

Efficient Window Aggregation with Stream Slicing

Jonas Traub
Research Associate (TU Berlin)



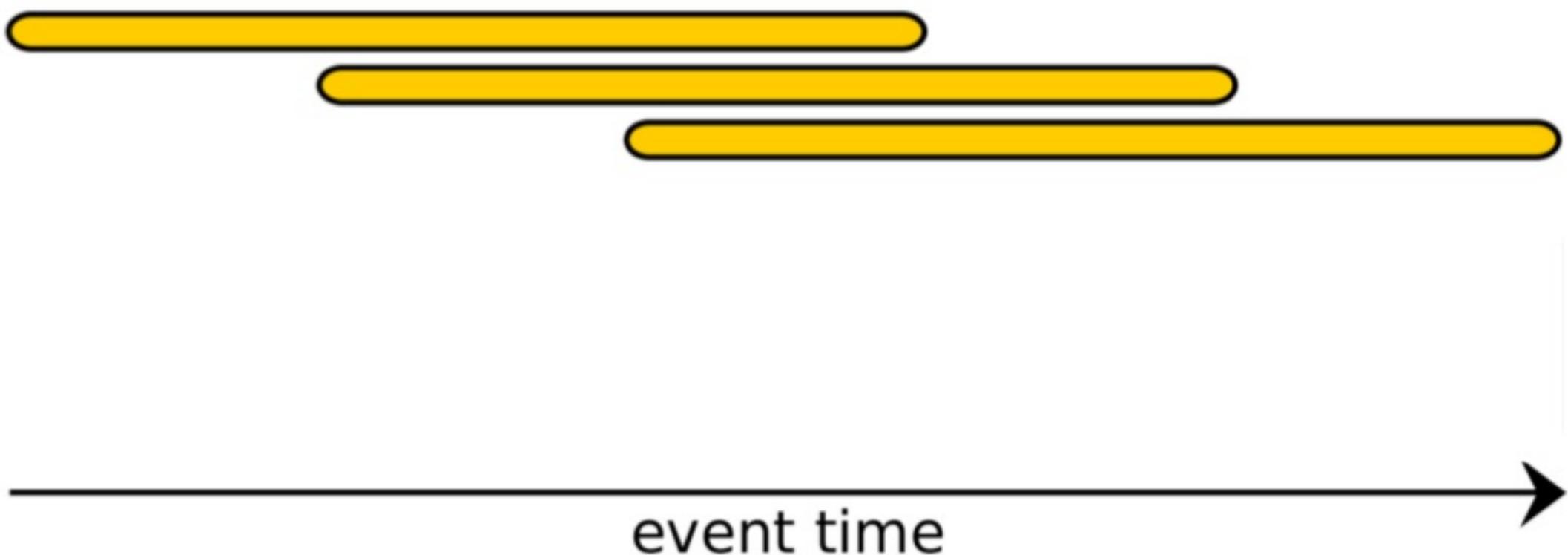
Philipp M. Grulich
Research Assistant (DFKI)



