*Welcome!*

# CS 744: DATAFLOW

Shivaram Venkataraman

Fall 2020

# ADMINISTRIVIA

- Assignment 2 grades are up! → *Canvas*

- Midterm grading in progress

- Course project proposal comments
  ↳ *Peer feedback*    *Thursday this week*
  ↳ *Instructor feedback*

- AEFIS feedback (next slide)

# AEFIS FEEDBACK

*Better organization*

Improve writing on the slides, speak slower

Get a better internet connection? Better microphone?

↳ *Let me know how this sounds?*
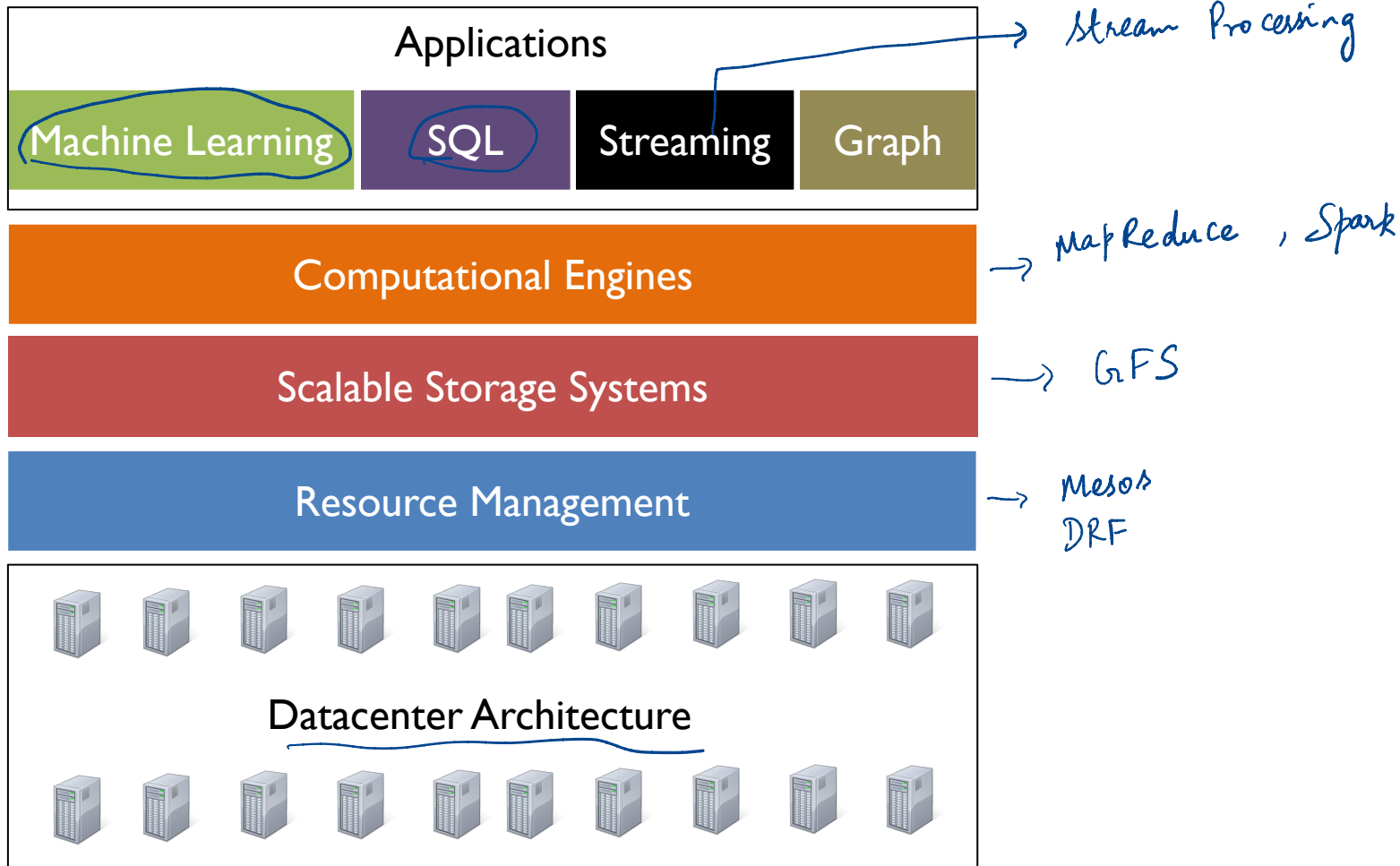
*Piazza*
{ More office hour slots

Discussion groups: same group each time? Also add prof. input

More time for Midterm exam, more guidance on deliverables

More homework/hands-on experience vs. too many evaluation components?

