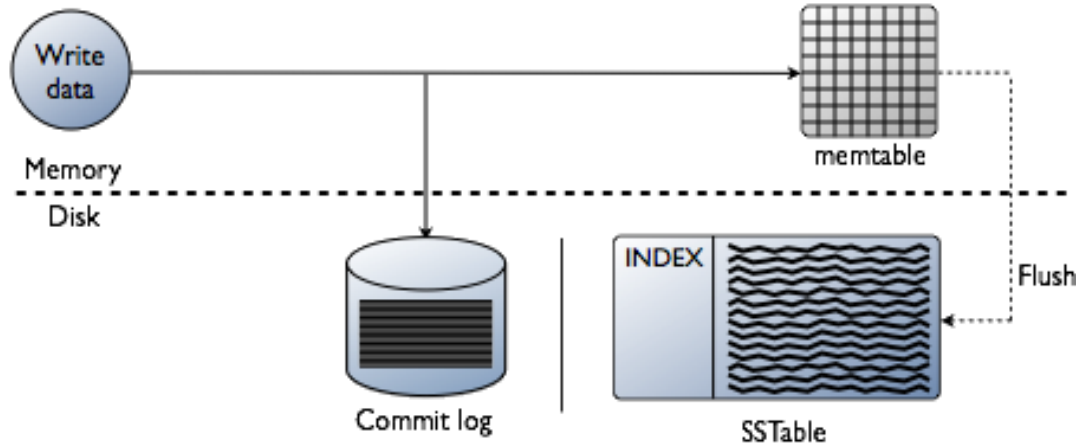
Turning the Database Inside Out



- Cf. Martin Kleppman's talks on re-designing data warehousing based on log-centric processing
- This view angle picks up some of these concepts
- Queryable State in Apache Flink = (Turning DB inside out)++



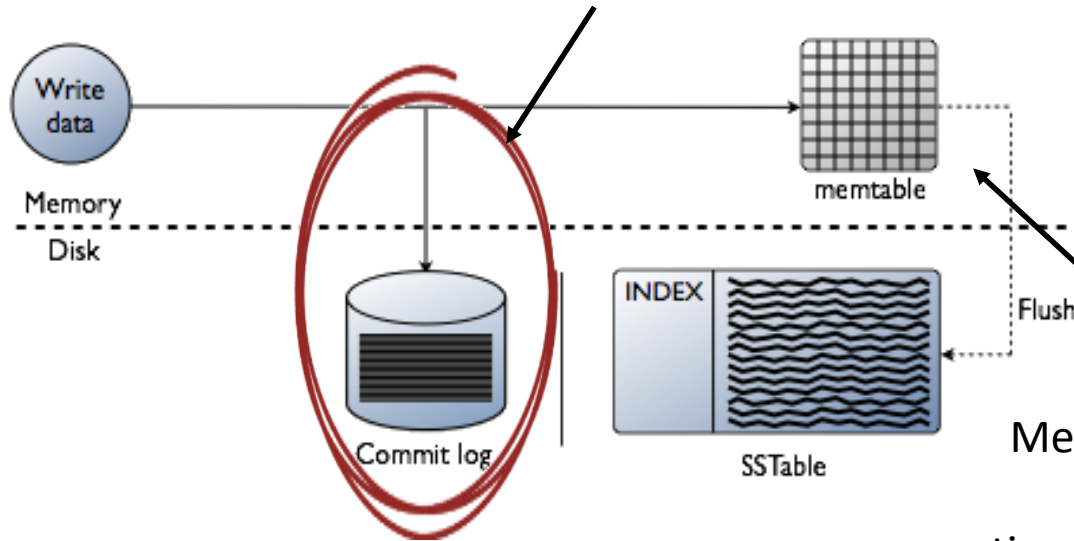
Write Path in Cassandra (simplified)



Write Path in Cassandra (simplified)



First step is durable write to the commit log
(in all databases that offer strong durability)