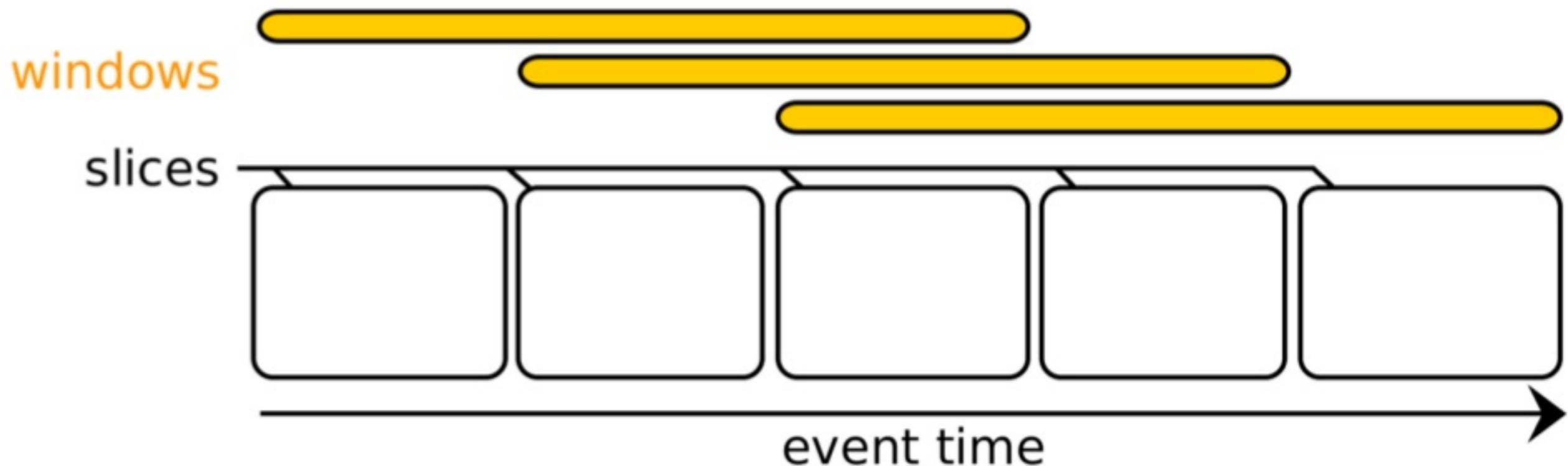
Stream Slicing Example

Stream Slicing Example

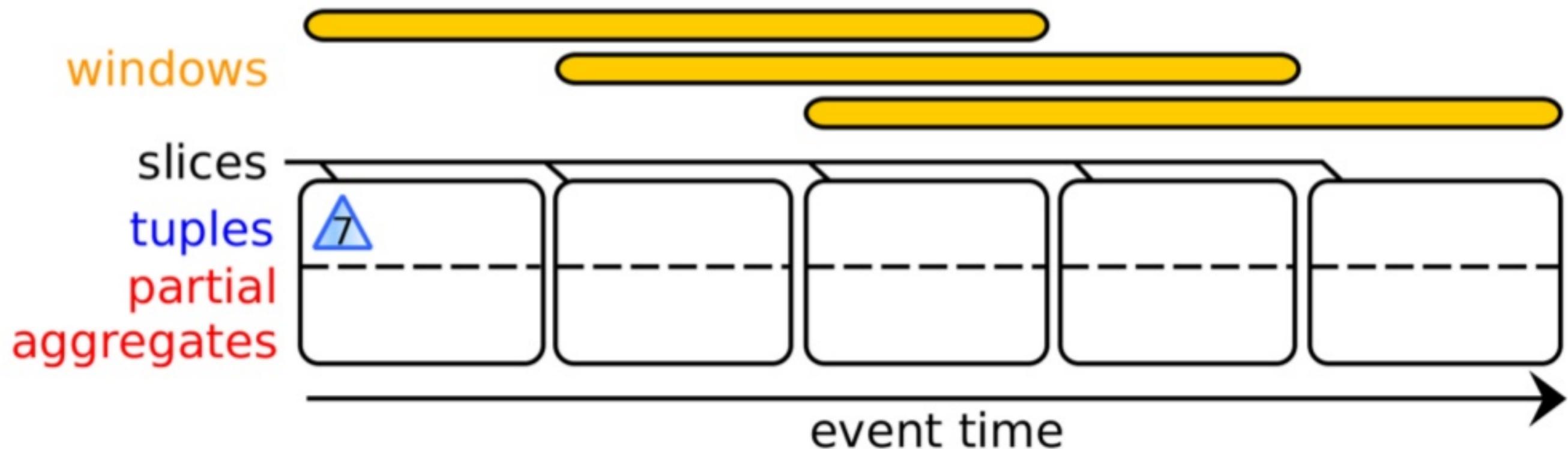
windows



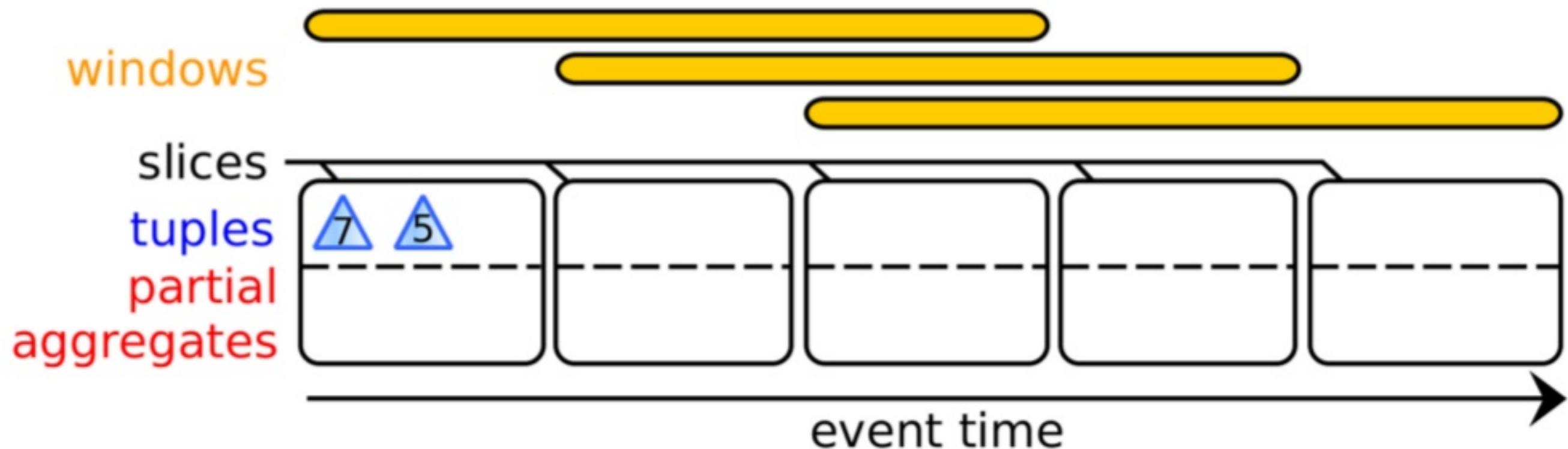
Stream Slicing Example



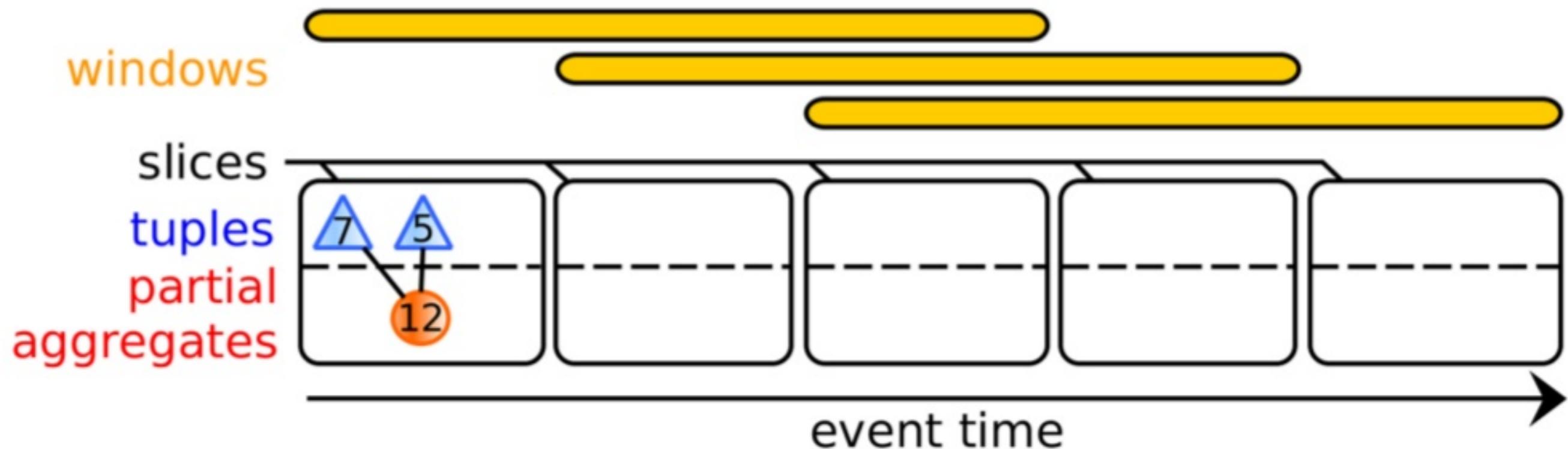
Stream Slicing Example



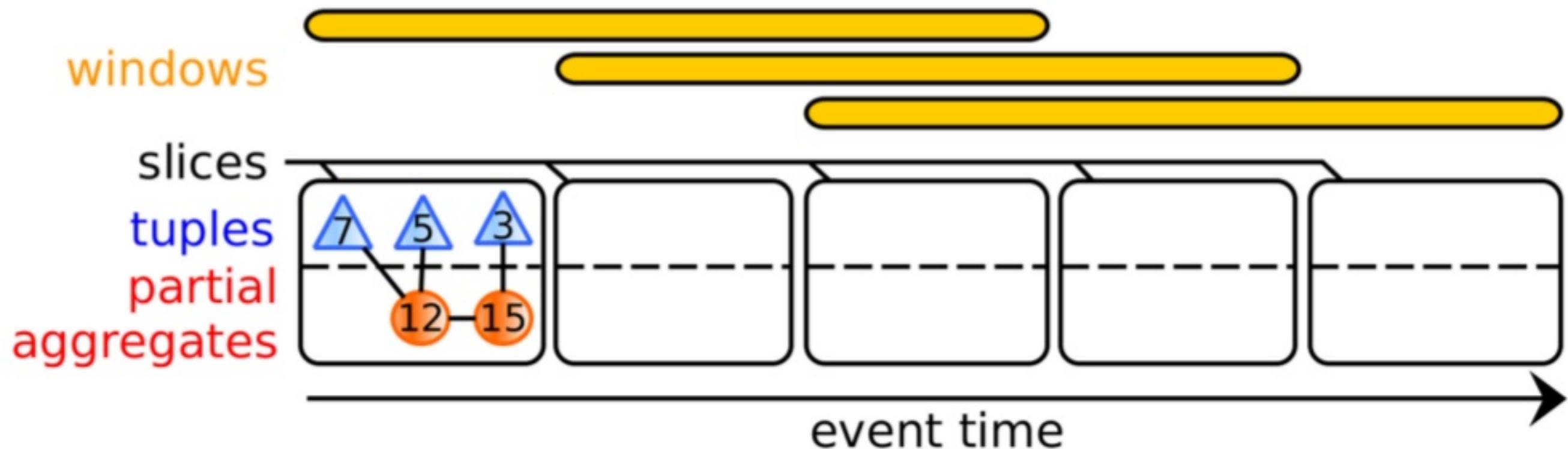
Stream Slicing Example



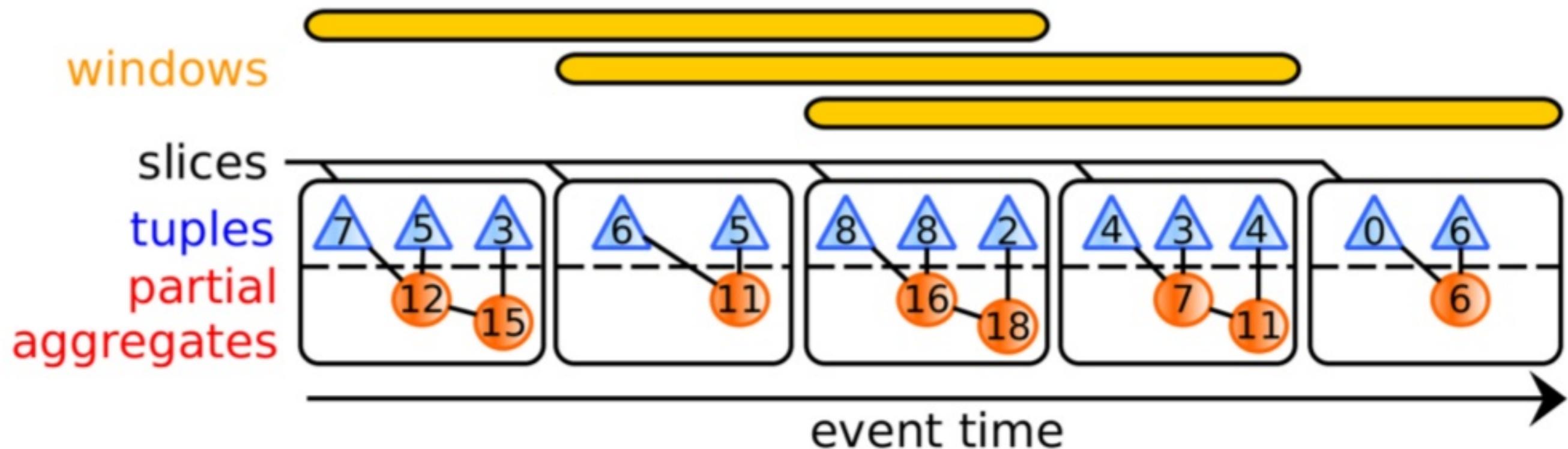
Stream Slicing Example



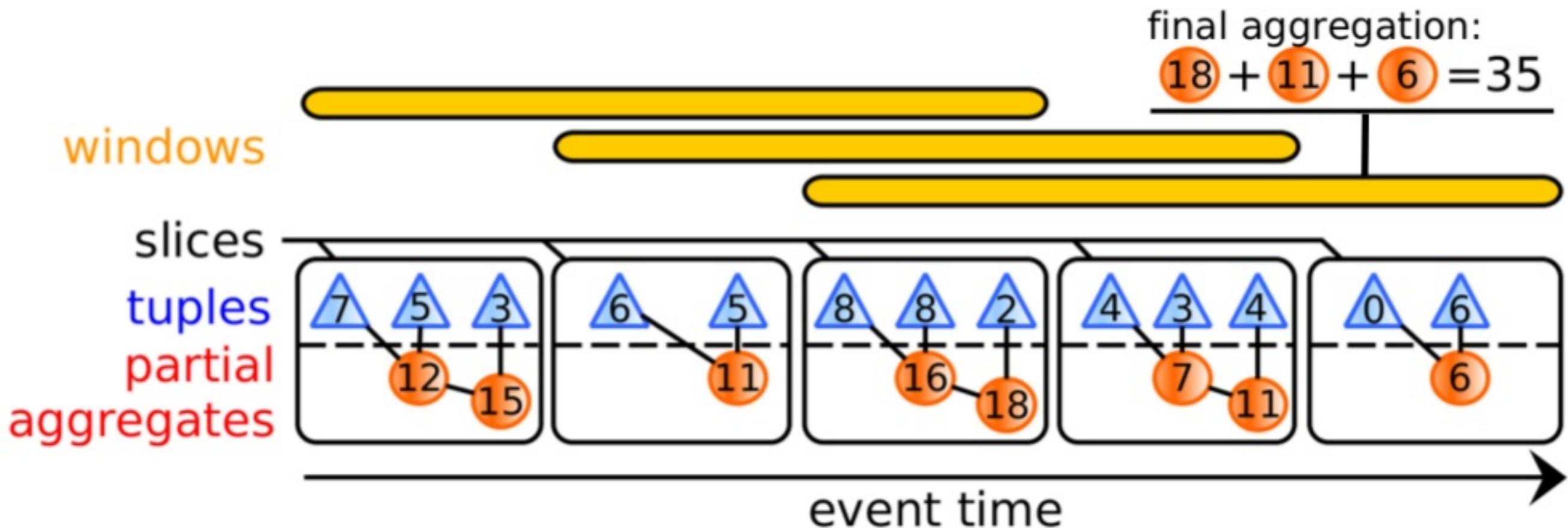
Stream Slicing Example



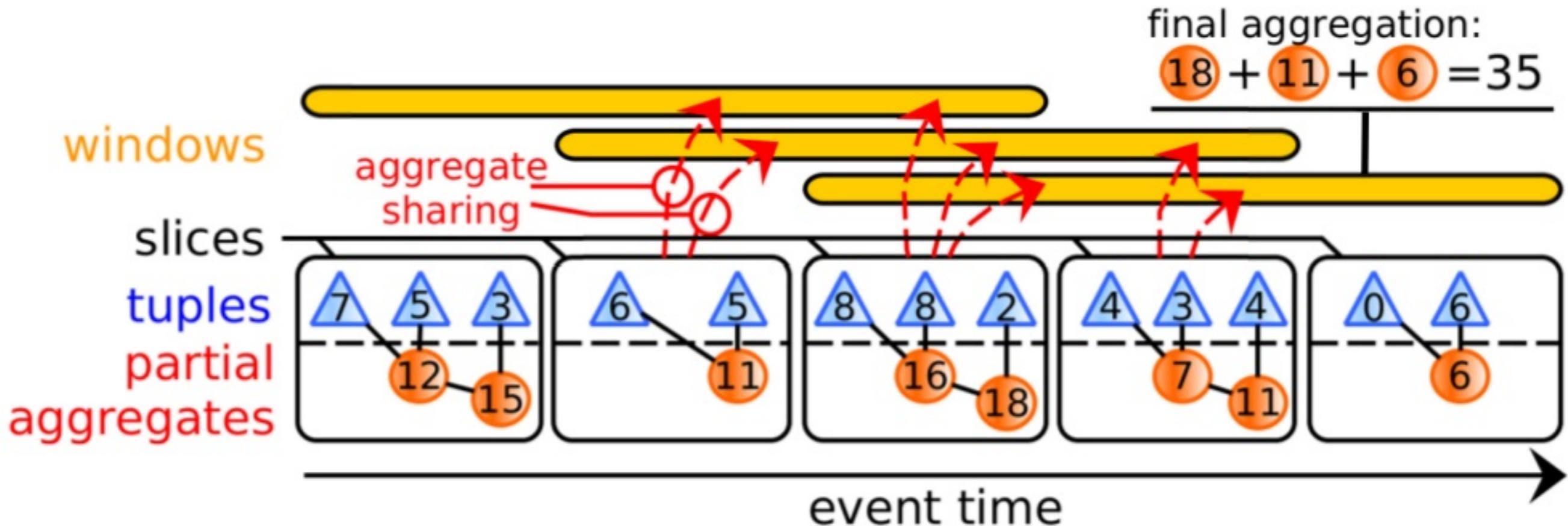
Stream Slicing Example



Stream Slicing Example



Stream Slicing Example



Stream Slicing Research

Cutty: Aggregate Sharing for User-Defined Windows

Paris Carbone[†]

