

Tutorial:

The Role of Event-Time Order in Data Streaming Analysis

Vincenzo Gulisano

Chalmers University of Technology
Gothenburg, Sweden
vincenzo.gulisano@chalmers.se

Bastian Havers

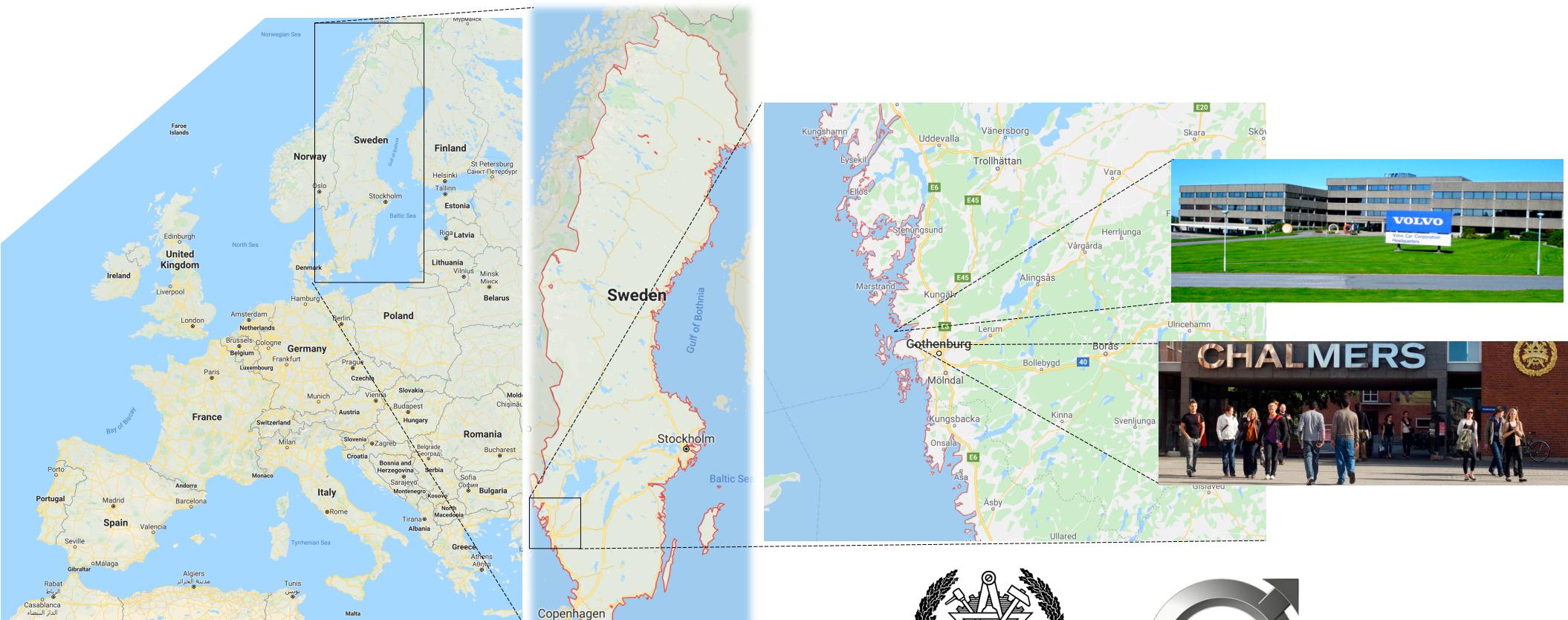
Chalmers University of Technology & Volvo Cars
Gothenburg, Sweden
havers@chalmers.se

Dimitris Palyvos-Giannas

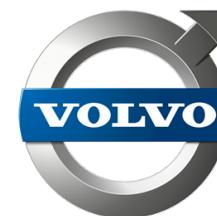
Chalmers University of Technology
Gothenburg, Sweden
palyvos@chalmers.se

Marina Papatriantafilou

Chalmers University of Technology
Gothenburg, Sweden
ptrianta@chalmers.se



Distributed Computing and Systems
Chalmers university of technology



Agenda

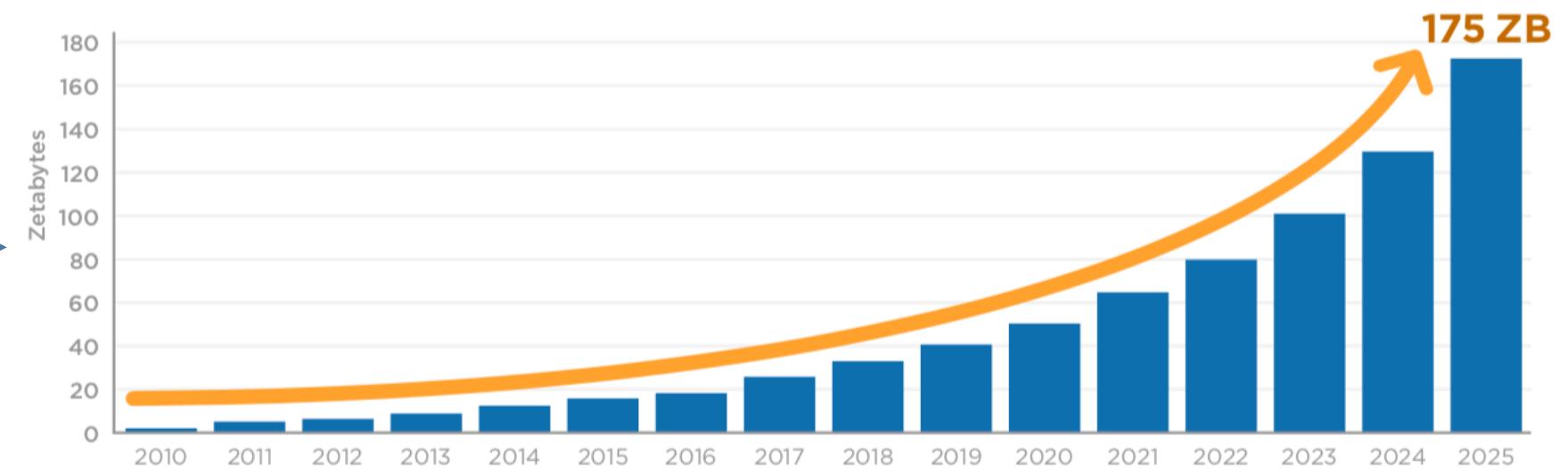
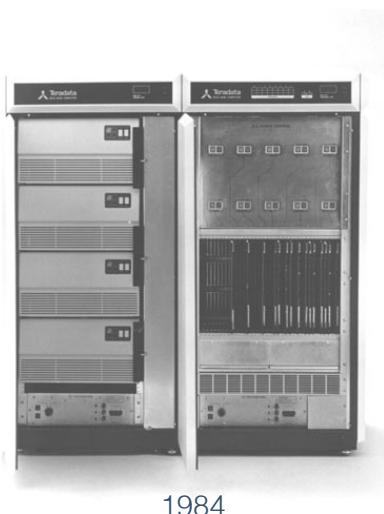
- Motivation, preliminaries and examples about data streaming and Stream Processing Engines
- Causes of out-of-order data and solutions enforcing total-ordering
- Pros/Cons of total-ordering
- Relaxation of total-ordering and the watermarks

https://github.com/vincenzo-gulisano/debs2020_tutorial_event_time

Agenda

- Motivation, preliminaries and examples about data streaming and Stream Processing Engines
- Causes of out-of-order data and solutions enforcing total-ordering
- Pros/Cons of total-ordering
- Relaxation of total-ordering and the watermarks

The era of Big Data



Source: Data Age 2025, sponsored by Seagate with data from IDC Global DataSphere, Nov 2018

Where does Big Data Originate?

219 billion photos
uploaded to Facebook

500 hours of video
uploaded to YouTube
every minute

1 trillion sensors
in the **Internet of Things (IoT)**
by 2030

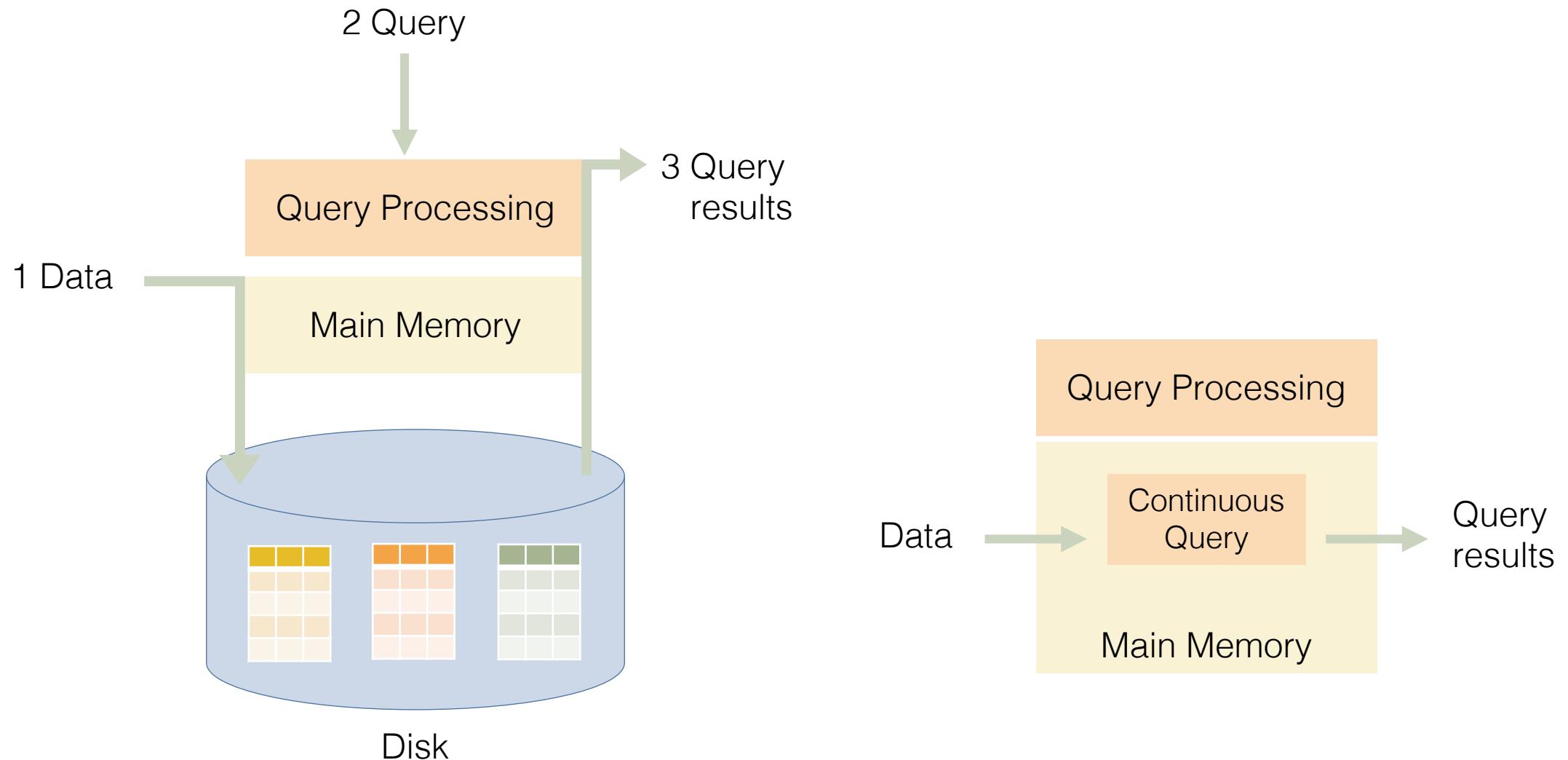


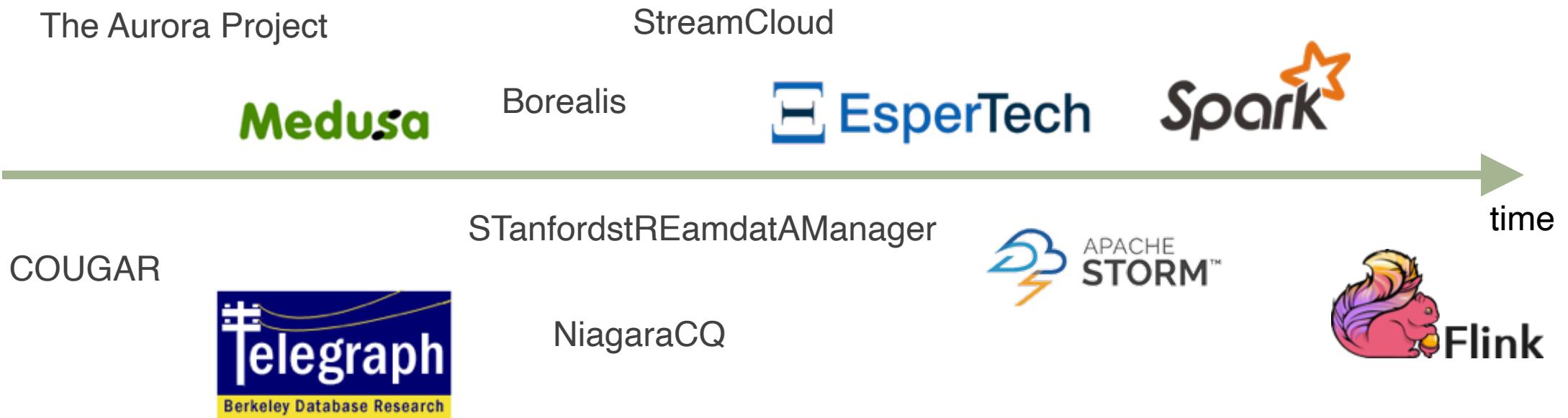
2.32 billion
Facebook users

1 online interaction
every 18 seconds
by 2025

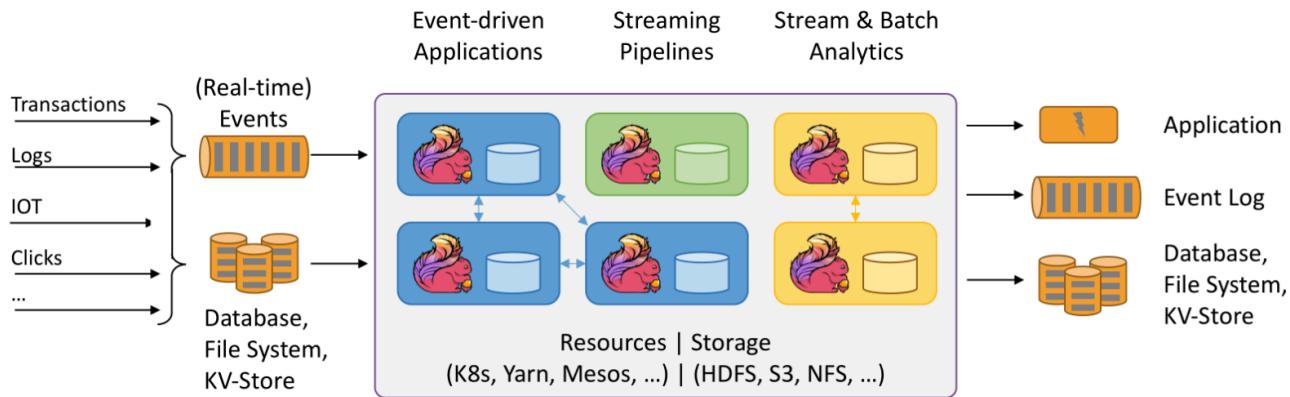


Database Management Systems (DBMSs) vs. Stream Processing Engines (SPEs)





Apache Flink® – Stateful Computations over Data Streams



All streaming use cases

- Event-driven Applications
- Stream & Batch Analytics
- Data Pipelines & ETL

[Learn more](#)

✓ Guaranteed correctness

- Exactly-once state consistency
- Event-time processing
- Sophisticated late data handling

[Learn more](#)

⬇ Layered APIs

- SQL on Stream & Batch Data
- DataStream API & DataSet API
- ProcessFunction (Time & State)

[Learn more](#)

⌚ Operational Focus

- Flexible deployment
- High-availability setup
- Savepoints

[Learn more](#)

⇄ Scales to any use case

- Scale-out architecture
- Support for very large state
- Incremental checkpointing

[Learn more](#)

⚡ Excellent Performance

- Low latency
- High throughput
- In-Memory computing

[Learn more](#)

Apache Flink

Software

Apache Flink is an open-source stream-processing framework developed by the Apache Software Foundation. The core of Apache Flink is a distributed streaming data-flow engine written in Java and Scala. Flink executes arbitrary dataflow programs in a data-parallel and pipelined manner. [Wikipedia](#)

Developed by: [Apache Software Foundation](#)

Stable release: 1.10.0 / [February 11, 2020](#); 2 months ago

Initial release: May 2011; 9 years ago

License: [Apache License 2.0](#)

Written in: Java, Scala

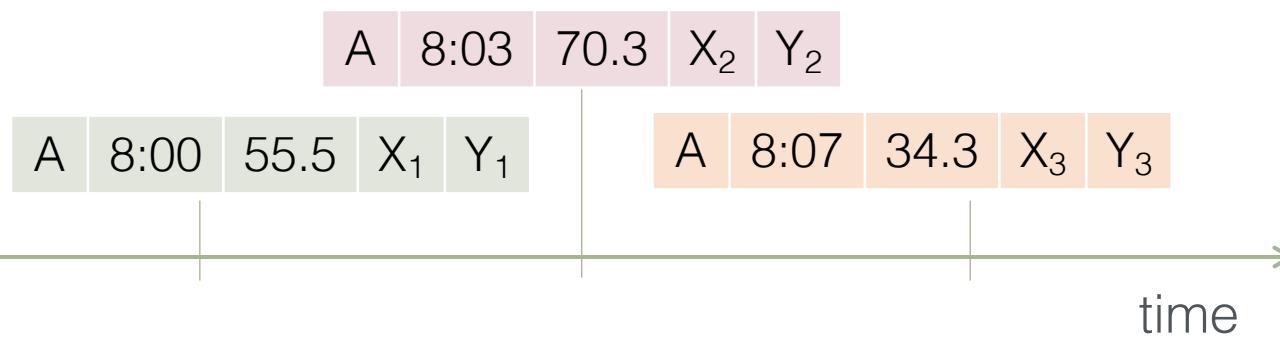
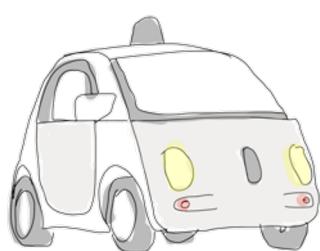
Data Stream:

unbounded sequence of tuples
sharing the same schema

Example: vehicles' speed reports

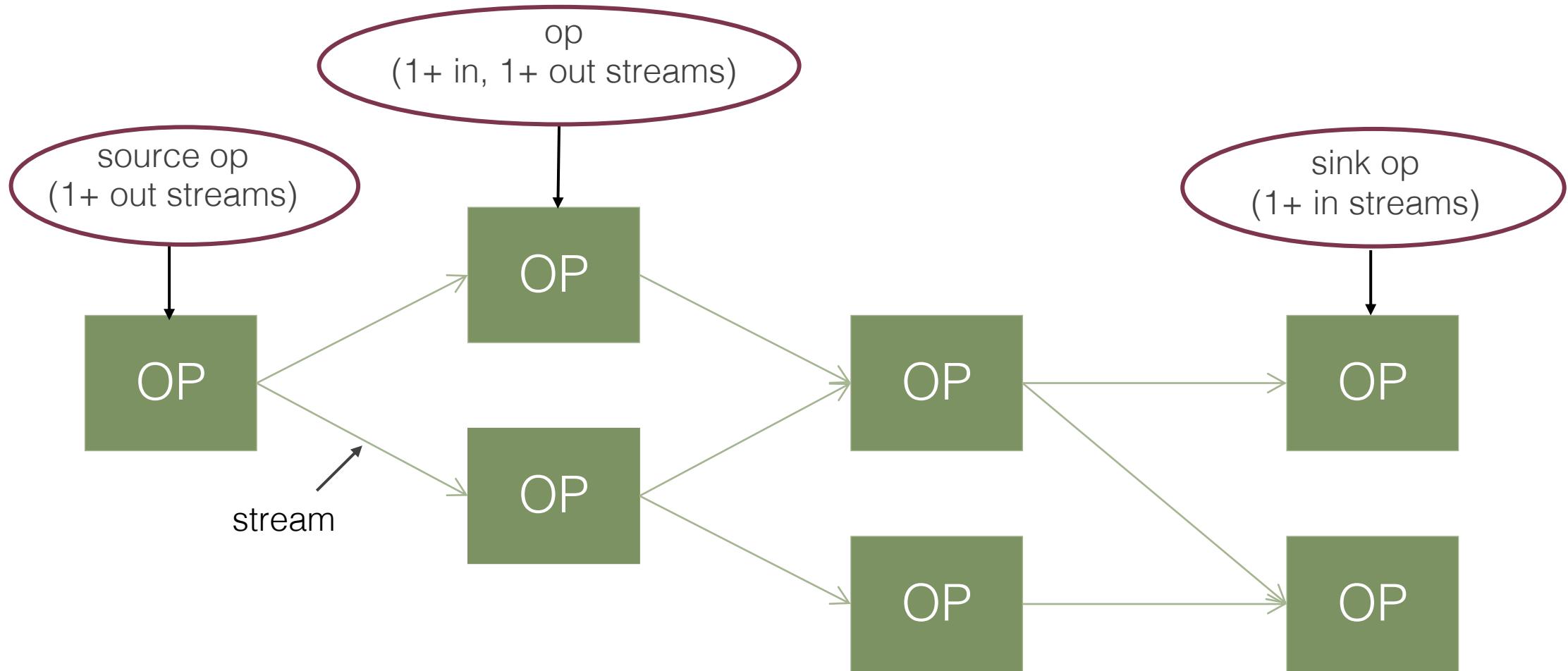
Field	Field
vehicle id	text
time (secs)	text
speed (Km/h)	double
X coordinate	double
Y coordinate	double

