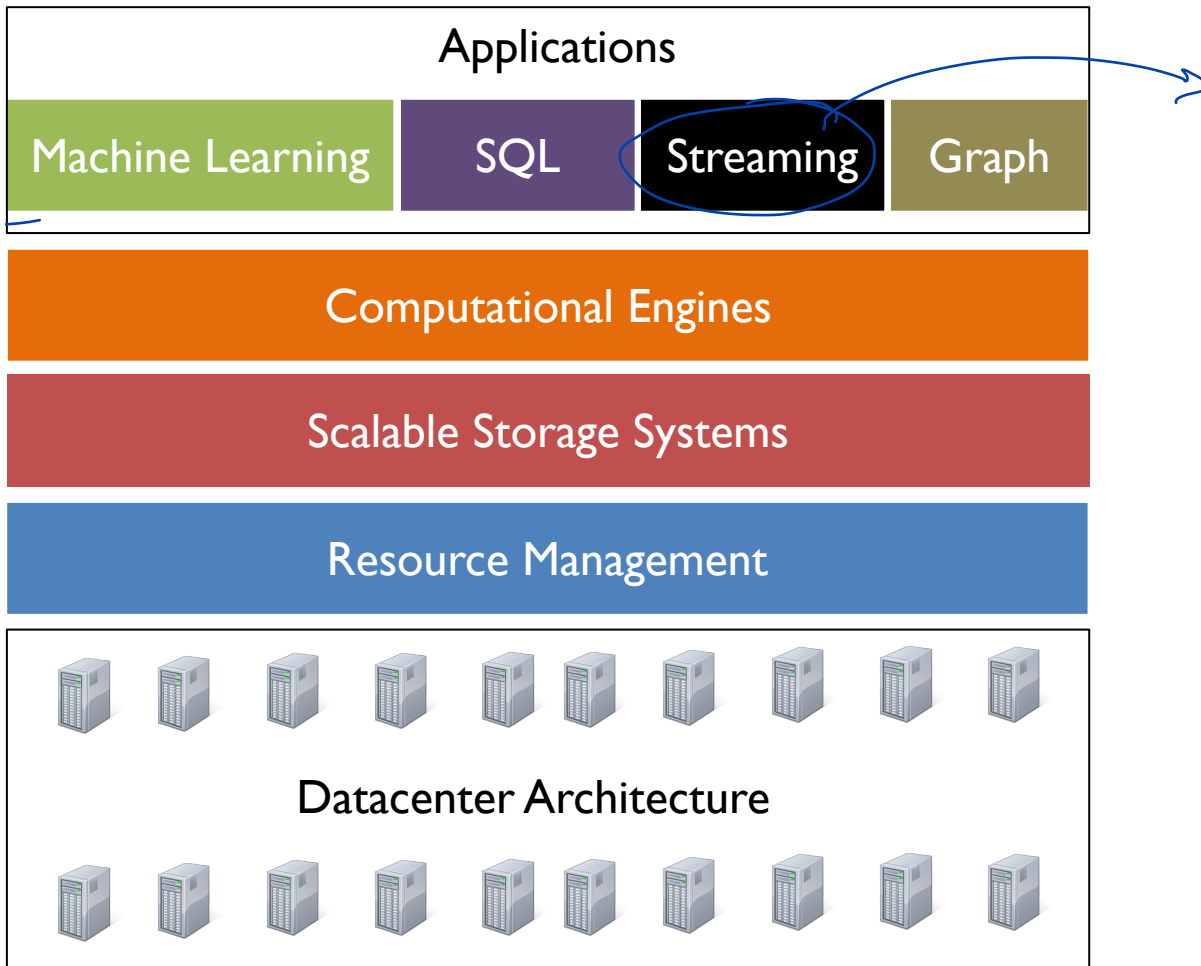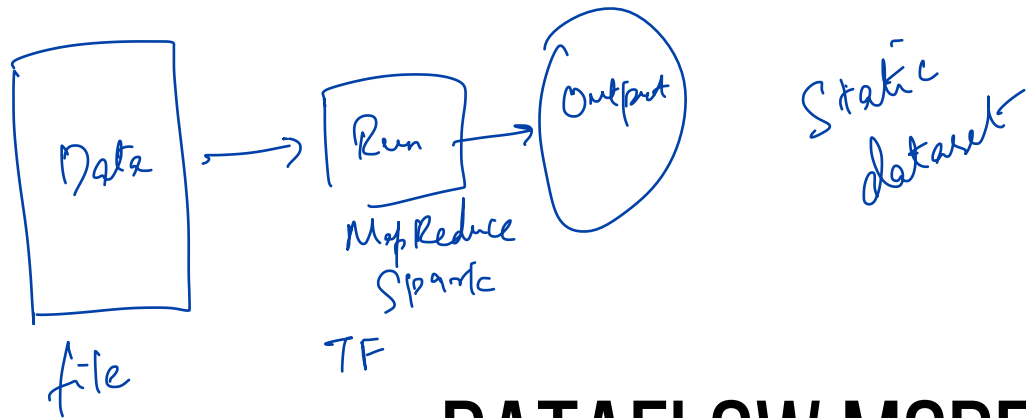# CS 744: DATAFLOW