Seif Haridi[†]

Jonas Traub[‡]

Volker Markl[‡]

Asterios Katsifodimos[‡]

Cutty: Aggregate Sharing for User-Defined Windows

Paris Carbone[†]

Jonas Traub[‡]

Asterios Katsifodimos[‡]

Seif Haridi[†]

Volker Markl[‡]

Scotty: Efficient Window Aggregation for out-of-order Stream Processing

Jonas Traub¹, Philipp M. Grulich², Alejandro Rodríguez Cuéllar¹, Sebastian Breß^{1,2}
Asterios Katsifodimos³, Tilmann Rabl^{1,2}, Volker Markl^{1,2}

Example Query:

```
.window(SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10)))  
.sum()
```

Example Query:

```
.window(SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10)))
.sum()
```

Processing with Buckets:

Events:

↓
Eventtime

Buckets:

Example Query:

```
.window(SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10)))  
.sum()
```

Processing with Buckets:



Example Query:

```
.window(SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10)))  
.sum()
```

Processing with Buckets:



Example Query:

```
.window(SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10)))  
.sum()
```

Processing with Buckets:



Example Query:

```
.window(SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10)))  
.sum()
```

Processing with Buckets:



Example Query:

```
.window(SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10)))  
.sum()
```

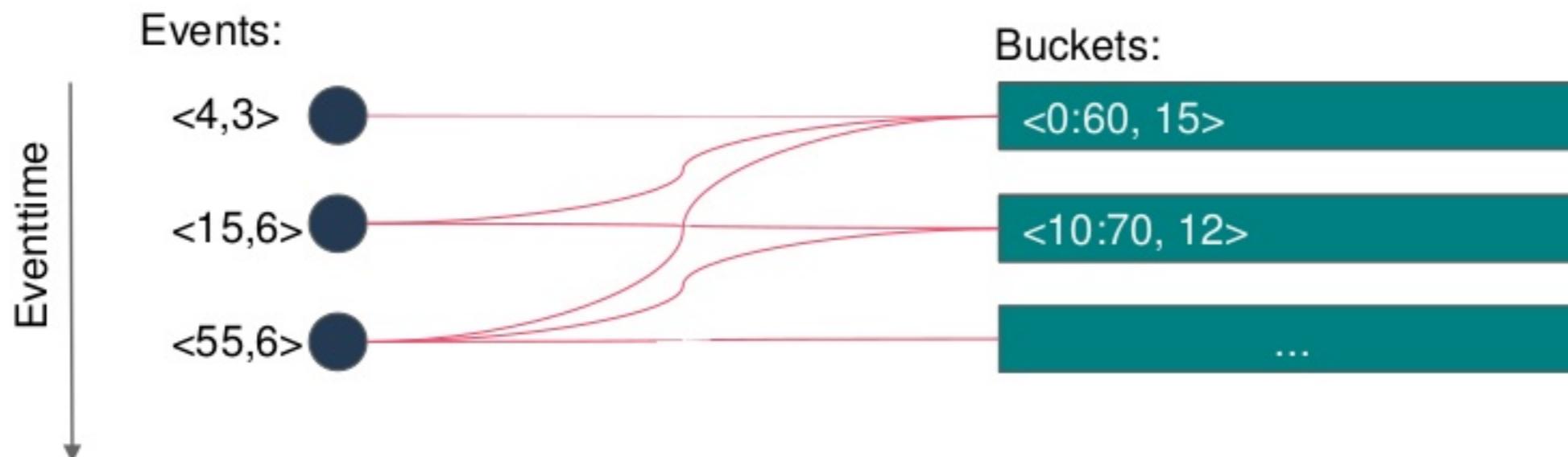
Processing with Buckets:



Example Query:

```
.window(SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10)))  
.sum()
```

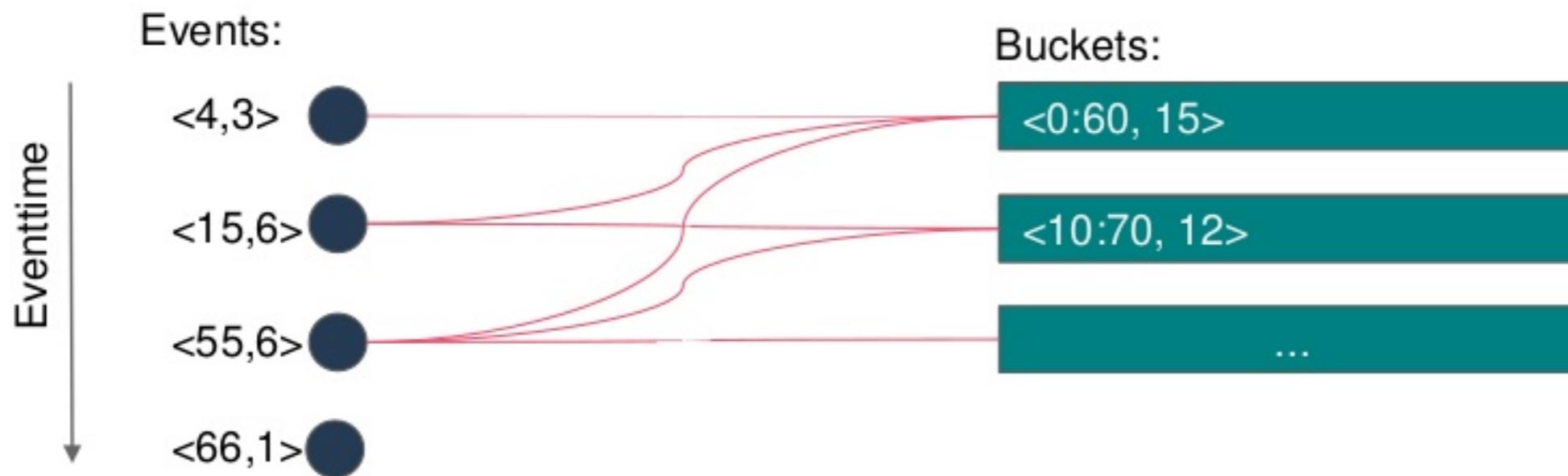
Processing with Buckets:



Example Query:

```
.window(SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10)))  
.sum()
```

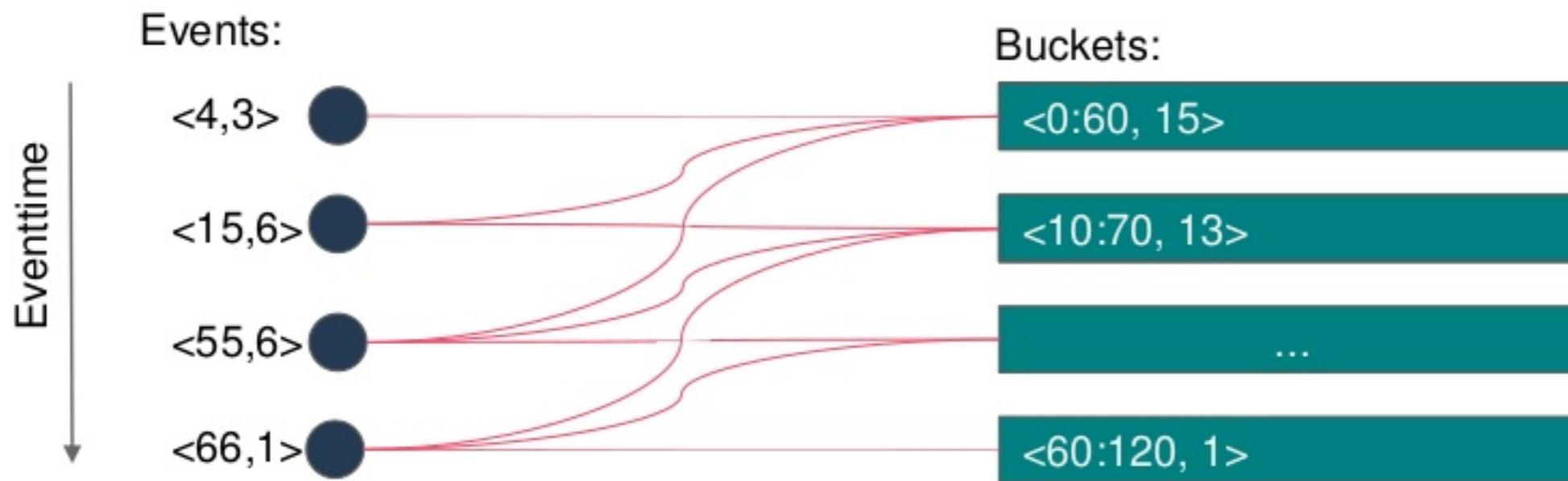
Processing with Buckets:



Example Query:

```
.window(SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10)))  
.sum()
```

Processing with Buckets:



Flink Windowing Bottlenecks

Number of Buckets = Window Length / Slide Length

Number of Buckets = Window Length / Slide Length

```
SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10)) --> 6 Buckets
```

Number of Buckets = Window Length / Slide Length

SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10))

--> 6 Buckets

SlidingEventTimeWindows.of(Time.day(1), Time.seconds(10))

--> 8640 Buckets

Number of Buckets = Window Length / Slide Length

SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10))

--> 6 Buckets

SlidingEventTimeWindows.of(Time.day(1), Time.seconds(10))

--> 8640 Buckets

Overlapping windows cause:

Number of Buckets = Window Length / Slide Length

SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10))

--> 6 Buckets

SlidingEventTimeWindows.of(Time.day(1), Time.seconds(10))

--> 8640 Buckets

Overlapping windows cause:

- Every event is assigned to many windows.