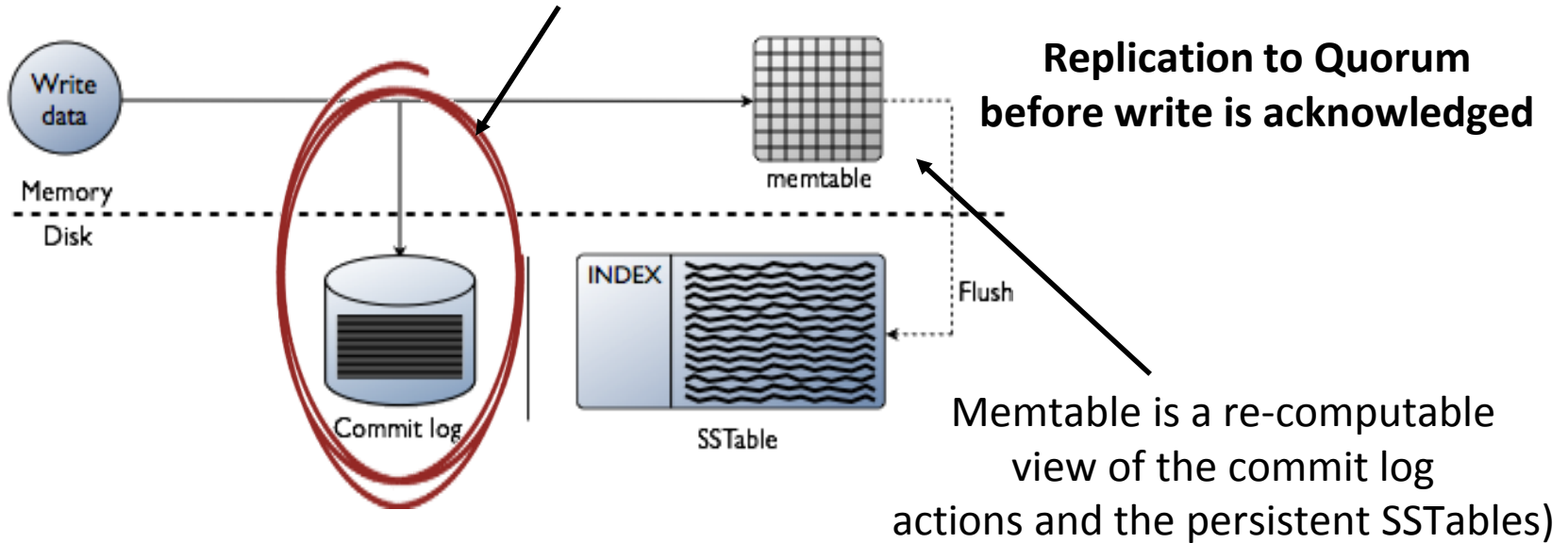


Memtable is a re-computable
view of the commit log
actions and the persistent SSTables)