Let's assume each source
(e.g., vehicle)
produces and delivers
a timestamp-sorted stream



Continuous query (or simply query):

Directed Acyclic Graph (DAG) of streams and operators



Data Streaming Operators

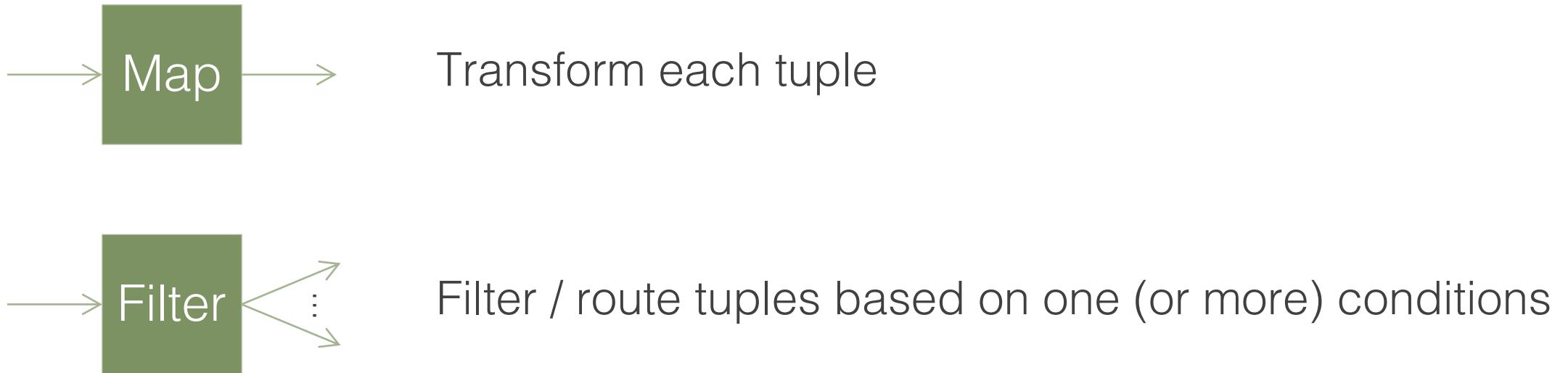
Two main types:

OP

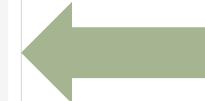
- Stateless operators
 - do not maintain any state
 - one-by-one processing
 - if they maintain some state, such state does not evolve depending on the tuples being processed
- Stateful operators
 - maintain a state that evolves depending on the tuples being processed
 - produce output tuples that depend on multiple input tuples

OP

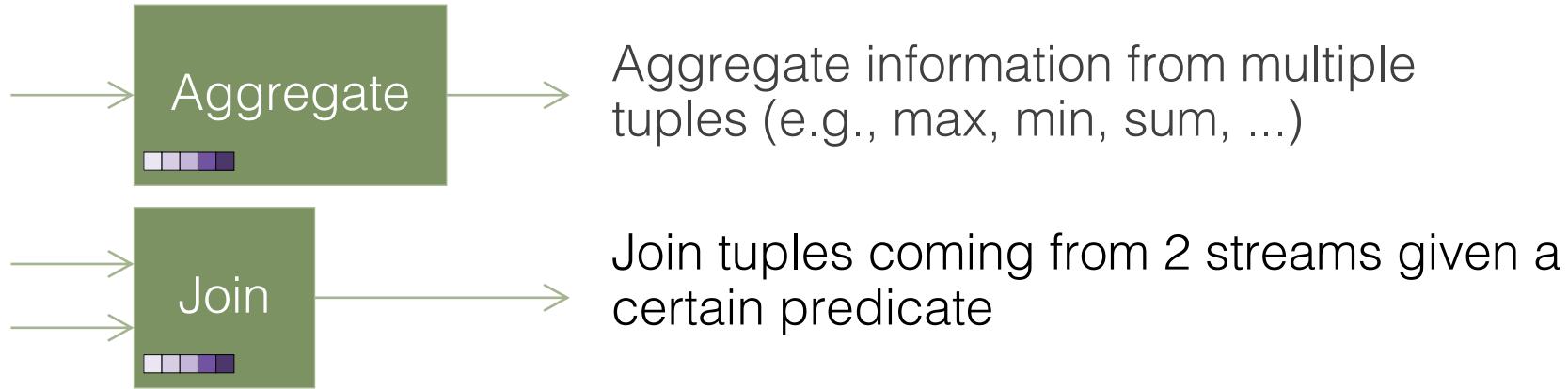
Stateless Operators



Map DataStream → DataStream	Takes one element and produces one element. A map function that doubles the values of the input stream: <pre>DataStream<Integer> dataStream = ...; dataStream.map(new MapFunction<Integer, Integer>() { @Override public Integer map(Integer value) throws Exception { return 2 * value; } });</pre>
FlatMap DataStream → DataStream	Takes one element and produces zero, one, or more elements. A flatmap function that splits sentences to words: <pre>dataStream.flatMap(new FlatMapFunction<String, String>() { @Override public void flatMap(String value, Collector<String> out) throws Exception { for(String word: value.split(" ")){ out.collect(word); } } });</pre>
Filter DataStream → DataStream	Evaluates a boolean function for each element and retains those for which the function returns true. A filter that filters out zero values: <pre>dataStream.filter(new FilterFunction<Integer>() { @Override public boolean filter(Integer value) throws Exception { return value != 0; } });</pre>



Stateful Operators



Aggregations

KeyedStream → DataStream

Rolling aggregations on a keyed data stream. The difference between min and minBy is that min returns the minimum value, whereas minBy returns the element that has the minimum value in this field (same for max and maxBy).

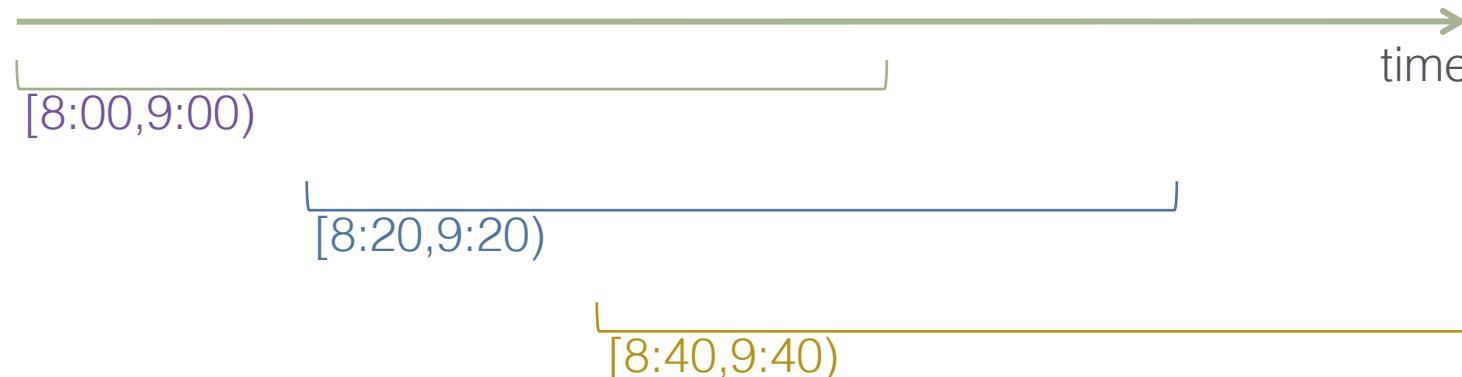
```
keyedStream.sum(0);
keyedStream.sum("key");
keyedStream.min(0);
keyedStream.min("key");
keyedStream.max(0);
keyedStream.max("key");
keyedStream.minBy(0);
keyedStream.minBy("key");
keyedStream.maxBy(0);
keyedStream.maxBy("key");
```

Windows and Stateful Analysis

Stateful operations are done over windows:

- Time-based (e.g., tuples in the last 10 minutes)
 - Tuple-based (e.g., given the last 50 tuples)
- How many tuple in a window?
→ Which time period does a window span?

Example of time-based window of size 1 hour and advance 20 minutes

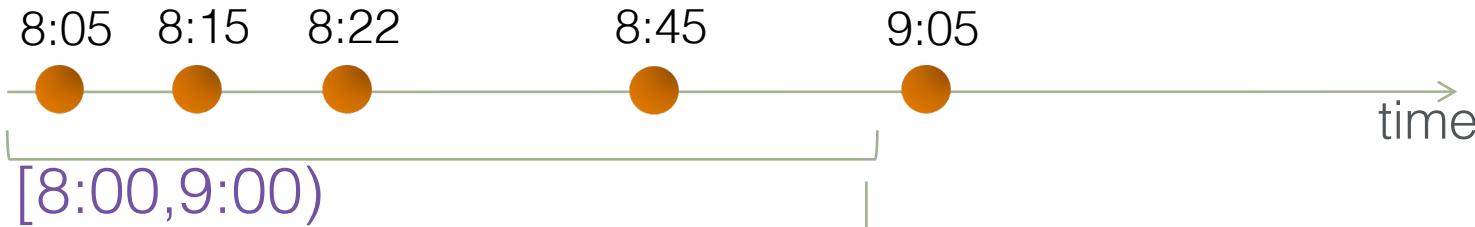


Time-based sliding window aggregation (count)

Counter: 1 Counter: 3

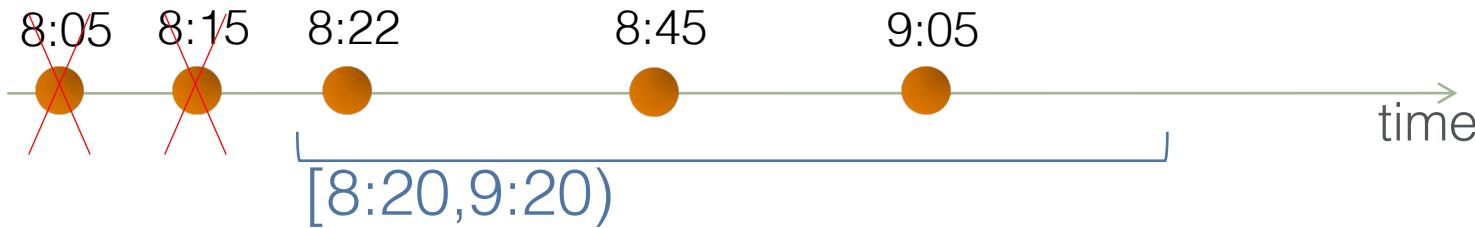
Counter: 2

Counter: 4

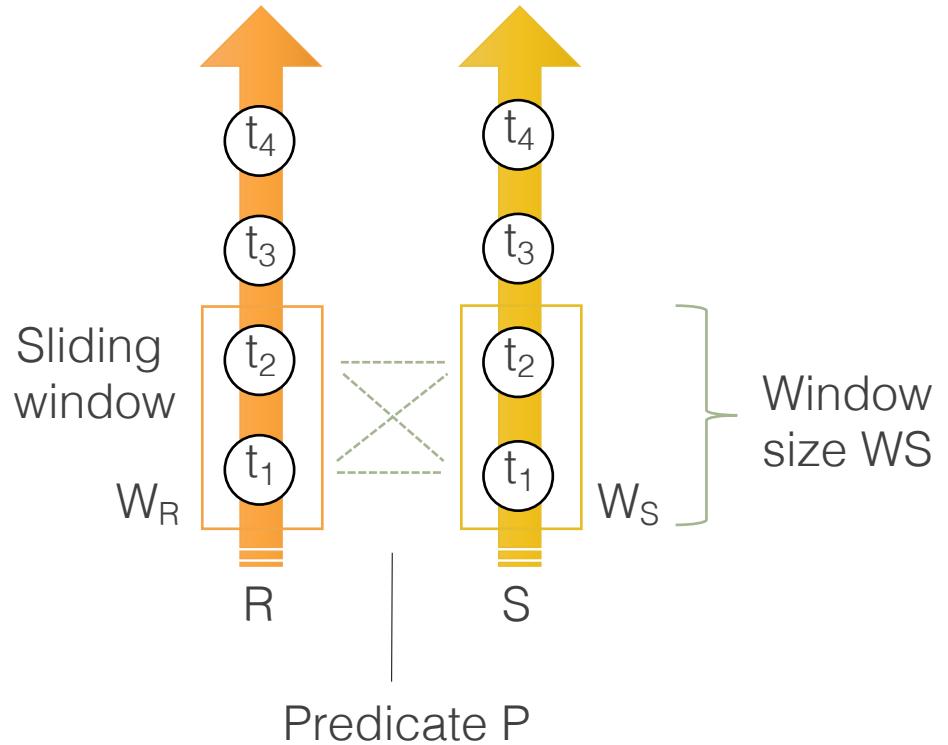


Output: 4

Counter: 3



Time-based sliding window joining



Windows and stateful analysis

Window
KeyedStream → WindowedStream

Windows can be defined on already partitioned KeyedStreams. Windows group the data in each key according to some characteristic (e.g., the data that arrived within the last 5 seconds). See [windows](#) for a complete description of windows.

```
dataStream.keyBy(0).window(TumblingEventTimeWindows.of(Time.seconds(5))); // Last  
5 seconds of data
```

Basic operators and user-defined operators

Besides a set of basic operators, SPEs usually allow the user to define ad-hoc operators (e.g., when existing aggregation are not enough)

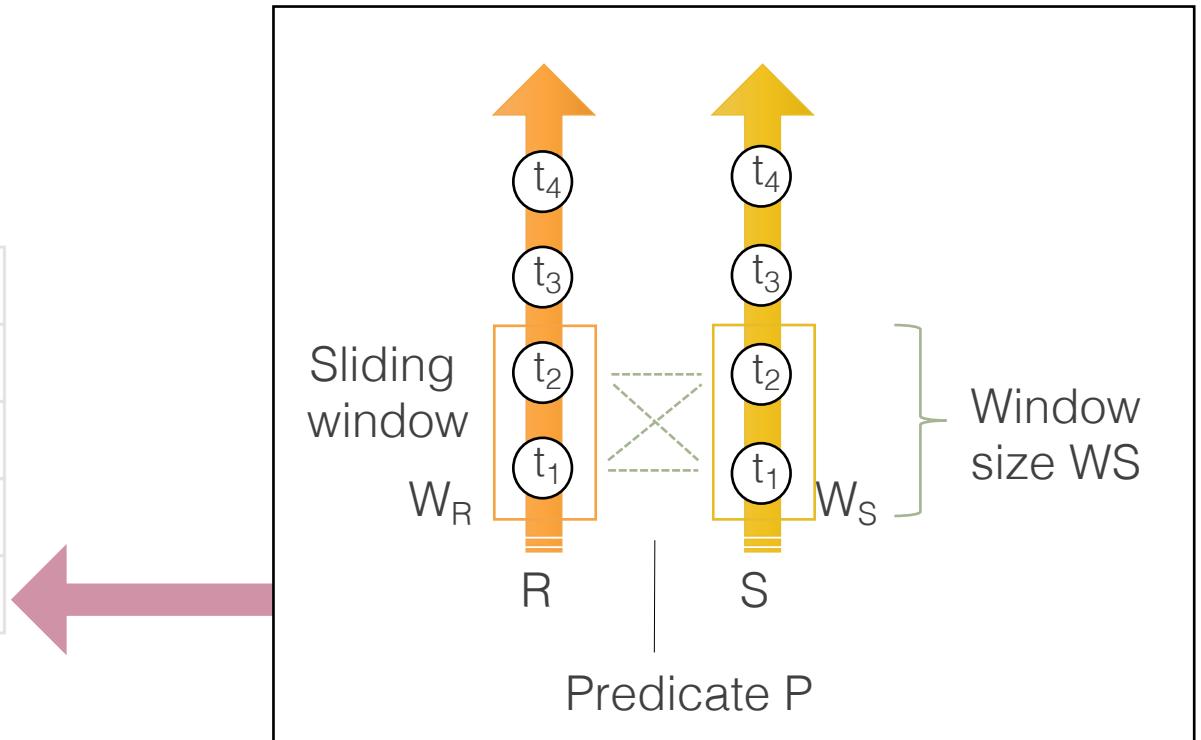
Window Apply WindowedStream → DataStream AllWindowedStream → DataStream	<p>Applies a general function to the window as a whole. Below is a function that manually sums the elements of a window.</p> <p>Note: If you are using a windowAll transformation, you need to use an AllWindowFunction instead.</p> <pre>windowedStream.apply (new WindowFunction<Tuple2<String, Integer>, Integer, Tuple, Window>() { public void apply (Tuple tuple, Window window, Iterable<Tuple2<String, Integer>> values, Collector<Integer> out) throws Exception { int sum = 0; for (value t: values) { sum += t.f1; } out.collect (new Integer(sum)); } });</pre>
--	--

SPEs and operators' variants

- Each SPE might define its own variants of certain streaming operators:

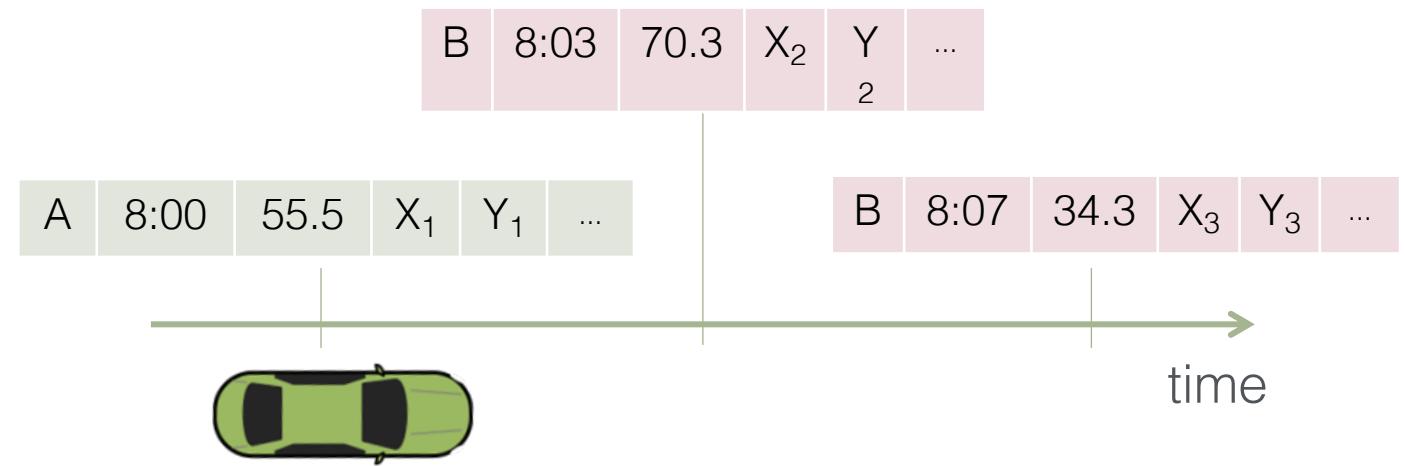
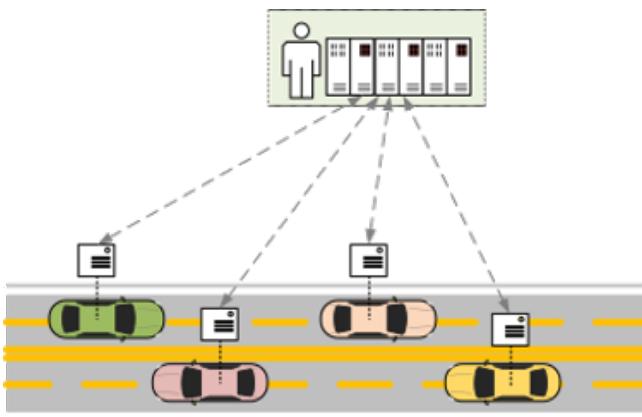
Joining

Window Join
Tumbling Window Join
Sliding Window Join
Session Window Join
Interval Join

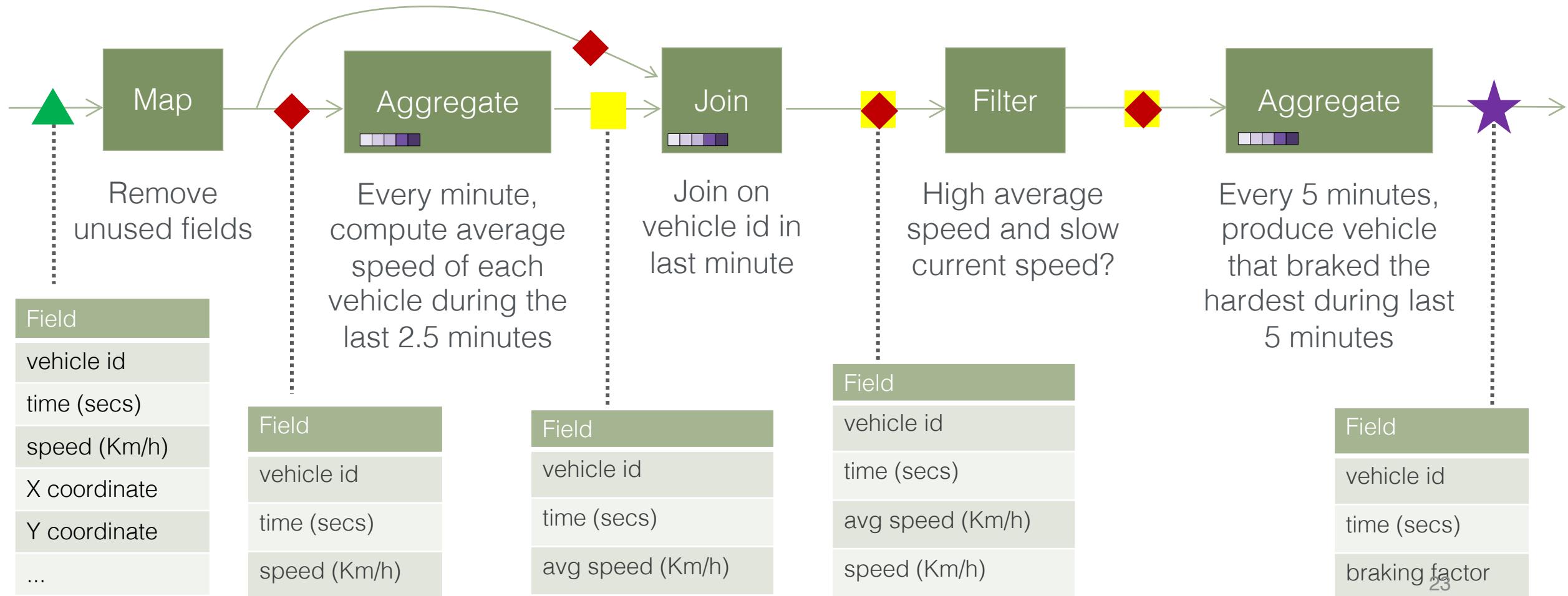


Sample Query

"every five minutes, of all vehicles that braked significantly,
find the one that braked the hardest"