Number of Buckets = Window Length / Slide Length

SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10))

--> 6 Buckets

SlidingEventTimeWindows.of(Time.day(1), Time.seconds(10))

--> 8640 Buckets

Overlapping windows cause:

- Every event is assigned to many windows.
- Repeated aggregations --> aggregation function is called on every window

Number of Buckets = Window Length / Slide Length

SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10))

--> 6 Buckets

SlidingEventTimeWindows.of(Time.day(1), Time.seconds(10))

--> 8640 Buckets

Overlapping windows cause:

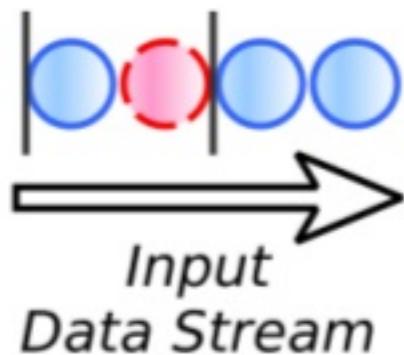
- Every event is assigned to many windows.
- Repeated aggregations --> aggregation function is called on every window
- High memory consumption --> especially for windows without incremental aggregation

Number of Buckets = Window Length / Slide Length

<i>SlidingEventTimeWindows.of(Time.minutes(1), Time.seconds(10))</i>	--> 6 Buckets
<i>SlidingEventTimeWindows.of(Time.day(1), Time.seconds(10))</i>	--> 8640 Buckets

Overlapping windows cause:

- Every event is assigned to many windows.
- Repeated aggregations --> aggregation function is called on every window
- High memory consumption --> especially for windows without incremental aggregation
- Check for merging windows

