

Good morning!

CS 744: DATAFLOW

Shivaram Venkataraman

Spring 2024

ADMINISTRIVIA

Grading In Progress

- Course project proposal

→ *steer your project*

- Assignment 2 ~ *90% graded*

- Midterm

MID-SEMESTER FEEDBACK

“...instead of having everyone submit full discussion answers regardless of what group they're in, maybe just have a single check box on the form...”

“I hope there can be more details about the diagrams in discussion and exams...”

....

“... I found midterm exercises were more tricky/challenging ...”

“... hard time in the midterm exam...”

“Its an early morning class”

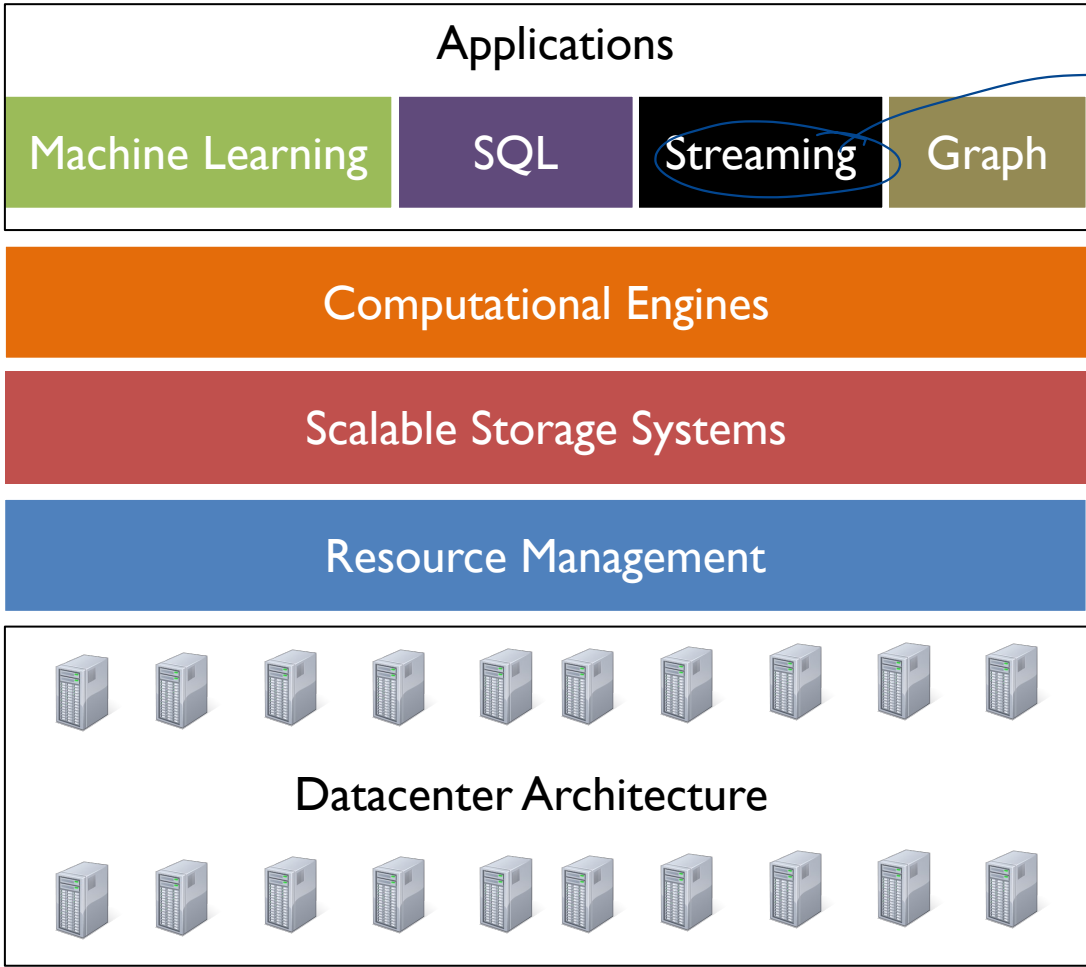
“The morning timings of the class”

“I want reading groups to talk and understand the paper better.”

“... If we had paper reading groups like Distributed Systems course...”

checkbox

After Spring
break



Properties
Dataflow model
↓
flink spark
 Streaming

Spark

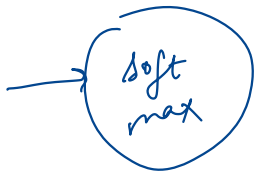


...