## Applications

| Machine Learning | SQL | Streaming | Graph |
|---|---|---|---|

→ Stream Processing

## Computational Engines

→ MapReduce , Spark

## Scalable Storage Systems

→ GFS

## Resource Management

→ Mesos
DRF

## Datacenter Architecture

operators or DAG of operators

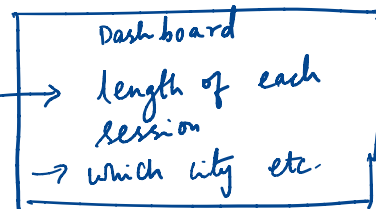Spark          SCOPE          PyTorch

# DATAFLOW MODEL (?)

# MOTIVATION

Streaming Video Provider

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

$\longrightarrow$ Mobile phone
Offline

Goals

- Easy to program
- Balance correctness, latency and cost

$\longrightarrow$ how accurate are your results

ESPN. com
- videos, each video has some ads

Dashboard
$\rightarrow$ length of each session
$\rightarrow$ which city etc.

Unbounded data, out of order
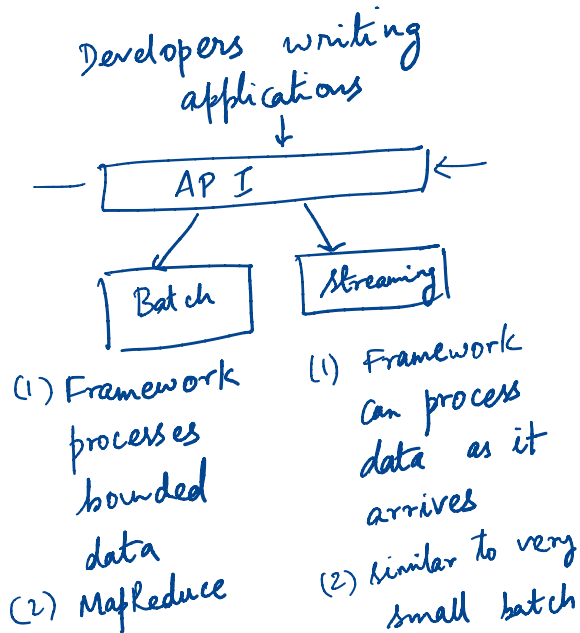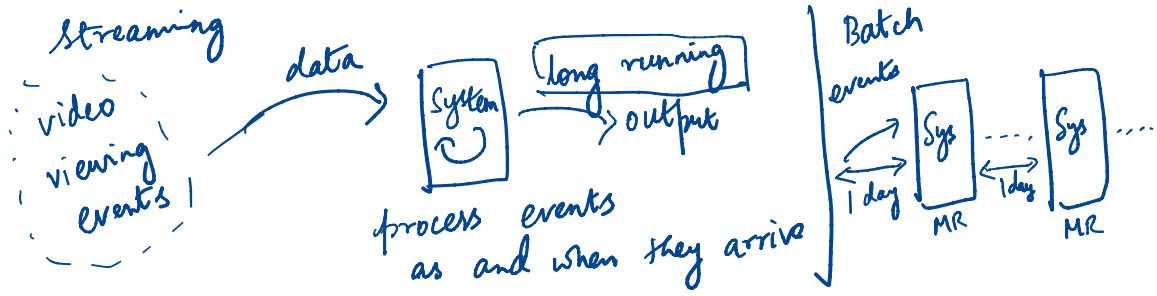how much delay till results are available