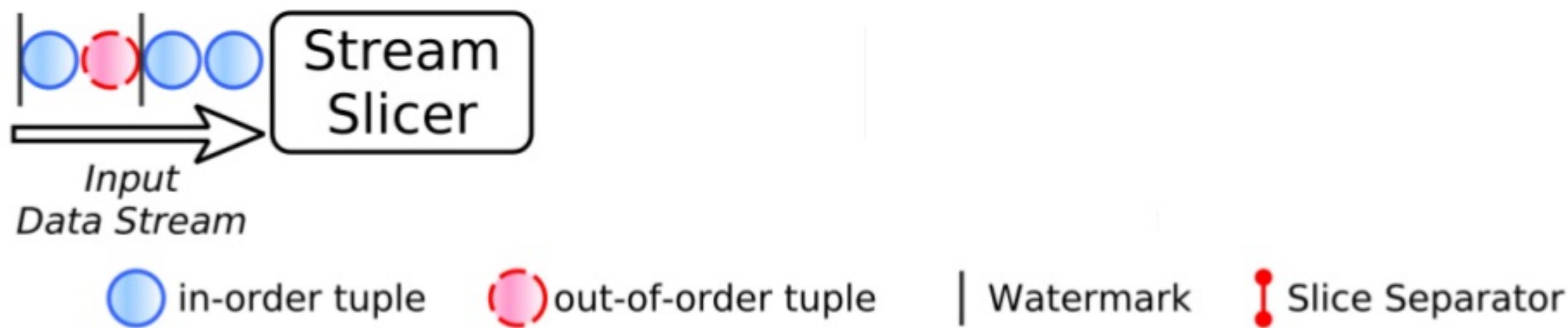


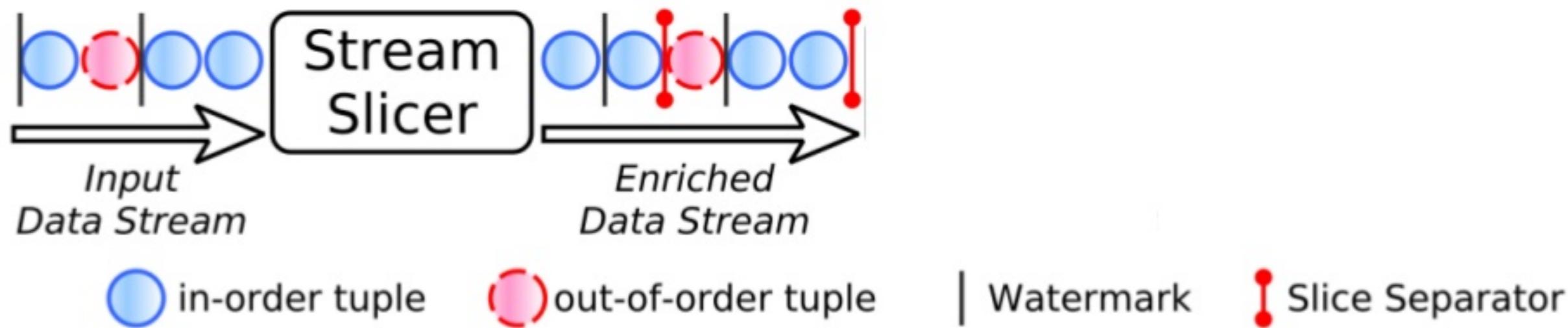
in-order tuple

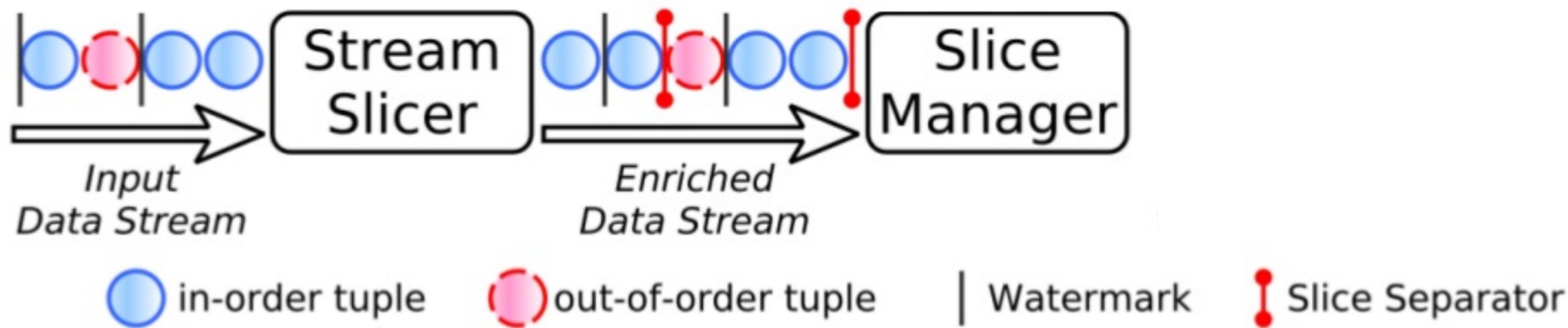
out-of-order tuple

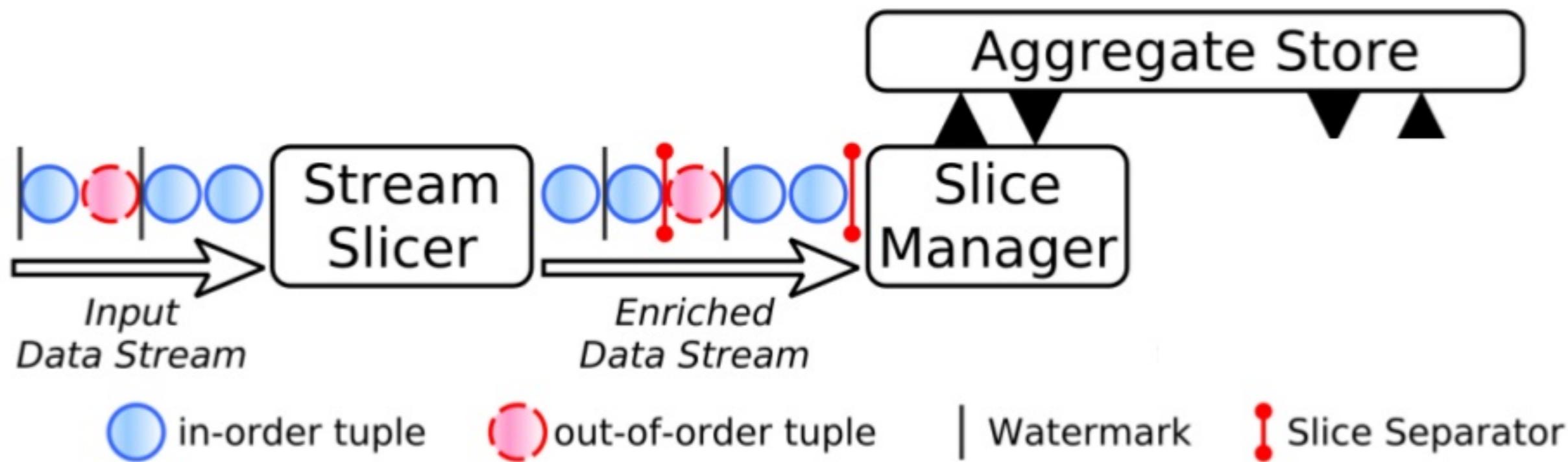
Watermark

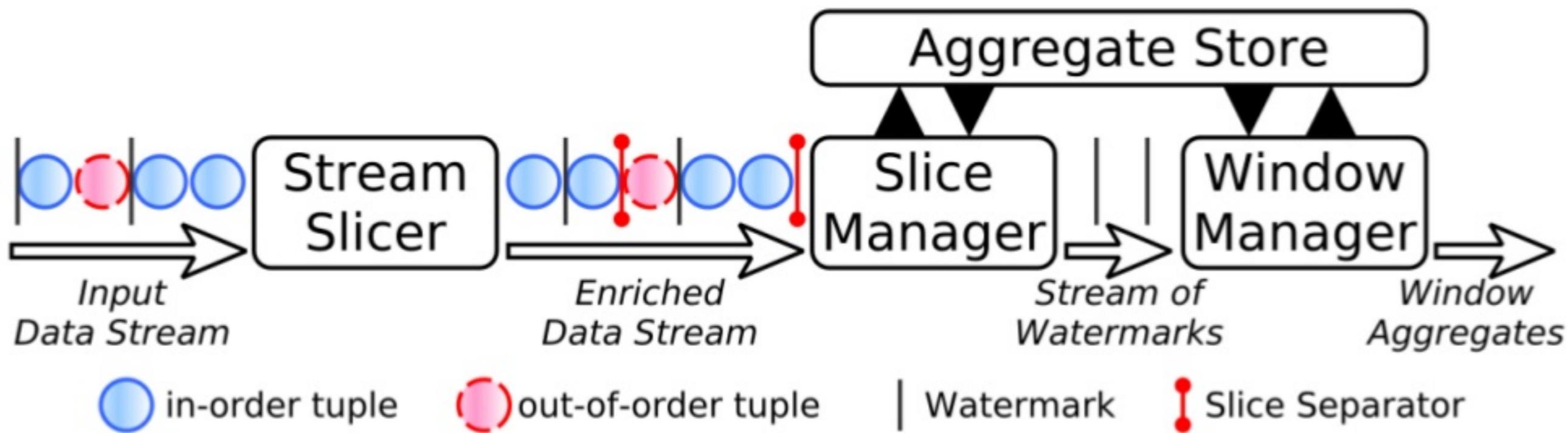
Slice Separator



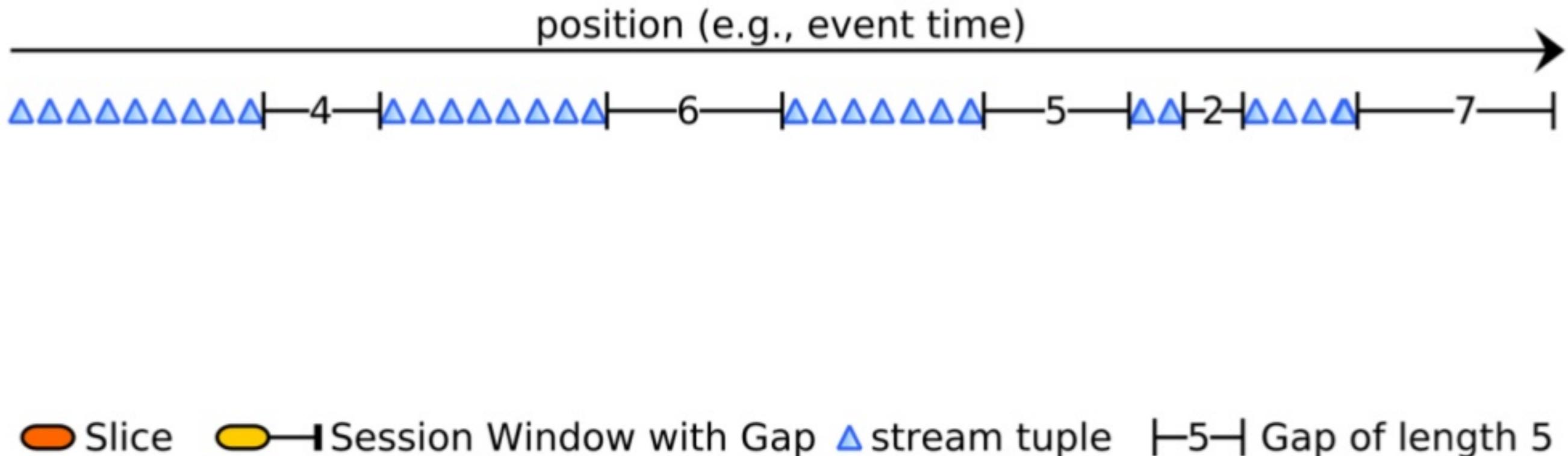