DATAFLOW MODEL (?)

PyTorch

data



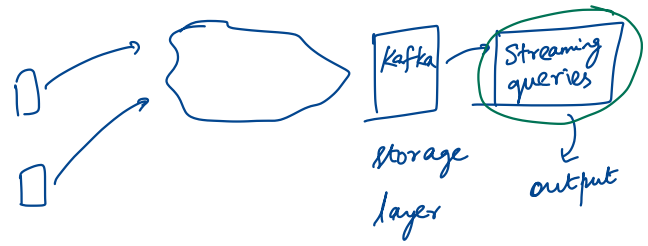
...

"Unbounded data"

- click logs on web pages
- IoT sensors
 - ↳ continuous operation
- Metrics from mobile phone apps

App logs
or
logs in
Data
Centers

MOTIVATION



Streaming Video Provider

- How much to bill each advertiser ?
- Need per-user, per-video viewing sessions
- Handle out of order data

unbounded
Data might arrive in
a different order than what was
generated
↳ Big challenge in
system design

Goals

- Easy to program ↳ similar to PyTorch, Spark
- Balance correctness, latency and cost

APPROACH

Separate user API from execution

Count views / video every minute, every 5 mins dashboard

Decompose queries into

Count / video - What is being computed

group every minute - Where in time is it computed

every 5 mins - When is it materialized

overwrite dashboard - How does it relate to earlier results

Select * from users | (map/reduce)

API for streaming queries

System to execute them

data has some ts associated with it

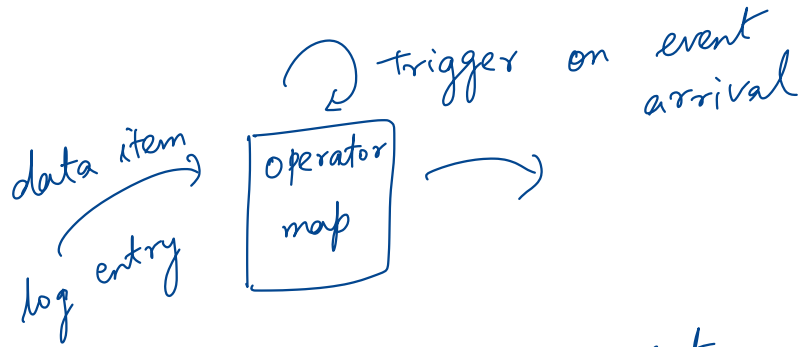
Triggers execution

Diff / overlap between results

STREAMING VS. BATCH

Streaming

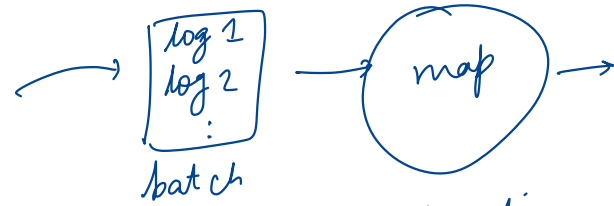
→ Flink, Timely Dataflow



- lower latency or outputs are produced quickly
- more "real time"

Batch

→ MapReduce, Spark, SCOPE



trigger at fixed time interval (5 mins)

- less real time
- more efficient → better output

TIMESTAMPS

Event time:

↳ Time at which event happened
associated by the sensor / device

Processing time:

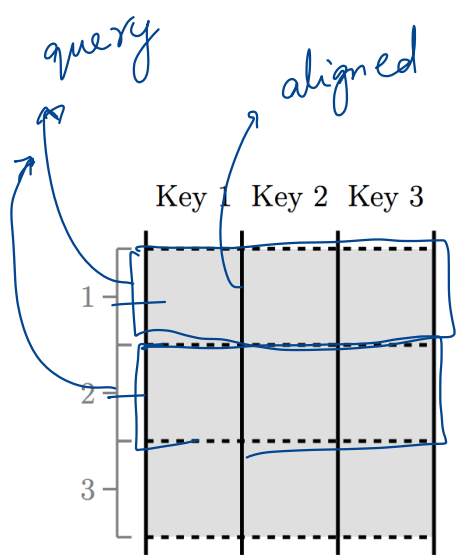
↳ Time at which event was processed
by stream processing system