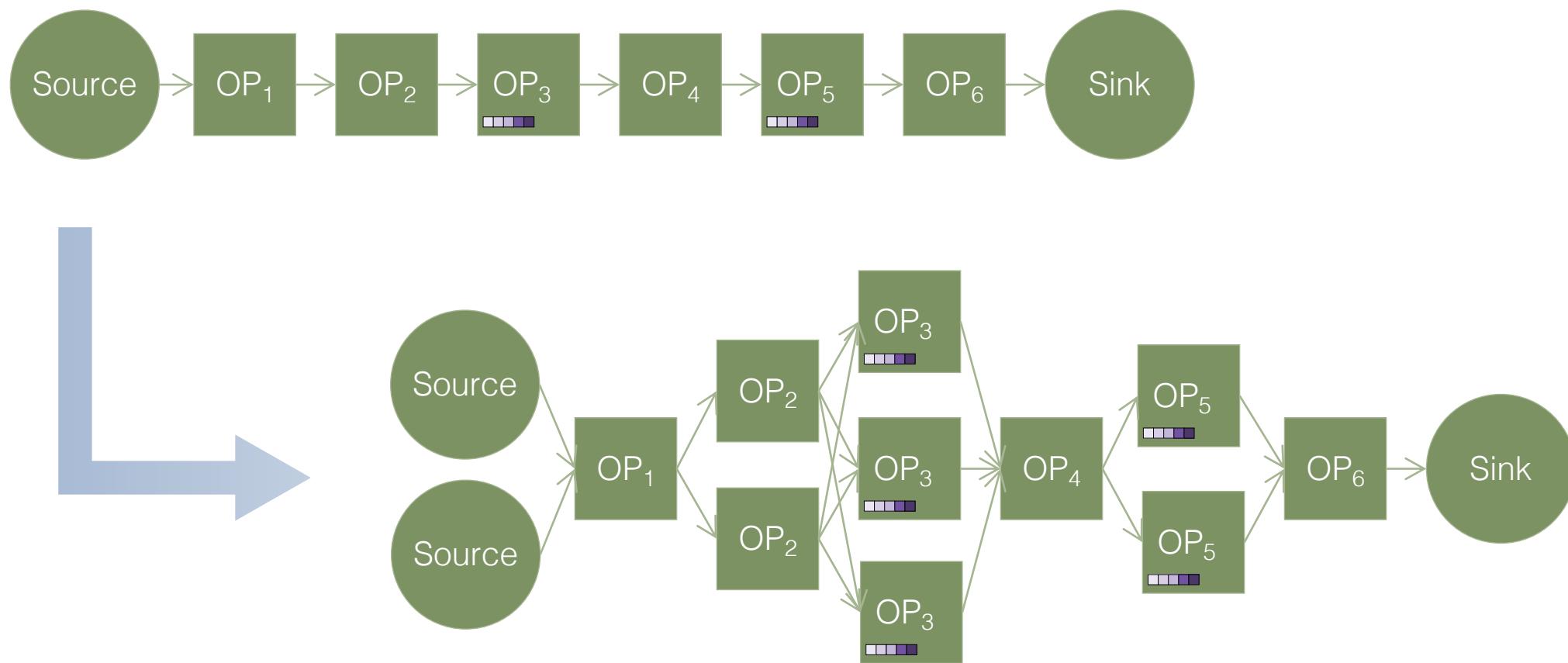
Sample Query



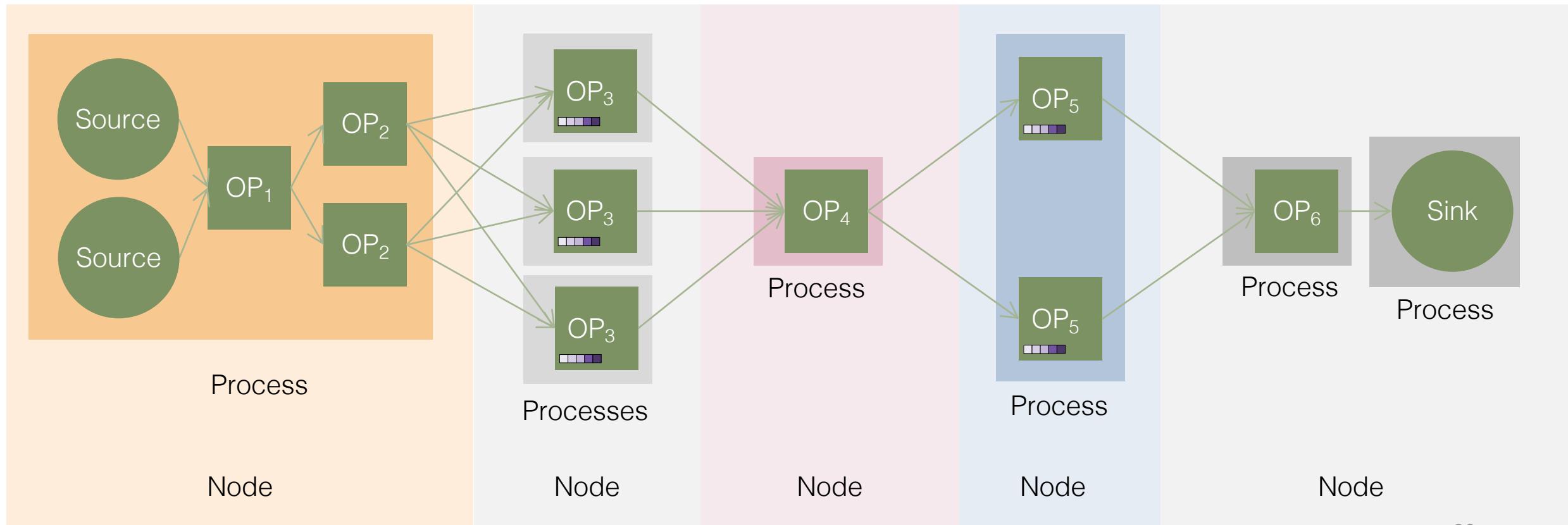
Agenda

- Motivation, preliminaries and examples about data streaming and Stream Processing Engines
- Causes of out-of-order data and solutions enforcing total-ordering
- Pros/Cons of total-ordering
- Relaxation of total-ordering and the watermarks

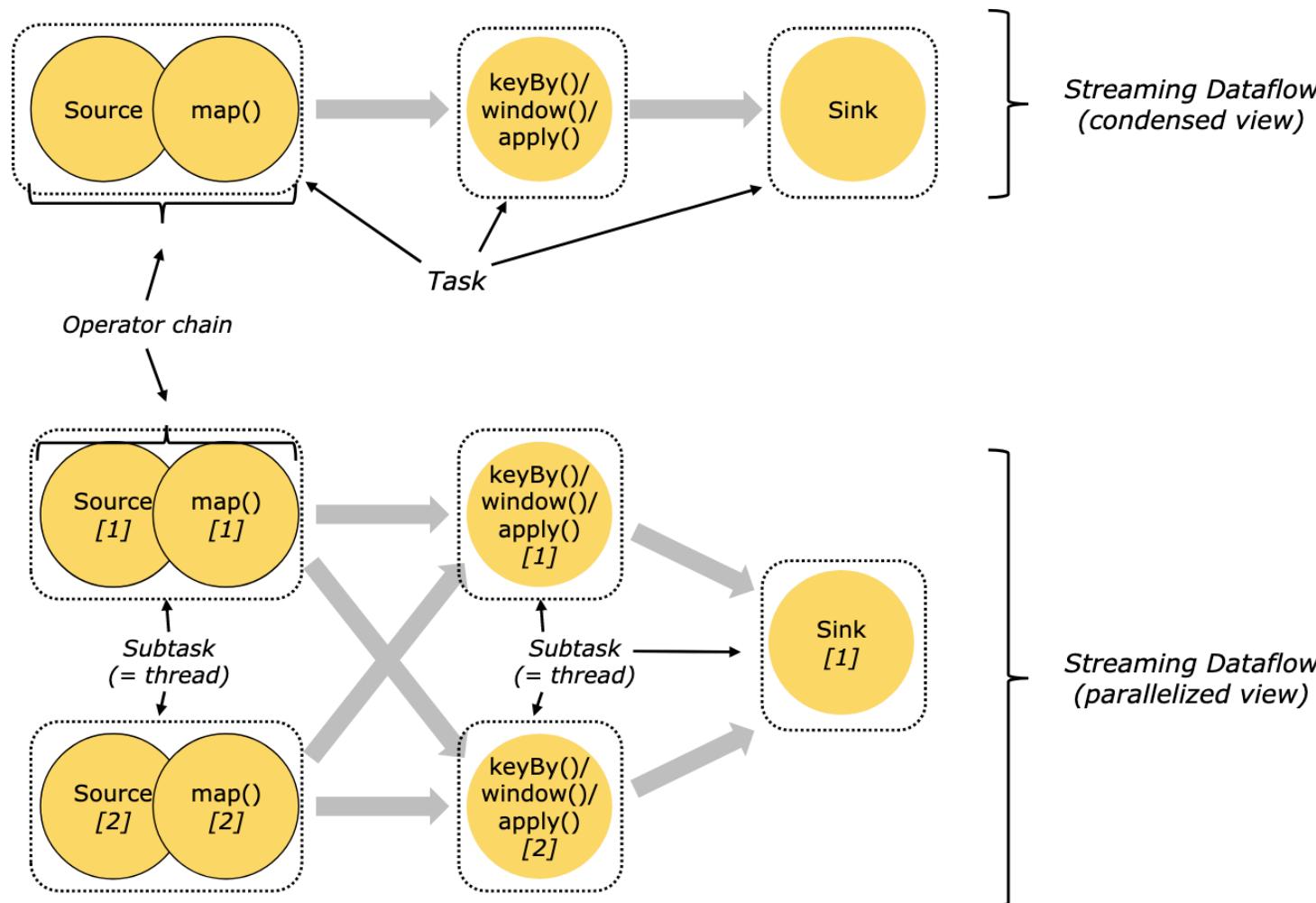
From an abstract query ... to streaming application run by an SPE



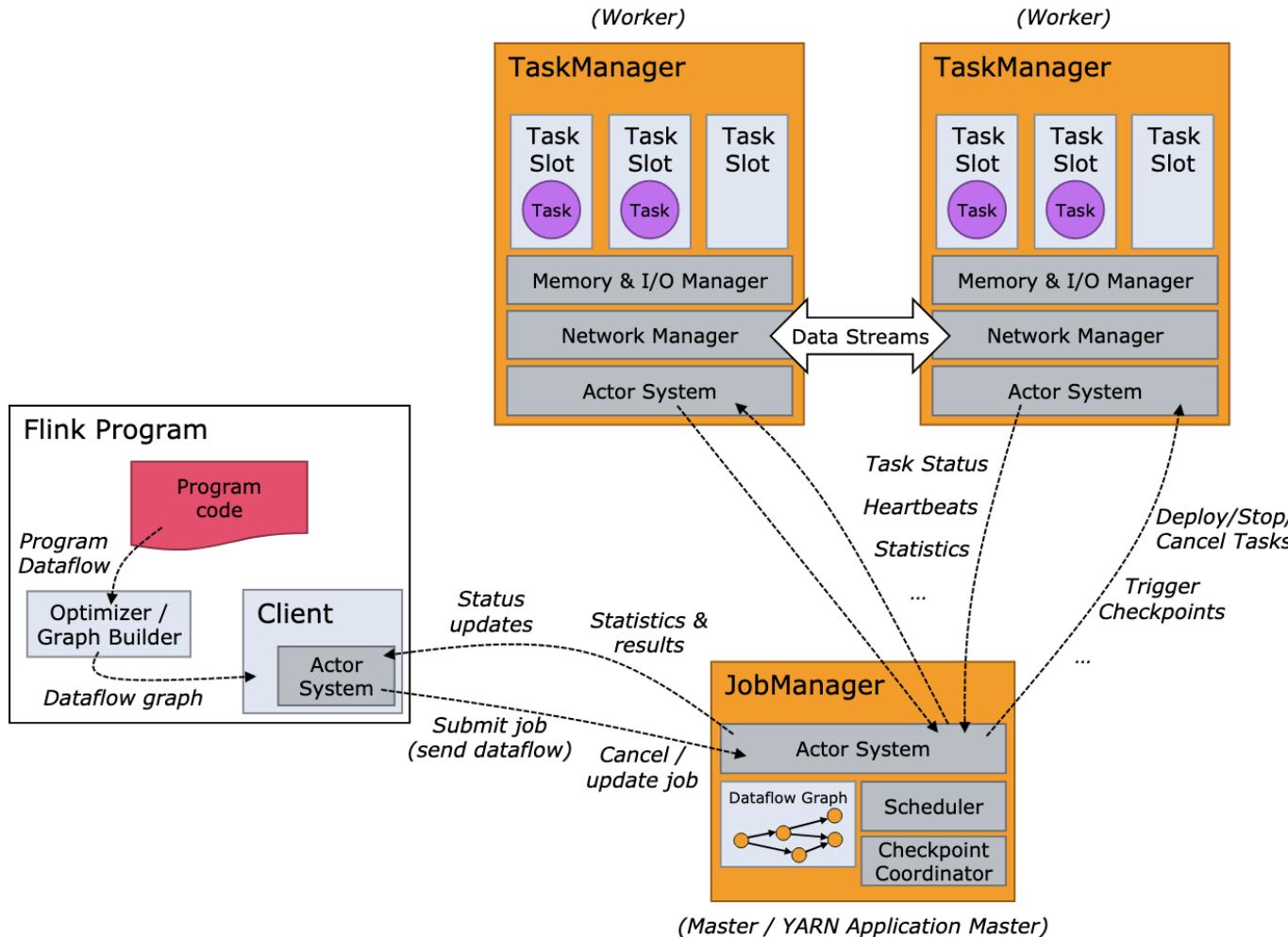
From an abstract query ... to streaming application run by an SPE



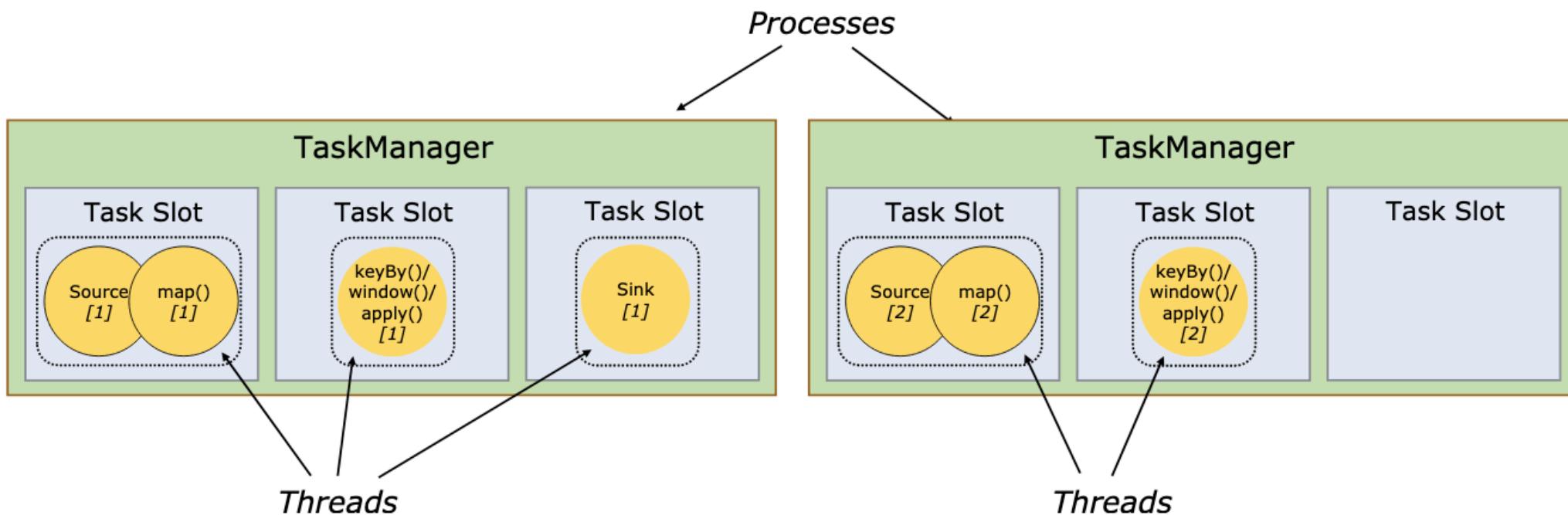
From an abstract query ...to streaming application run by Flink - 1/3



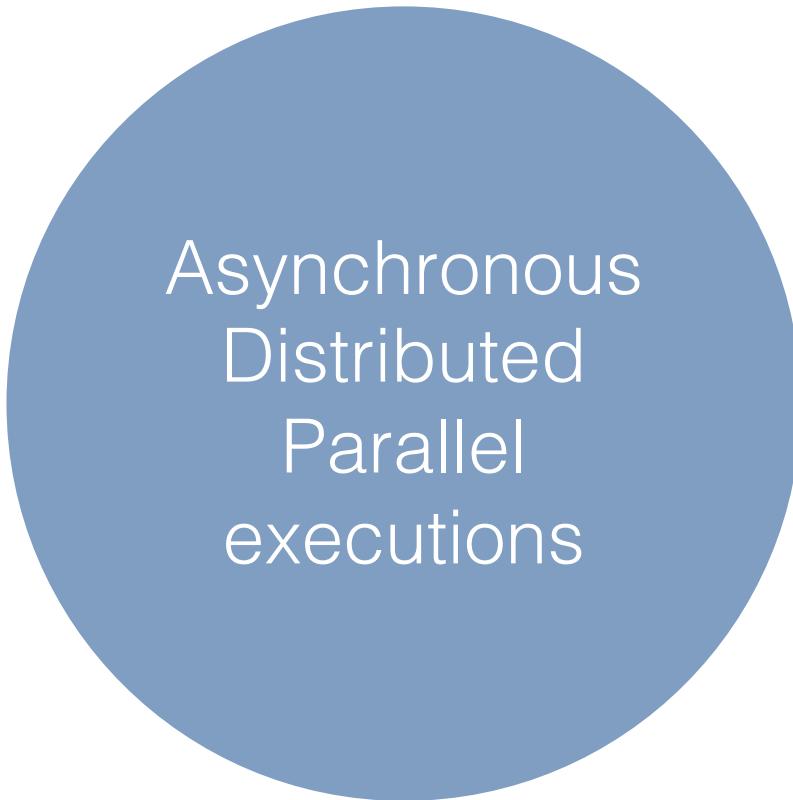
From an abstract query ...to streaming application run by Flink - 1/3



From an abstract query ...to streaming application run by Flink - 1/3



Causes of out-of-order data:



Data sources that produce out-of-order data

- Discussed in many related works (e.g., Babu, Shivnath, Utkarsh Srivastava, and Jennifer Widom. "Exploiting k-constraints to reduce memory overhead in continuous queries over data streams." ACM Transactions on Database Systems (TODS) 29.3 (2004): 545-580.)
- Battery-operated devices, unreliable wireless networks, ...