Shivaram Venkataraman

Fall 2019
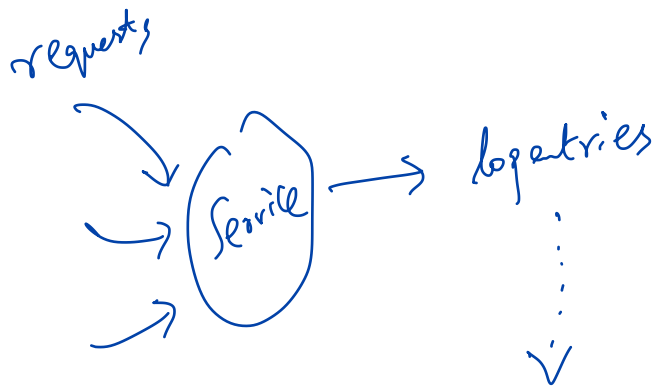
# ADMINISTRIVIA

- Assignment 2 grades up

- Midterm grading

- Course project proposal comments


- AEFIS feedback → Slides

  → Discussion


- No Class next Tuesday?

# DATAFLOW MODEL (?)

Data → Run → Output

file

MapReduce
Spark
TF

Static dataset

requests → Service → logentries
⋮
↓

stream dataset
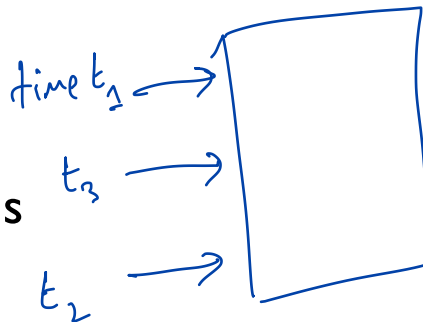
user 1 : < videos
              user has watched

Correctness vs. latency

# MOTIVATION

Streaming Video Provider

    - How much to bill each advertiser ?

    - Need per-user, per-video viewing sessions

    - Handle out of order data

Goals

    - Easy to program

    - Balance correctness, latency and cost

# APPROACH

API Design
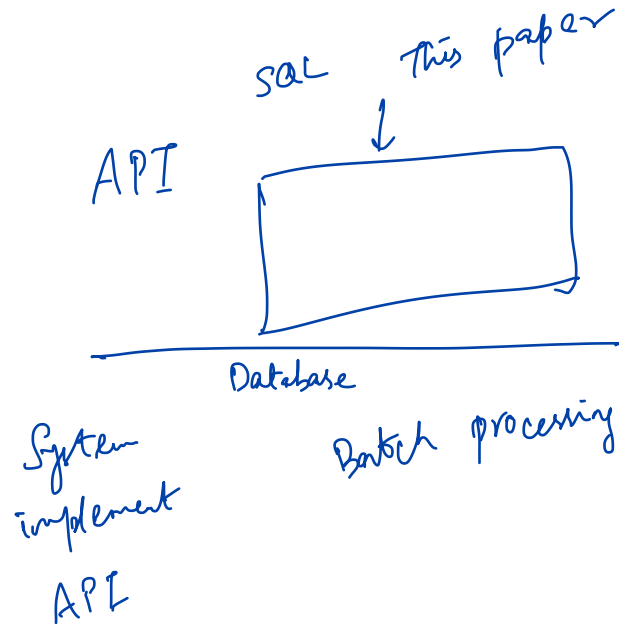
    Separate user-facing model from execution