# APPROACH

API Design → Dataflow Model

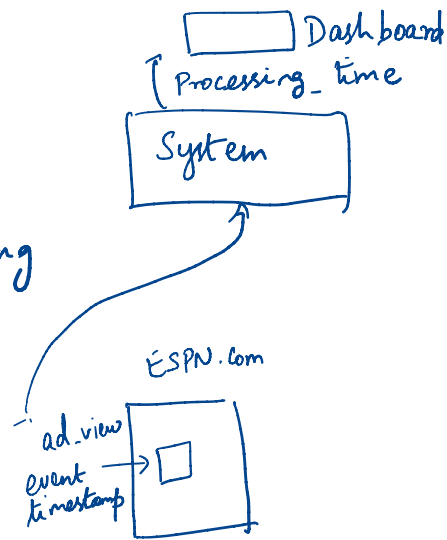  Separate user-facing model from execution

  Decompose queries into

    - What is being computed

    - Where in time is it computed

    - When is it materialized → Output

    - How does it relate to earlier results

Developers writing
applications
↓
AP I

Batch          Streaming

(1) Framework        (1) Framework
processes            can process
bounded              data as it
data                 arrives

(2) MapReduce        (2) similar to very
                         small batch

streaming

video
viewing
events

→ data → System → long running → output

process events
as and when they arrive

Batch
events

Sys ←1 day→ Sys ····· Sys ·····
     MR          MR

# TERMINOLOGY

Unbounded/bounded data $\longrightarrow$ Data is constantly arriving

Streaming/Batch execution
$\quad\quad\quad \rightarrow$ See previous slide

Dashboard
Processing_time

System

ESPN.com

ad_view
event
timestamp

Timestamps

Event time: Time when event occurs wrt user/input
e.g, time at ad was viewed in video

Processing time: Time at which an event is processed
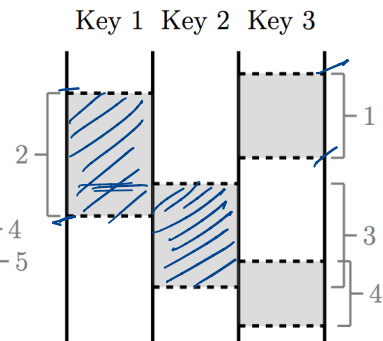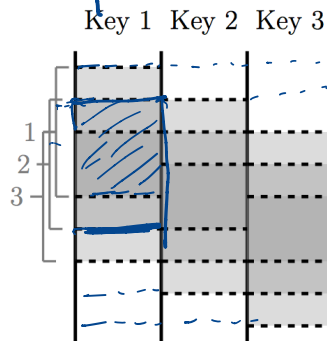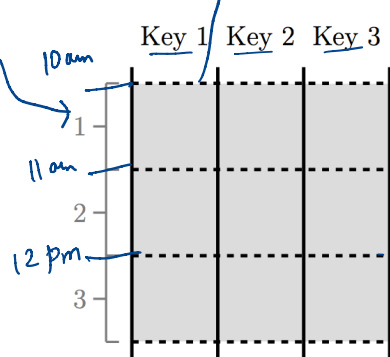e.g, time at which ad_view event is
processed to update the dashboard.

# WINDOWING

logical
constructs

windows
are aligned
across keys

10 am to 11 am
10:30 am to 11:30 am
⋮

window - Id

Range of
event - timestamp

Do not
overlap
with each
other

|  | Key 1 | Key 2 | Key 3 |
|---|---|---|---|
| 10am | | | |
| 1 | | | |
| 11am | | | |
| 2 | | | |
| 12 pm | | | |
| 3 | | | |

Fixed
or
Tumbling
windows

|  | Key 1 | Key 2 | Key 3 |
|---|---|---|---|
| 1 | | | |
| 2 | | | |
| 3 | | | |
| 4 | | | |
| 5 | | | |

Sliding

overlap between
consecutive
windows

|  | Key 1 | Key 2 | Key 3 |
|---|---|---|---|
| 1 | | | |
| 2 | | | |
| 3 | | | |
| 4 | | | |