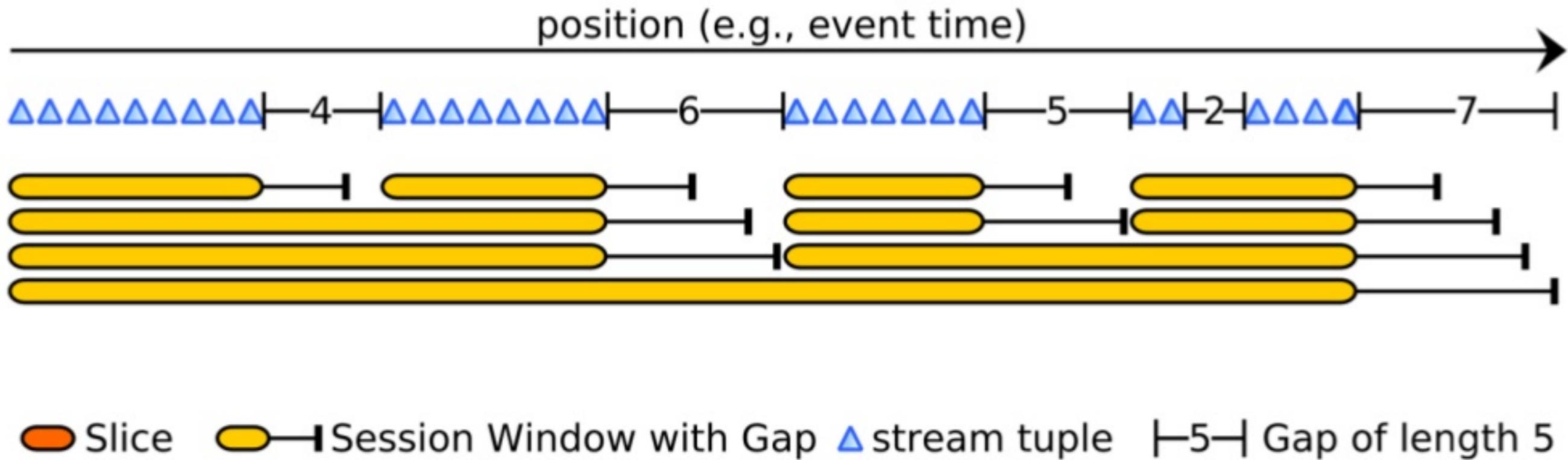




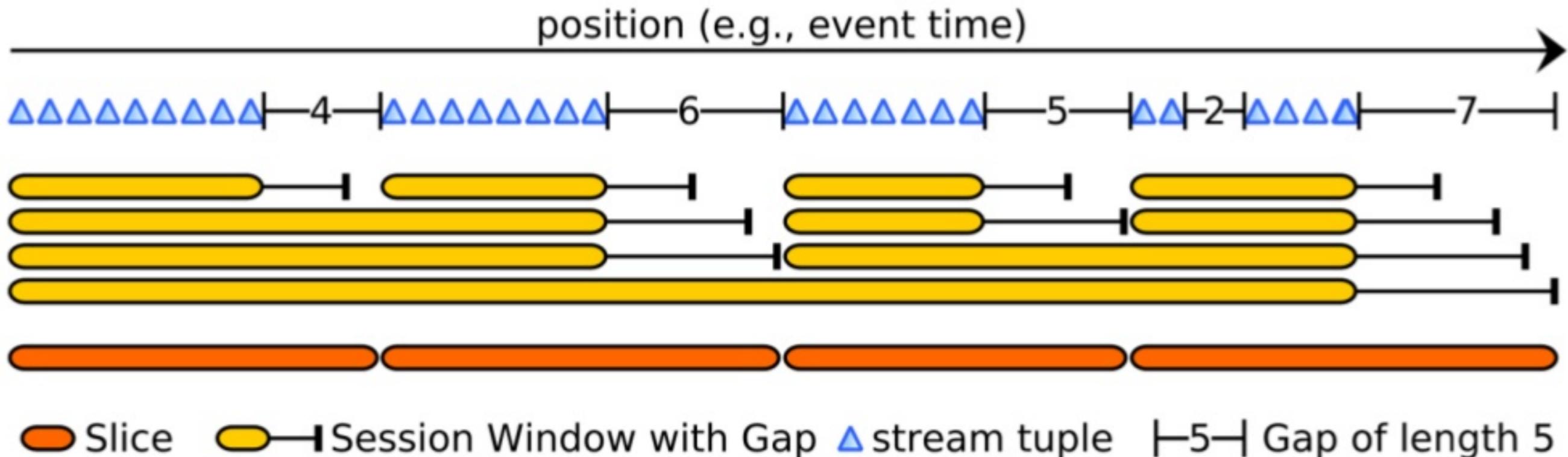
Session Window Aggregate Sharing

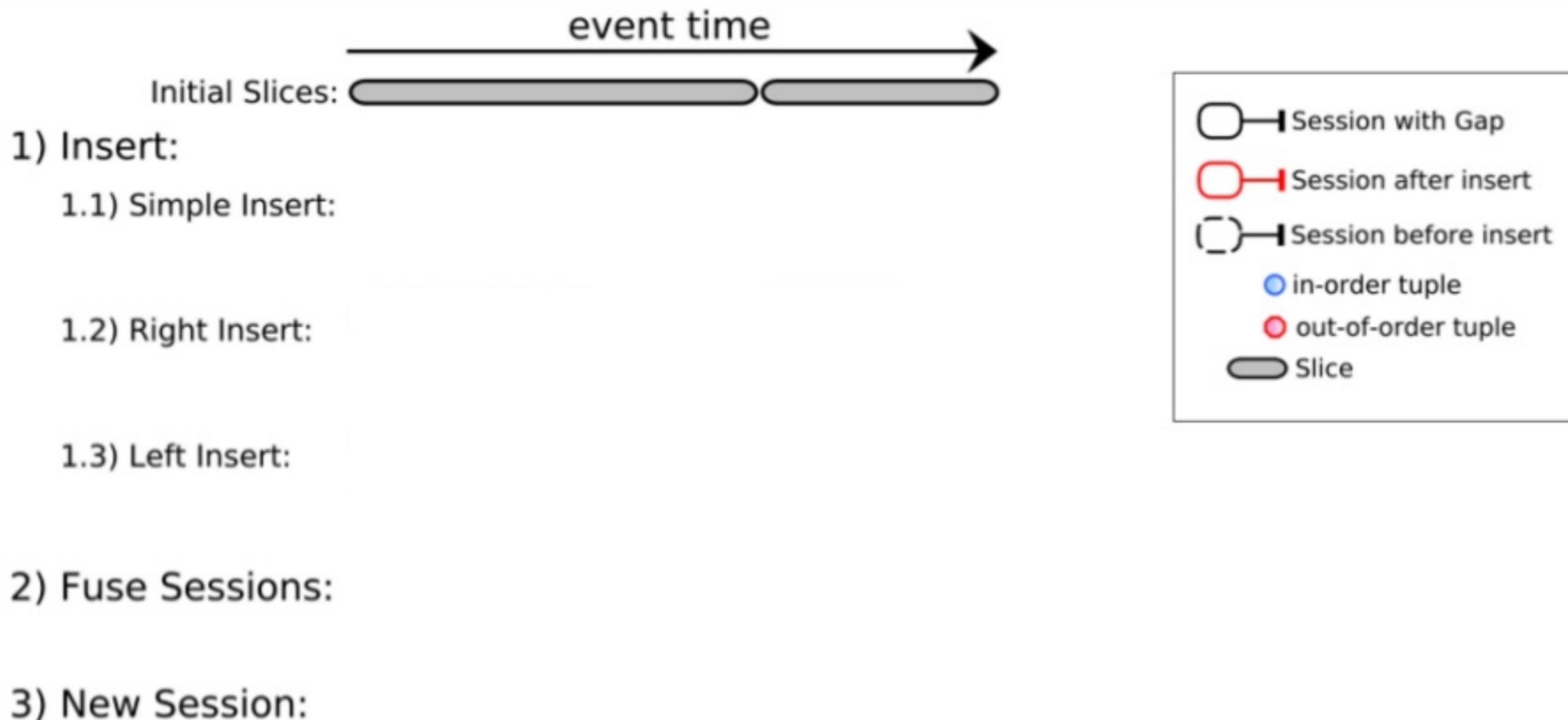


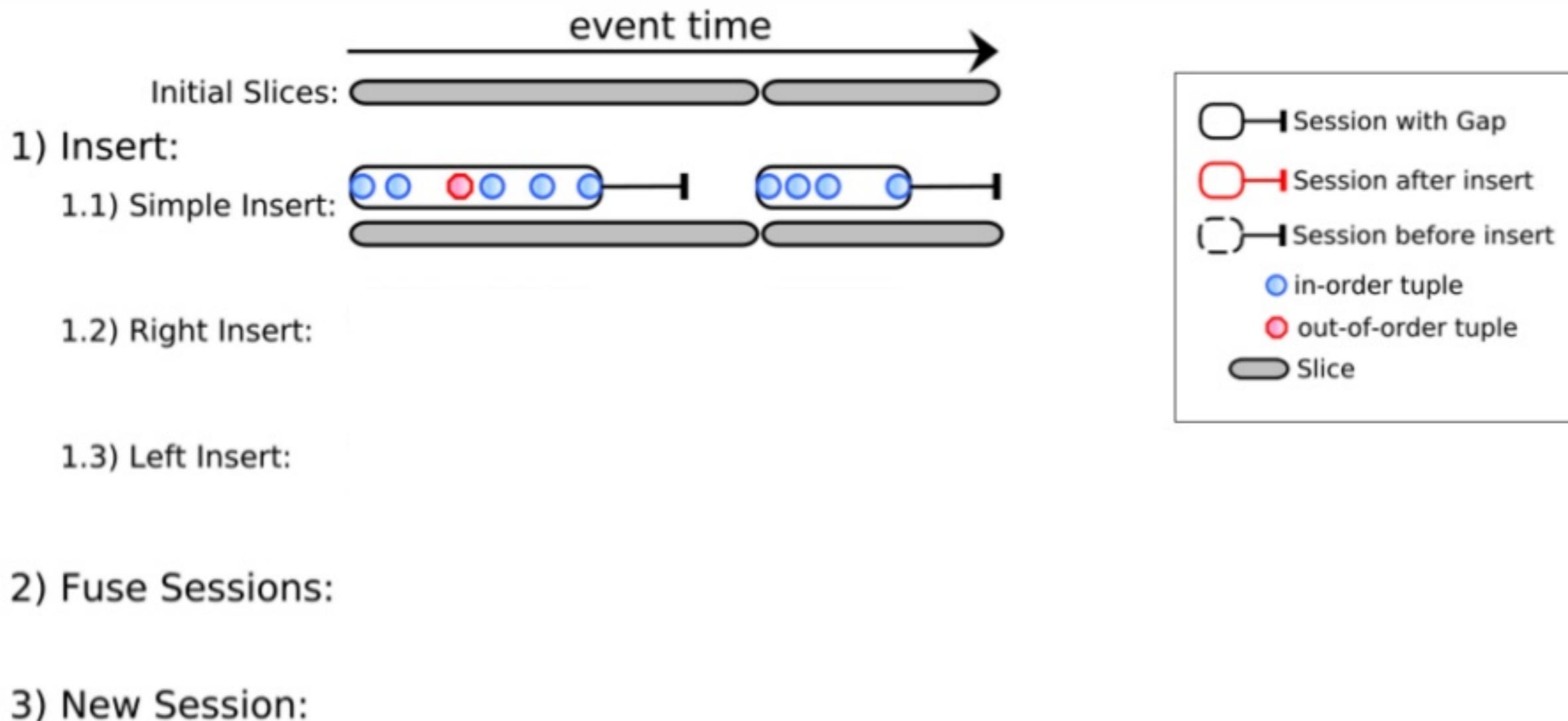
Session Window Aggregate Sharing

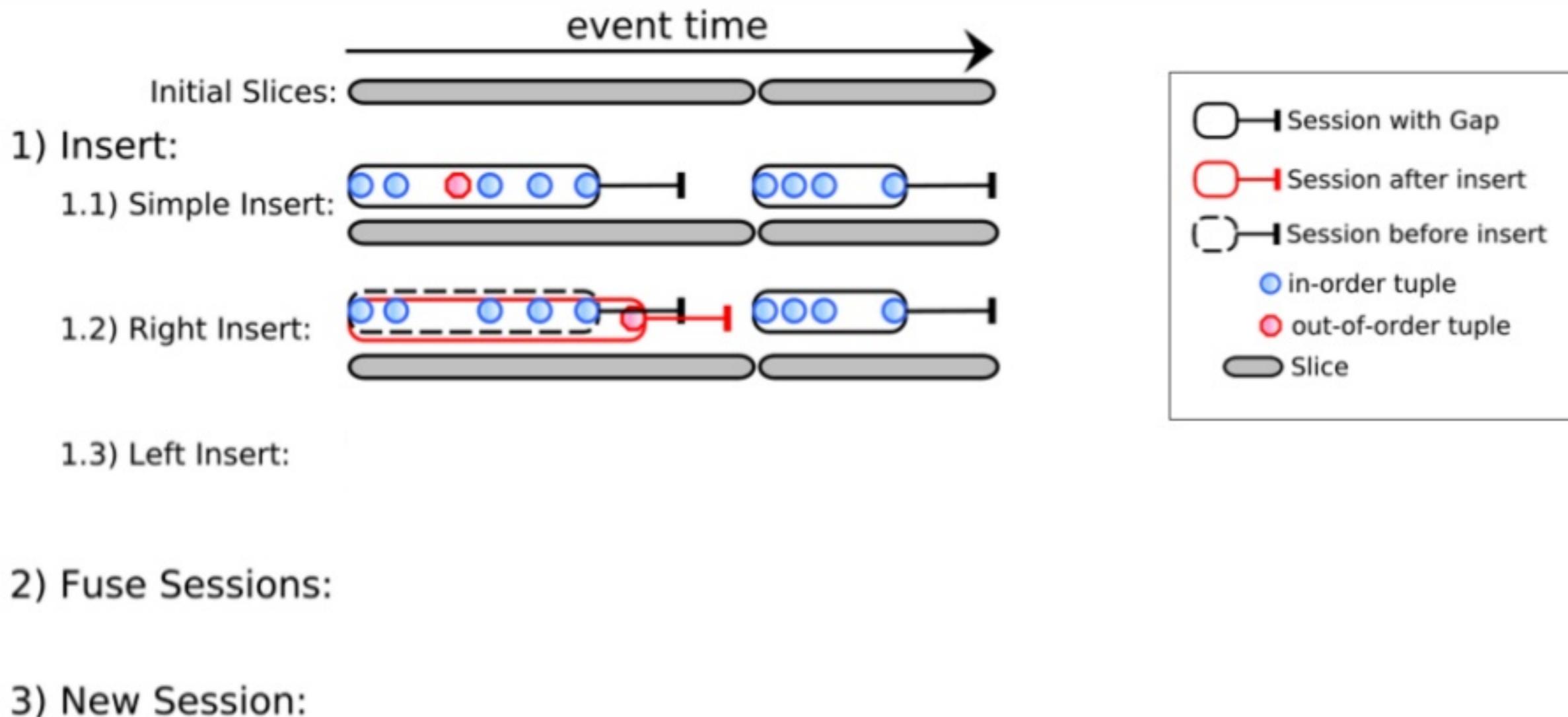


Session Window Aggregate Sharing