Processing time \geq event time

WINDOWING

→ streaming late 90s

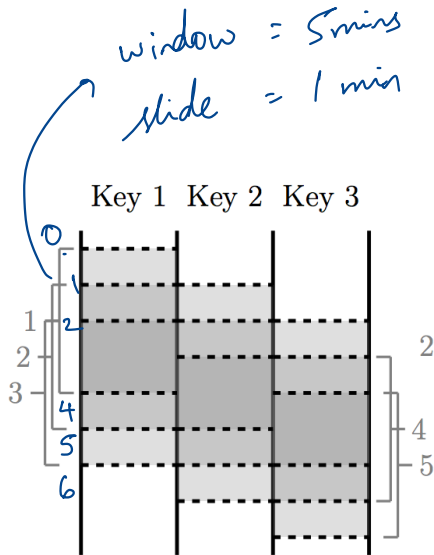
- event
time
stamp
1 min
window



Fixed

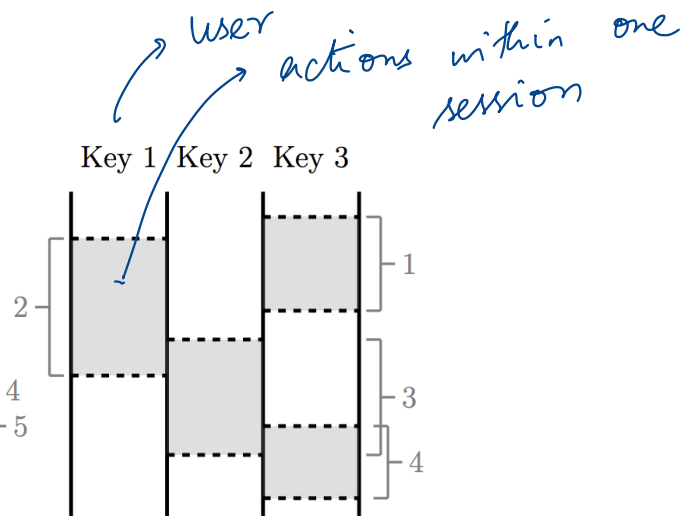
or

Tumbling
windows



Sliding

windows



Sessions

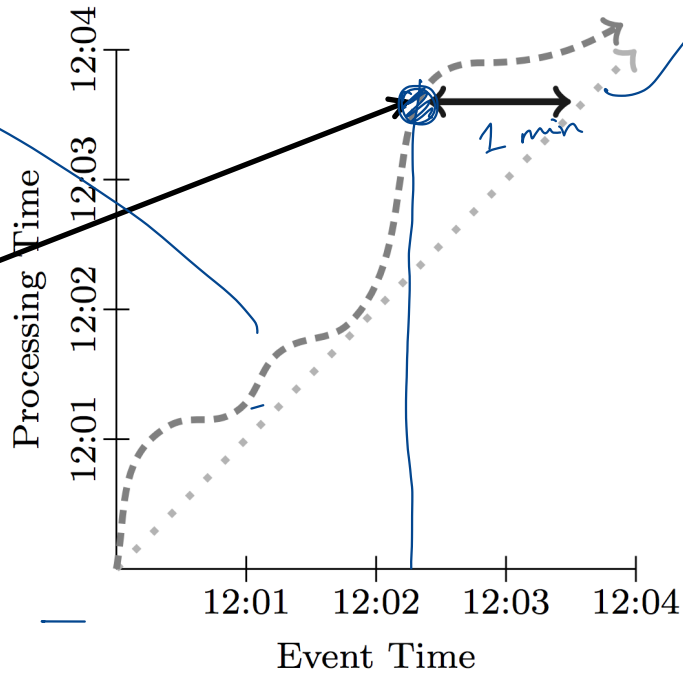
any activity
within a time
window

WATERMARK OR SKEW

App specific heuristics.

network latency ~ 1min

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



Proc = event time

gap between processing time between event time

- Actual watermark: ----->
- Ideal watermark:>
- Event Time Skew: <----->

API

ParDo:

GroupByKey:

