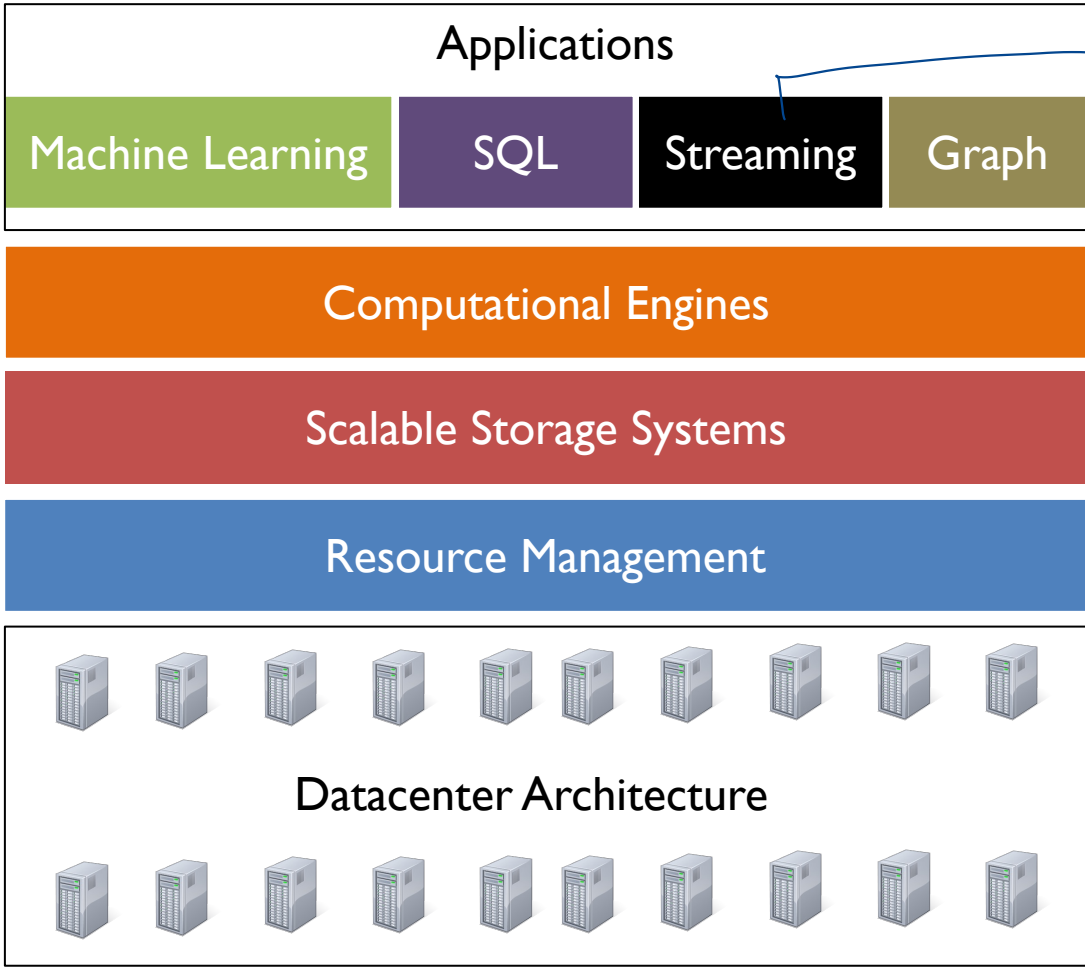Welcome back!

# CS 744: SPARK STREAMING

Shivaram Venkataraman

Spring 2024

# ADMINISTRIVIA

- Course Projects feedback → *Canvas*

- Midterm grades – this week?

- Cloudlab reservations → *Course Project*
  - Per-user from now
    ↳ *only you will be able to use that reservation*