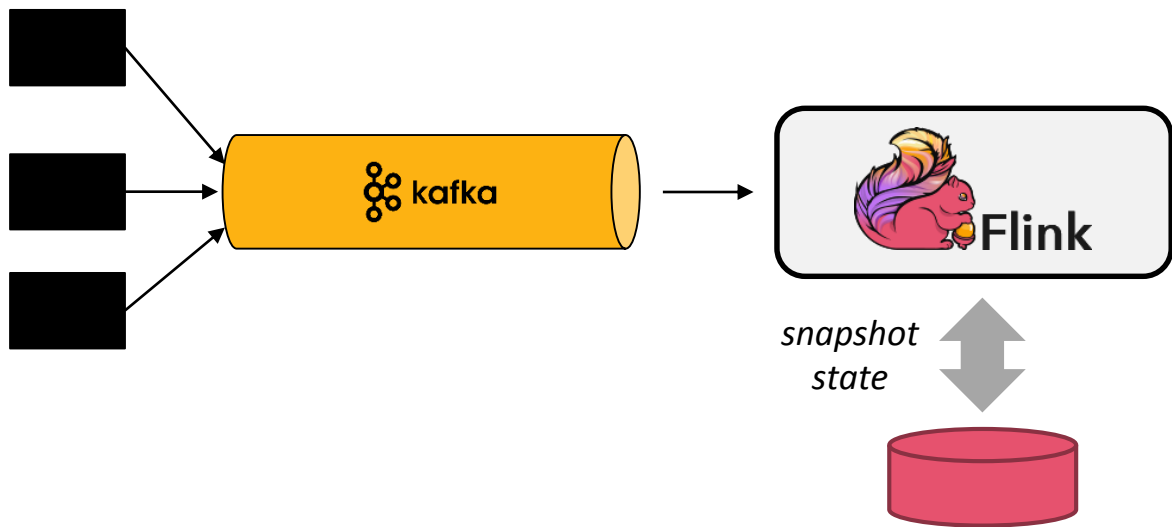
Write Path in Cassandra (simplified)



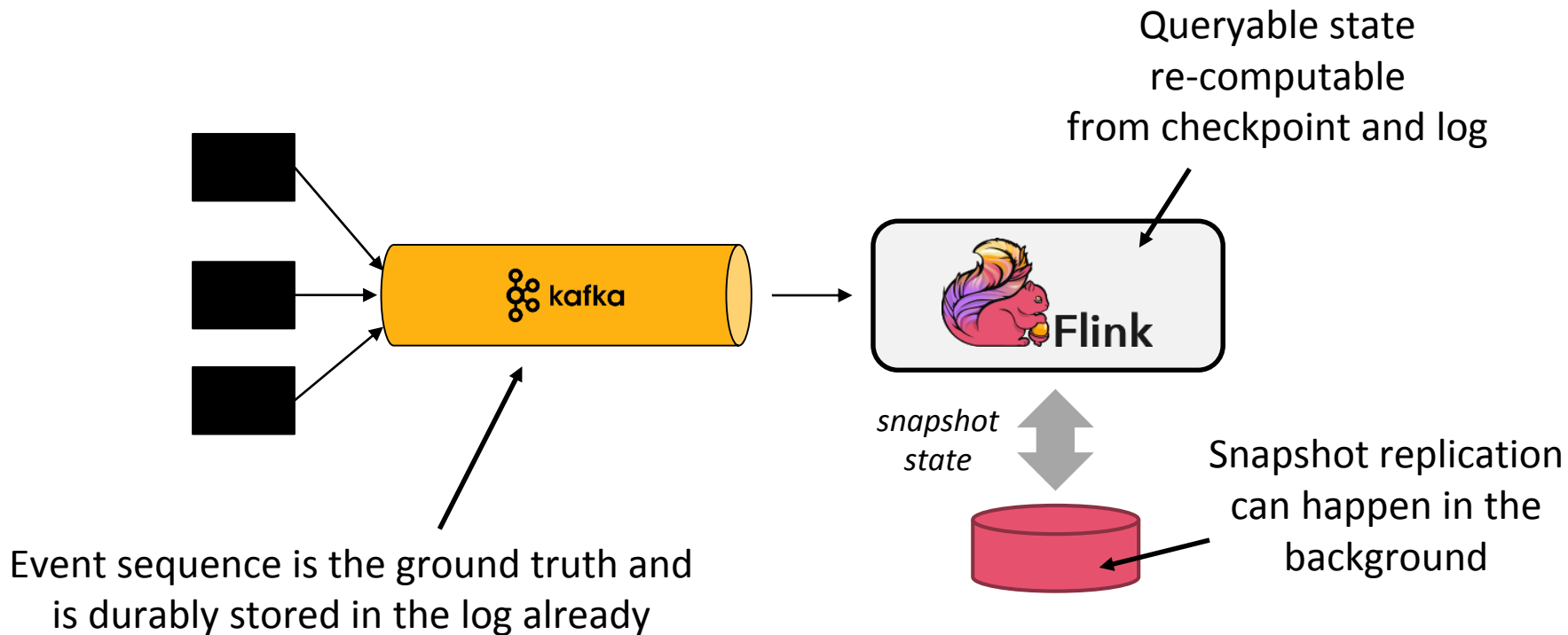
First step is durable write to the commit log
(in all databases that offer strong durability)