1 trillion sensors

in the Internet of Things (IoT)

by 2030

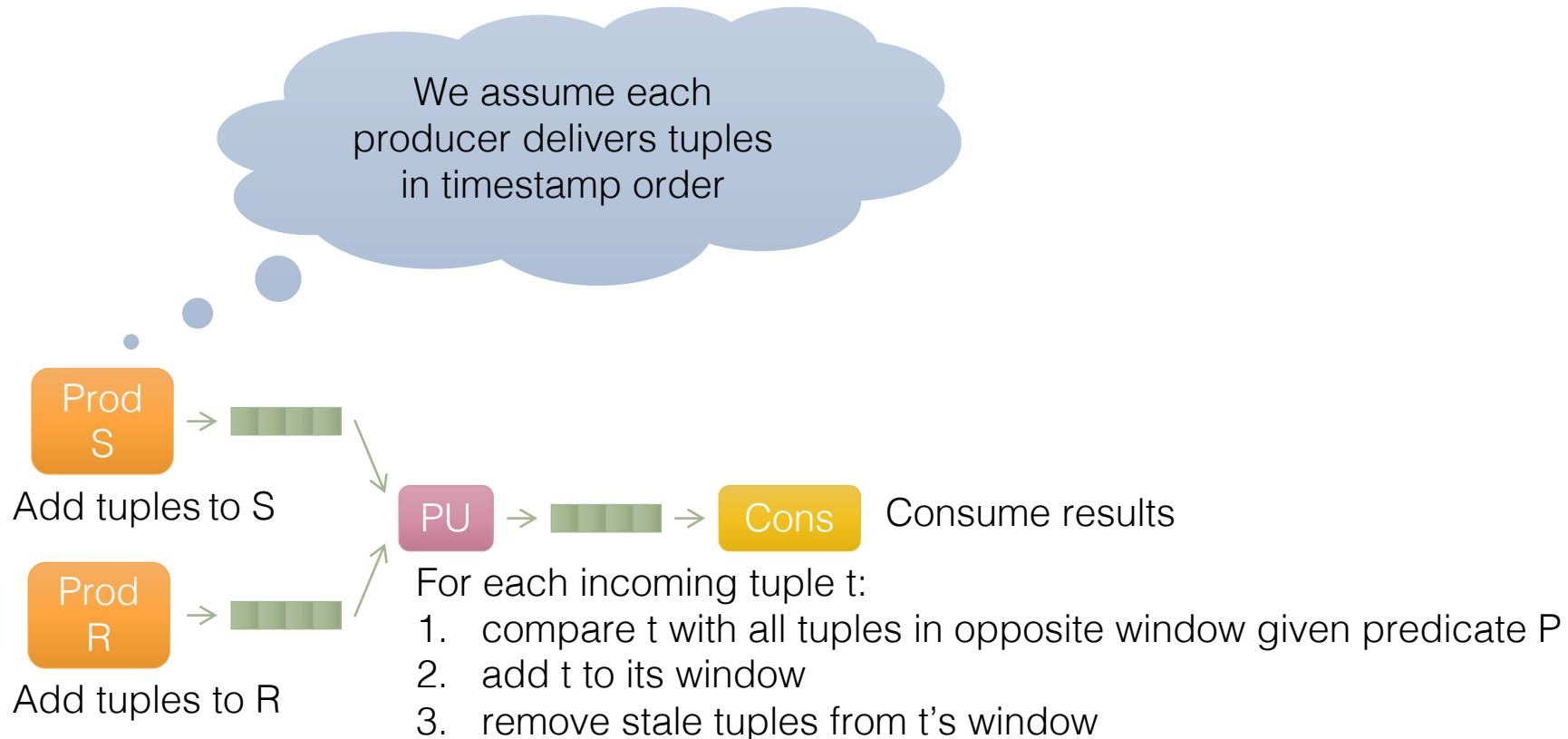


Causes of out-of-order data

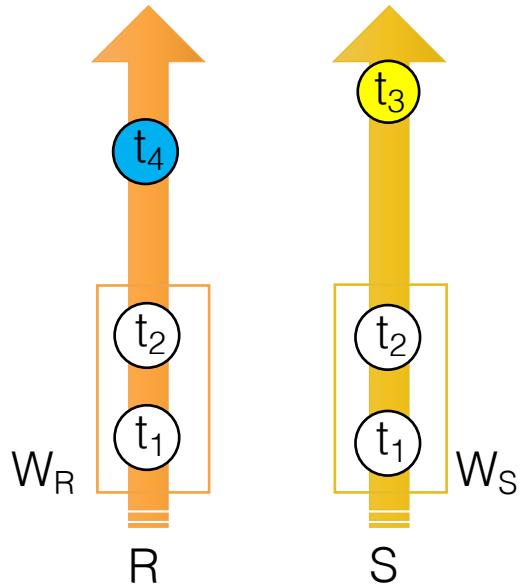
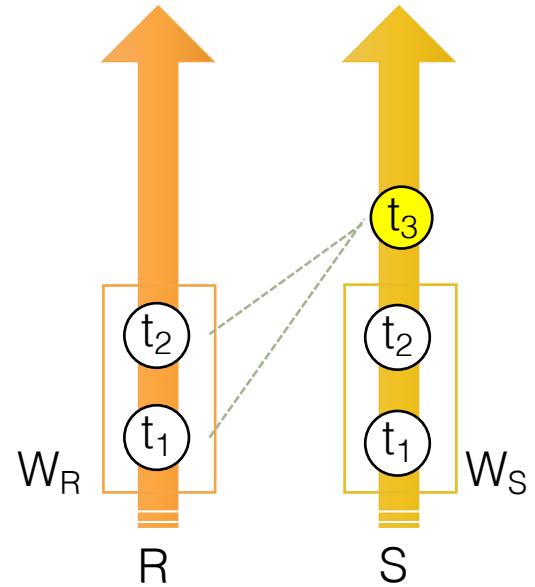
Sources
themselves

Asynchronous
Distributed
Parallel
executions

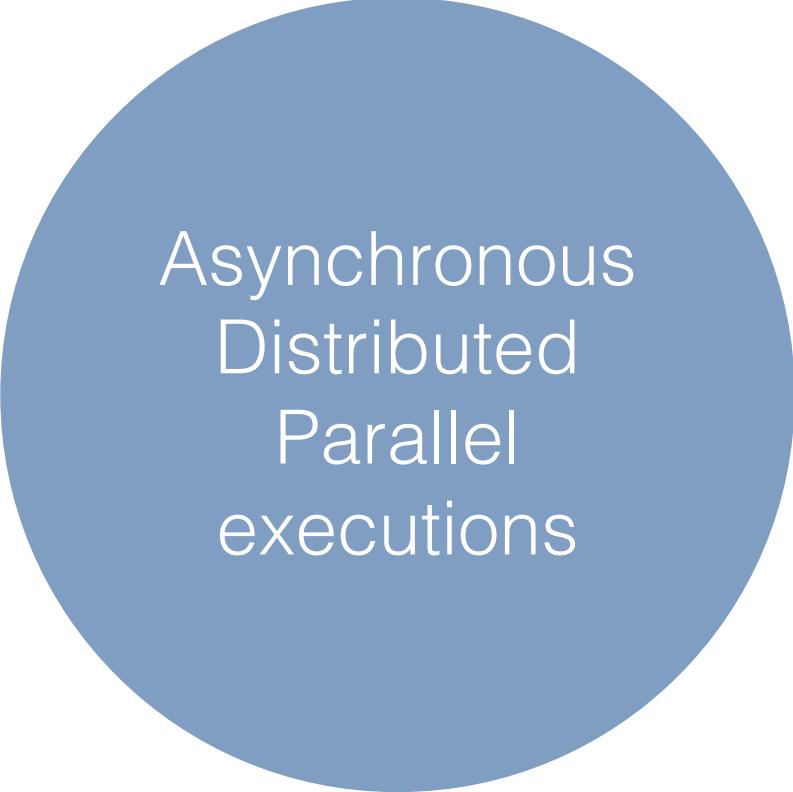
The 3-step procedure (sequential stream join)



The 3-step procedure, is it enough?



Causes of out-of-order data

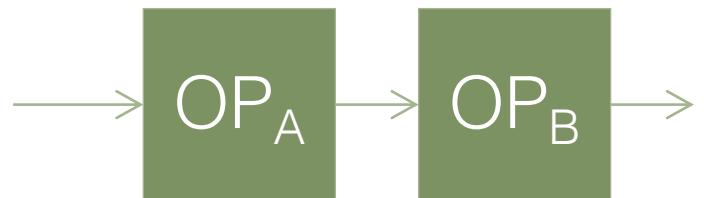


Asynchronous
Distributed
Parallel
executions

Any operator fed data from multiple logical / physical stream can potentially observe out-of-order data

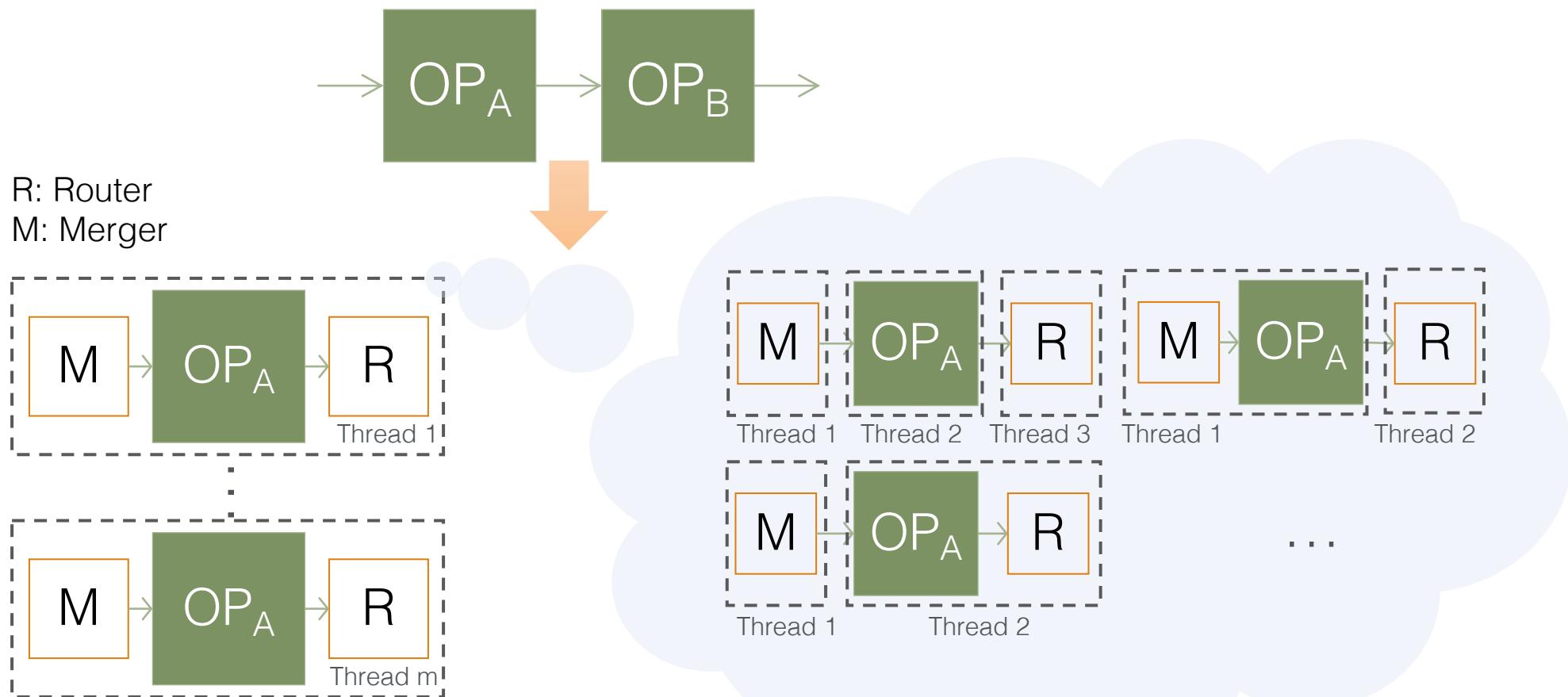
Parallel execution

- General approach



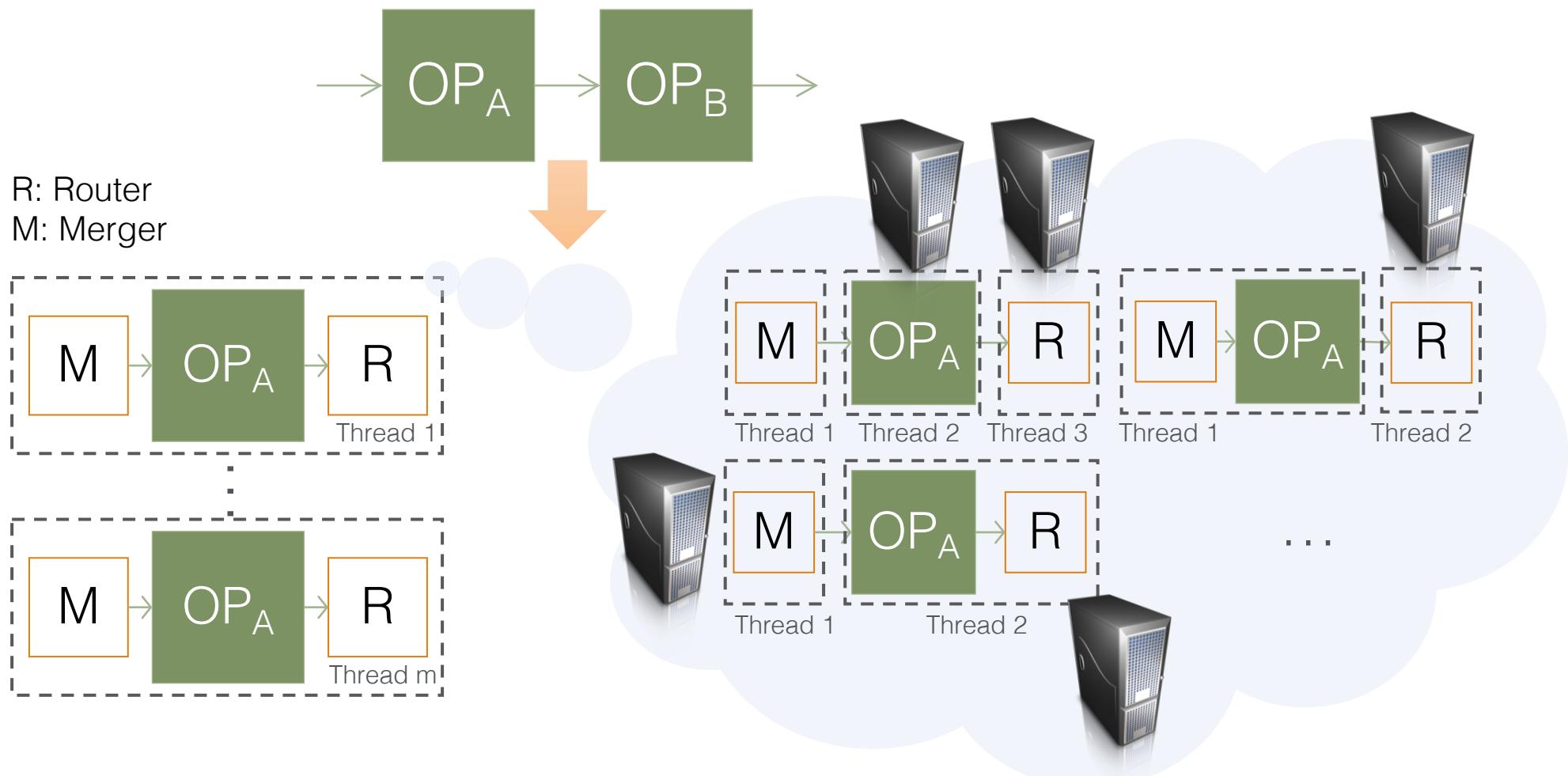
Parallel execution

- General approach



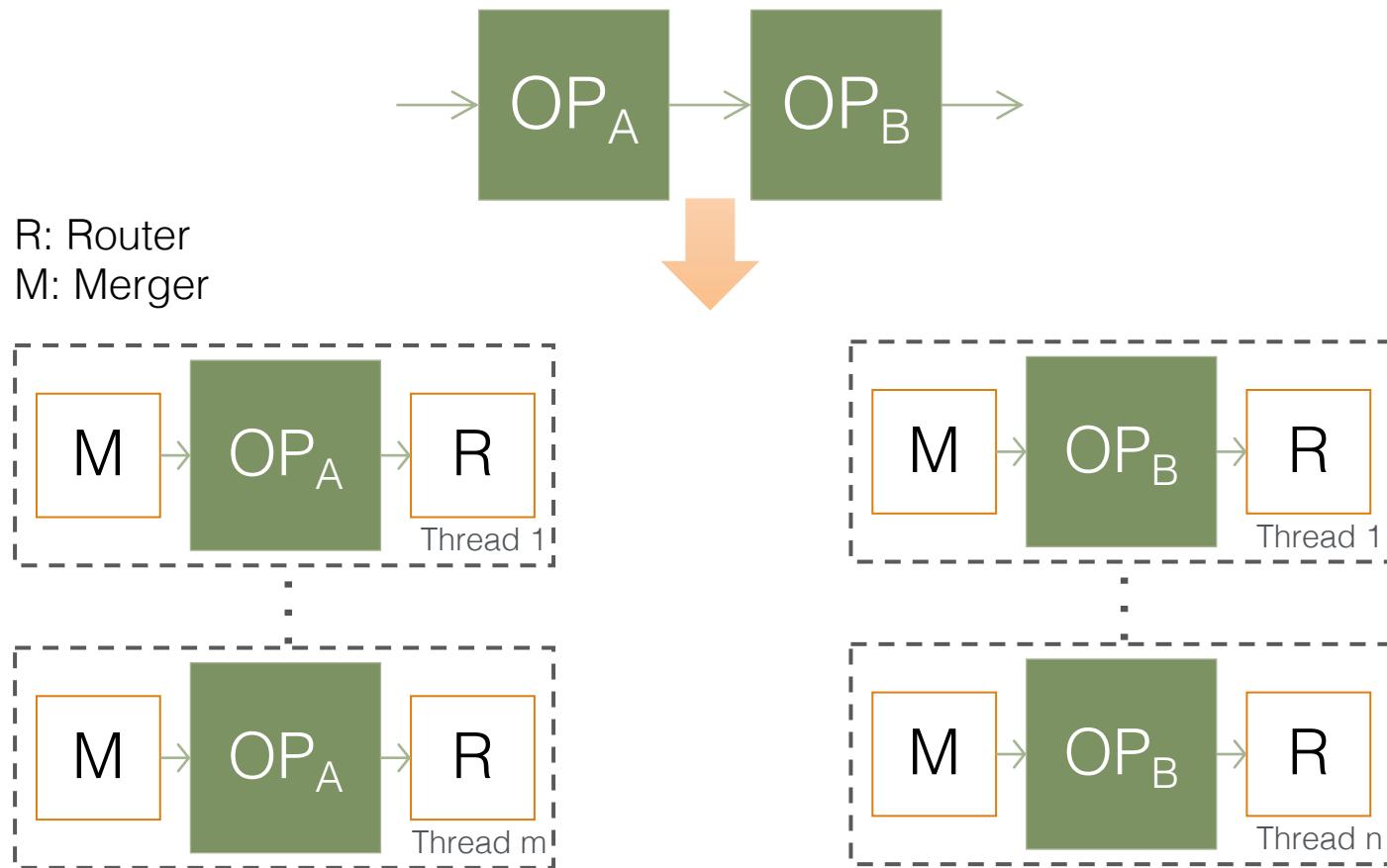
Parallel execution

- General approach



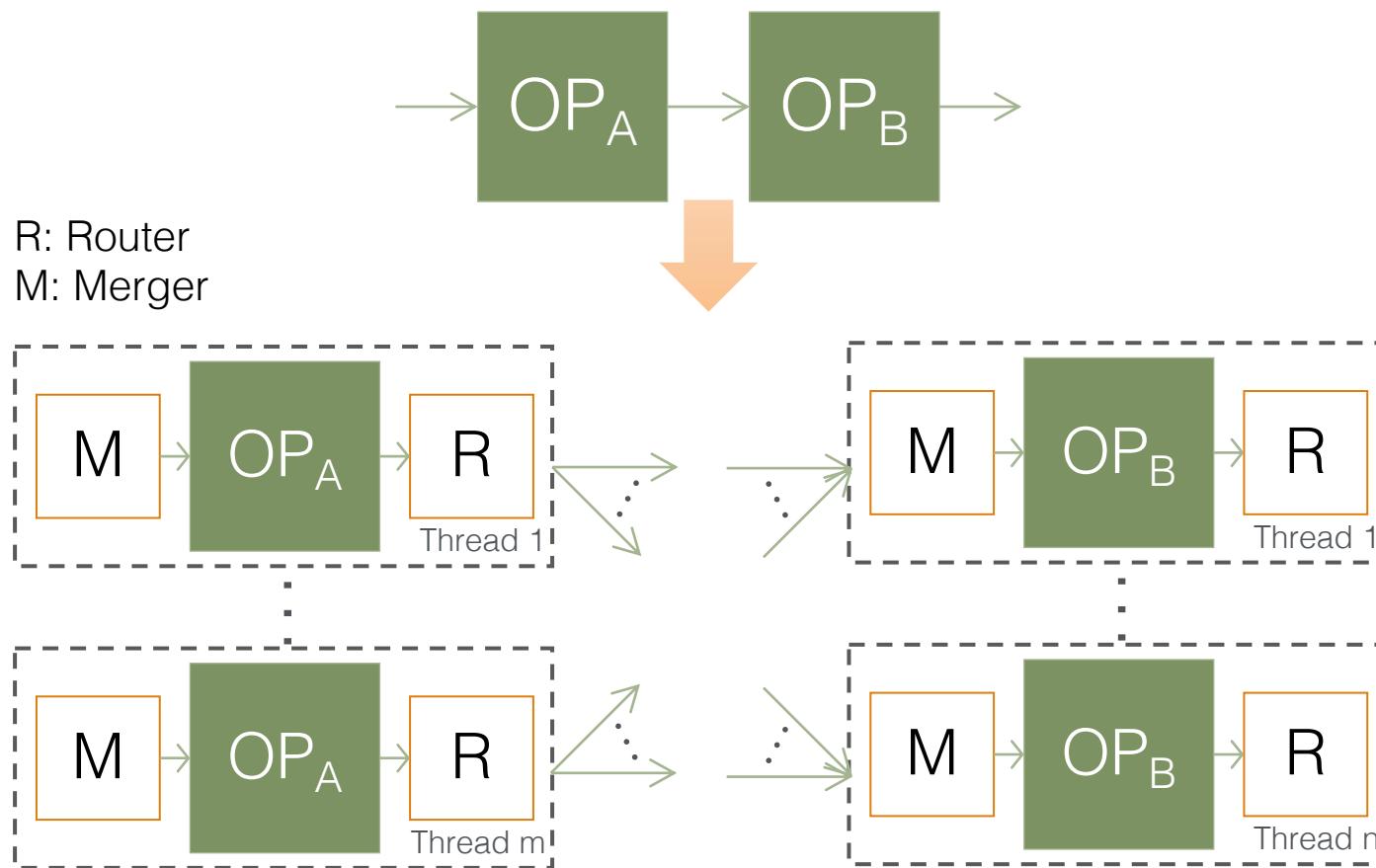
Parallel execution

- General approach



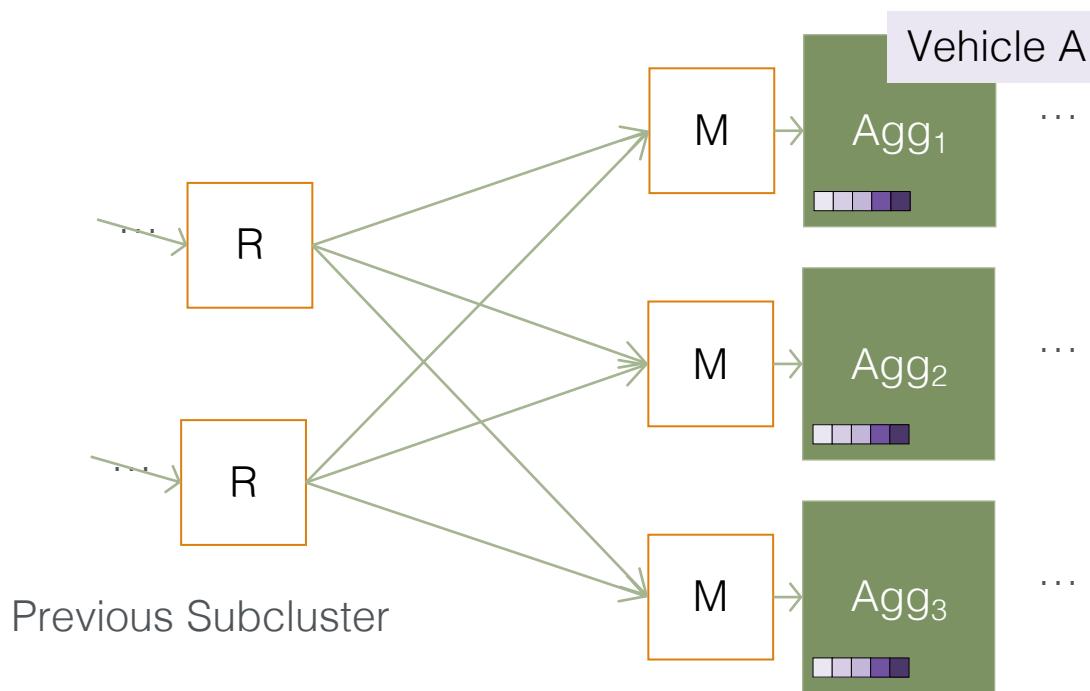
Parallel execution

- General approach



Parallel execution

- Stateful operators: Semantic awareness
 - Aggregate: count within last hour, group-by vehicle id

