good morning !!

# CS 744: FLINK

Shivaram Venkataraman
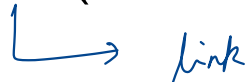
Spring 2024

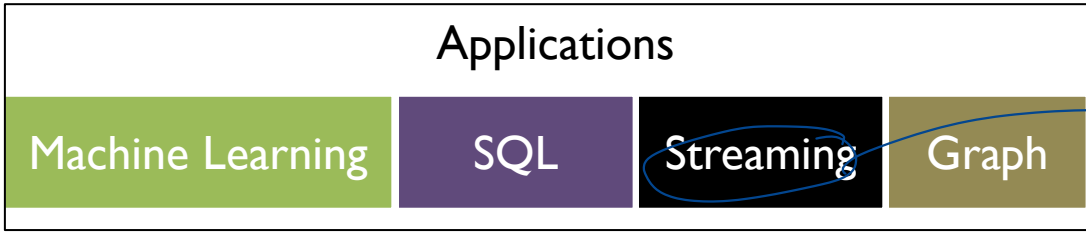# ADMINISTRIVIA

Grading

- Assignment 2 grading → *released*

- Course Project Proposal feedback → *today / tomorrow*

- Midterm → *next week*


Resources for Course Projects

    - Cloudlab (Reservations?) → *Survey ?*

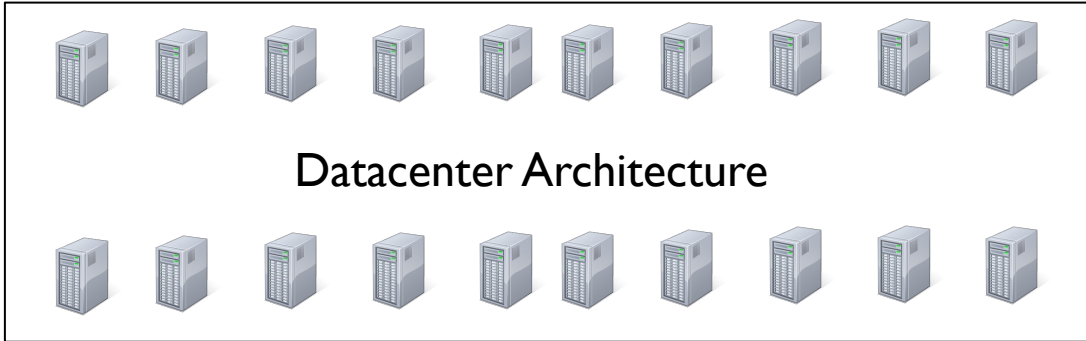    - GCP credits (Email Tzu-Tao and me)

        → *link*

# DASHBOARDS

*update    regularly*

## Sales Dashboard

| Total Sales | Number of Deals | Avg Deal Size | Rev. per Salesperson |
|---|---|---|---|
| **$3,256.8M** | **17,164** | **$189,545** | **$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**
$147,043   $336,519

*per hour / per day etc.*

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

### Sales Team Performance

| Sales Team | Salesperson |
|---|---|
| Enterprise | Susan Olson |
|  | Steve Watkins |
|  | Raymond Hawk.. |
|  | Robert Hudson |
|  | Preston Rose |
|  | Sophia Willis |
|  | Dan Rivera |
|  | Sarah Stephens |
|  | Ross Spencer |
|  | Ellie Price |

# STREAMING COMPUTATION

① MR → wait for
all mappers
before reducers

Here:
we are
doing them
concurrently

② FT is different

③ Bottlenecks

event TS

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

(phone2, 2:05, open)

(phone4, 2:03, close)

→ Intermediate
data disk. Streaming
them here!

Putting it
in window

(2:00, open)

(2:00, close)

windows on event ts

→ trigger

| open | 1:00 | 5 |
| open | 2:00 | 1 | +1 |

| close | 1:00 | 5 |
| close | 2:00 | 2 | +1 |

Database /
dashboard
etc.