    Decompose queries into

        - What is being computed

        - Where in time is it computed

        - When is it materialized

        - How does it relate to earlier results

*Data*

*Scheduling*

*Output*

*whether you are incrementing or accumulating*

*SQL*  *This paper*

*API*

*Database*

*System implement*

*API*

*Batch processing*

# TERMINOLOGY

All of them have
proc. time
~ 199

**Unbounded/bounded data**

$100$ →
$101$

Day 1 → MR → Output

- - 199

200 →

Day 2 → MR → Output

201 →

    **Streaming/Batch execution**

**Timestamps**

Streaming
System

internally
updates →

Output

**Event time:** Time when event
occured wrt usual input

**Processing time:**
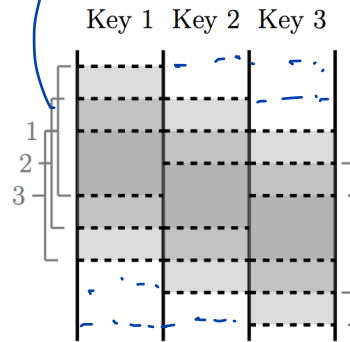Time at which event
is processed

# WINDOWING



Fixed

Sliding

Sessions

duration

no overlap

align
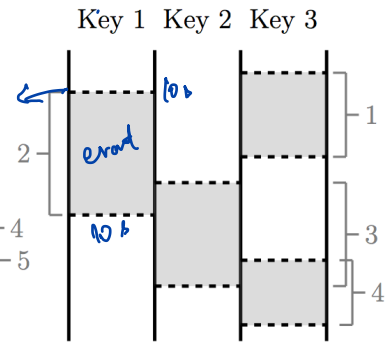
Can overlap
slide
duration

specific for key
unaligned

window size
↓ Developer

# WATERMARK OR SKEW



System has processed all events up to 12:02:30

Actual watermark:
Ideal watermark:
Event Time Skew:

Processing time lags event time

→ Time out
→ Count

event = proc. time

# API

ParDo:  map / flat Map          Input  k, v  $\rightarrow$  $(k_1, v_1)$
                                                             $(k_2, v_2)$
                                                                 $\vdots$

GroupByKey:  reduce
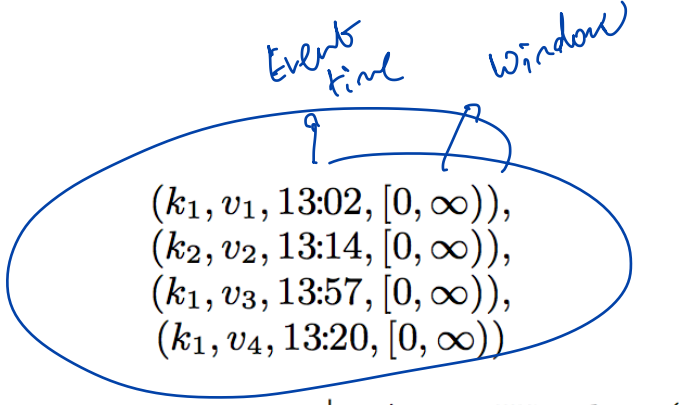
$k \rightarrow (v_1, v_2, v_3 \cdots)$

