

Hello!

# CS 744: FLINK

Shivaram Venkataraman

Fall 2022

# ADMINISTRIVIA

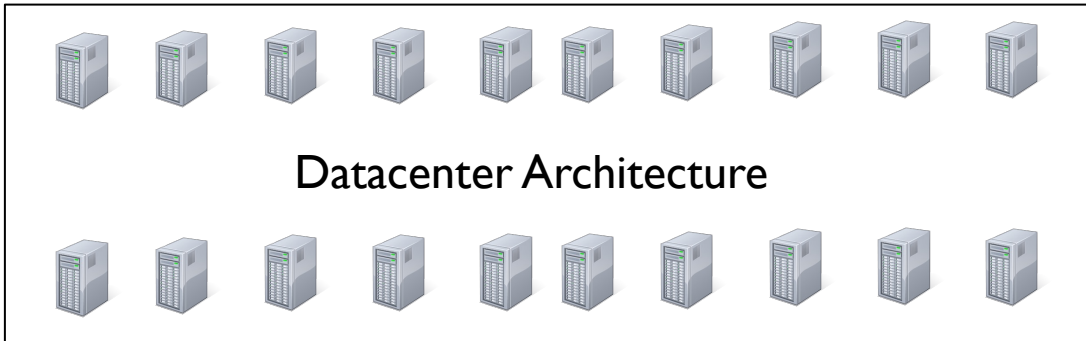
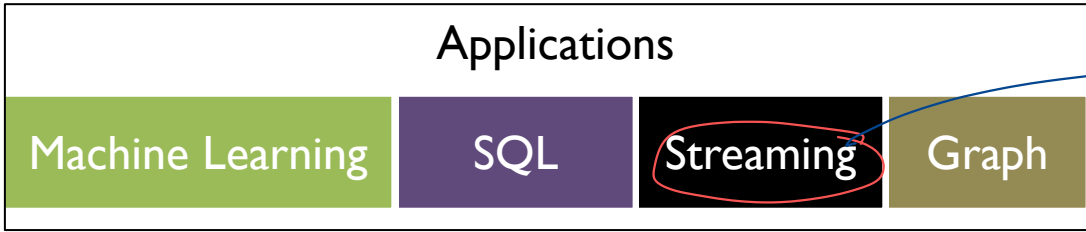
## In Progress

- Course Project Proposal feedback *~ end of this week*
- Midterm, [Assignment 2 grading]

## Resources for Course Projects

- Cloumlab reservations (Check Piazza)
- GCP credits (Email Roger and me)

*Google Cloud \$50*



API requirements / features semantics

# DASHBOARDS

*Input* unbounded stream of data  
*Users* express queries. *Output* ↻

## Sales Dashboard

Total Sales  
**\$3,256.8M**

Number of Deals  
**17,164**

Avg Deal Size  
**\$189,545**

Rev. per Salesperson  
**\$20.5M**

Week of Date Closed

December 6, 200 December 25, 20



Region

(All)

Country

(All)