Sessions

session
length

Not aligned
across all the
keys

# WATERMARK OR SKEW

- Watermark is not easy to know
- Heuristics
  - After 10 mins most devices send events

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

12:03:30

catch up

- Processing time lags event time
- Event time skew



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

No gap between event-t & processing time

# API

ParDo:      *Map in MapReduce or FlatMap in Spark*
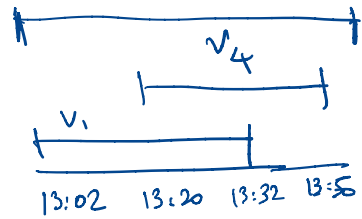
GroupByKey:   *Reduce in MapReduce*

Windowing
  AssignWindow    → *Buckets tuple into a window*

  MergeWindow    → *Merge buckets based on strategy*
                     *( sessions)*

# EXAMPLE

*Assign tuples to Sessions*

↓ ↓ ↳ *event - timestamp*

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

$\downarrow$ *AssignWindows(*
  *Sessions(30m))*

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

*add 30m to event timestamp*

$\downarrow$ *DropTimestamps*

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

GroupByKey

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

$\downarrow$ *MergeWindows(*
  *Sessions(30m))*

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

*overlap and merges them*

$\downarrow$ *GroupAlsoByWindow*

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

$\downarrow$ *ExpandToElements*

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

$v_4$
$v_1$
13:02  13:20  13:32  13:50

# TRIGGERS AND INCREMENTAL PROCESSING

Windowing: where in event time data are grouped

Triggering: when in processing time groups are emitted

Strategies

    Discarding $\equiv$ 6

    Accumulating $\equiv$ 11

    Accumulating & Retracting

$\equiv$ -5, 11

    ↓         ↓

  retracting   Accumulating

---

v1, 2

v1, 3

⋮

v1 6

Counter, sum of all views for a video
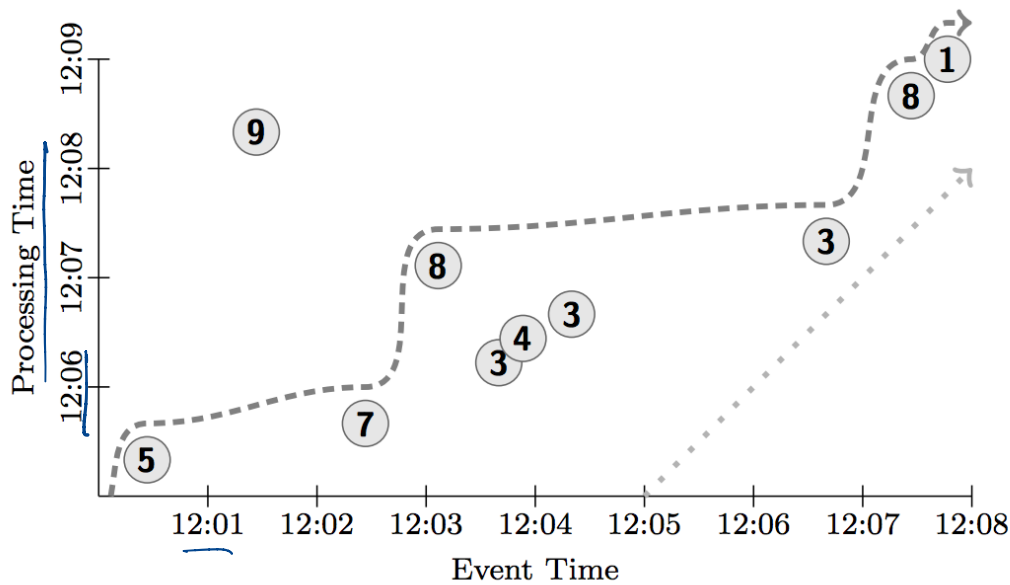
2+3 → Output = 5