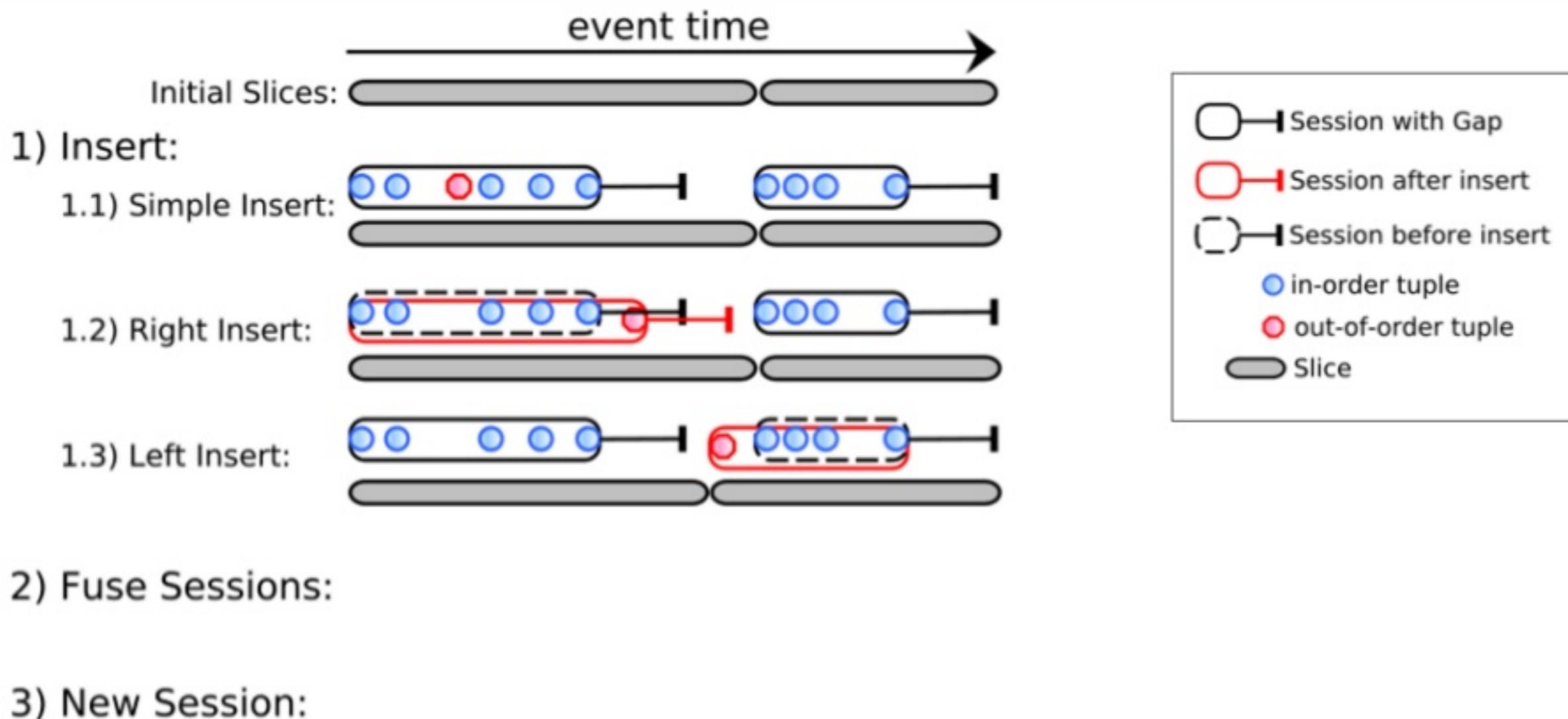




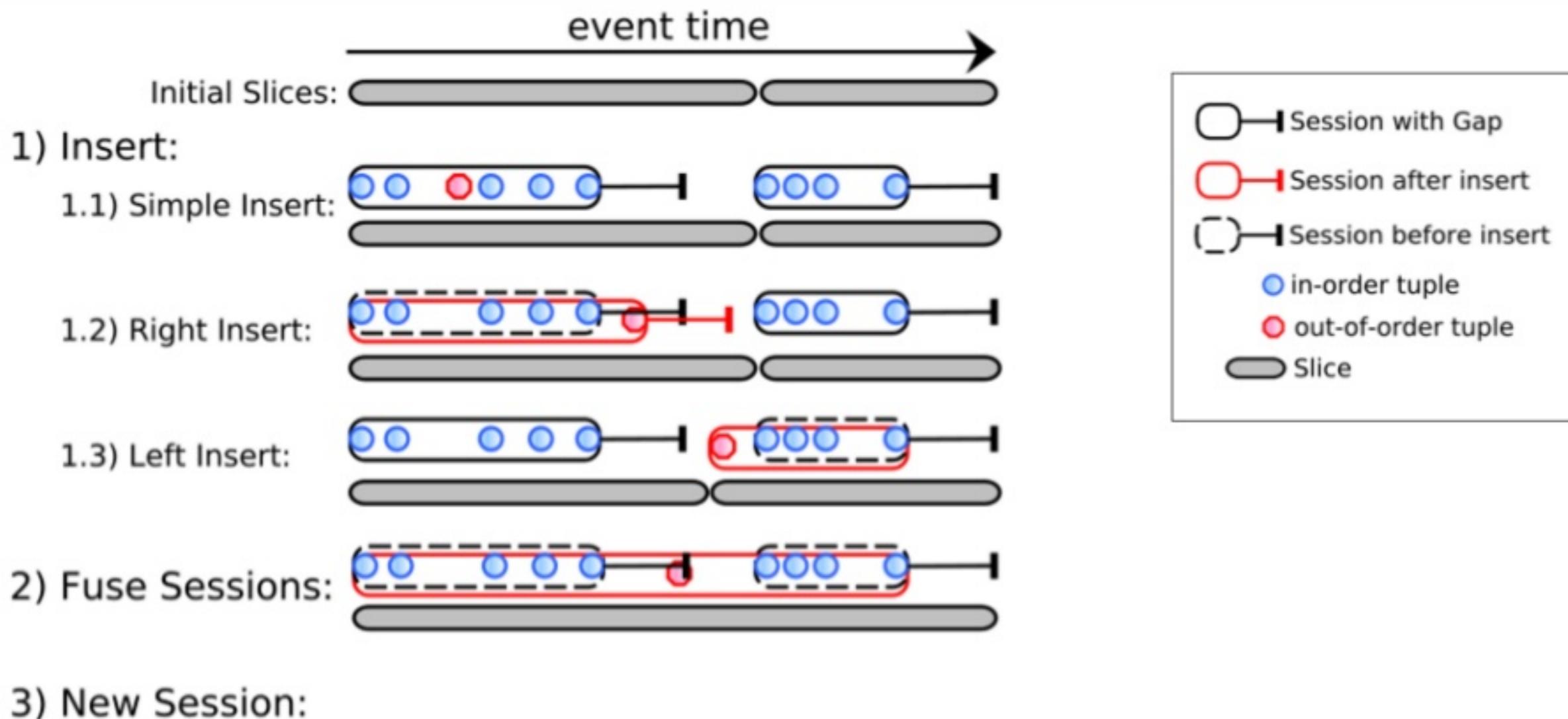




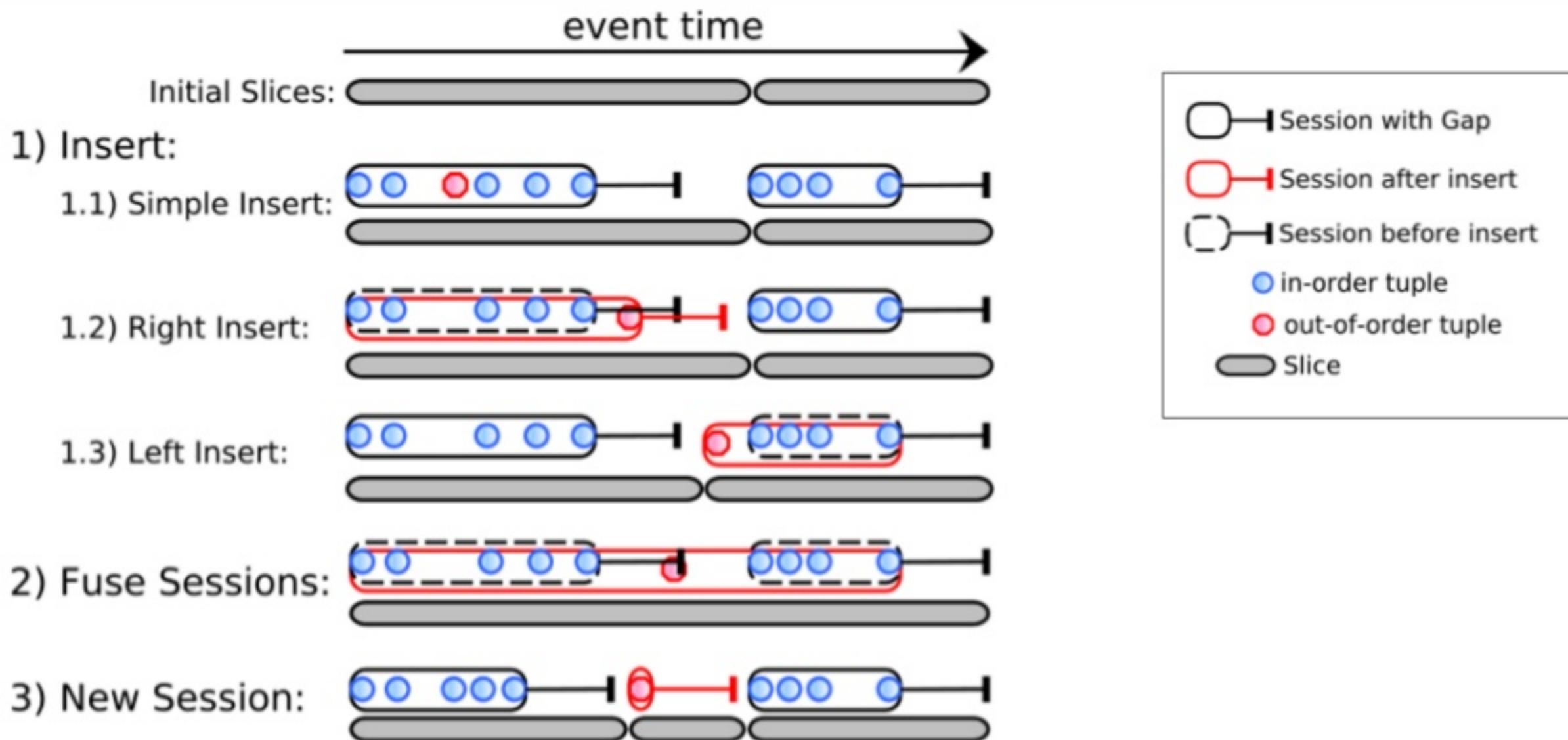
FLINK FORWARD Out-of-Order Processing and Sessions



FLINK FORWARD Out-of-Order Processing and Sessions



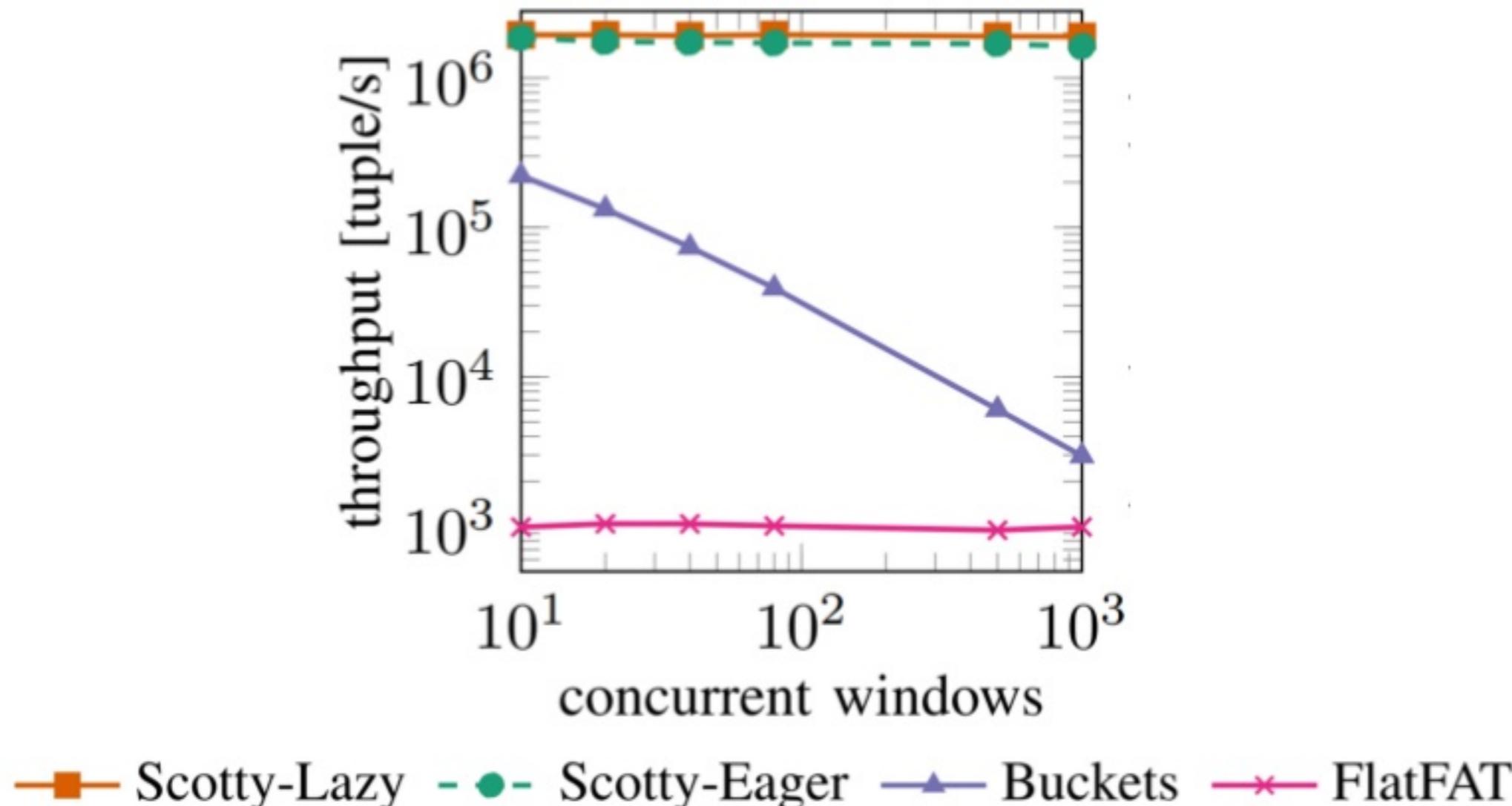
FLINK FORWARD Out-of-Order Processing and Sessions