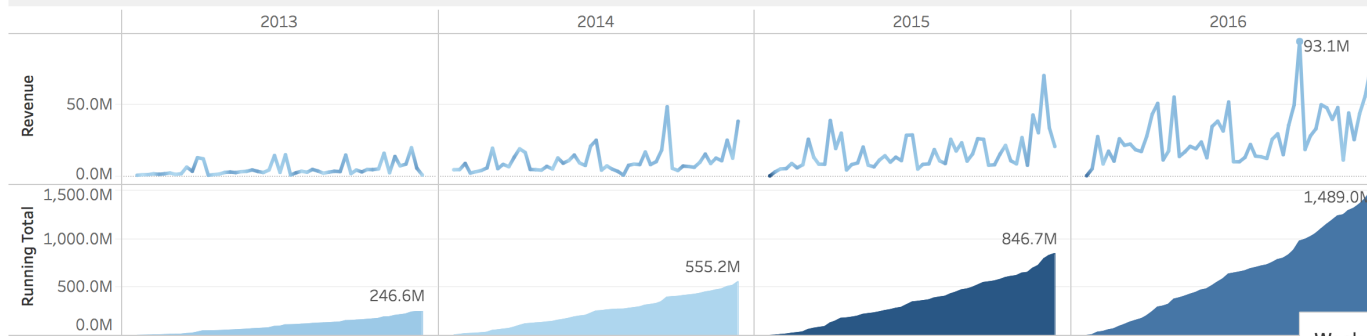
Sales Team

- (All)
- Small and Midmarket
- Enterprise

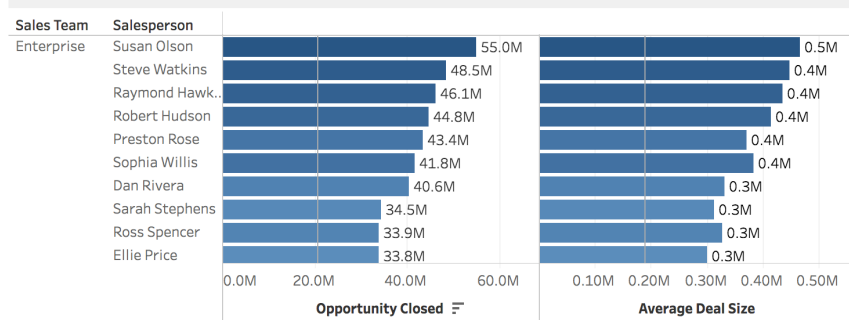
Avg Deal Size/Salesperson



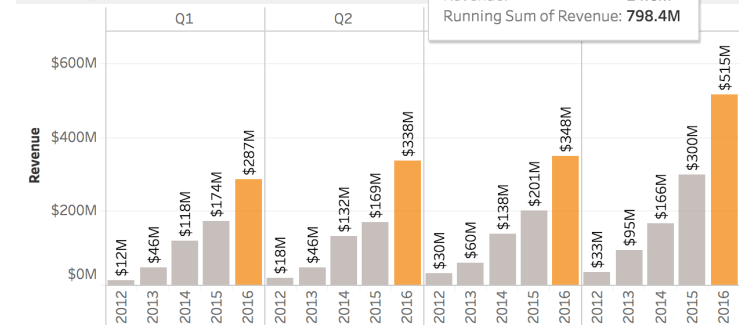
### Revenue Over Time



### Sales Team Performance



### Revenue by Quarter



**Week of September 4, 2016**  
 Revenue: **14.6M**  
 Running Sum of Revenue: **798.4M**

# STREAMING COMPUTATION

event time

(phone1, 2:01, open)  
(phone1, 2:03, close)  
...



(phone2, 2:05, open)



(phone4, 2:03, close)



Mappers  
split by (action, hour)

hash (open)

[2:00, open]

(2:00, close)

hashmap

open	1:00	5
open	2:00	1

close	1:00	5
close	2:00	2

grouping events by action  
- Counting how many  
window of 1 hr



Output Sink

Reducers  
count by (action, hour)

extracting the hour

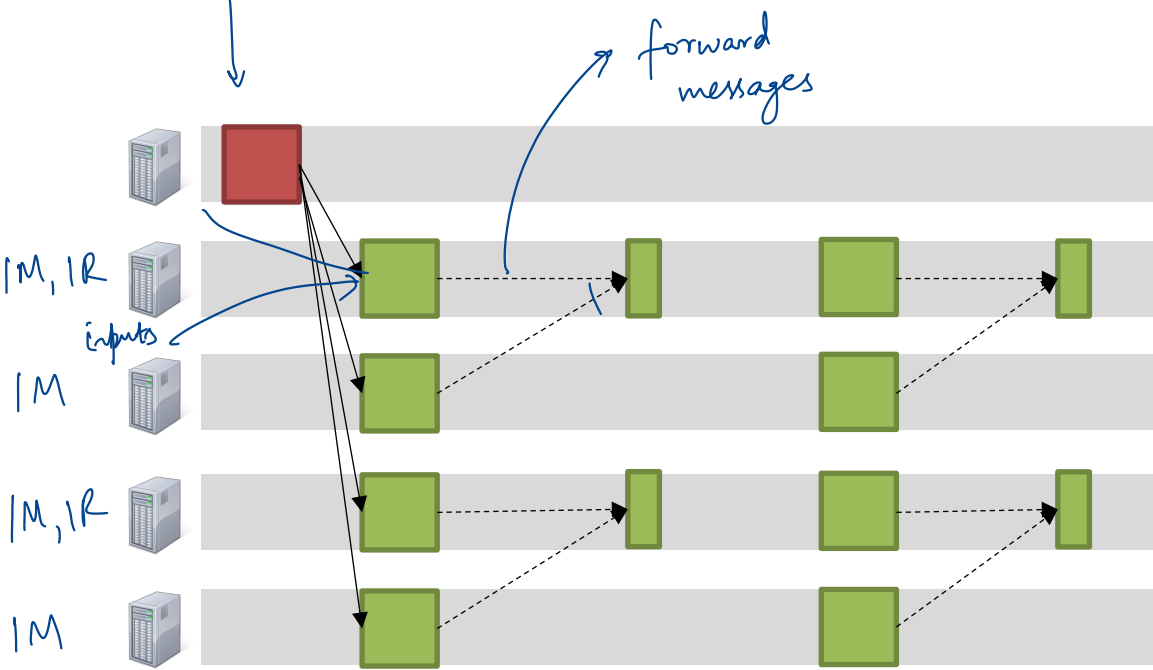
window based  
on hour  
and then count  
how many in each  
window