Windowing

    AssignWindow  $\rightarrow$  Bucket tuples into window

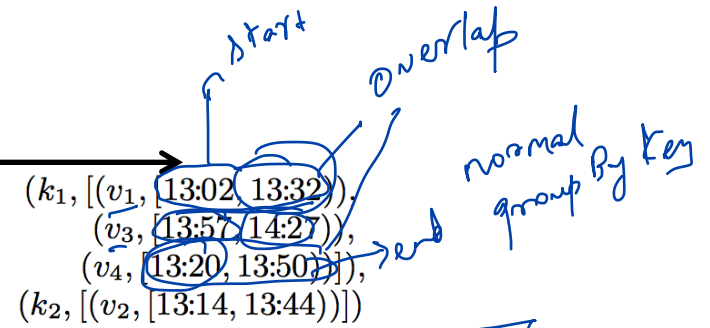    MergeWindow  $\rightarrow$  Merge buckets based on strategy

# EXAMPLE

*Events time*    *Window*

$$(k_1, v_1, 13{:}02, [0, \infty)),$$
$$(k_2, v_2, 13{:}14, [0, \infty)),$$
$$(k_1, v_3, 13{:}57, [0, \infty)),$$
$$(k_1, v_4, 13{:}20, [0, \infty))$$

↓ $AssignWindows($
$Sessions(30m))$

$$(k_1, v_1, 13{:}02, [13{:}02, 13{:}32)),$$
$$(k_2, v_2, 13{:}14, [13{:}14, 13{:}44)),$$
$$(k_1, v_3, 13{:}57, [13{:}57, 14{:}27)),$$
$$(k_1, v_4, 13{:}20, [13{:}20, 13{:}50))$$

↓ $DropTimestamps$

$$(k_1, v_1, [13{:}02, 13{:}32)),$$
$$(k_2, v_2, [13{:}14, 13{:}44)),$$
$$(k_1, v_3, [13{:}57, 14{:}27)),$$
$$(k_1, v_4, [13{:}20, 13{:}50))$$

GroupByKey

*start*   *Overlap*   *normal group By Key*

$$(k_1, [(v_1, [13{:}02, 13{:}32)),$$
$$(v_3, [13{:}57, 14{:}27)),$$
$$(v_4, [13{:}20, 13{:}50))]),$$
$$(k_2, [(v_2, [13{:}14, 13{:}44))])$$

*end*

↓ $MergeWindows($
$Sessions(30m))$

$$(k_1, [(v_1, [\mathbf{13{:}02}, \mathbf{13{:}50})),$$
$$(v_3, [13{:}57, 14{:}27)),$$
$$(v_4, [\mathbf{13{:}02}, \mathbf{13{:}50}))]),$$
$$(k_2, [(v_2, [13{:}14, 13{:}44))])$$

↓ $GroupAlsoByWindow$

$$(k_1, [((\mathbf{v_1}, \mathbf{v_4}], [13{:}02, 13{:}50)),$$
$$([\mathbf{v_3}], [13{:}57, 14{:}27))]),$$
$$(k_2, [([\mathbf{v_2}], [13{:}14, 13{:}44))])$$

↓ $ExpandToElements$   *timestamp*

$$(k_1, [v_1, v_4], \mathbf{13{:}50}, [13{:}02, 13{:}50)),$$
$$(k_1, [v_3], \mathbf{14{:}27}, [13{:}57, 14{:}27)),$$
$$(k_2, [v_2], \mathbf{13{:}44}, [13{:}14, 13{:}44))$$

*query*

# TRIGGERS AND INCREMENTAL PROCESSING

Windowing: where in event time data are grouped

Triggering: when in processing time groups are emitted
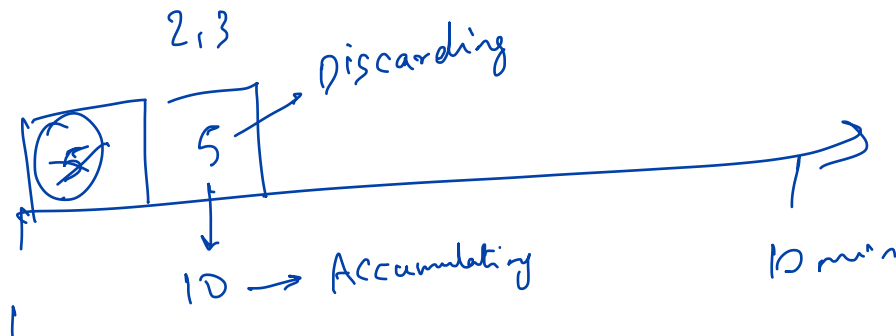
Timestamp
10 min window
every 1 min