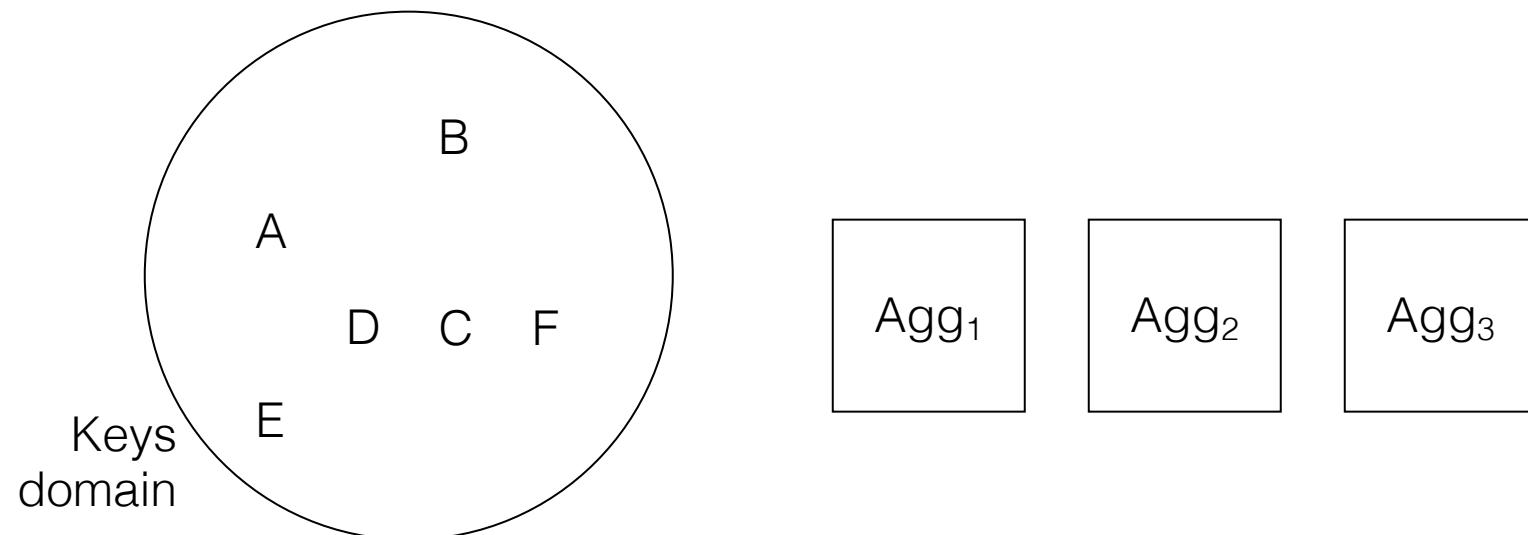


Parallel execution

- Depending on the stateful operator semantic:
 - Partition input stream into **keys**
 - Each key is processed by 1 thread
- $\# \text{ keys} >> \# \text{ threads/nodes}$

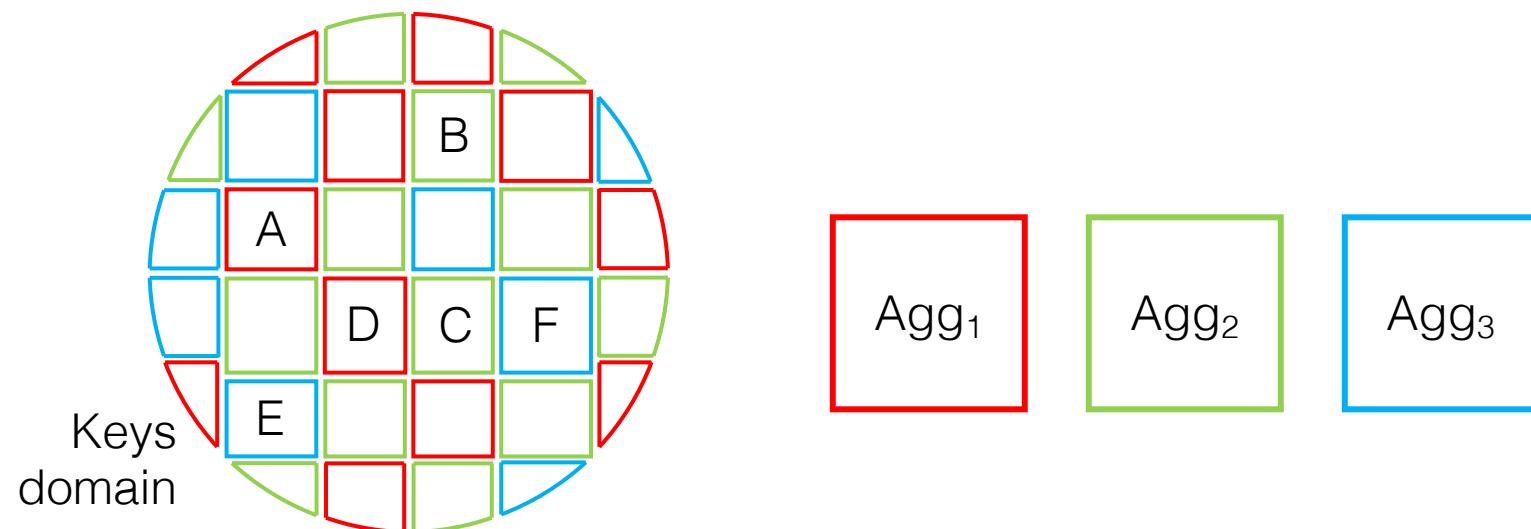
Parallel execution

- Depending on the stateful operator semantic:
 - Partition input stream into **keys**
 - Each key is processed by 1 thread
- # keys >> # threads/nodes



Parallel execution

- Depending on the stateful operator semantics:
 - Partition input stream into **keys**
 - Each key is processed by 1 thread
- # keys >> # threads/nodes



Aggregations

KeyedStream → DataStream

Rolling aggregations on a keyed data stream. The difference between min and minBy is that min returns the minimum value, whereas minBy returns the element that has the minimum value in this field (same for max and maxBy).

```
keyedStream.sum(0);
keyedStream.sum("key");
keyedStream.min(0);
keyedStream.min("key");
keyedStream.max(0);
keyedStream.max("key");
keyedStream.minBy(0);
keyedStream.minBy("key");
keyedStream.maxBy(0);
keyedStream.maxBy("key");
```

Window

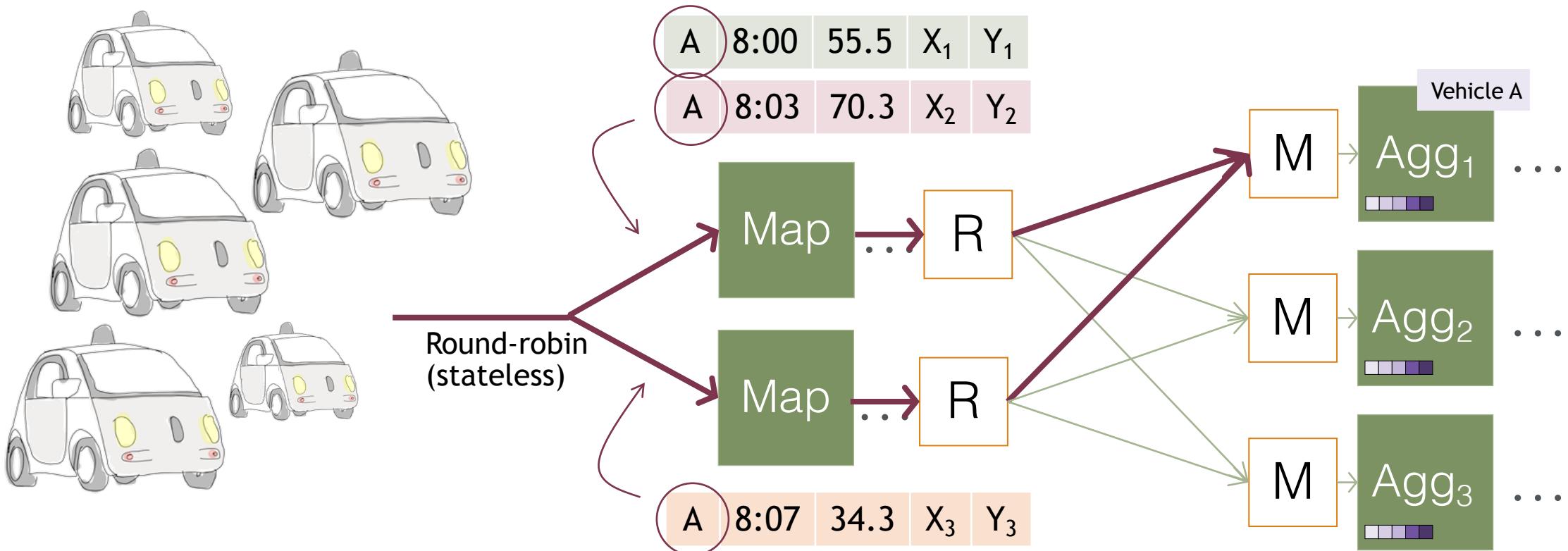
KeyedStream → WindowedStream

Windows can be defined on already partitioned KeyedStreams. Windows group the data in each key according to some characteristic (e.g., the data that arrived within the last 5 seconds). See [windows](#) for a complete description of windows.

```
dataStream.keyBy(0).window(TumblingEventTimeWindows.of(Time.seconds(5))); // Last
5 seconds of data
```

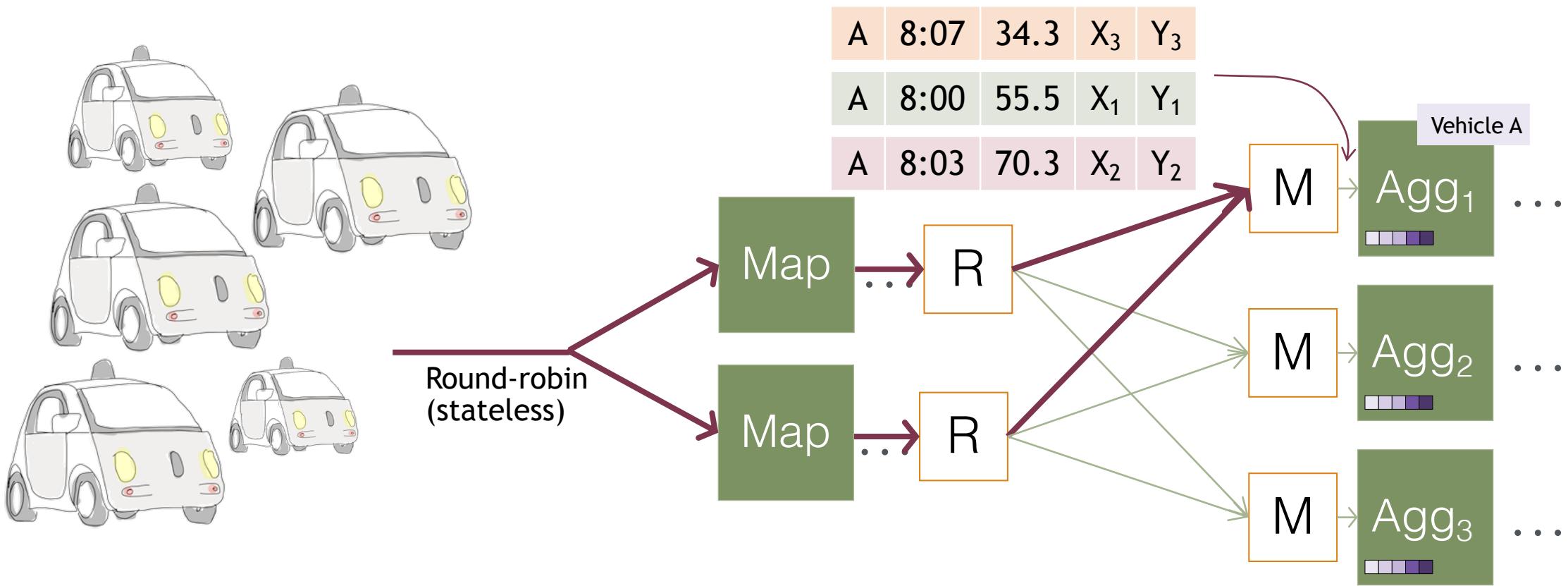
Parallel execution

- Stateful operators: Semantic awareness
 - Aggregate: count within last hour, group-by vehicle id

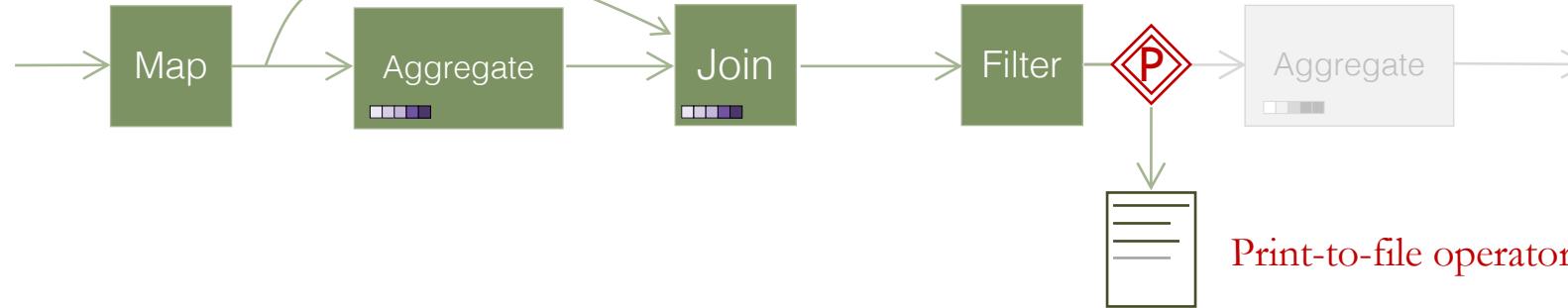


Parallel execution

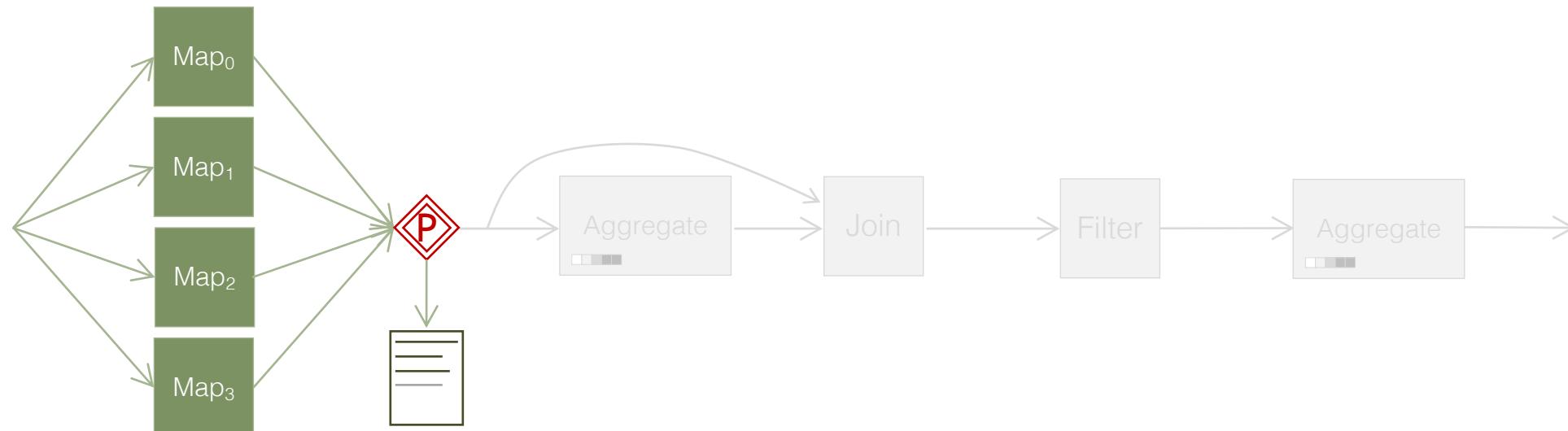
- Stateful operators: Semantic awareness
 - Aggregate: count within last hour, group-by vehicle id



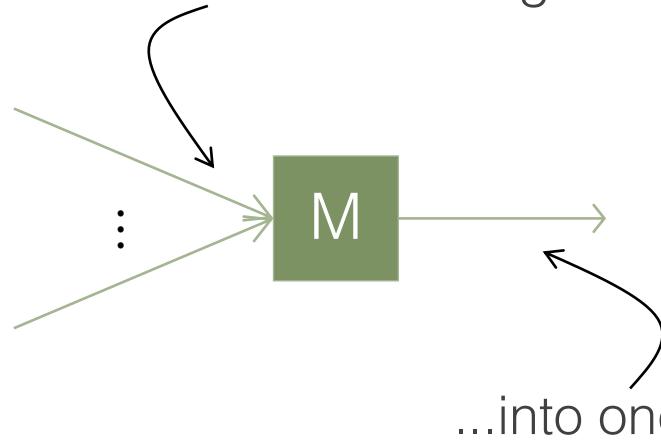
Inherent disorder



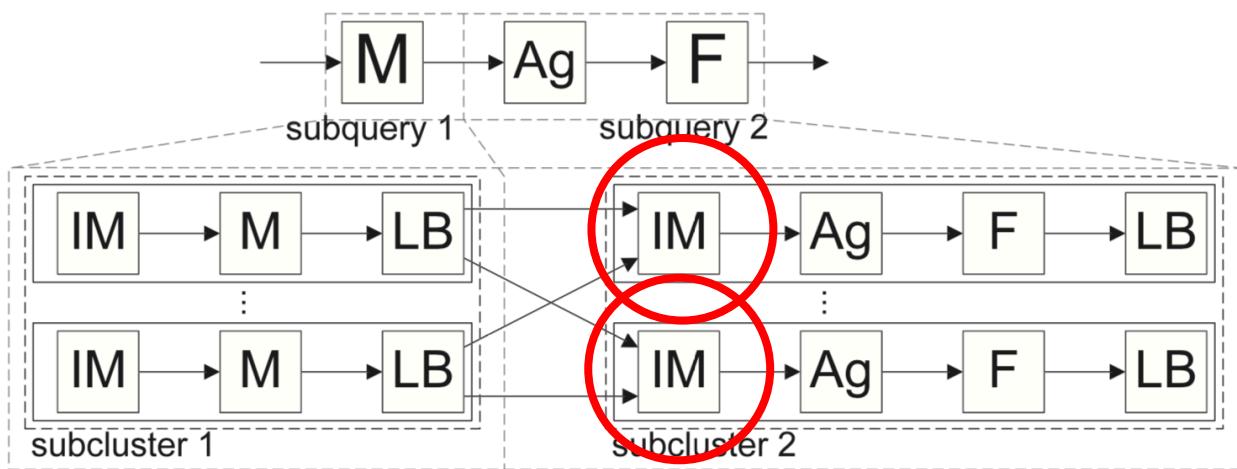
Disorder from parallelism



how to merge several timestamp-sorted streams...