① Reducers  
don't wait  
for maps to  
end

② Tasks are  
long running  
→ allocate resources  
for all tasks

# FLINK: COMPUTATION MODEL

Query → Streaming Word Count



operators implement Map / Reduce / Windows

Long-lived operators

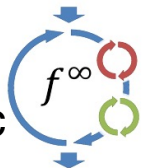
Mutable State

Operators internally maintain state

Google MillWheel



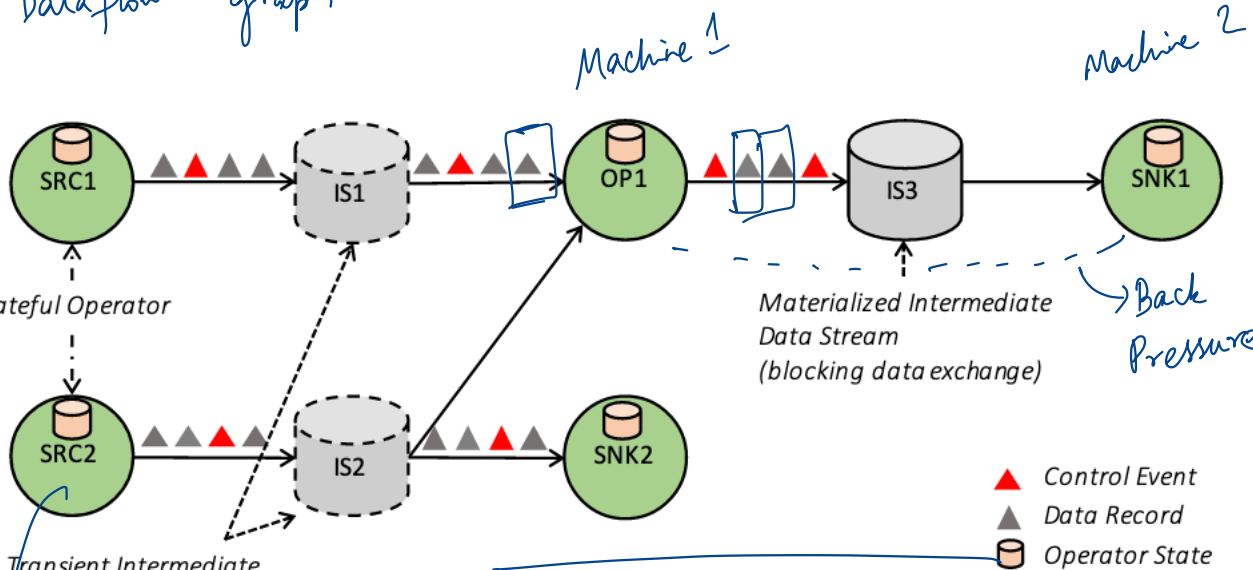
Streaming DBs: Borealis, Flux etc



Naiad

# INTERMEDIATE DATA STREAMS

Data flow graph



- Operators can be partitioned

- Streams can be 1-1 or 1-many

- Pipelined stream  
- Comp. & comm. pipelined

- Store intermediate tuples
- Back Pressure
- Increase num partition / machines (slowest operators)

operators

# STATEFUL OPERATORS

## Examples?

- All windowing operators
- Aggregation / maintaining summaries

Stateless

- `map()` : input → output
- `filter()` : input, filter criteria

## Challenge

How to ensure fault tolerance?