Windowing

AssignWindow

MergeWindow

EXAMPLE

Input

Key value event ts window
 $(k_1, v_1, 13:02, [0, \infty))$,
 $(k_2, v_2, 13:14, [0, \infty))$,
 $(k_1, v_3, 13:57, [0, \infty))$,
 $(k_1, v_4, 13:20, [0, \infty))$

operator
 $\text{AssignWindows}(\text{Sessions}(30m))$

$(k_1, v_1, 13:02, [13:02, 13:32))$, \rightarrow event-ts + 30 mins
 $(k_2, v_2, 13:14, [13:14, 13:44))$,
 $(k_1, v_3, 13:57, [13:57, 14:27))$,
 $(k_1, v_4, 13:20, [13:20, 13:50))$

\downarrow DropTimestamps

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

map

GroupByKey

key values
same window w

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

\downarrow MergeWindows(
 $\text{Sessions}(30m)$)

$(k_1, [(v_1, [13:02, 13:50)),$ \rightarrow merge overlapping windows
 $(v_3, [13:57, 14:27)),$
 $(v_4, [13:02, 13:50))])$,
 $(k_2, [(v_2, [13:14, 13:44))])$

\downarrow GroupAlsoByWindow

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

\downarrow ExpandToElements

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

TRIGGERS AND INCREMENTAL PROCESSING

Windowing: where in event time are data grouped

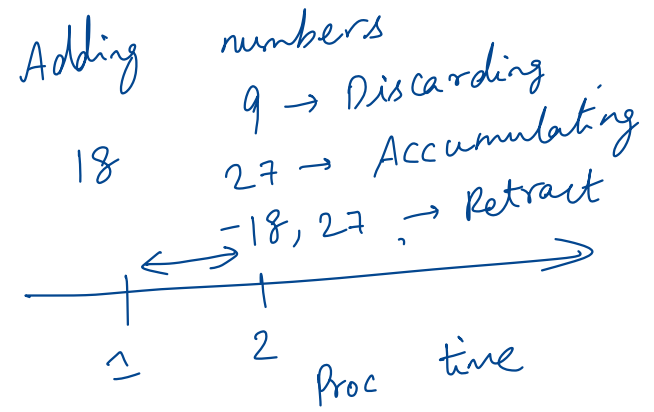
Triggering: when in processing time are groups emitted