MySQL

Output Sink

Mappers
split by (action, hour)

Reducers
count by (action, hour)

④ Input source

↳ mutable state that is
present in operators

# FLINK: COMPUTATION MODEL

Started at beginning

Long-lived operators

Mutable State

Google
MillWheel

Flink

Streaming DBs:
Borealis, Flux etc $f^\infty$ Naiad

Manager/Driver      ⟶   Control Message

Task/Computation     ⤑   Network Transfer

# INTERMEDIATE DATA STREAMS



10 tuples /sec

5 tuples /sec

Buffer

Stateful Operator

network

Materialized Intermediate
Data Stream
(blocking data exchange)

Similar
to MR

Transient Intermediate
Data Stream (pipelined data exchange)

▲ Control Event
▲ Data Record
🛢 Operator State

Data Stream

Transient
  ↳ Pipeline

Comm. with

Computation

  ↳ 1-1 or 1-many
         or
   many -1 dependency

Trigger
— Buffer full
— Timeout [20 ms]

Back pressure
— signal that goes from
  dest to src to say
  slow down

# STATEFUL OPERATORS

Examples?

→ Windowing operations

→ Aggregation

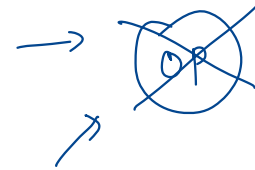Challenge

How to ensure fault tolerance?

Explicitly register local variables

StateBackends that are automatically saved/recovered

Stateless operators
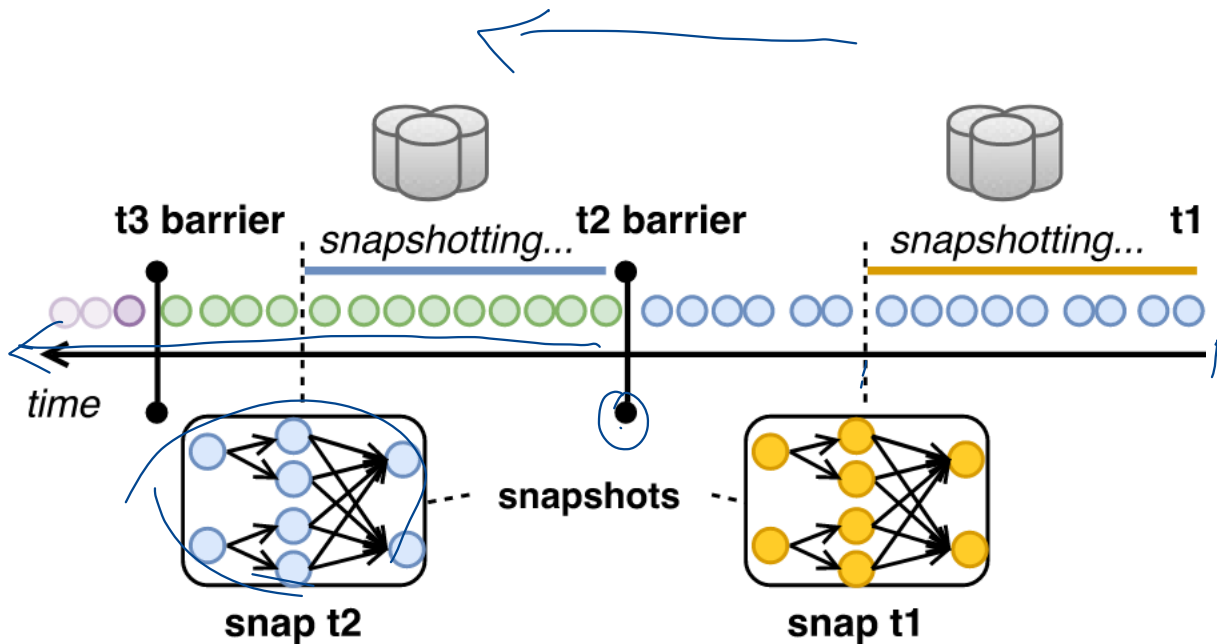
→ map() → input tuple
↓
output tuple

→ filter()

→ ⊘ → what do we do?