Explicitly register local variables → annotations

StateBackends that are automatically saved/recovered

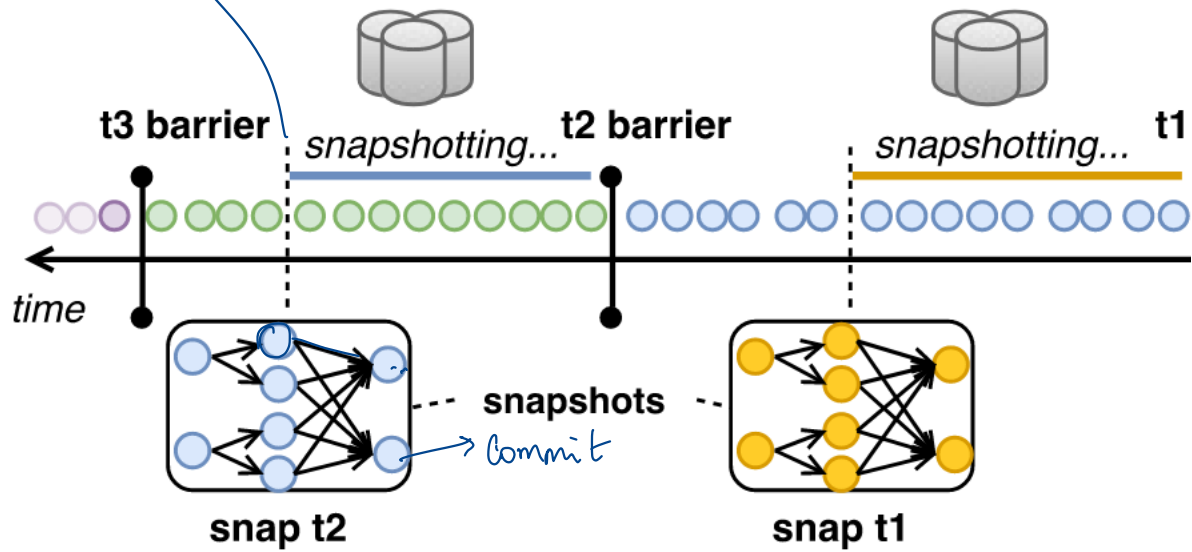


→ need to restore state when op fails.



# FAULT TOLERANCE: CHECKPOINTING

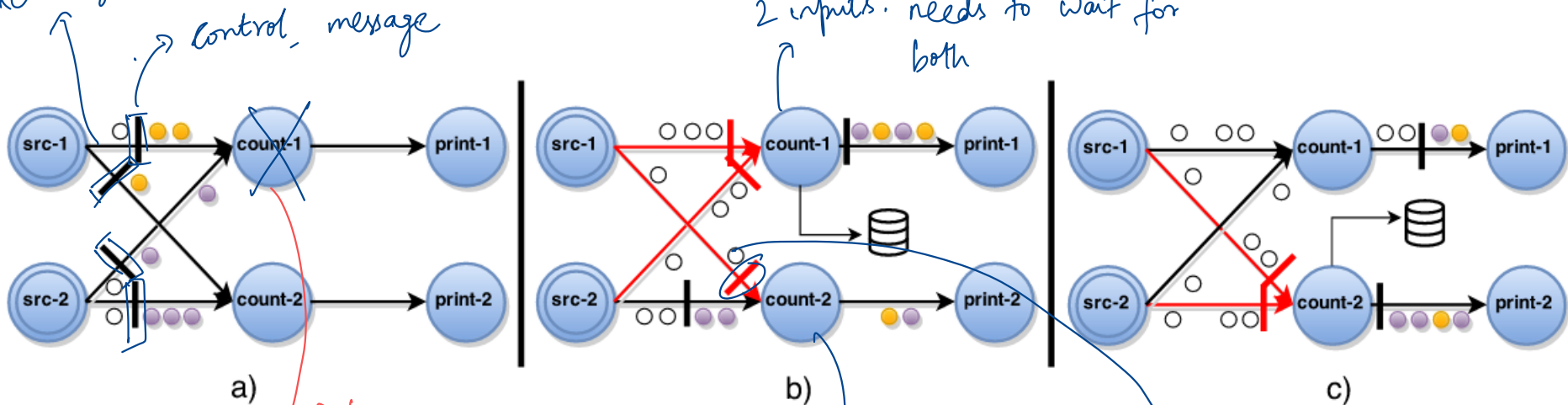
start snapshot → Control msg to first op  
→ Fwd to downstream