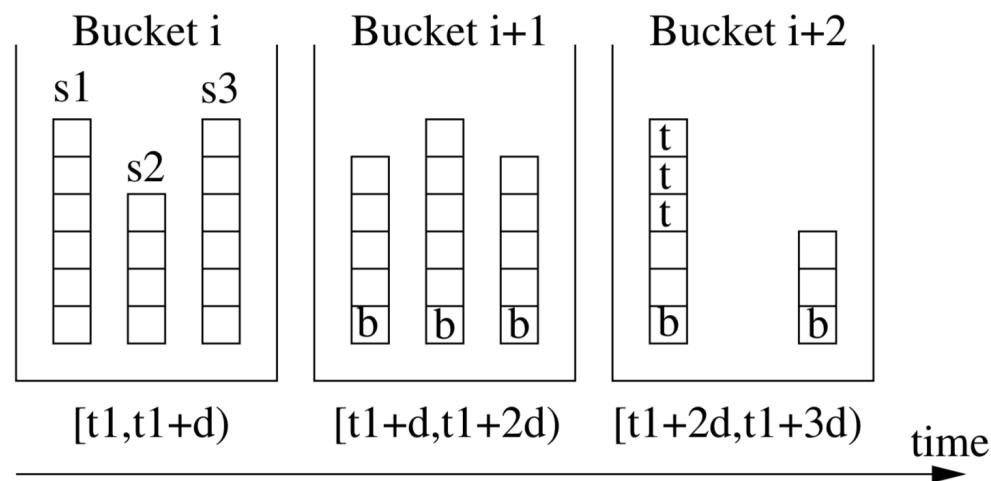
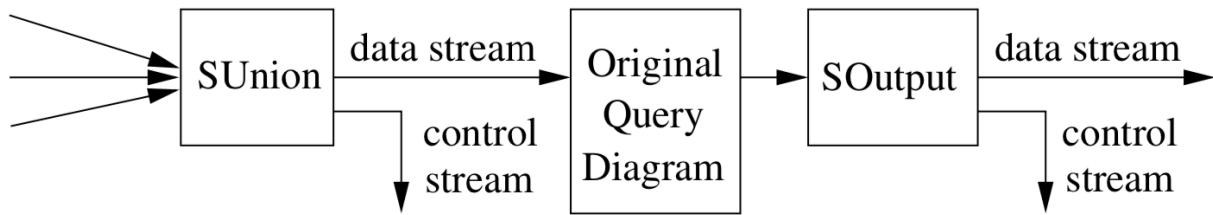
...into one timestamp-sorted stream?



Input Merger

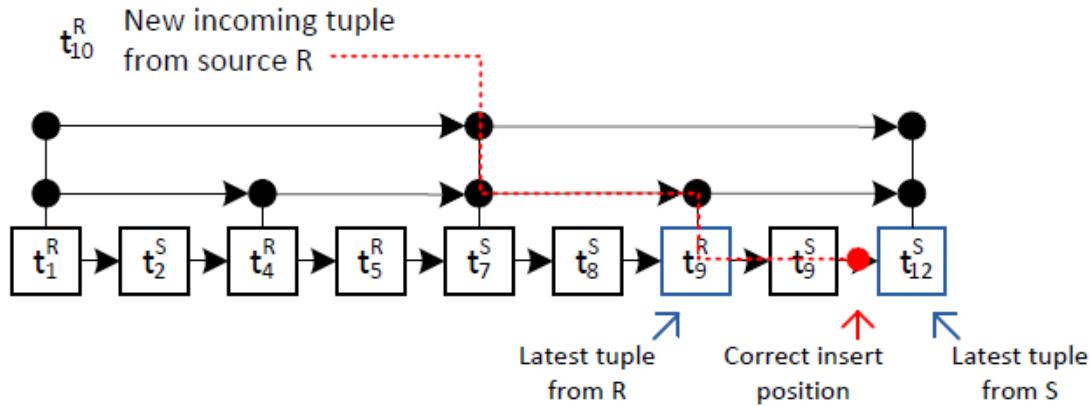
Upon: arrival of t from stream i :

- 9: $\text{buffer}[i].enqueue(t)$
- 10: **if** $\forall i$, $\text{buffer}[i].nonEmpty()$ **then**
- 11: $t_o = \text{earliestTuple}(\text{buffer})$
- 12: **if** $\neg \text{isDummy}(t_o)$ **then**
- 13: $\text{forward}(t_o)$

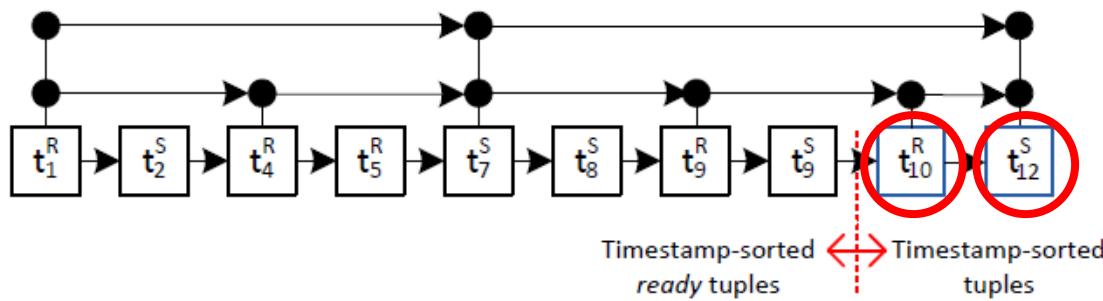


Balazinska, Magdalena, et al. "Fault-tolerance in the Borealis distributed stream processing system." *Proceedings of the 2005 ACM SIGMOD international conference on Management of data.* 2005.

a) Insertion of a new tuple



b) Distinguishing *ready* from non *ready* tuples



Which one to choose?

- Different options have different costs
 - might require different types of data structures
 - might require “special tuples” (more on this in the following part of the tutorial)
- When the price of merge-sorting is paid by who provides the tuples rather than who receives them, the system can scale better

Tutorial:

The Role of Event-Time Order in Data Streaming Analysis

Vincenzo Gulisano

Chalmers University of Technology
Gothenburg, Sweden
vincenzo.gulisano@chalmers.se

Bastian Havers

Chalmers University of Technology & Volvo Cars
Gothenburg, Sweden
havers@chalmers.se

Dimitris Palyvos-Giannas

Chalmers University of Technology
Gothenburg, Sweden
palyvos@chalmers.se

Marina Papatriantafilou

Chalmers University of Technology
Gothenburg, Sweden
ptrianta@chalmers.se

Tutorial:

The Role of Event-Time Order in Data Streaming Analysis

Vincenzo Gulisano

Chalmers University of Technology
Gothenburg, Sweden
vincenzo.gulisano@chalmers.se

Bastian Havers

Chalmers University of Technology & Volvo Cars
Gothenburg, Sweden
havers@chalmers.se

Dimitris Palyvos-Giannas

Chalmers University of Technology
Gothenburg, Sweden
palyvos@chalmers.se

Marina Papatriantafilou

Chalmers University of Technology
Gothenburg, Sweden
ptrianta@chalmers.se

Agenda

- Motivation, preliminaries and examples about data streaming and Stream Processing Engines
- Causes of out-of-order data and solutions enforcing total-ordering
- Pros/Cons of total-ordering
- Relaxation of total-ordering and the watermarks

Agenda

- Motivation, preliminaries and examples about data streaming and Stream Processing Engines
- Causes of out-of-order data and solutions enforcing total-ordering
- Pros/Cons of total-ordering
- Relaxation of total-ordering and the watermarks

Pros/Cons of total-ordering

- Cons:
 - Expensive (computation- and latency-wise)
 - An “overkill” for certain applications (more of this in the following slides)
- Pros:
 - Determinism
 - Synchronization
 - Eager purging of stale state

Pros/Cons of total-ordering

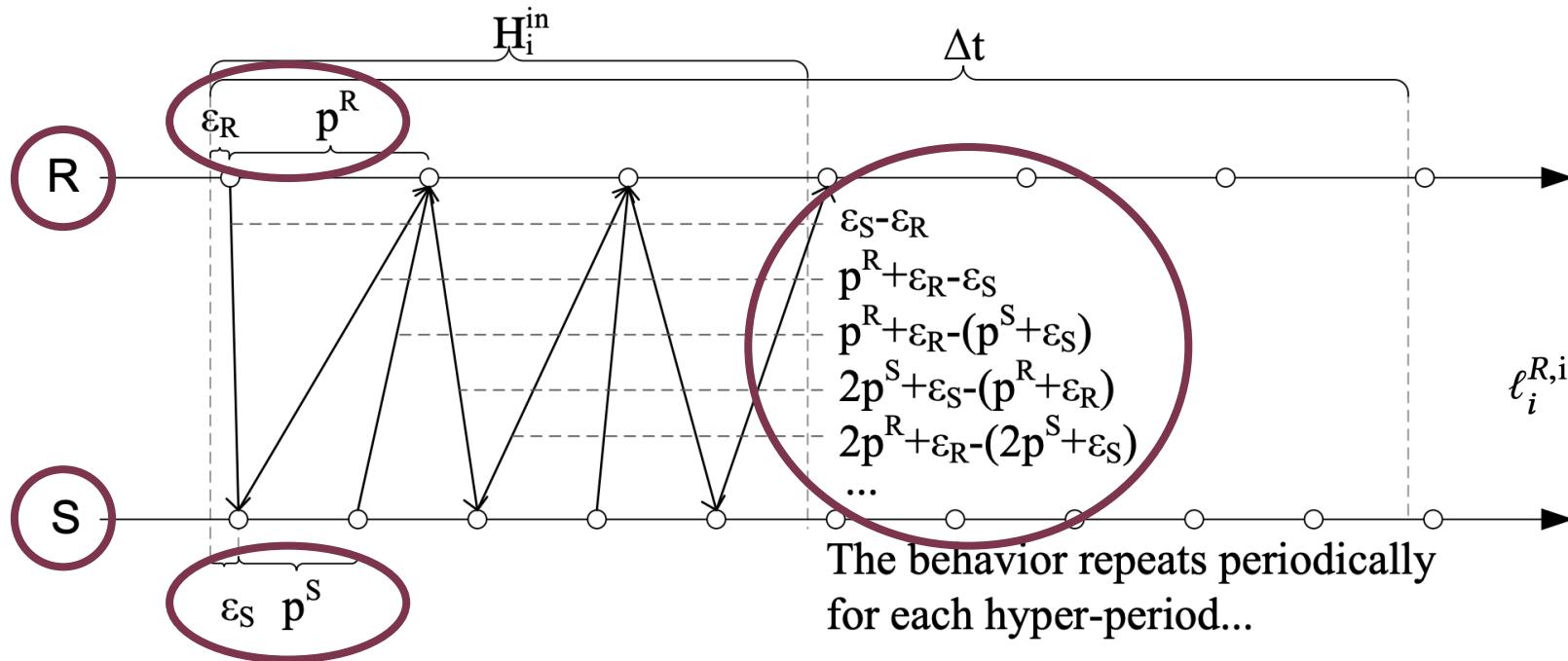
- Cons:
 - Expensive (computation- and latency-wise)
 - An “overkill” for certain applications (more of this in the following slides)
- Pros:
 - Determinism
 - Synchronization
 - Eager purging of stale state

Cost

- We need to temporary maintain tuples
 - Linear in the number of tuples we receive, which depends on the streams' rate
- We need to sort tuples... $O(n \log(n))$
 - (n is number of sources or tuples, depending on the case)
- We need data from all sources, the processing latency depends on the slowest source.
 - The latency overhead can be estimated based on the sources' rates¹

¹Gulisano, Vincenzo, et al. "Performance modeling of stream joins." *Proceedings of the 11th ACM International Conference on Distributed and Event-based Systems*. 2017.

Estimating the latency overhead



$$H_i^{\text{in}} = \text{LCM} (p_i^R, p_i^S)$$

$$\ell_i^{R,\text{in}} = \sum_{m=0}^{H_i^{\text{in}} r_i - 1} \left(p_i^S \left[\frac{mp_i^R + \varepsilon_R}{p_i^S} \right] + \varepsilon_S - (mp_i^R + \varepsilon_R) \right)$$

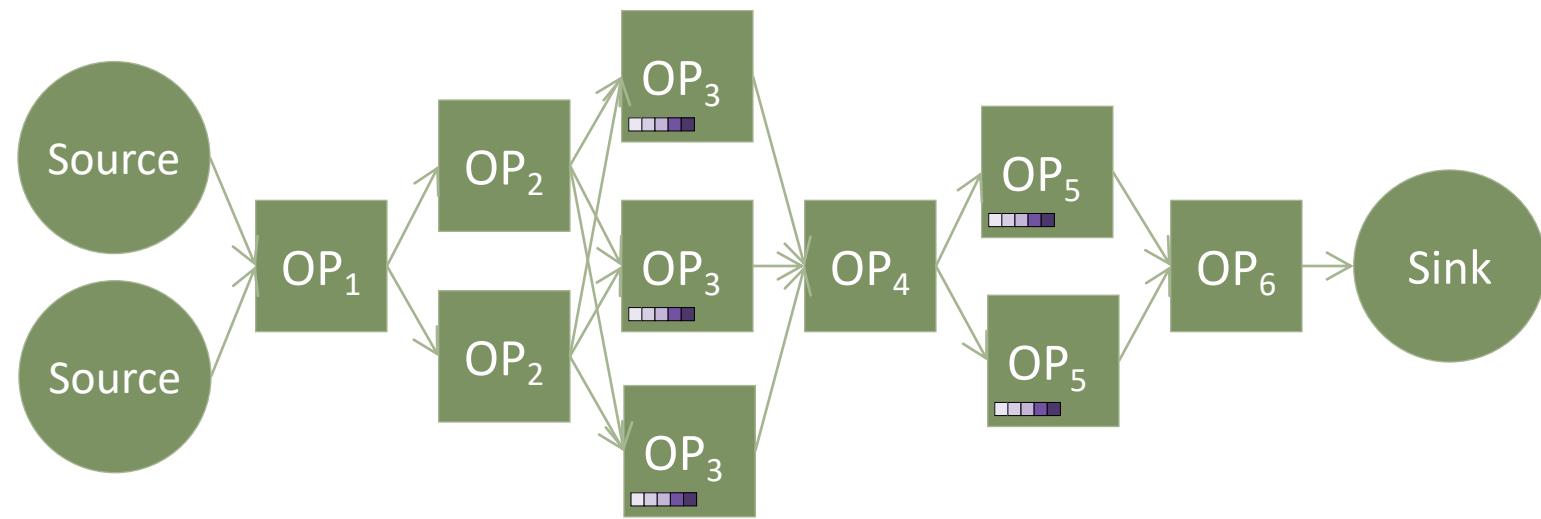
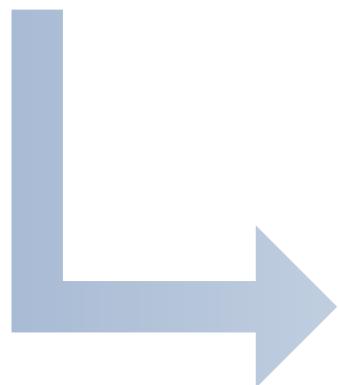
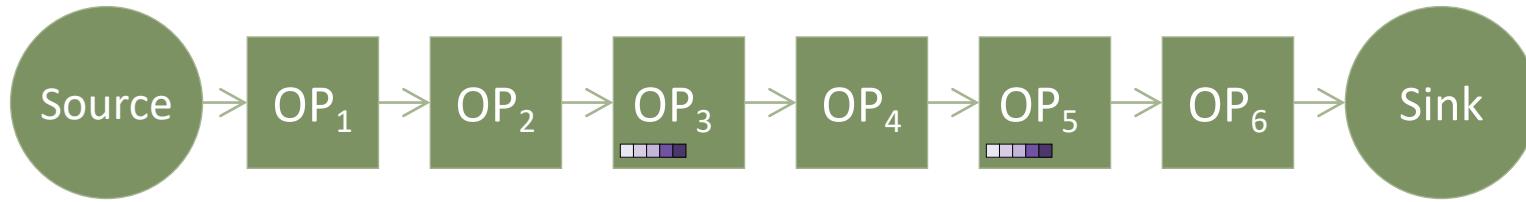
$$\ell_i^{\text{in}} = \frac{\ell_i^{R,\text{in}} + \ell_i^{S,\text{in}}}{H_i^{\text{in}} (r_i + s_i)}.$$

¹Gulisano, Vincenzo, et al. "Performance modeling of stream joins." Proceedings of the 11th ACM International Conference on Distributed and Event-based Systems. 2017.

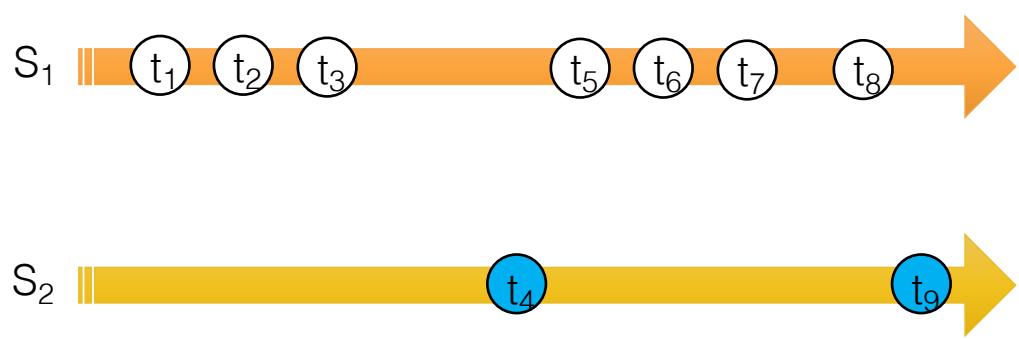
Pros/Cons of total-ordering

- Cons:
 - Expensive (computation- and latency-wise)
 - An “overkill” for certain applications (more of this in the following slides)
- Pros:
 - Determinism
 - Synchronization
 - Eager purging of stale state

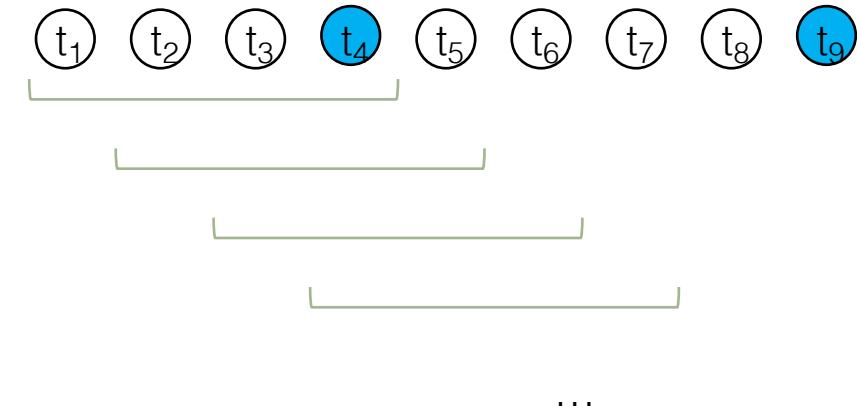
Determinism



Determinism



Tuple-based window, size: 4 / advance: 1



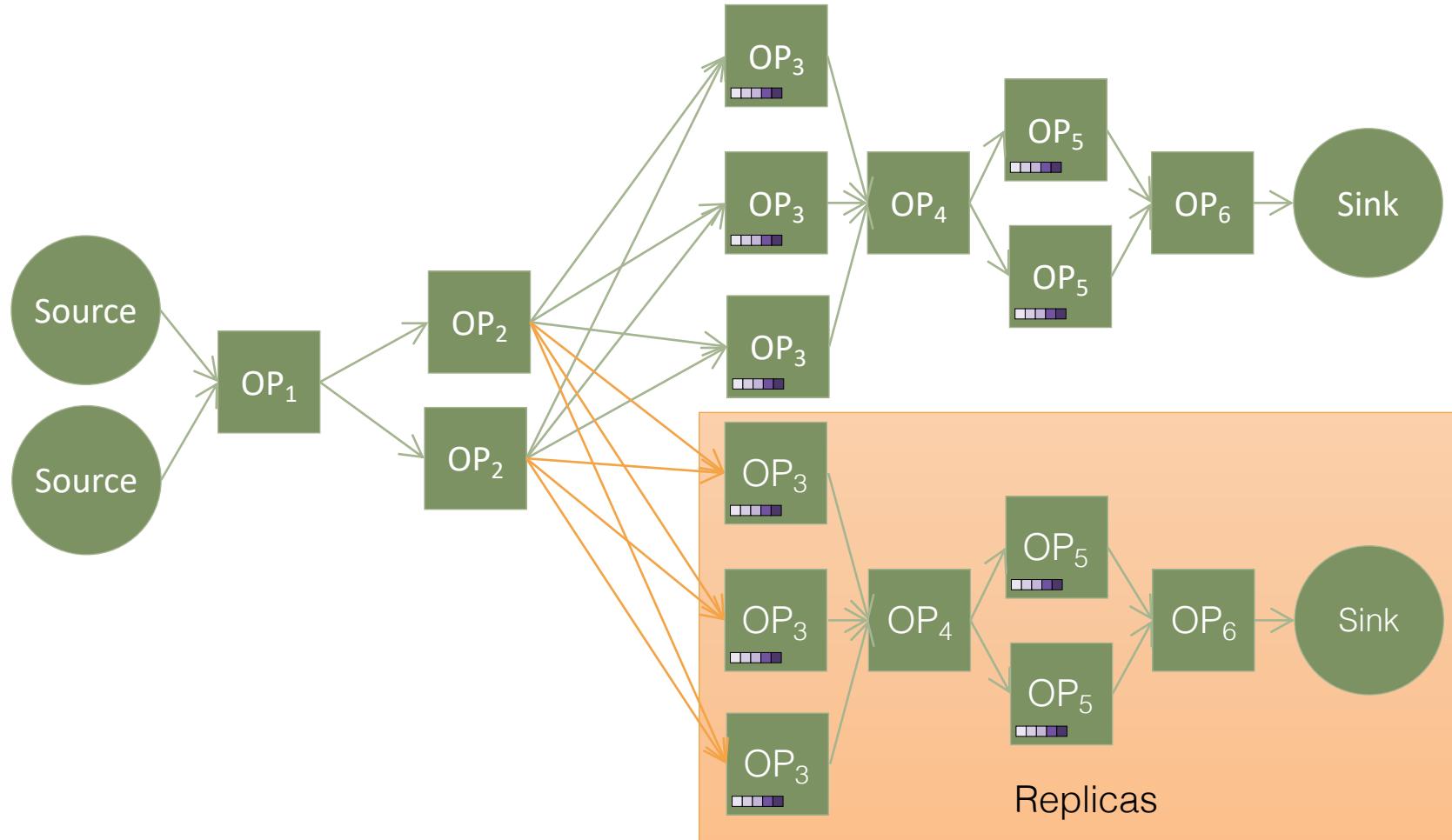
Hwang, Jeong-Hyon, Ugur Cetintemel, and Stan Zdonik. "Fast and reliable stream processing over wide area networks." 2007 IEEE 23rd International Conference on Data Engineering Workshop. IEEE, 2007.