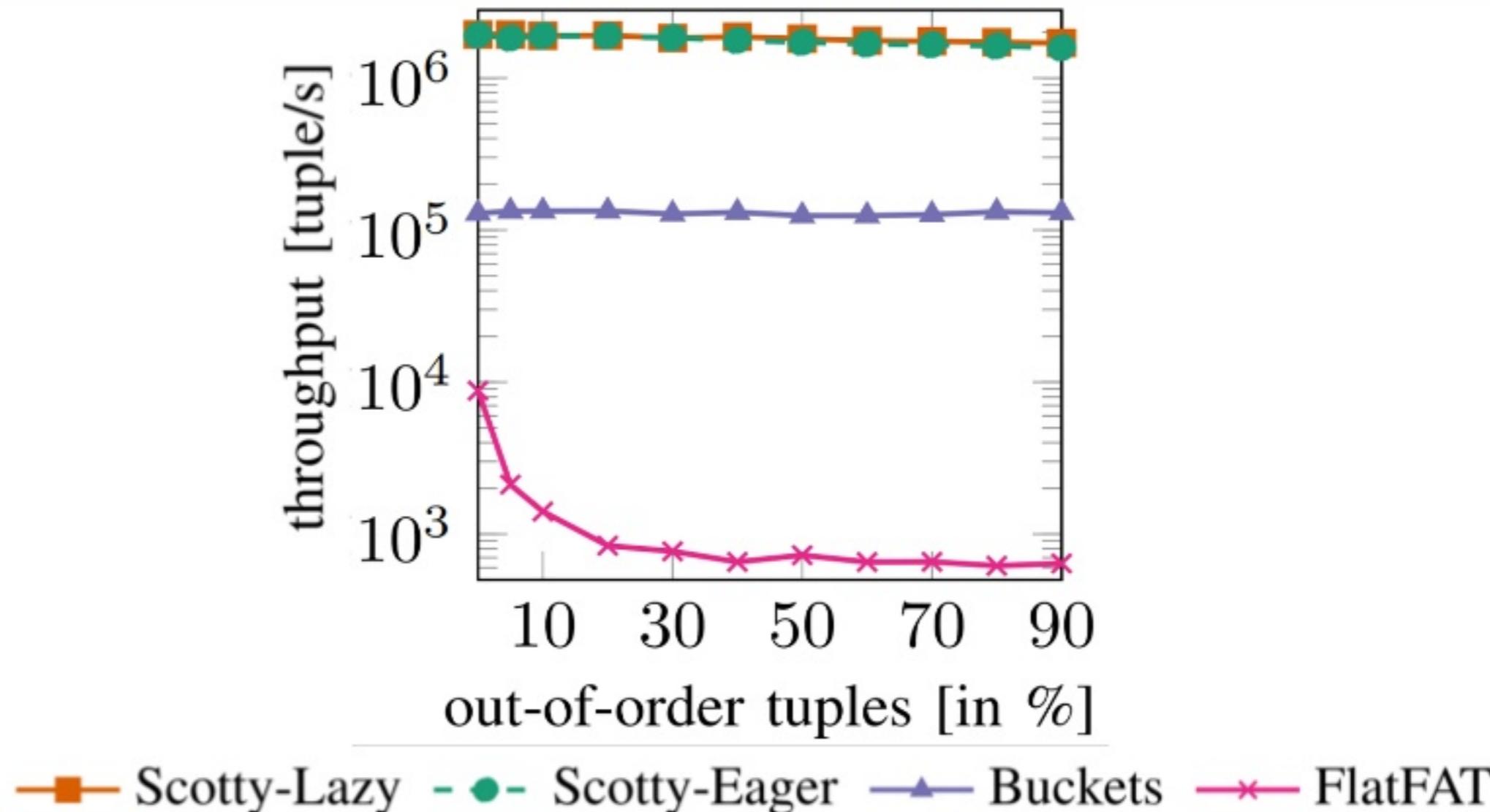
Multi-Window Processing (Example: Fitness Tracker)

```
[...].window(  
    // Daily report:  
    TumblingEventTimeWindows.of(Time.days(1)),  
    // Monitoring dashboard (last hour):  
    SlidingEventTimeWindows.of(Time.hours(1), Time.seconds(1)),  
    // Activity periods:  
    EventTimeSessionWindows.of(Time.minutes(1)))  
.sum()
```

Stream Slicing Performance



Stream Slicing Performance



Runtime-Dynamic Windows

Event Stream:



Window Definition Stream:



<WindowDefinition>

Event Stream:



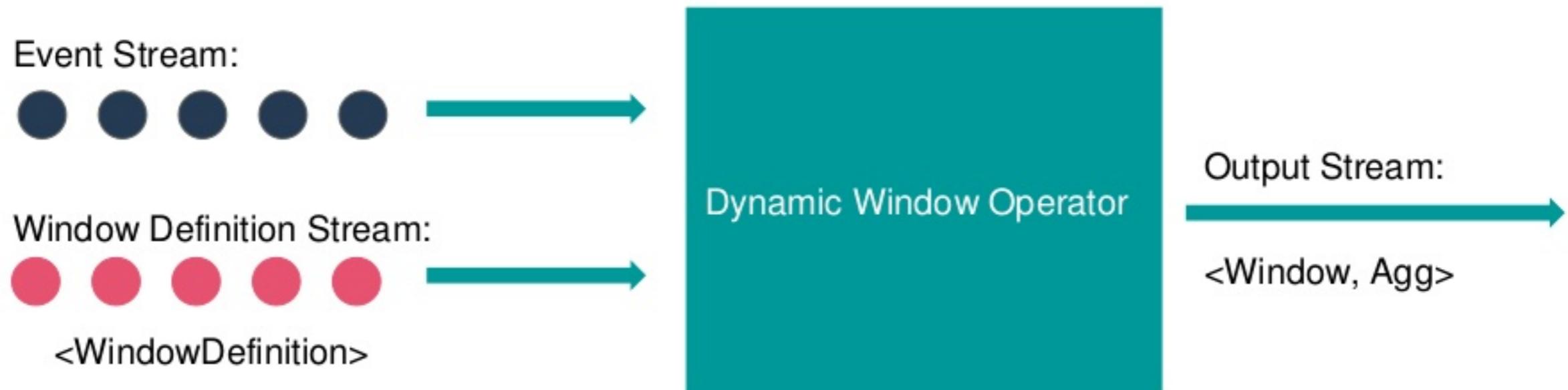
Window Definition Stream:



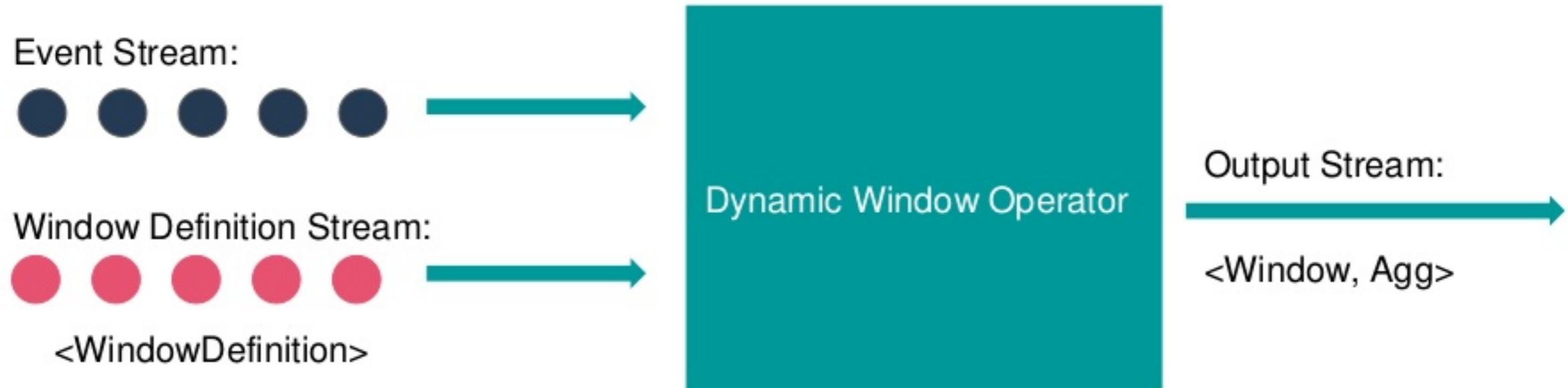
<WindowDefinition>



Runtime-Dynamic Windows



Runtime-Dynamic Windows



```
.dynamicWindow(windowDefinitionStream)  
.sum()
```

Research



Production

Research

- Implement complete fault-tolerance and state-management

Production

Research

- Implement complete fault-tolerance and state-management
- State migration

Production

Research

- Implement complete fault-tolerance and state-management
- State migration
 - Hard limitation: Aggregated buckets in state snapshots cannot be migrated

Production

Research

- Implement complete fault-tolerance and state-management
- State migration
 - Hard limitation: Aggregated buckets in state snapshots cannot be migrated
- Sophisticated testing

Production

Research

- Implement complete fault-tolerance and state-management
- State migration
 - Hard limitation: Aggregated buckets in state snapshots cannot be migrated
- Sophisticated testing

How to expose multi-windows and dynamic-windows to users?

Production

Scotty Features:

- stream slicing
- pre-aggregation
- aggregate sharing
- out-of-order processing
- session window support
- multi-window support
- runtime-dynamic window support

Let's bring it to production!

JIRA: [FLINK-7001]



This talk is supported by the European Union Horizon 2020 Projects Proteus (687691), Streamline (688191), SAGE (671500), and E2Data (780245) and by the German Ministry for Education and Research as Berlin Big Data Center (01IS14013A) and Software Campus (01IS12056).