Data Streams : Control messages

- Cannot replay all the data from beginning
- Input is "replayable"
  - Apache Kafka
- Fault Tolerance :
  - "exactly once semantics"
  - at least once
  - at most once

# ASYNCHRONOUS BARRIER SNAPSHOTTING → arXiv



In-order delivery

Control message

2 inputs. needs to wait for both

a)

b)

c)

failure

reset all operators to snapshot!

only saw 1 barrier so far

no tuple after barrier should be in snapshot

<span style="color: yellow;">●</span> Preshot records src-1	Operator snapshot
<span style="color: purple;">●</span> Preshot records src-2	Snapshot barrier
<span style="color: grey;">○</span> Postshot records	Blocked channel

# WATERMARKS, WINDOWS

Implements similar model as Dataflow

“Watermarks originate at the sources of a topology”  $\longrightarrow$  External to Flink  
Propagate through the other operators of dataflow

Windows based on event-time, processing time, ingest time(?)

Same  
idea

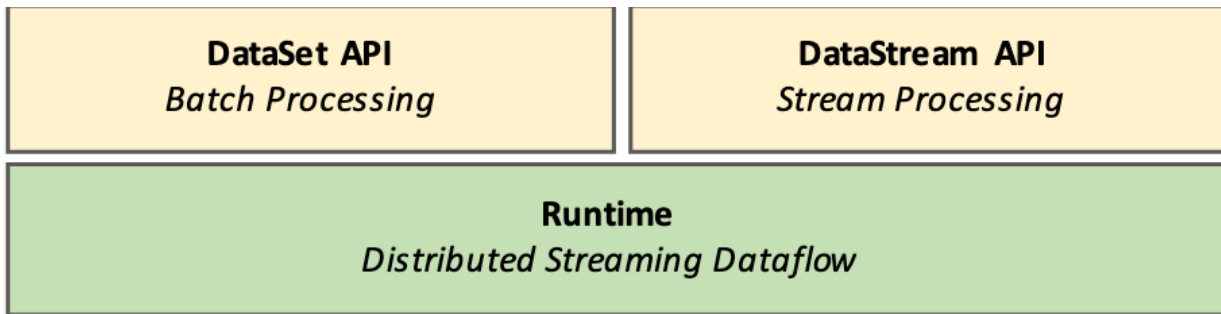
```
stream  
  .window(SlidingTimeWindows.of(  
    Time.of(6, SECONDS), Time.of(2, SECONDS))  
  ).trigger(EventTimeTrigger.create())
```

# COMBINING BATCH, STREAMING

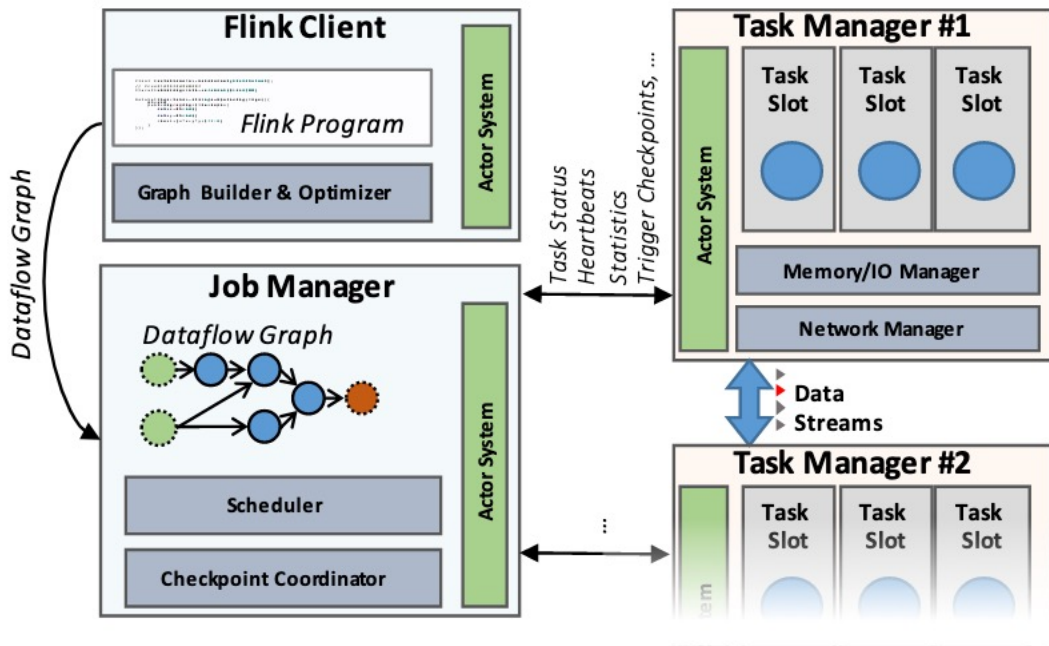
Blocked DataStreams → Intermediate data in MR