→

→ annotate some variables in operator → Periodically state variables as state variables checkpoint variables

# FAULT TOLERANCE: CHECKPOINTING



**t3 barrier** *snapshotting...* **t2 barrier** *snapshotting...* **t1**

time

**snapshots**

**snap t2**        **snap t1**

Periodically take checkpoints

Restart the operator and replay tuples that came after checkpoint.

Guarantee that we want:
→ Exactly once semantics
↳ useful

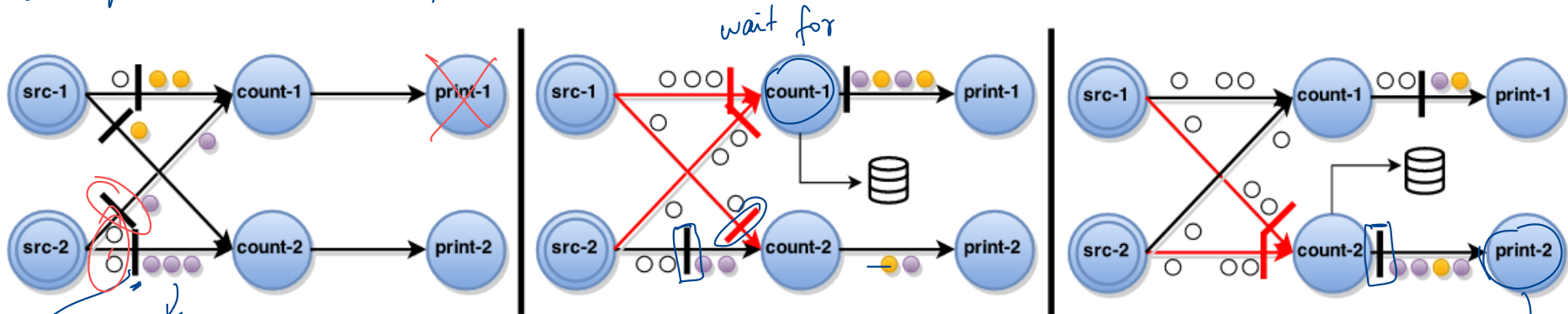each input's effect is seen exactly once in the output!

Input is replayable!

# ASYNCHRONOUS BARRIER SNAPSHOTTING

all operators have processed upto some .tuple → consistent snapshot

wait for



a)  b)  c)

check point

data records

reset all operator to check point state

ckpt is complete

| | |
|---|---|
| 🟠 | Preshot records src-1 |
| 🟣 | Preshot records src-2 |
| ⚪ | Postshot records |

| | |
|---|---|
| 🗄 | Operator snapshot |
| ▮ | Snapshot barrier |
| → | Blocked channel |

① Control message to ckpt
state is saved when
message is recvd.

② checkpointing takes time / slowdown
during ckpt

# WATERMARKS, WINDOWS

Implements similar model as Dataflow

"Watermarks originate at the sources of a topology"