Gulisano, Vincenzo, Yiannis Nikolakopoulos, Marina Papatriantafilou, and Philippas Tsigas. "Scalejoin: A deterministic, disjoint-parallel and skew-resilient stream join." IEEE Transactions on Big Data (2016).

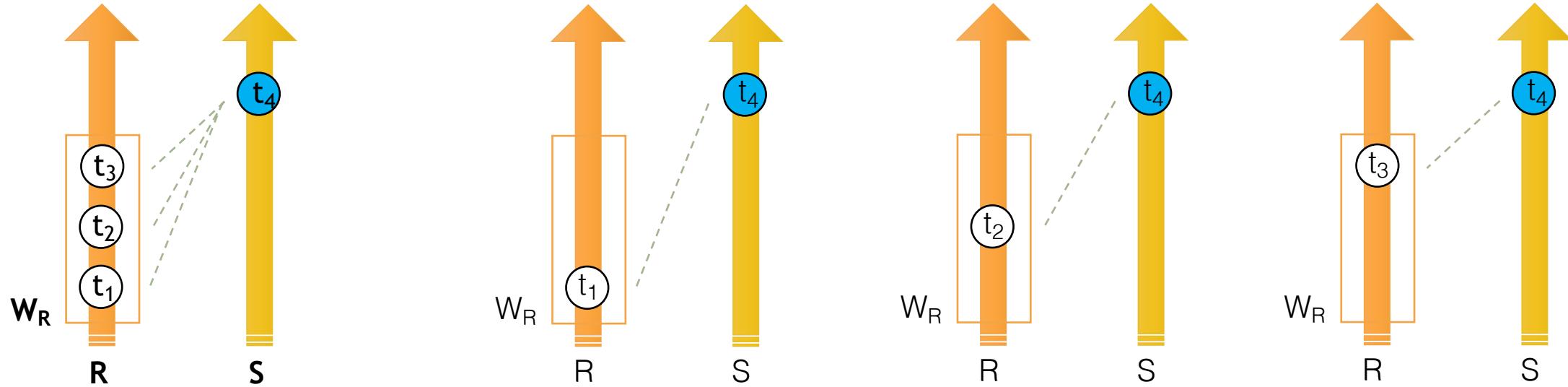
Pros/Cons of total-ordering

- Cons:
 - Expensive (computation- and latency-wise)
 - An “overkill” for certain applications (more of this in the following slides)
- Pros:
 - Determinism
 - Synchronization
 - Eager purging of stale state

Synchronization

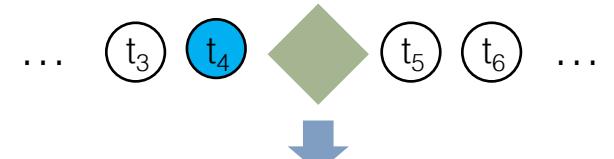
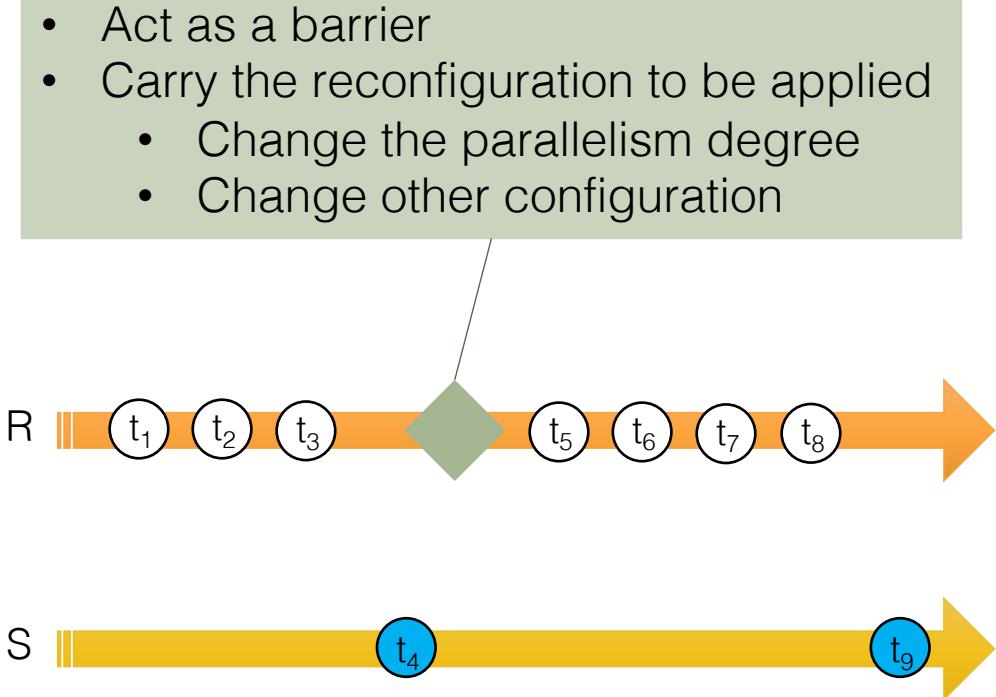


Synchronization



Gulisano, Vincenzo, Yiannis Nikolopoulos, Marina Papatriantafilou, and Philippas Tsigas. "Scalejoin: A deterministic, disjoint-parallel and skew-resilient stream join." *IEEE Transactions on Big Data* (2016).

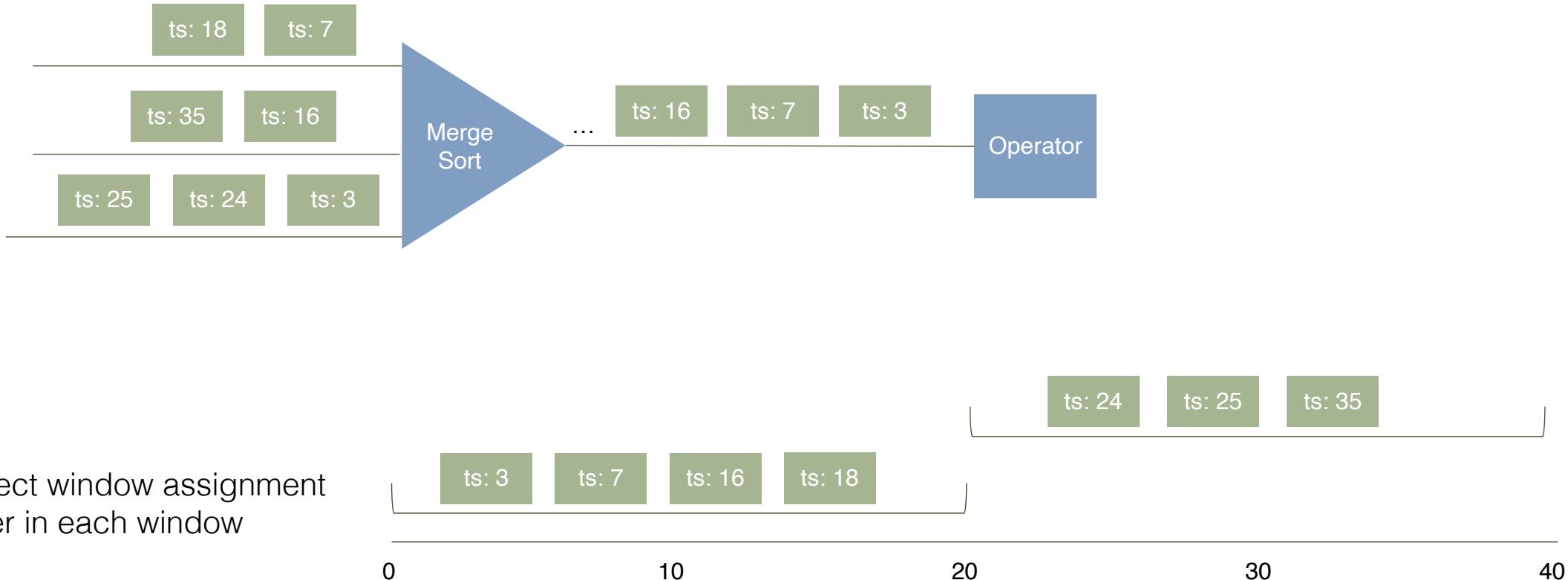
Synchronization



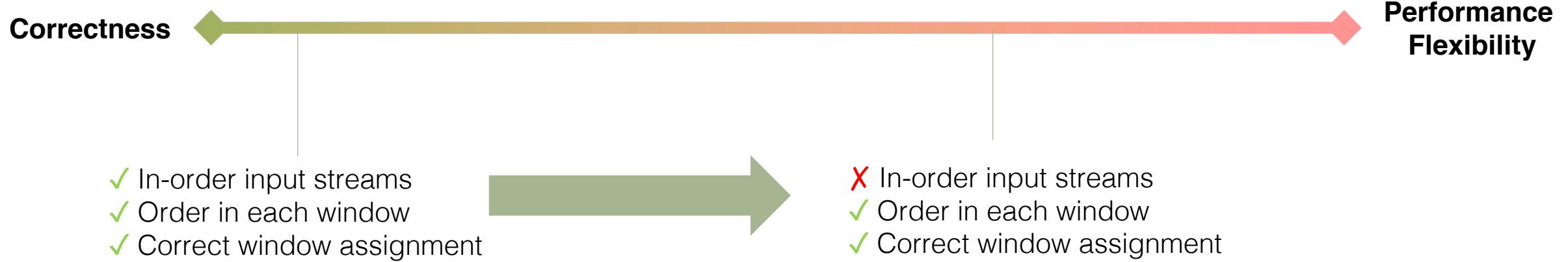
Agenda

- Motivation, preliminaries and examples about data streaming and Stream Processing Engines
- Causes of out-of-order data and solutions enforcing total-ordering
- Pros/Cons of total-ordering
- Relaxation of total-ordering and the watermarks

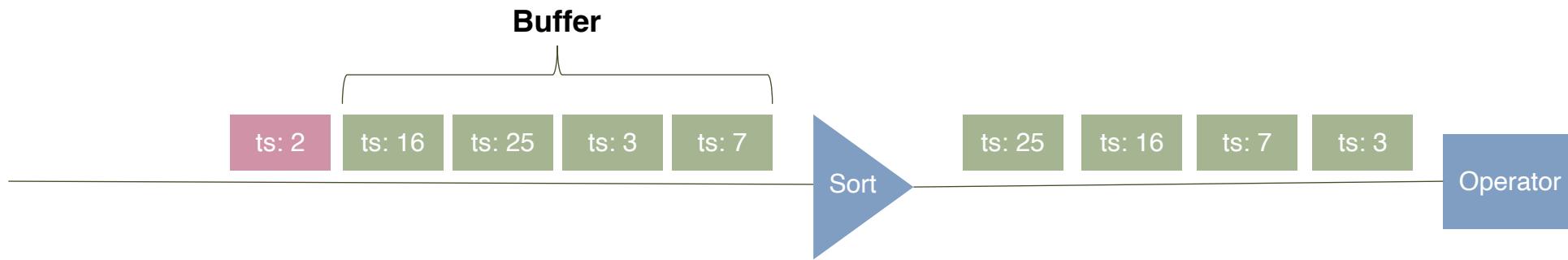
Total-Ordering Recap



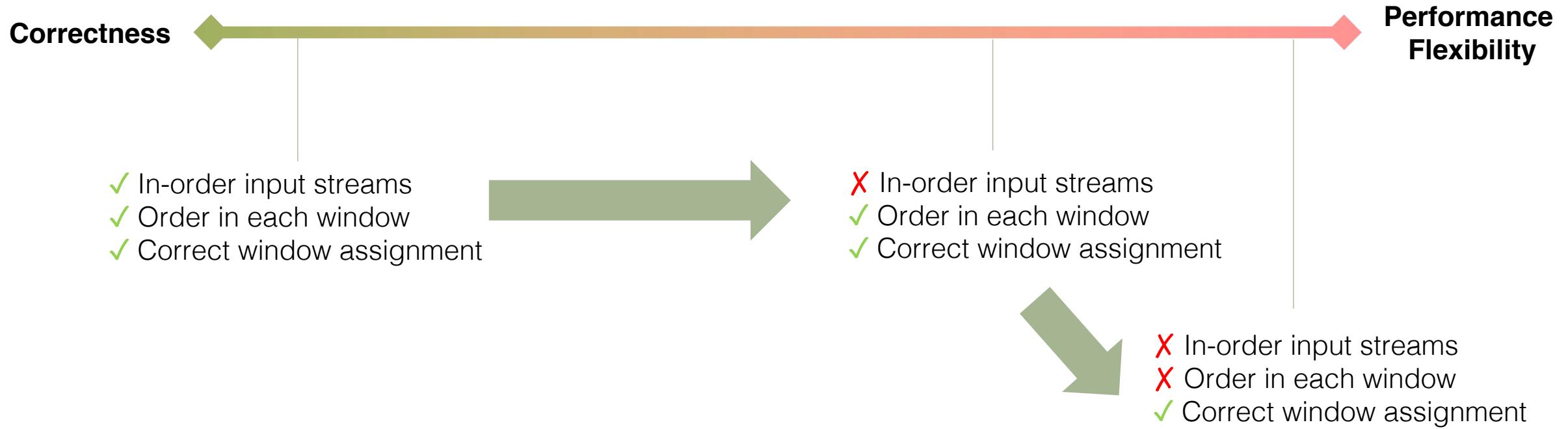
Relaxing Correctness



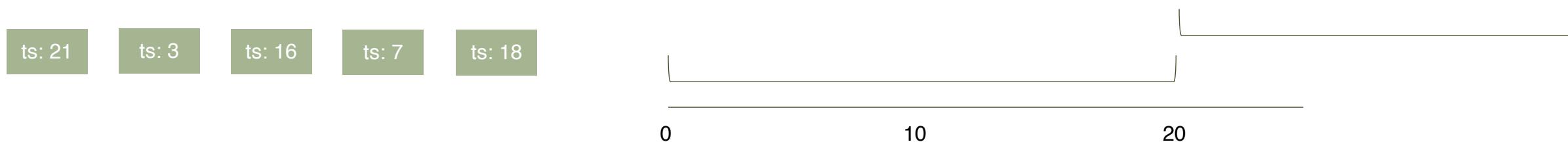
Buffering Approaches



Relaxing Correctness

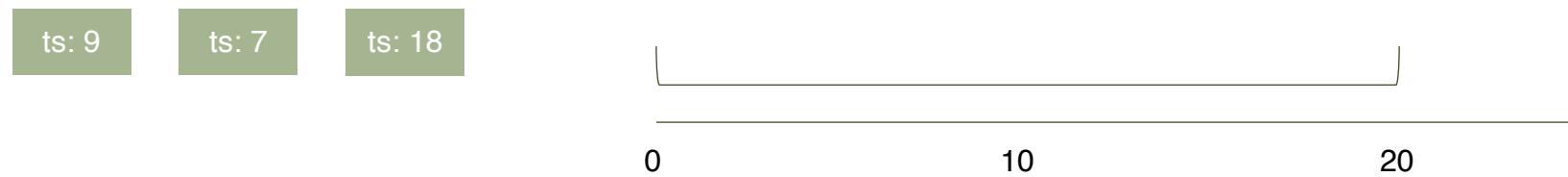


Disorder + Correct Window Assignment



1. When to **create** each window?
2. When to **close** each window (and produce result)?

Creating Windows



Closing Windows



Operator needs some guarantee
it will not receive tuples with $ts < W$



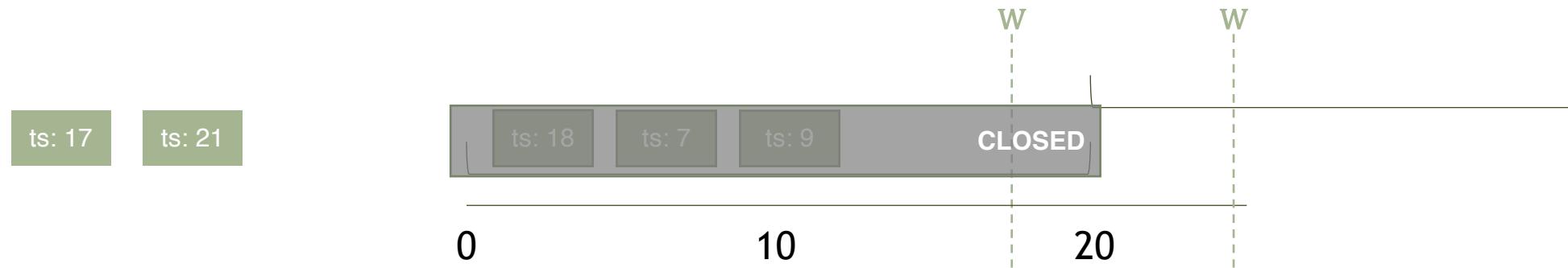
Safely close all windows where
 $right_boundary \leq W$

Watermarks

Assume we can compute a monotonic function

$$F: O \rightarrow E$$

that returns $W \in E$ the earliest event time of any tuple
that can arrive at operator O in the future

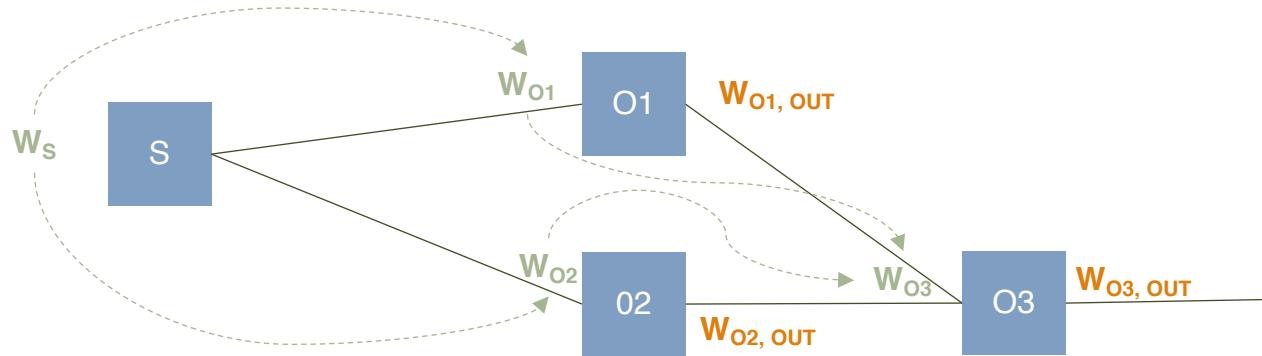


We call the value of this monotonic function F^* the (low) **watermark** of operator O !

- ✓ Monotonicity → no tuples with $ts < W$ will arrive in the future.
- ✓ Solves problem of safely closing windows!

* The watermark of operator O is function of O and all its upstream peers, but we omit the latter for brevity.

Watermarks in Practice

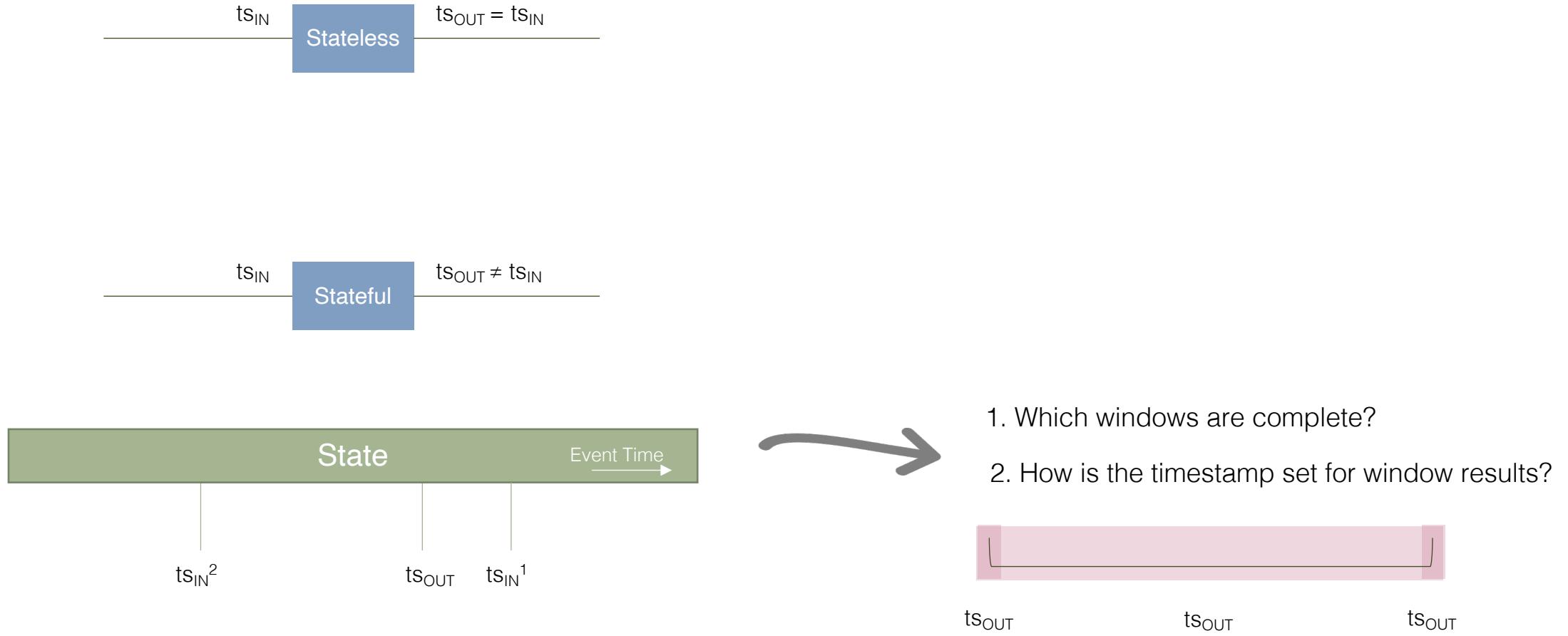


Input Watermark: Earliest ts that O can receive.