Strategies

Discarding

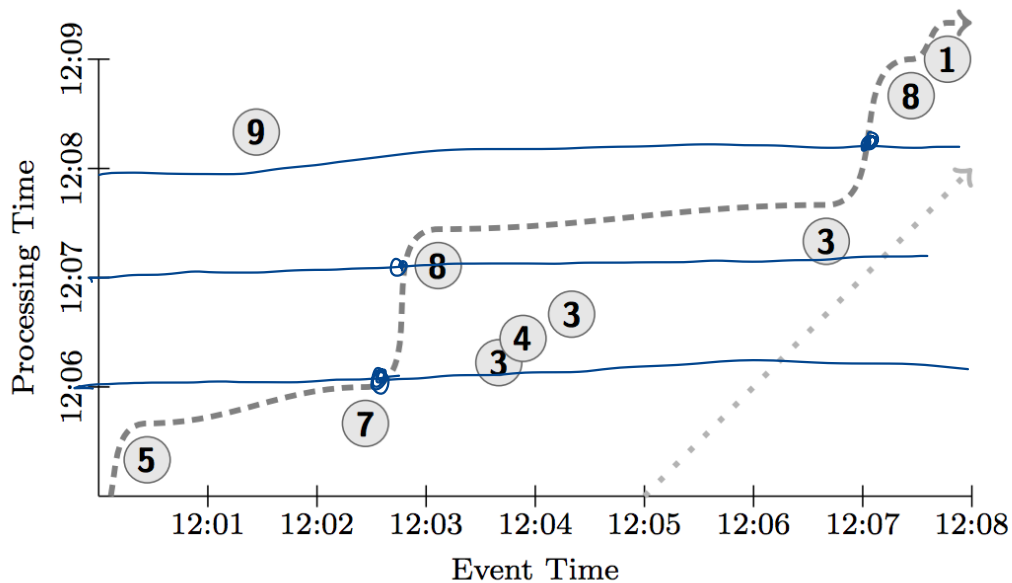
Accumulating

Accumulating & Retracting



RUNNING EXAMPLE

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



Actual watermark: ----->
Ideal watermark: >

GLOBAL WINDOWS, ACCUMULATE

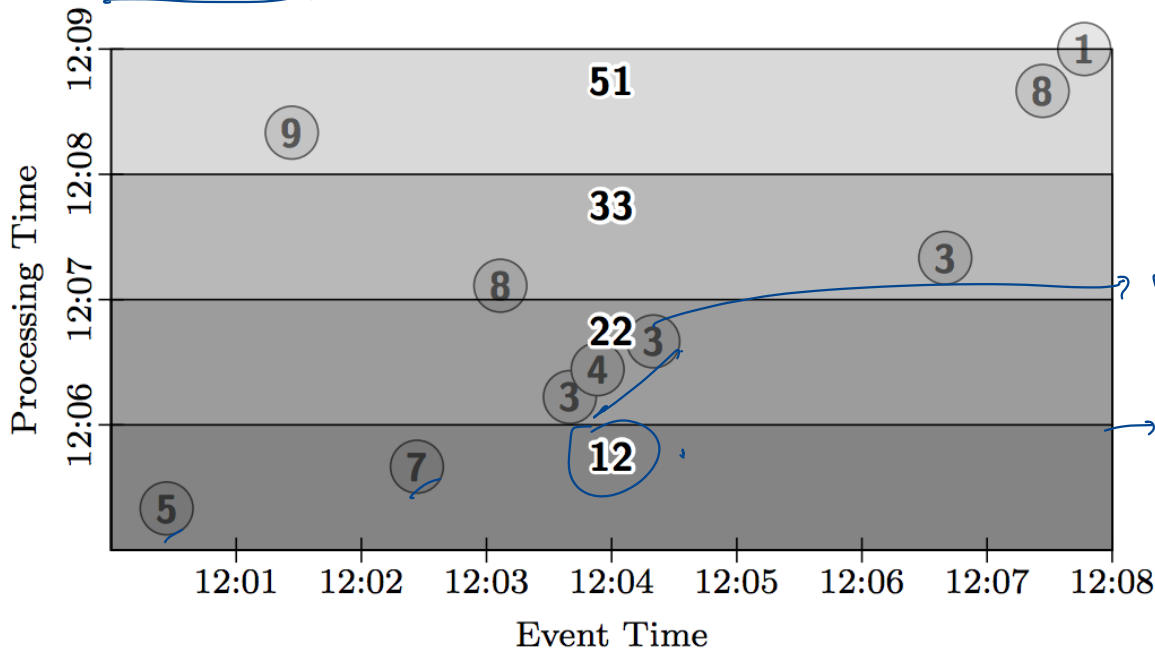
```
PCollection<KV<String, Integer>> output = input
```

```
.apply(Window.trigger(Repeat(AtPeriod(1, MINUTE))))
```

```
.accumulating()
```

```
.apply(Sum.integersPerKey());
```

tumbling every window 1 min



GLOBAL WINDOWS, COUNT, DISCARDING

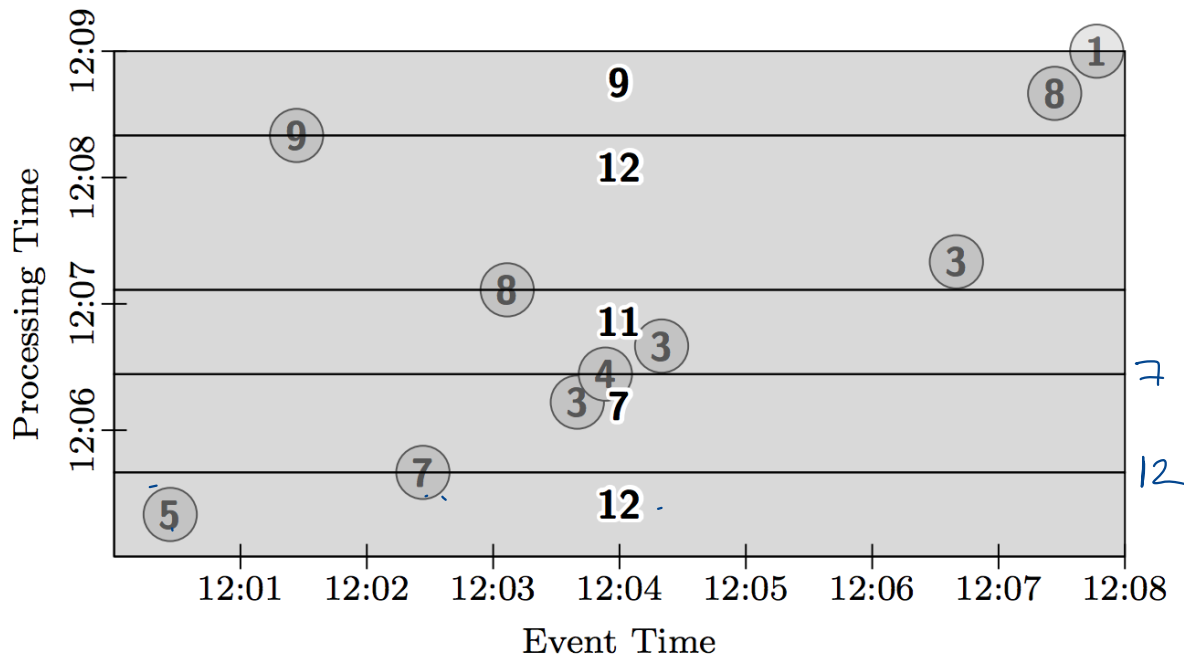
```
PCollection<KV<String, Integer>> output = input
```

```
.apply(Window.trigger(Repeat(AtCount(2))))
```

```
.discarding()
```

```
.apply(Sum.integersPerKey());
```