Durability of Queryable state



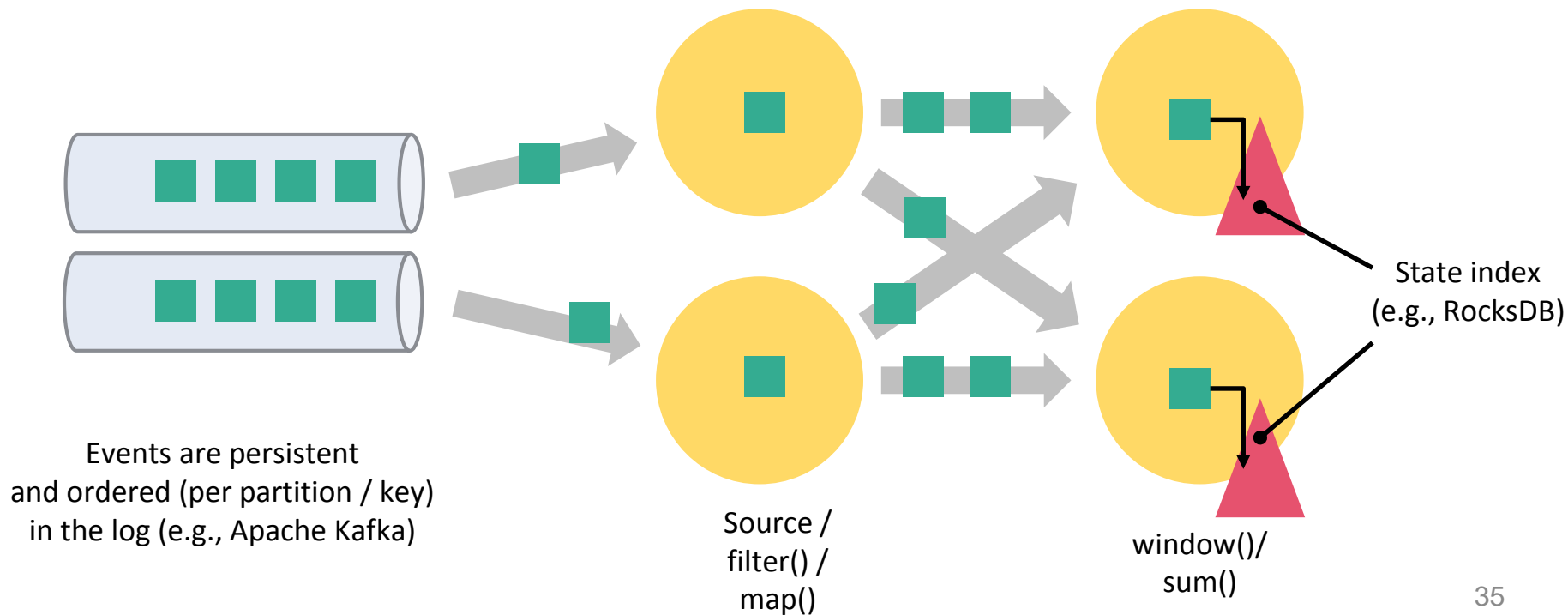
Durability of Queryable state