Output =

# RUNNING EXAMPLE

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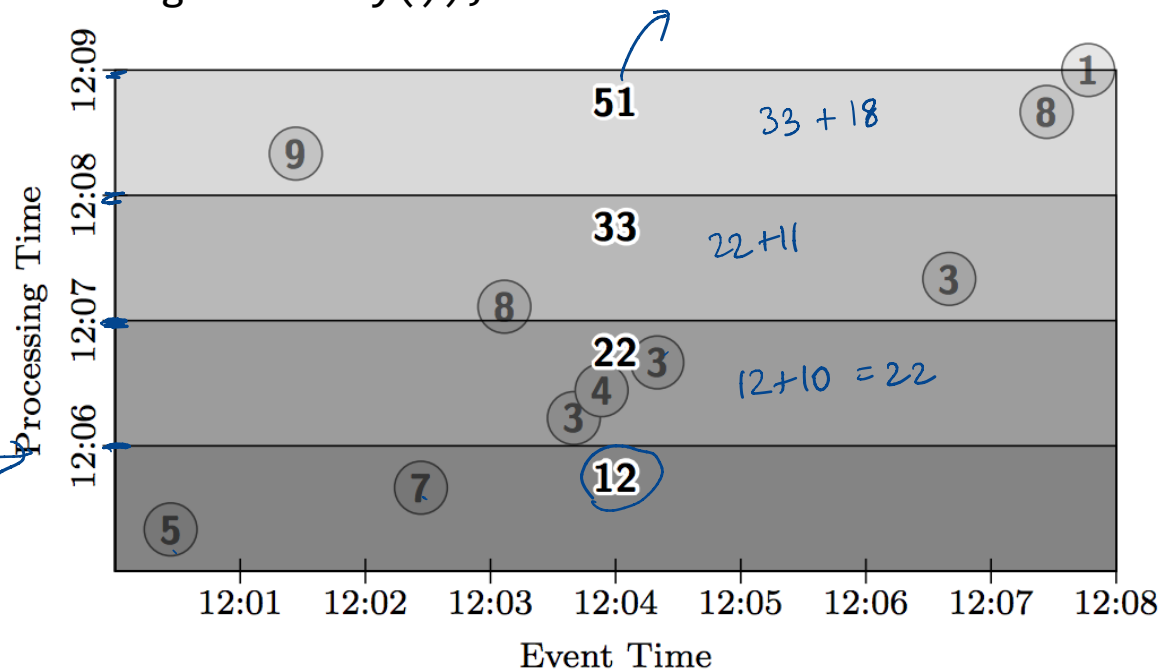
# GLOBAL WINDOWS, ACCUMULATE

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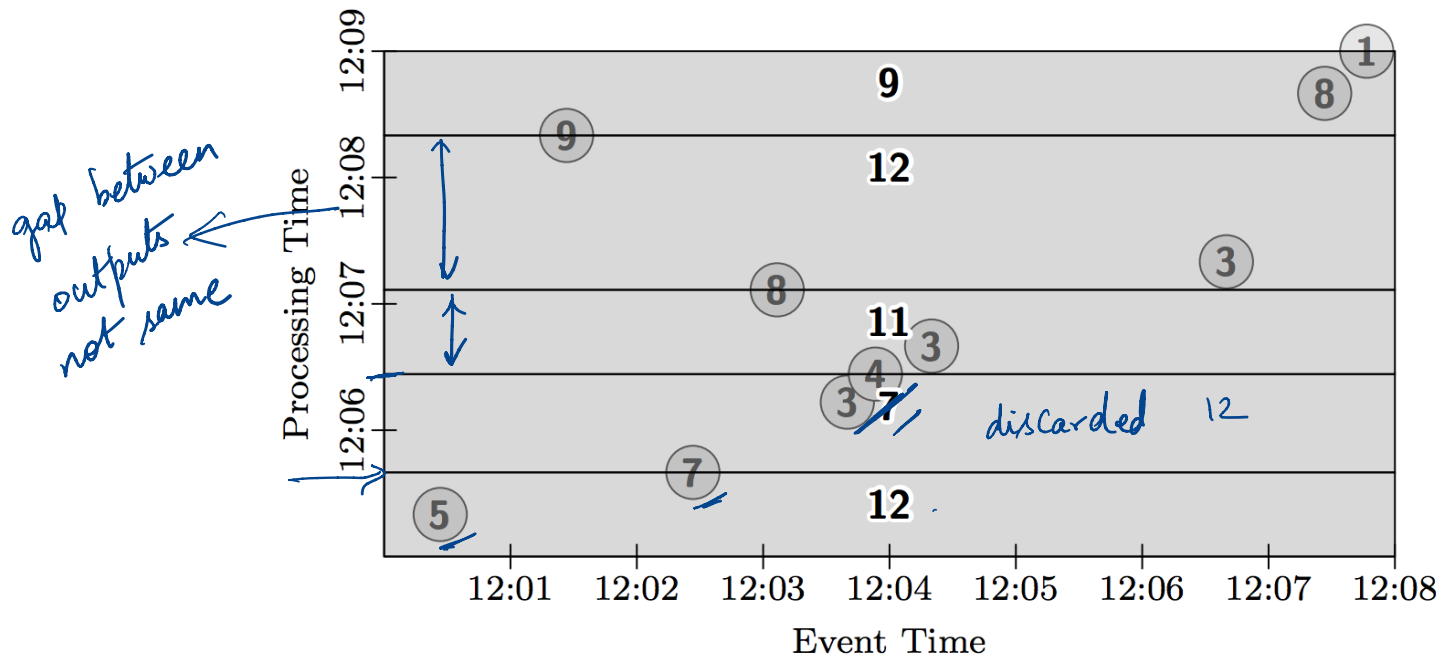
# GLOBAL WINDOWS, COUNT, DISCARDING