two keys



FIXED WINDOWS, MICRO BATCH

→ triggered every 1 min

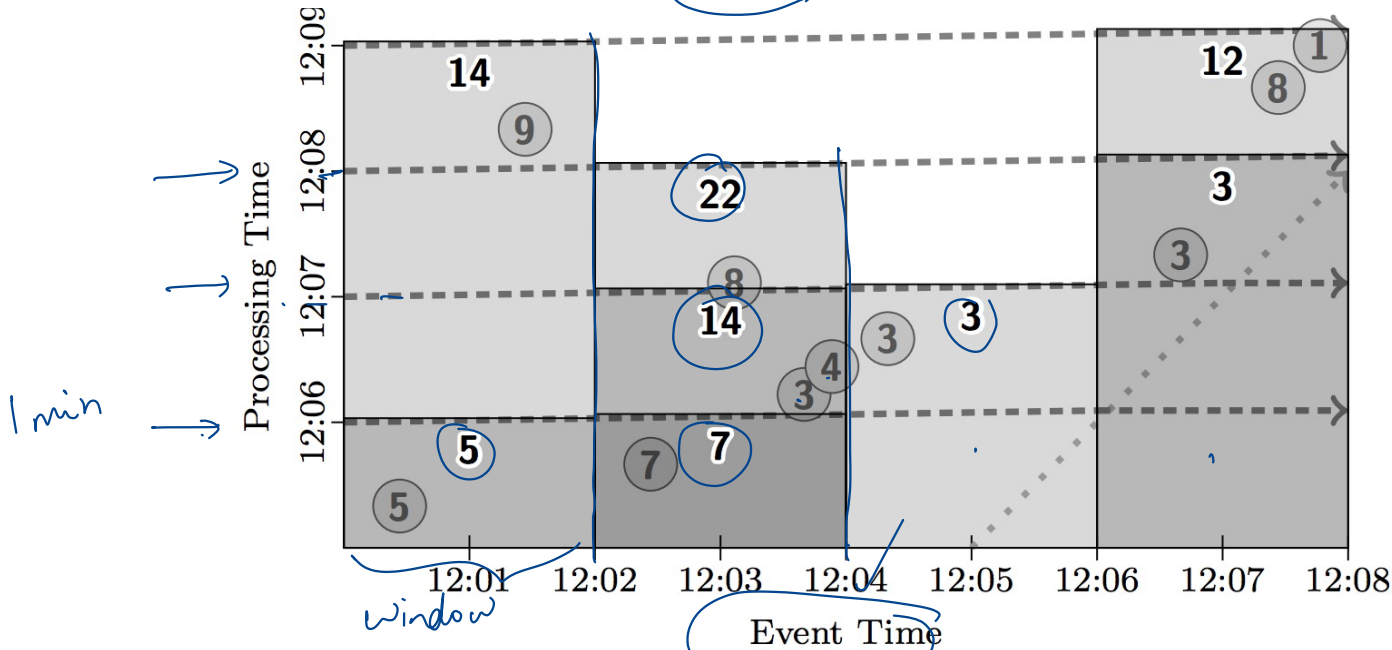
PCollection<KV<String, Integer>> output = input

```
.apply(Window.into(FixedWindows.of(2, MINUTES)))
```

```
.trigger(Repeat(AtWatermark()))
```

```
.accumulating()
```