Strategies

Discarding

Accumulating

Accumulating & Retracting

State: 2,3
5

2,3

5 → Discarding

10 → Accumulating
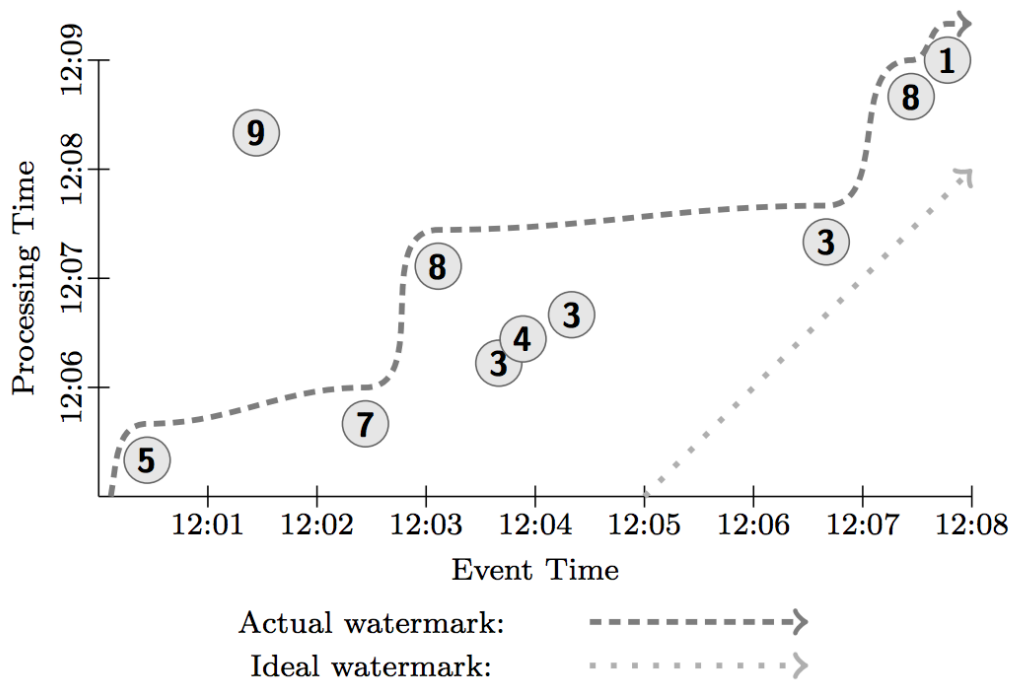
10 min

old state, updates
↓ new state

−5, 10

# RUNNING EXAMPLE

```
PCollection<KV<String, Integer>> input = IO.read(...);
PCollection<KV<String, Integer>> output =
        input.apply(Sum.integersPerKey());
```
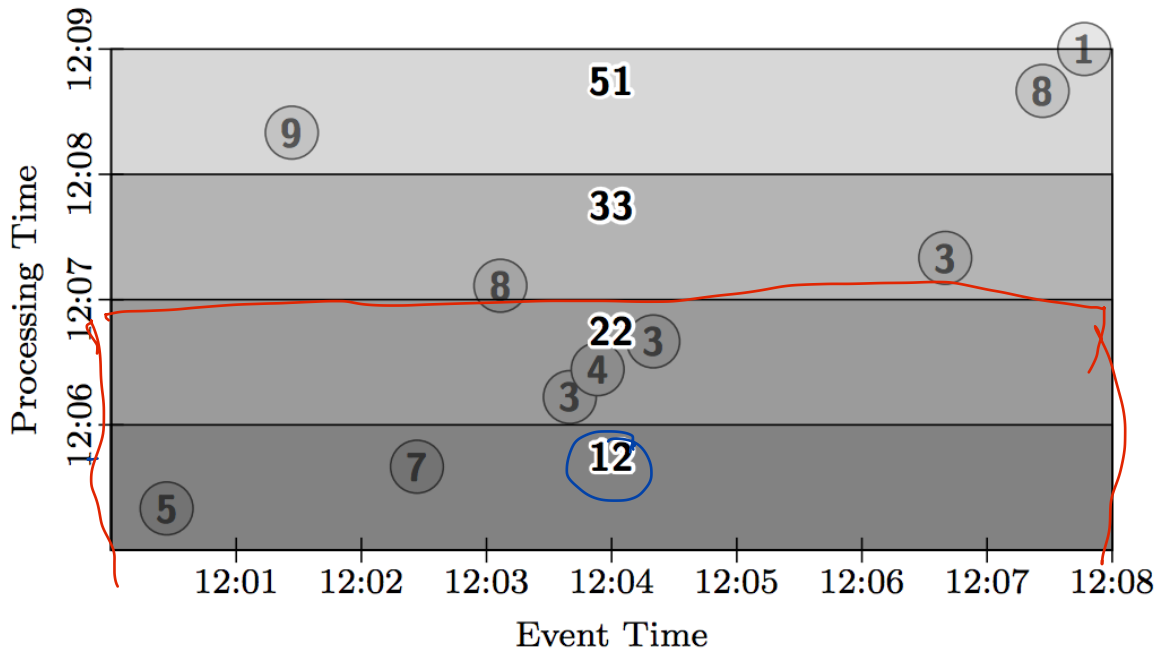