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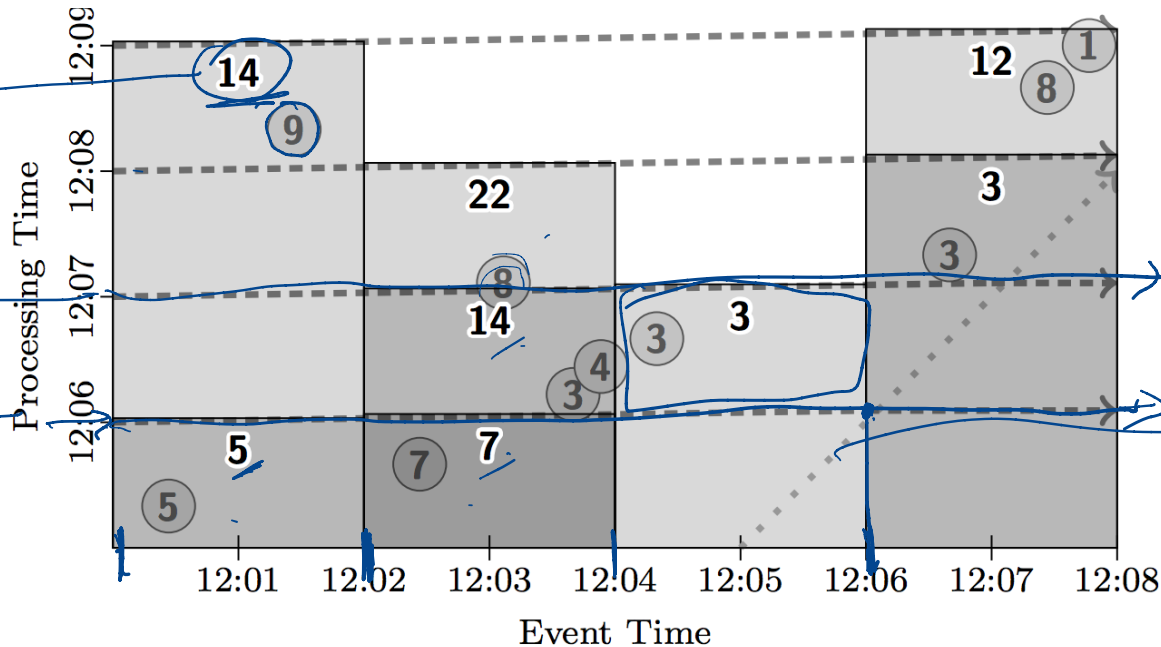
# FIXED WINDOWS, MICRO BATCH