Turn off periodic snapshots → Restart job

Blocking operators (e.g., sort) → New operators batch specific



# OVERALL ARCHITECTURE



# SUMMARY

Stream processing → Increasingly important workload trend

Flink: Distributed streaming dataflow to run streaming, batch, iterative

Distributed streaming dataflow

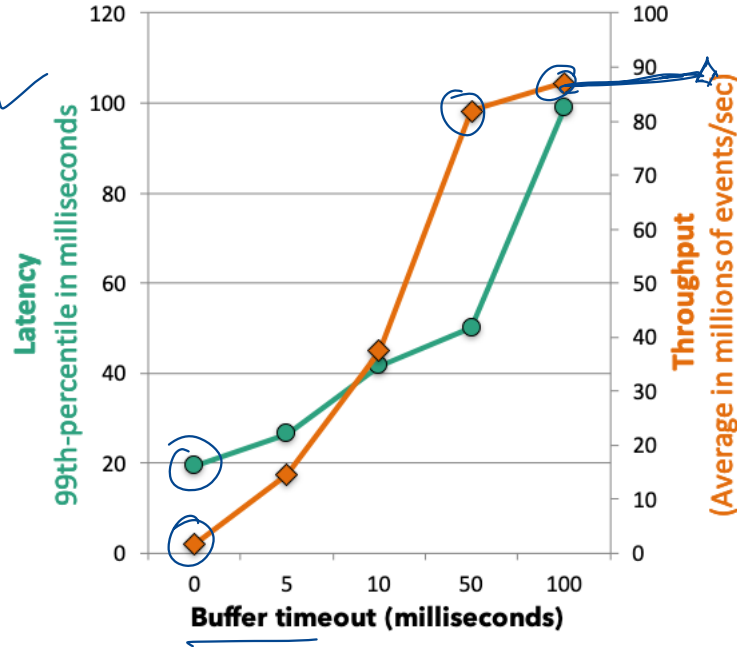
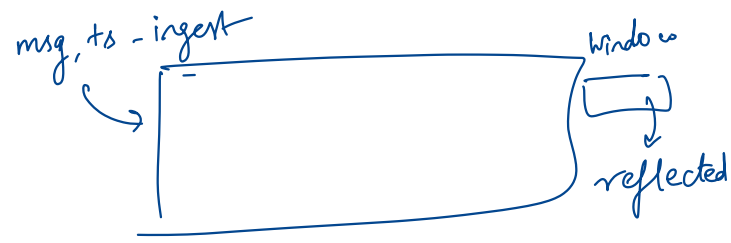
- Stateful operators
- Checkpointing based FT

# DISCUSSION

<https://forms.gle/idnKJWoBFHpRNC2j7>

# Data Streams

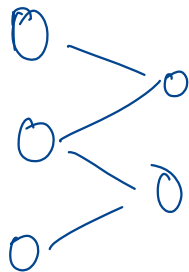
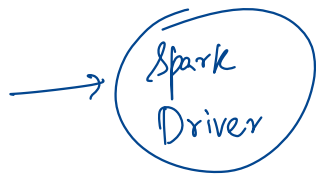
- Small buffer timeout
- low hput X
- low 99<sup>th</sup> -ile latency ✓
- Very large timeout
- saturate at some hput
- latency might keep going up?





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?

Batch size  $bs = ls$



}  $\Rightarrow$  batch size is small  
lots of tasks  
high overhead for launching tasks

- Tasks are stateless. Windowing needs state!

- Caching is less useful? Data changes every time

# SUMMARY

Next class: Spark Streaming