

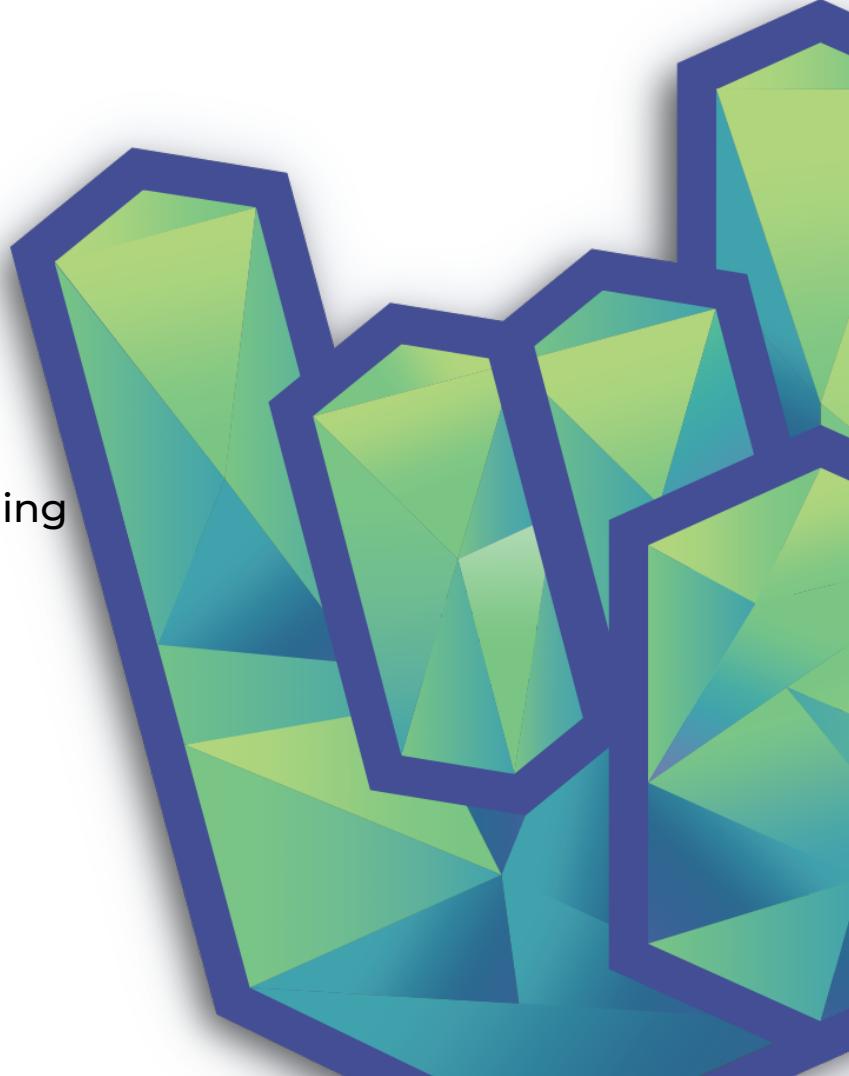
Event Time Windows



Objective

Handle records by event time

Learn window functions on Structured Streaming



Event Time

The moment when the record was *generated*

Set by the data generation system

Usually a column in the dataset

Different from *processing time* = the time the record arrives at Spark

Window Functions

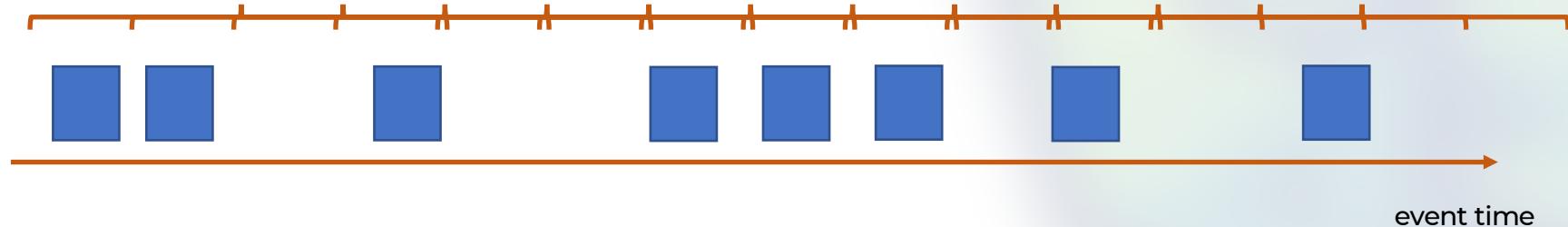
Aggregations on time-based groups

Essential concepts:

- window durations
- window sliding intervals