*outside of Flink*

Propagate through the other operators of dataflow

*same as in Dataflow paper.*

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

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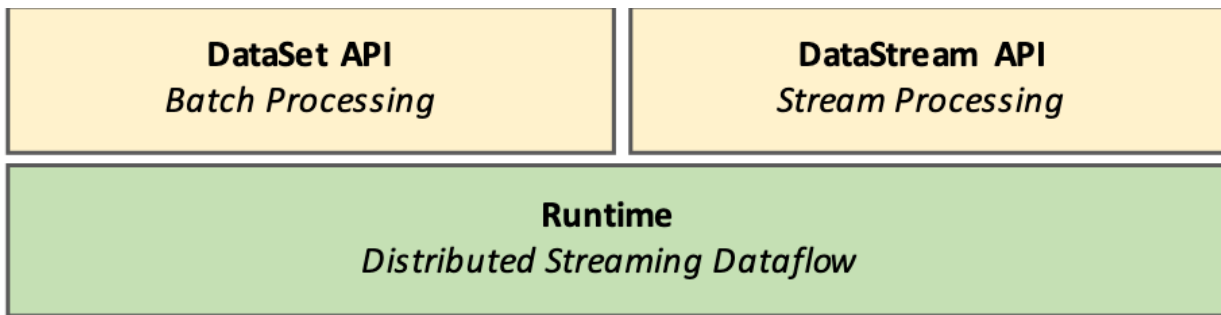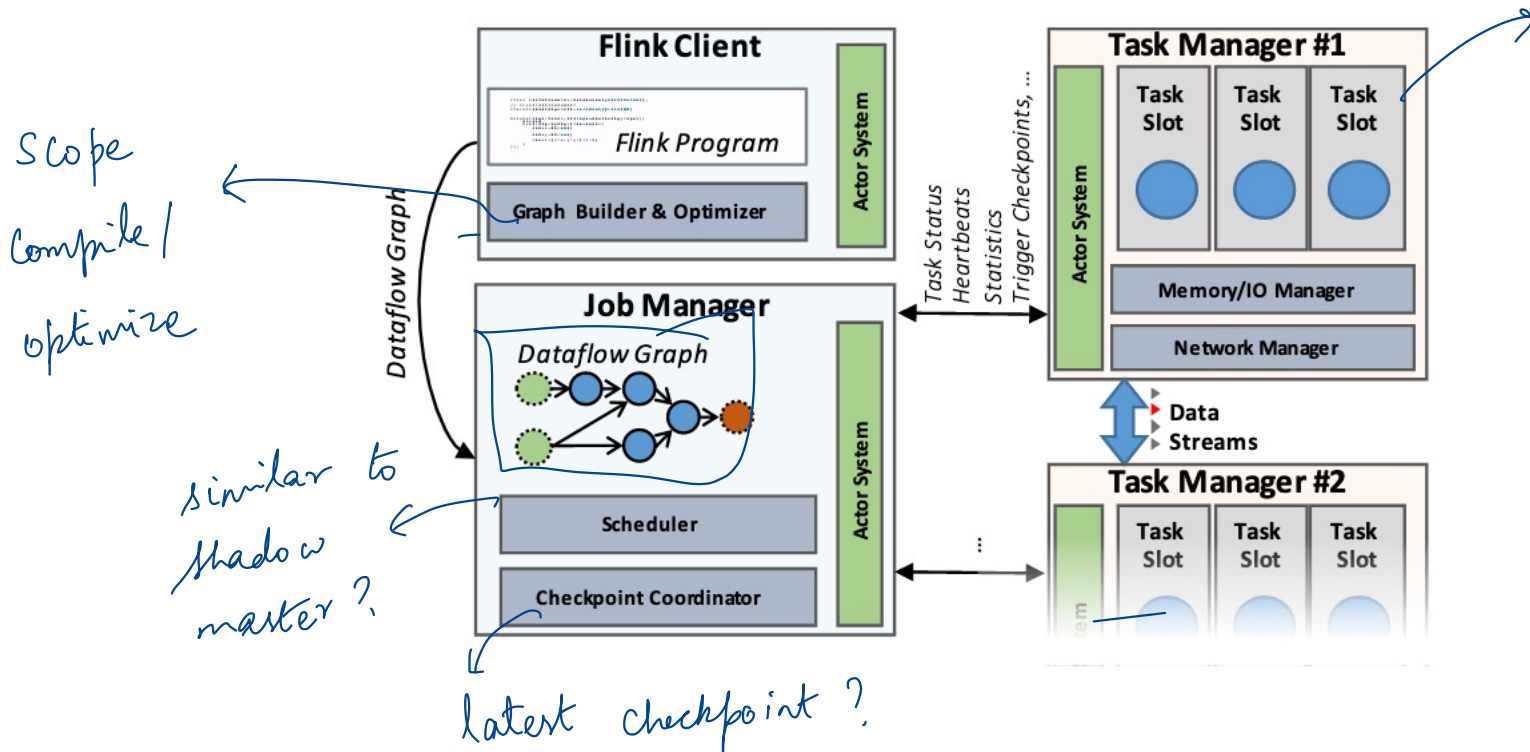
# COMBINING BATCH, STREAMING

Blocked DataStreams → *intermediate data to disk*

Turn off periodic snapshots

Blocking operators (e.g., sort)
↳ *batch specific operators*

| DataSet API | DataStream API |
|---|---|
| *Batch Processing* | *Stream Processing* |

| Runtime |
|---|
| *Distributed Streaming Dataflow* |

# OVERALL ARCHITECTURE



Flink Client
- Flink Program
- Graph Builder & Optimizer
- Actor System

Job Manager
- Dataflow Graph
- Scheduler
- Checkpoint Coordinator
- Actor System

Task Manager #1
- Task Slot
- Task Slot
- Task Slot
- Actor System
- Memory/IO Manager
- Network Manager

Task Manager #2
- Task Slot
- Task Slot
- Task Slot

Data Streams

Task Status
Heartbeats
Statistics
Trigger Checkpoints, ...