Performance of Flink's State



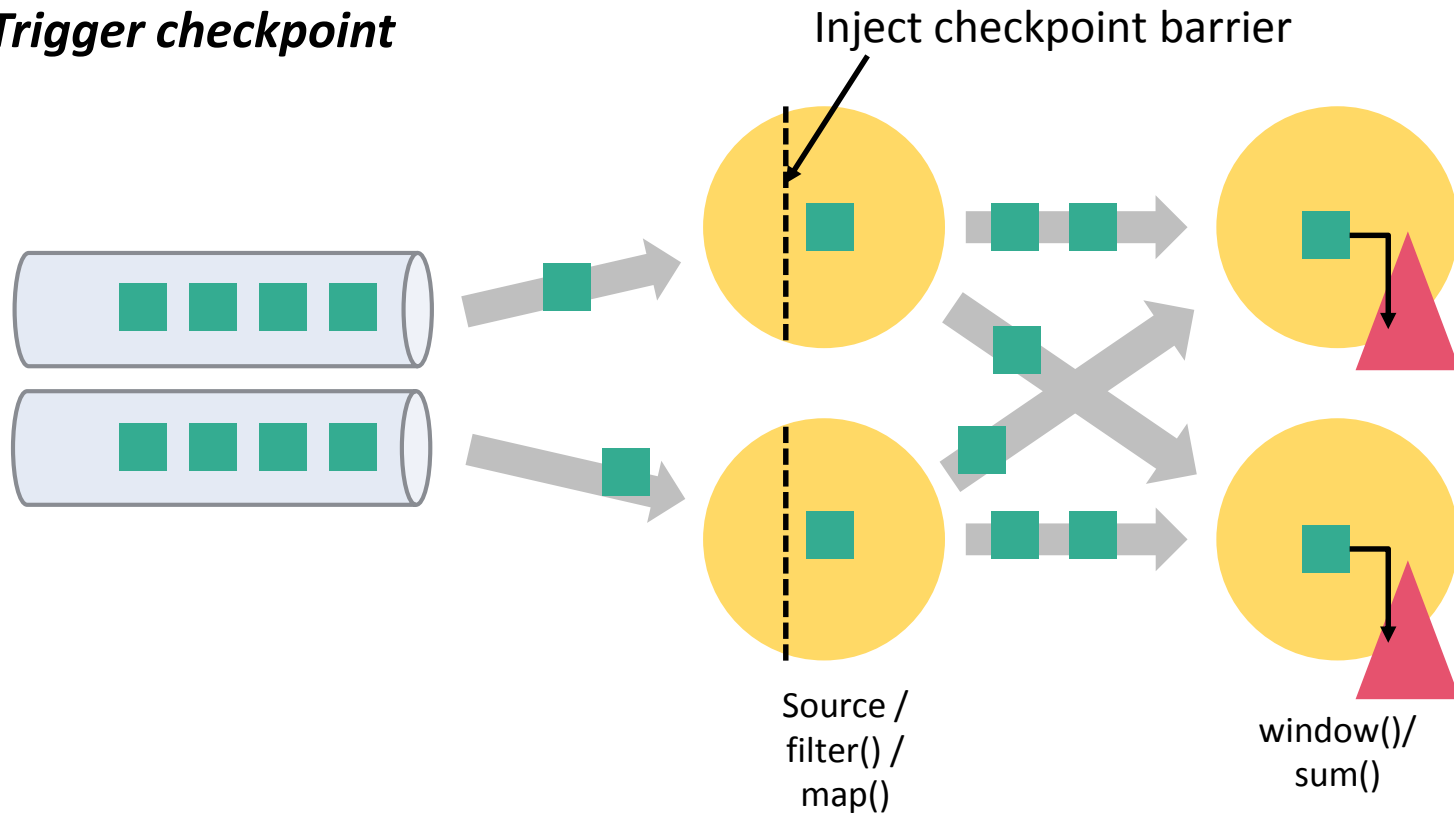
Events flow without replication or synchronous writes