As opposed to DStreams

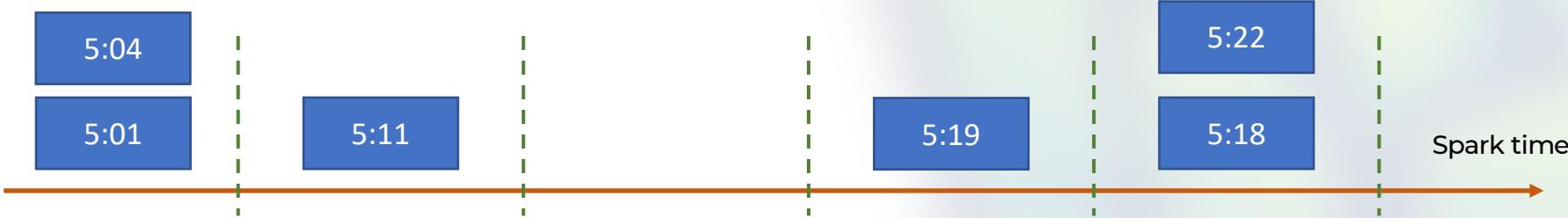
- records are not necessarily taken between "now" and a certain past date
- we can control output modes



Window Functions

Assume we have a count by window, in complete mode

- batch time = 10 minutes
- window duration = 20 minutes
- window sliding interval = 10 minutes



window	count
4:50 – 5:10	2
5:00 – 5:20	2

window	count
4:50 – 5:10	2
5:00 – 5:20	3

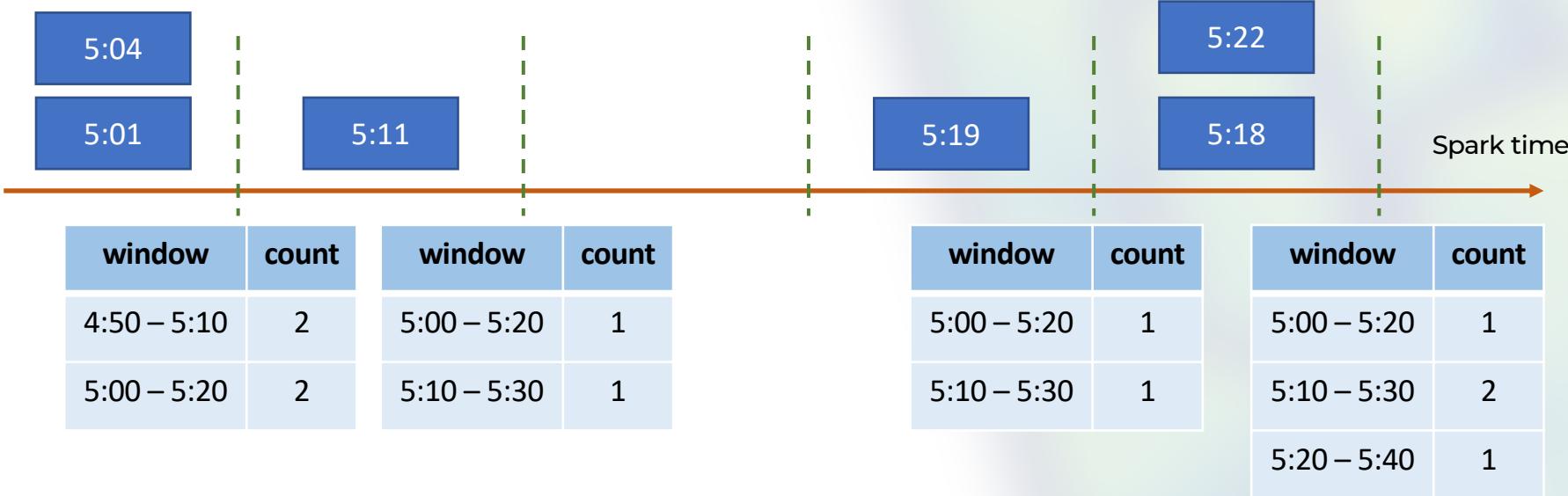
window	count
4:50 – 5:10	2
5:00 – 5:20	4
5:10 – 5:30	2