# GLOBAL WINDOWS, ACCUMULATE

```
PCollection<KV<String, Integer>> output = input
    .apply(Window.trigger(Repeat(AtPeriod(1, MINUTE)))
           .accumulating())
    .apply(Sum.integersPerKey());
```
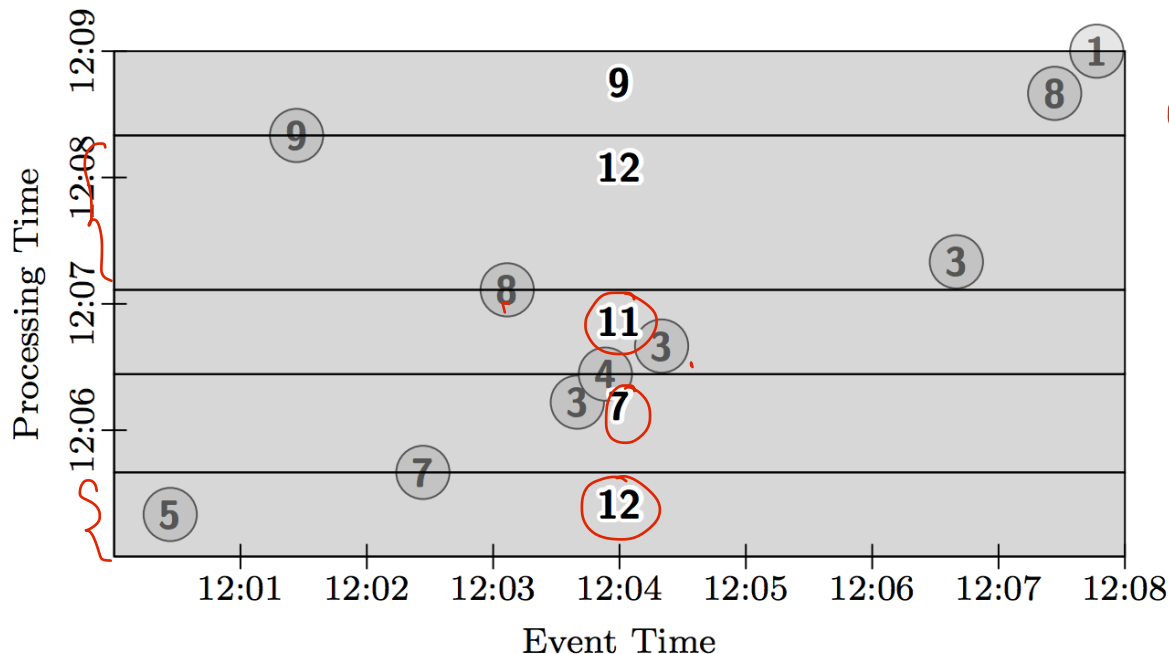
# GLOBAL WINDOWS, COUNT, DISCARDING

```
PCollection<KV<String, Integer>> output = input
    .apply(Window.trigger(Repeat(AtCount(2)))
           .discarding())
    .apply(Sum.integersPerKey());
```