Performance of Flink's State



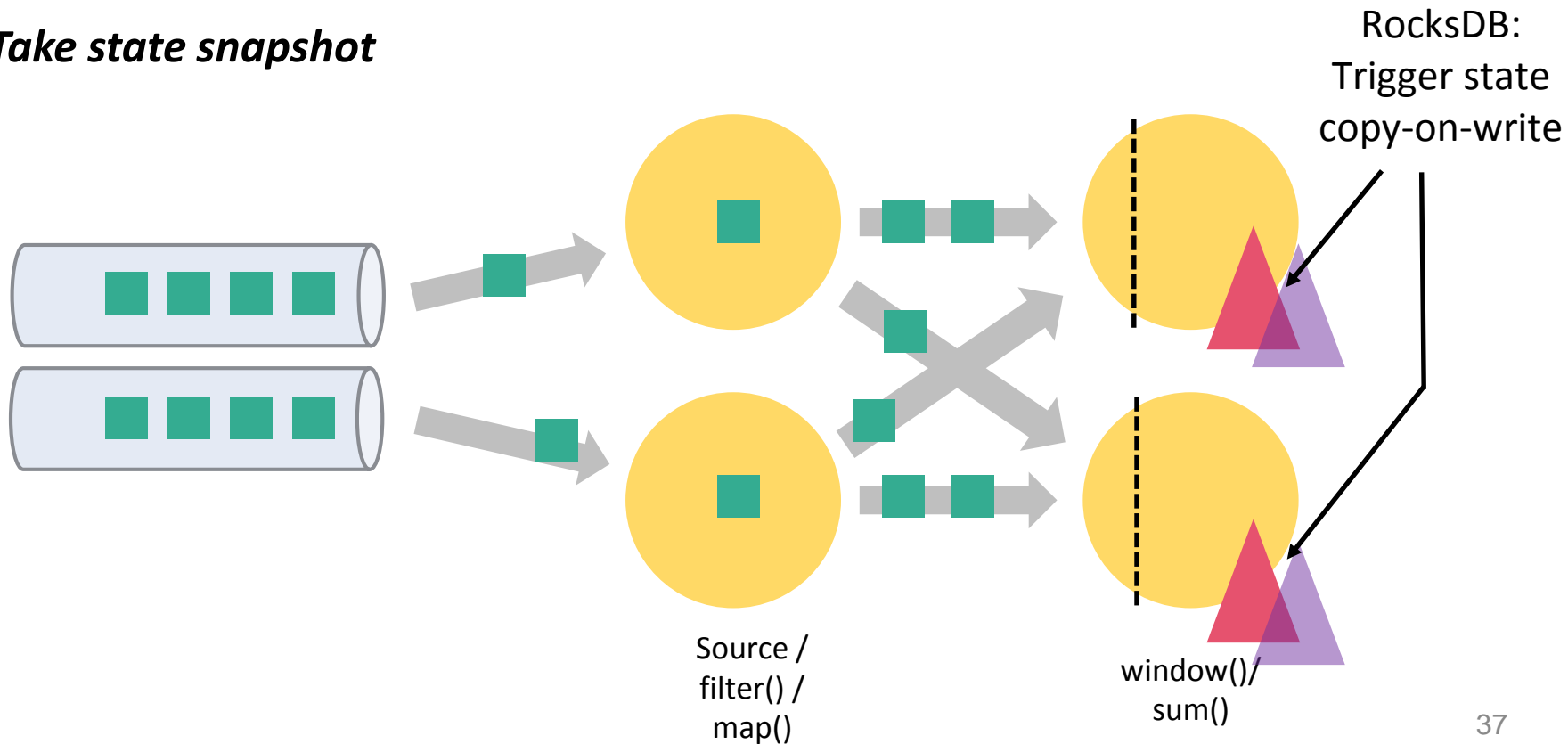
Trigger checkpoint



Performance of Flink's State



Take state snapshot



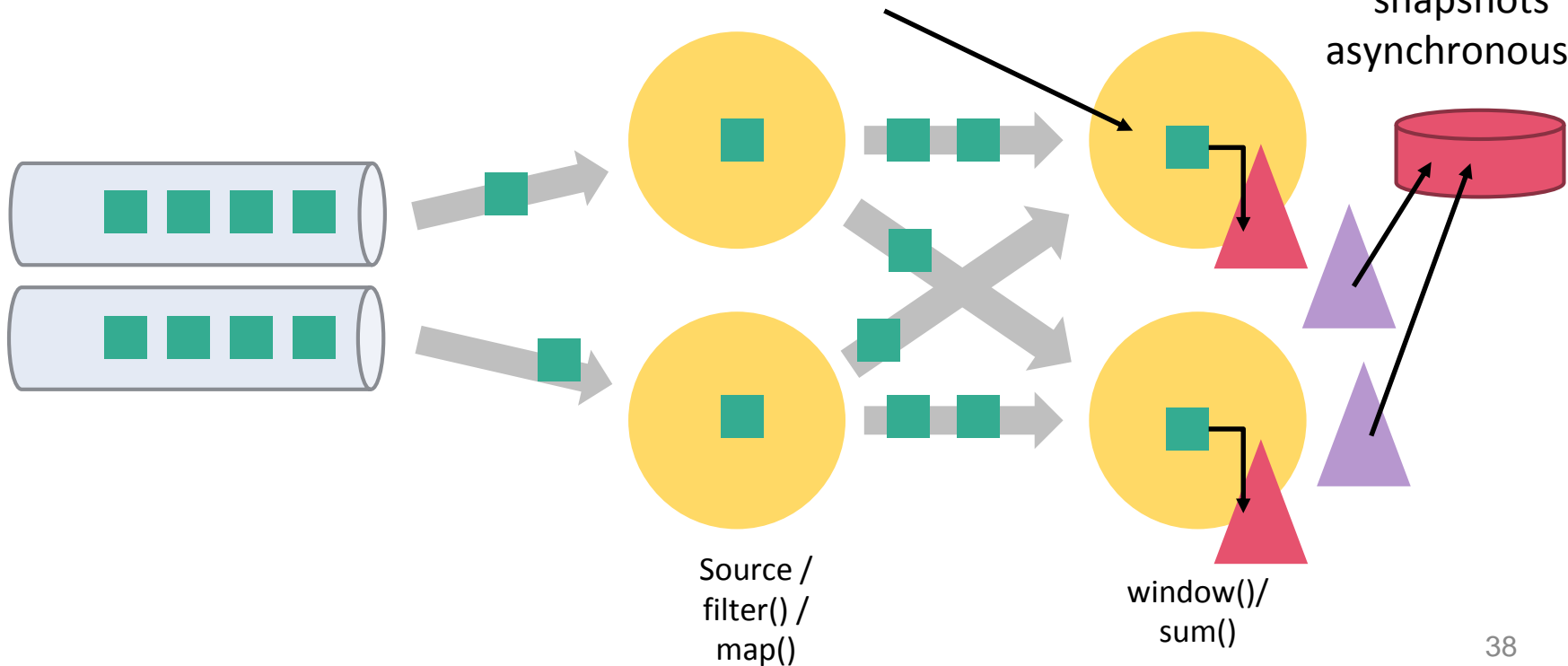
Performance of Flink's State



Persist state snapshots

Processing pipeline continues

Durably persist snapshots asynchronously





Conclusion

Takeaways



- Streaming applications are often not bound by the stream processor itself. **Cross system interaction is frequently biggest bottleneck**
- **Queryable state mitigates a big bottleneck:** Communication with external key/value stores to publish realtime results
- **Apache Flink's** sophisticated support for state makes this possible

Takeaways



Performance of Queryable State

- Data persistence is fast with logs (*Apache Kafka*)
 - Append only, and streaming replication
- Computed state is fast with local data structures and no synchronous replication (*Apache Flink*)
- Flink's checkpoint method makes computed state persistent with low overhead

Go Flink!