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

# SUMMARY/LESSONS

Design for unbounded data: Don't rely on completeness

Be flexible, diverse use cases
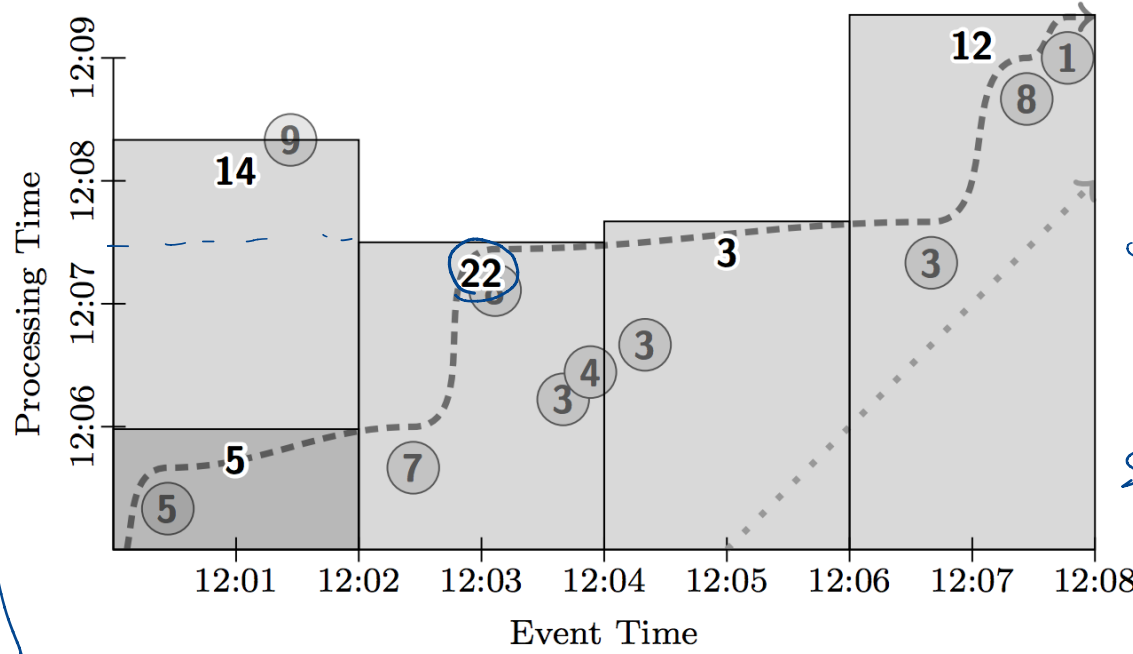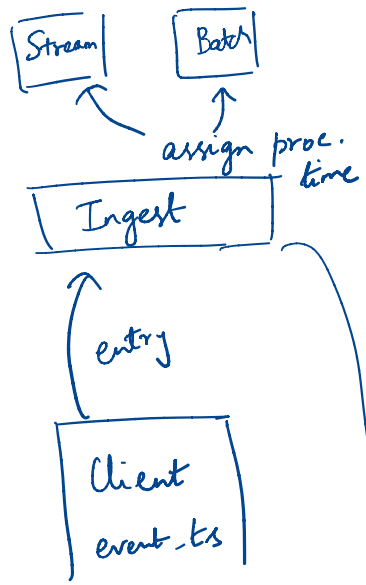
- Billing

- Recommendation

- Anomaly detection

Windowing, Trigger API to simplify programming on unbounded data

# DISCUSSION

https://forms.gle/jwHjTBbR49vyQASq6

Fixed windows streaming
Assume watermark is given

(1) Window fires every time
watermark pass
⇒ Worse latency
⇒ fewer outputs

Stream    Batch

assign proc. time

Ingest

entry

Client
event_ts

Written to System
Apache Kafka
persist disk

Pub-Sub

Ingest
time

proc_t

↳ Update query

(2) Micro batch
partial
results

Streaming ⇒
buffer events
until watermark

Processing Time

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?

# NEXT STEPS

Next class: Naiad

Course project proposal peer feedback