# FIXED WINDOWS, MICRO BATCH

Output

```
PCollection<KV<String, Integer>> output = input
    .apply(Window.into(FixedWindows.of(2, MINUTES))
        .trigger(Repeat(AtWatermark())))
        .accumulating())
```

12:00 - 12:02    5

12:02 - 12:04    14 22

12:04 - 12:06    3

# LESSONS / EXPERIENCES

Don't rely on completeness

Be flexible, diverse use cases
- Billing
- Recommendation
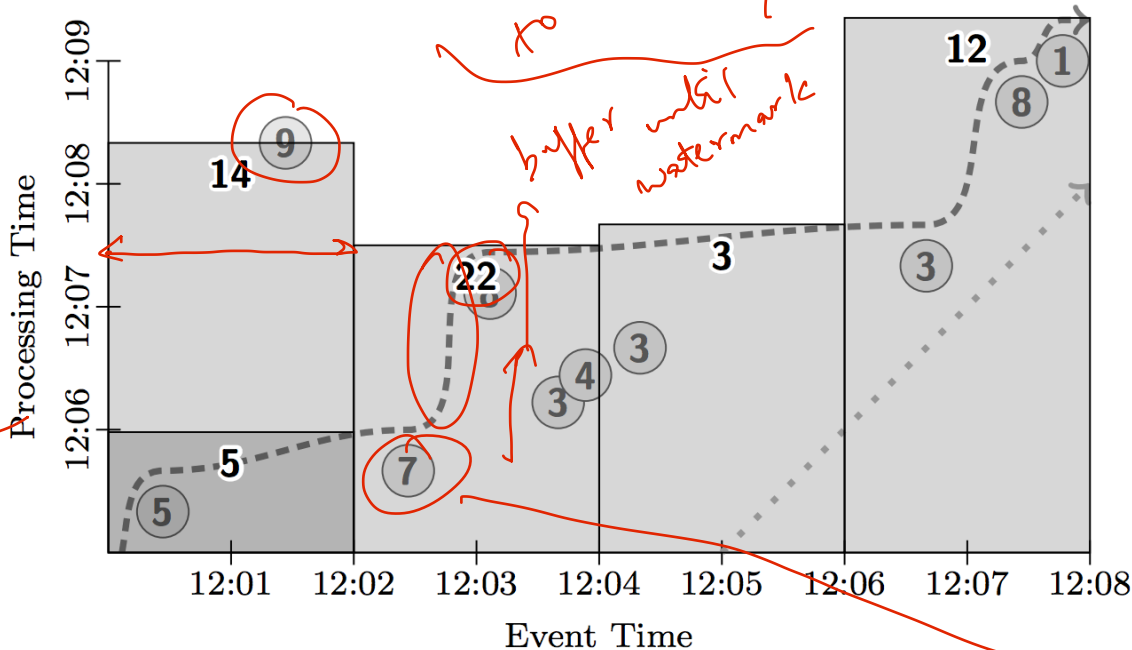- Anomaly detection

Support analysis in context of events

# DISCUSSION

https://forms.gle/s7T2r67BDvkGQhmN9

Computation overheads?

Number of windows to update

update lover for streaming

no

buffer until watermark

Batch
22
→ 12:08

Number of outputs is more

lower latery watermark
22 → 12:07:30

No partial results until watermark

is delayed till 12:07:30

14    9

22

3    4    3    3

5

5

7

12    1
8

3

Processing Time

12:09
12:08
12:07
12:06

12:01    12:02    12:03    12:04    12:05    12:06    12:07    12:08

Event Time

Consider you are implementing a micro-batch streaming API on top of Apache Spark. What are some of the bottlenecks/challenges you might have in building such a system?