Dataflow Graph

Scope compile / optimize

similar to shadow master ?

latest checkpoint ?

# SUMMARY

Stream processing → Increasingly important workload trend

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

Distribtuted streaming dataflow
- Stateful operators
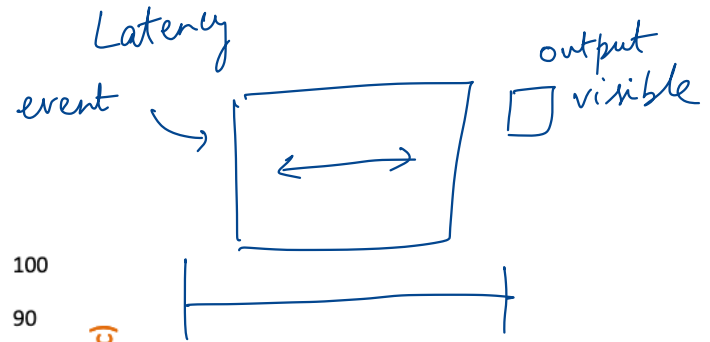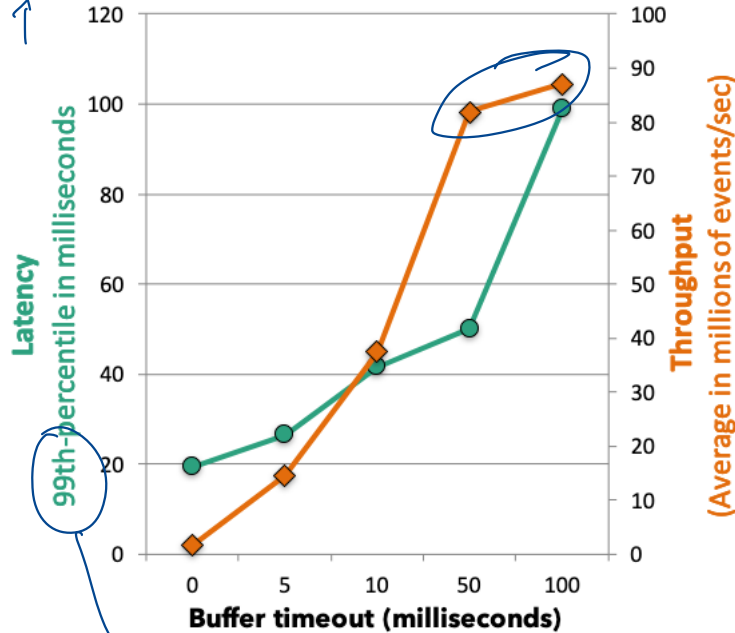- Checkpointing based FT

# DISCUSSION

https://forms.gle/j9Z7rm4qQpogbz5W8

① Timeout ↑
 - higher latency ↑
 - higher tput ↑

② Tput wins flatten
 out

 - You might
 buffer size
 limit before?



Latency
event →
output
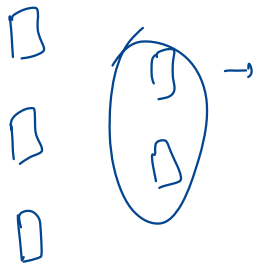☐ visible

worst case behavior

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?

→ Events to be processed fast
    → Tasks launched when prev. stage finishes ??

→ Windows ?
    → output for window → data arrives out of order ??

→ Data shuffling? → each batch is new data
       → state across batches? → memory ?

# SUMMARY

Next week: Spring break!!

Next class: Spark Streaming