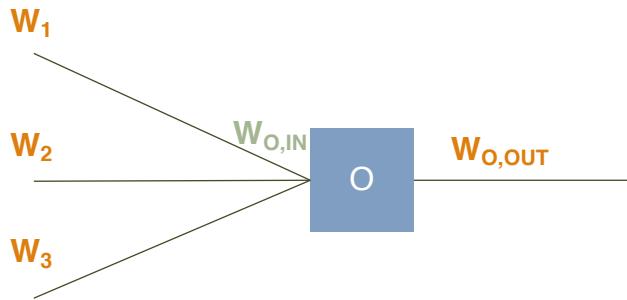
Output Watermark: Earliest ts that O can emit.

- ✓ Watermarks are **generated** at the sources.
- ✓ They (conceptually) **flow** through the pipeline.
- ✓ They propagate **regardless of data filtering** → all operators have up-to-date view of time

Input & Output Timestamps

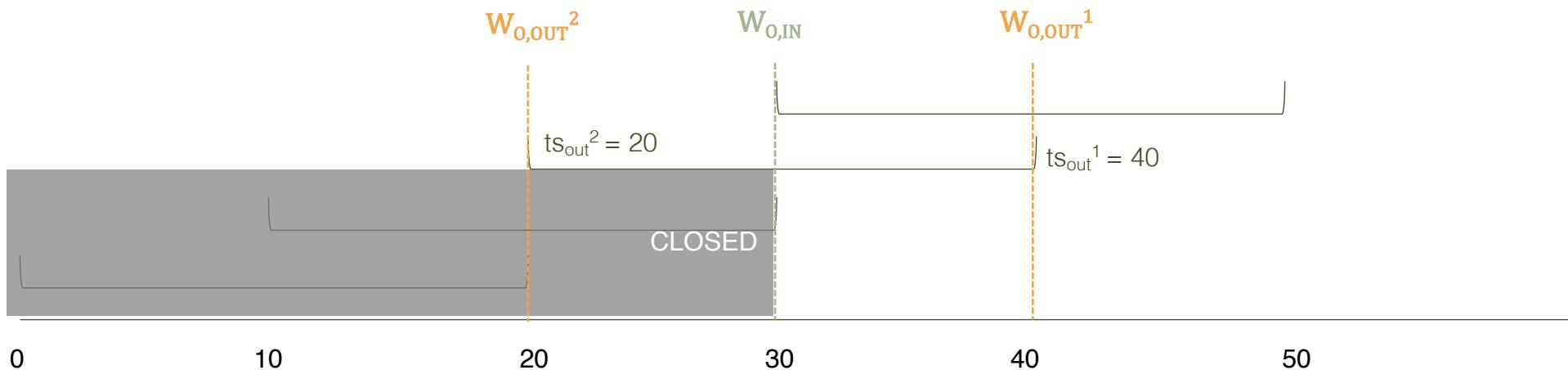


Computing Watermarks

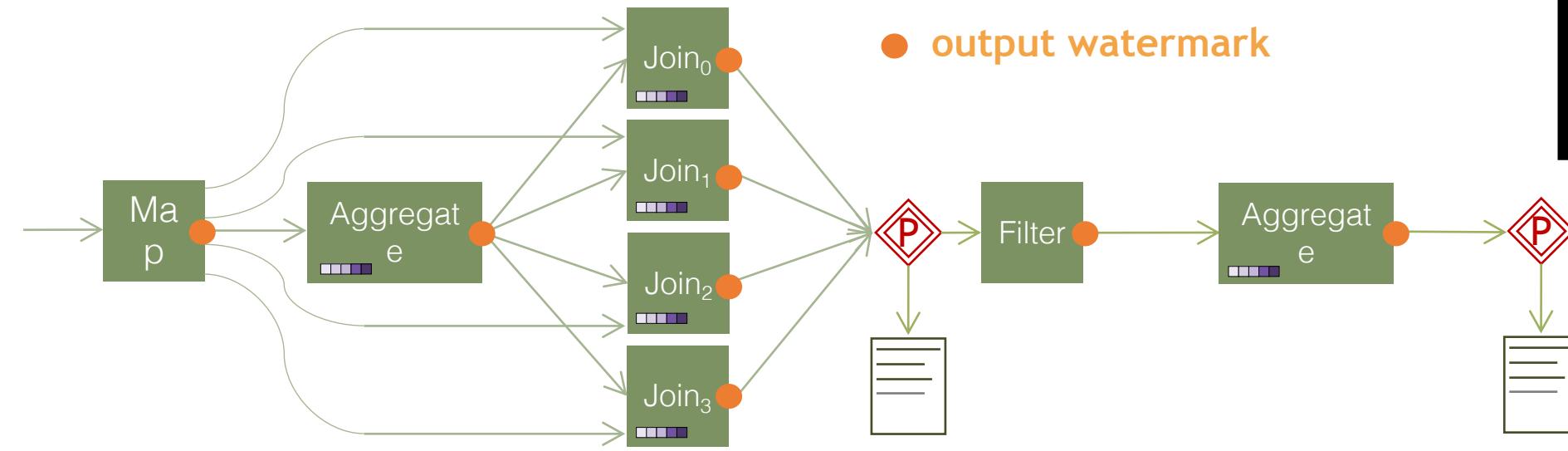


$$W_{O,IN} = \min(W_1, W_2, W_3)$$

$$W_{O,OUT} = \begin{cases} W_{O,IN} & \text{if } O \text{ stateless} \\ g(\text{state}_O, \text{semantics}_O) & \text{otherwise} \end{cases}$$

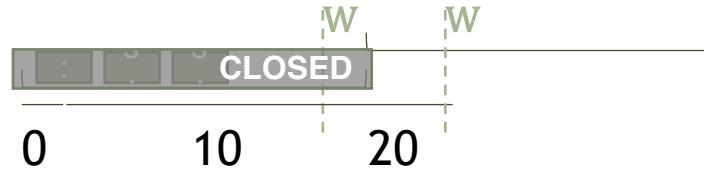


Flink Example

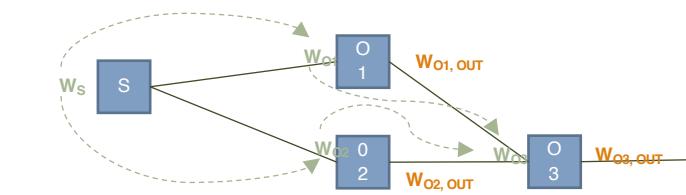


● output watermark

1. Result Correctness from correct window assignment



2. Watermark propagation



Generating Watermarks

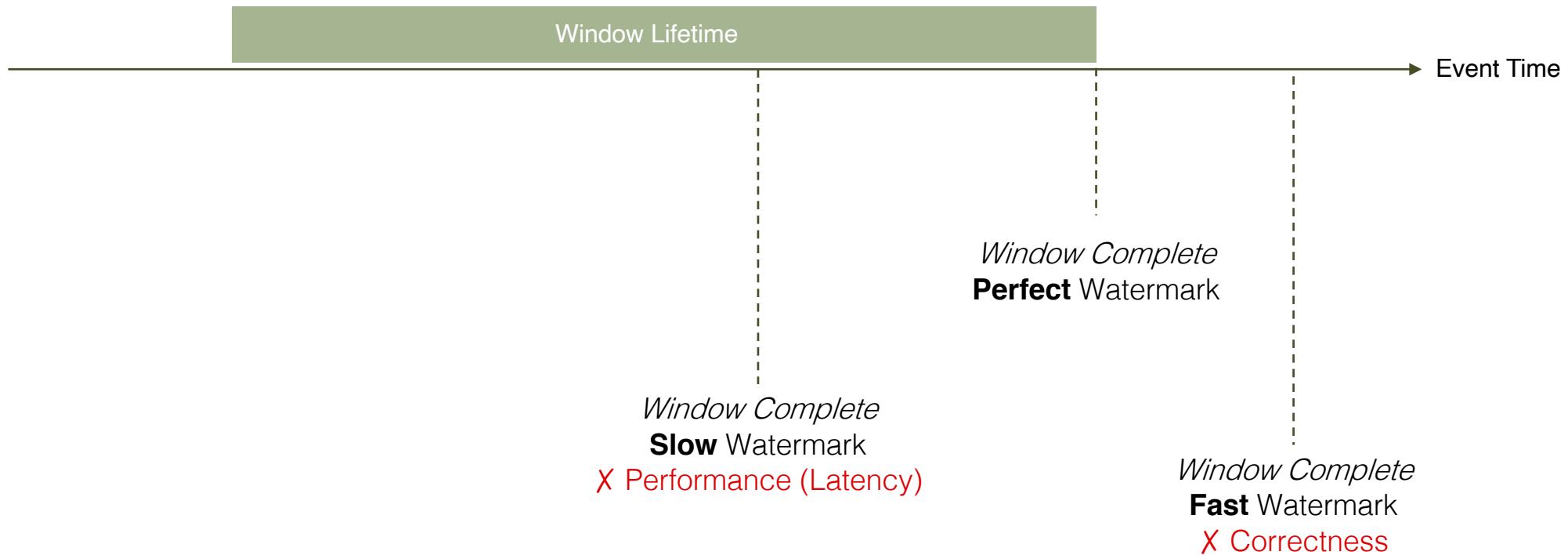
- **Perfect Watermarks**

- Sorted data or very predictable data sources.
- Determinism without sorting for order-independent window functions.
- Disorder inside windows is still possible!
- Sorting possible if needed, but not imposed.

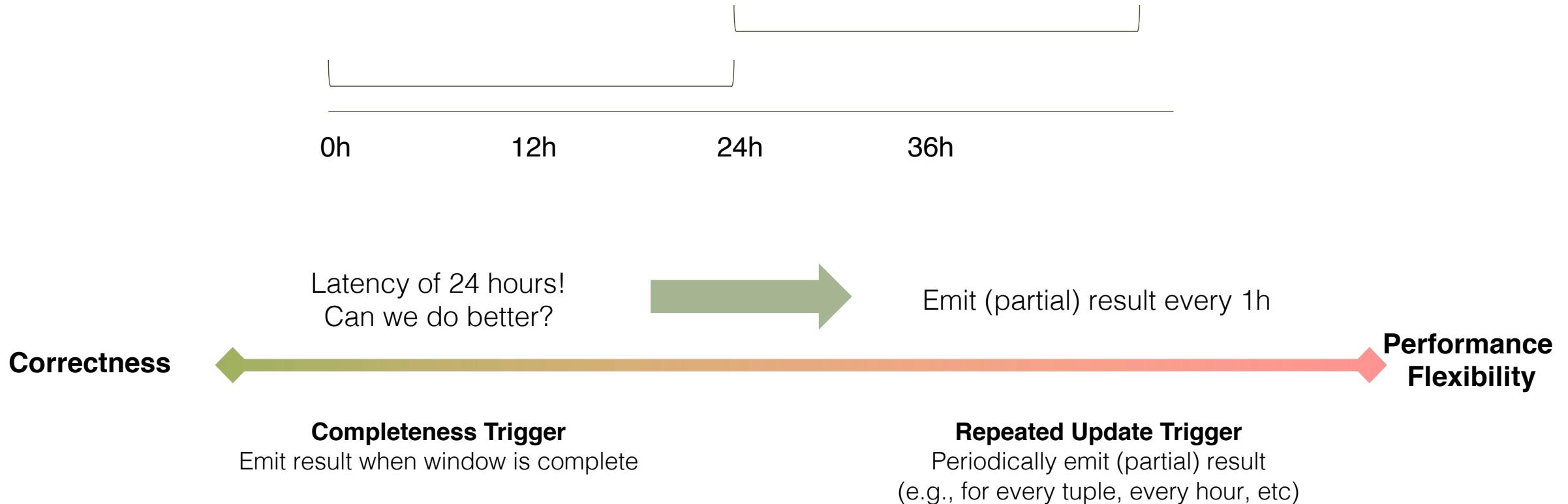
- **Heuristic Watermarks**

- When impossible to perfectly predict data (e.g., distributed sources).
- Best-effort prediction of event-time progress.
- Possibility for late data.
- More knowledge about internals of sources → less mispredictions (late data).

Fast and Slow Watermarks



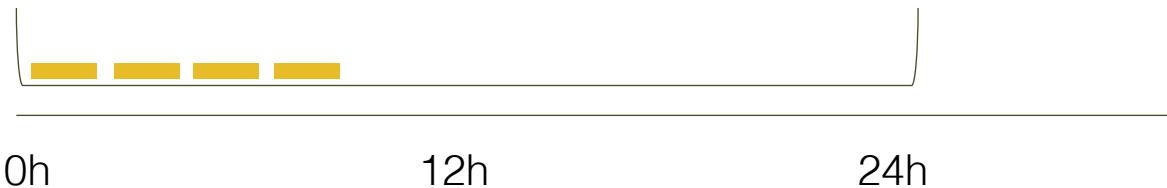
Triggering



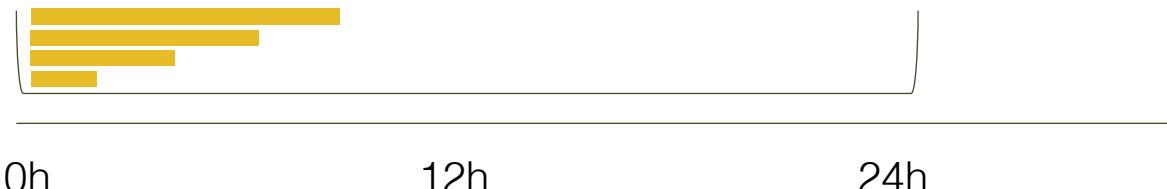
Repeated Update Trigger Results

Repeated Update Trigger → Every 1h

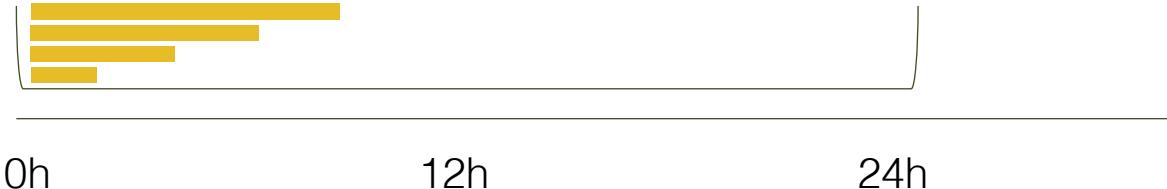
Discarding



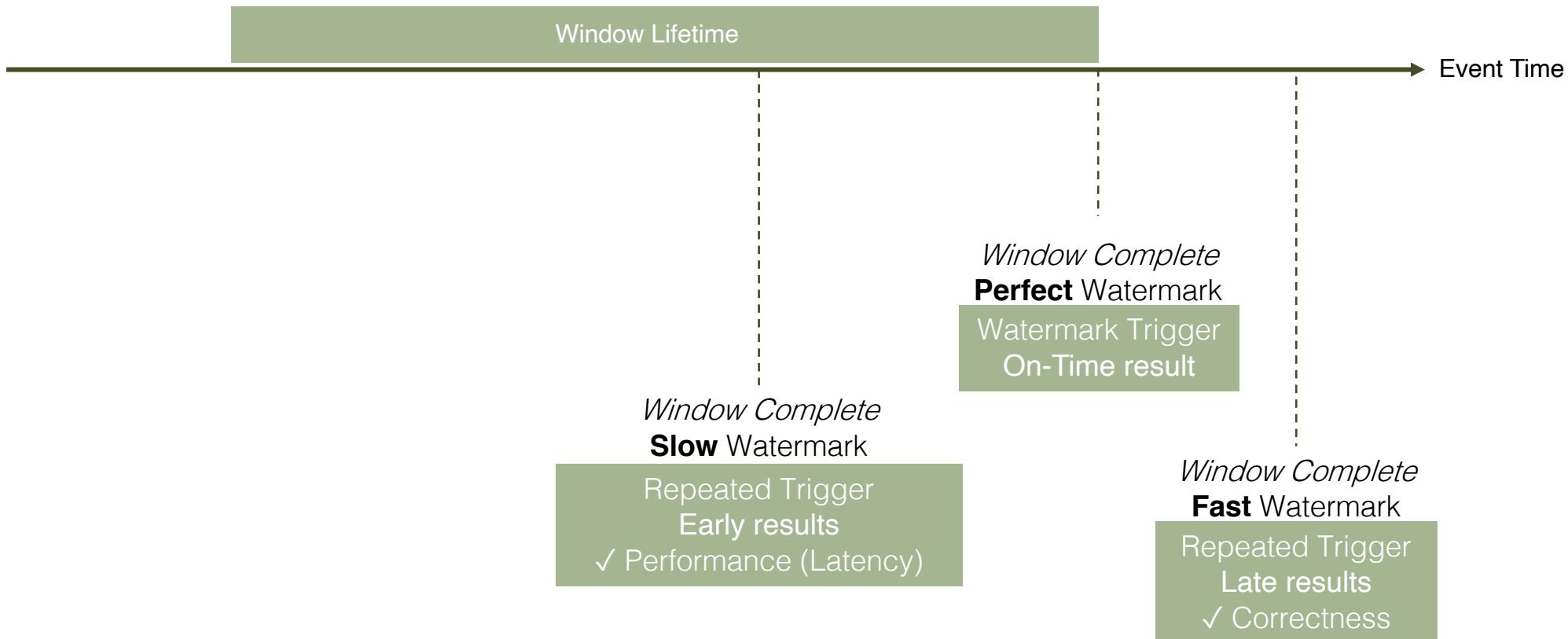
Accumulating



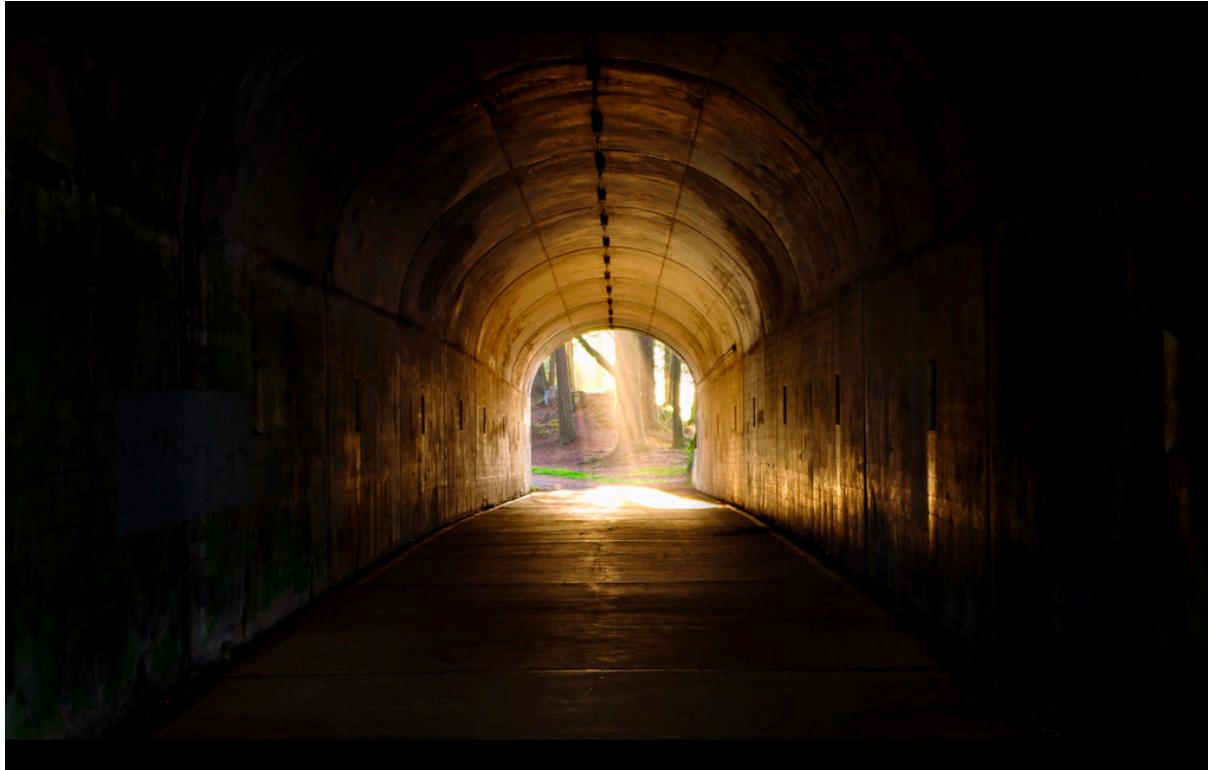
Accumulating +
Retracting



Putting It All Together



... the light at the end of the tunnel ...



- Motivation, preliminaries and examples about data streaming and Stream Processing Engines
- Causes of out-of-order data and solutions enforcing total-ordering
- Pros/Cons of total-ordering
- Relaxation of total-ordering and the watermarks

To summarize:

Event time advances based on:	Cons	Pros
Tuples themselves	<ul style="list-style-type: none">• Costly merge-sorting• Coupled processing / output of tuples	<ul style="list-style-type: none">• Determinism, for order sensitive / insensitive functions
Watermarks	<ul style="list-style-type: none">• Would require special support for order-sensitive functions• Latency depends on frequency of watermarks	<ul style="list-style-type: none">• Decoupled processing / output of tuples• Propagation of time passing even in the absence of tuples

Tutorial:

The Role of Event-Time Order in Data Streaming Analysis

Vincenzo Gulisano

Chalmers University of Technology
Gothenburg, Sweden
vincenzo.gulisano@chalmers.se

Bastian Havers

Chalmers University of Technology & Volvo Cars
Gothenburg, Sweden
havers@chalmers.se

Dimitris Palyvos-Giannas

Chalmers University of Technology
Gothenburg, Sweden
palyvos@chalmers.se

Marina Papatriantafilou

Chalmers University of Technology
Gothenburg, Sweden
ptrianta@chalmers.se