watermark



SUMMARY/LESSONS

Design for unbounded data: Don't rely on completeness

Be flexible, diverse use cases

- Billing
- Recommendation
- Anomaly detection

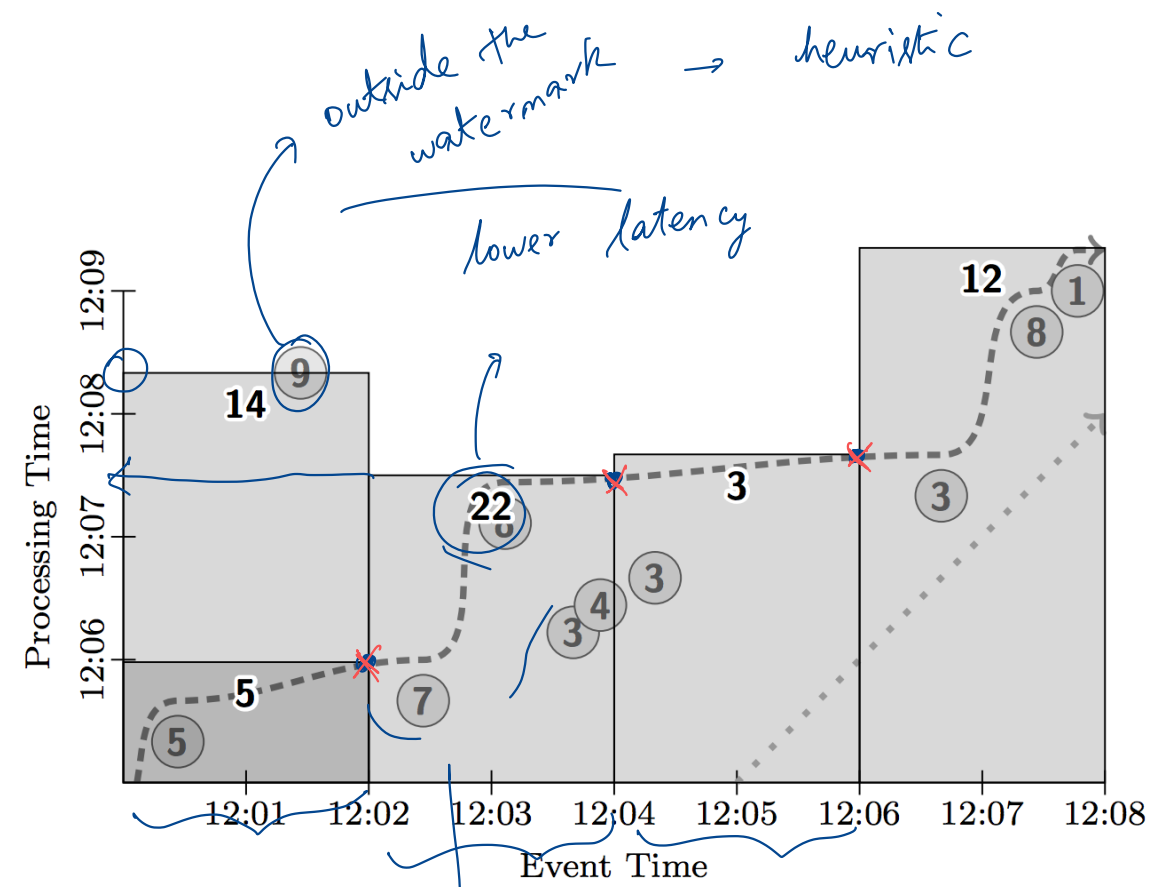
Windowing, Trigger API to simplify programming on unbounded data



DISCUSSION

<https://forms.gle/TB5kz2cH3uYc6rjv6>

Trigger a window when watermark crosses the window



Value 22 appear at 12:07:30
 At 12:08 with mini batch

with micro-batch, we get partial results at 1 min interval

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?

NEXT STEPS

Next class: Apache Flink