window	count
4:50 – 5:10	2
5:00 – 5:20	5
5:10 – 5:30	4
5:20 – 5:40	1

Window Functions

Assume we have a count by window, in append mode

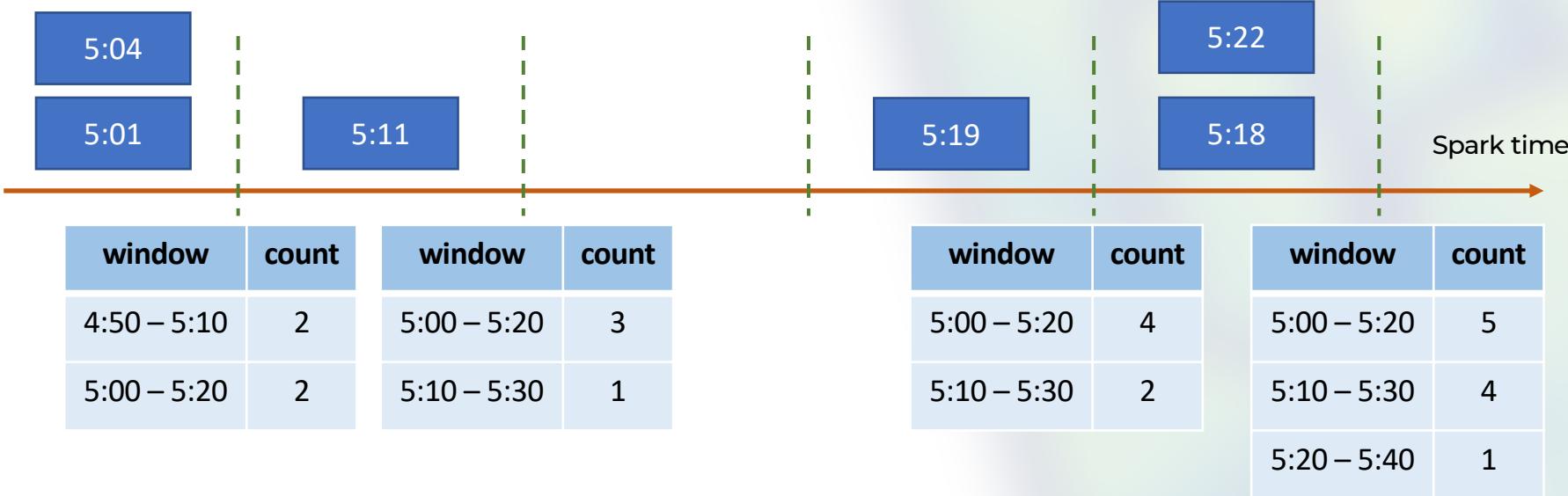
- batch time = 10 minutes
- window duration = 20 minutes
- window sliding interval = 10 minutes



Window Functions

Assume we have a count by window, in update mode

- batch time = 10 minutes
- window duration = 20 minutes
- window sliding interval = 10 minutes



Takeaways

Can group streamed data into time-based groups



Window duration and sliding interval must be a multiple of the batch interval

Output mode will influence results

Spark rocks