Applications

Machine Learning | SQL | Streaming | Graph

Computational Engines

Scalable Storage Systems

Resource Management

Datacenter Architecture

Dataflow model
↳ how to express streaming queries

Apache Flink
↳ realizes continuous operator

# DASHBOARDS

data continuously over time. → low latency
→ out of order delivery



## 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)
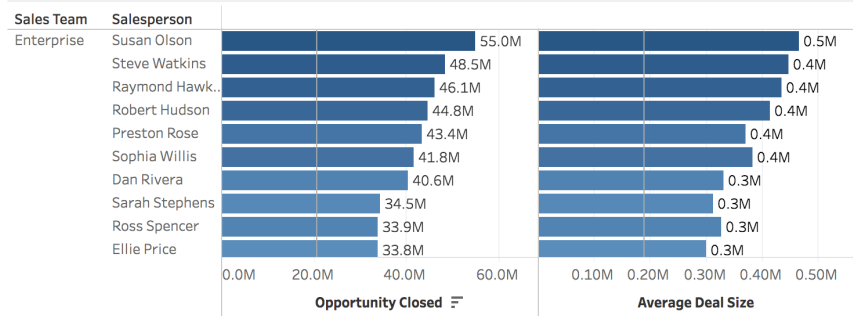
**Sales Team**
◉ (All)
○ Small and Midmarket
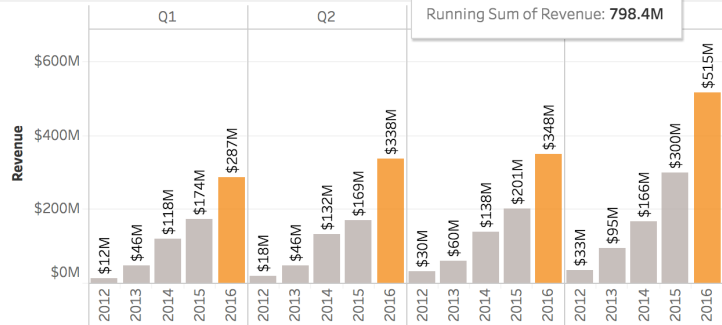○ Enterprise

**Avg Deal Size/Salesperson**
$147,043   $336,519

### Revenue Over Time

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

### Sales Team Performance

| Sales Team | Salesperson | Opportunity Closed | Average Deal Size |
|---|---|---|---|
| Enterprise | Susan Olson | 55.0M | 0.5M |
| | Steve Watkins | 48.5M | 0.4M |
| | Raymond Hawk.. | 46.1M | 0.4M |
| | Robert Hudson | 44.8M | 0.4M |
| | Preston Rose | 43.4M | 0.4M |
| | Sophia Willis | 41.8M | 0.4M |
| | Dan Rivera | 40.6M | 0.3M |
| | Sarah Stephens | 34.5M | 0.3M |
| | Ross Spencer | 33.9M | 0.3M |
| | Ellie Price | 33.8M | 0.3M |

### Revenue by Quarter

# CONTINUOUS OPERATOR MODEL

read log entries

group by userid

map

mutable state
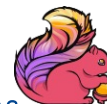
Long-lived operators

Mutable State

Distributed Checkpoints for Fault Recovery

Stragglers ?

Algorithm for checkpointing

Flink

$f^\infty$  Naiad
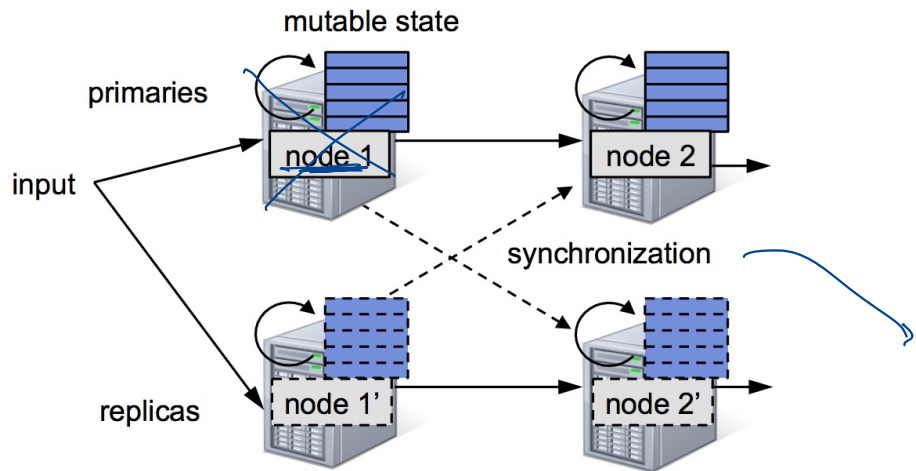
■ Driver  ⟶ Control Message

■ Task  ⤍ Network Transfer

# CONTINUOUS OPERATORS



mutable state

primaries

input

node 1     node 2

synchronization

replicas

node 1'     node 2'

replicate operators across nodes

↳ more resources to
support this scheme

same order of events
to the replica operator

↳ overhead during
processing

# SPARK STREAMING: GOALS

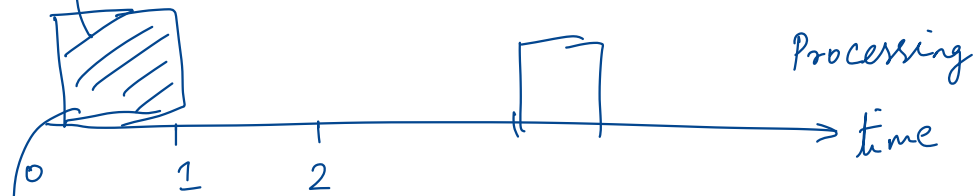1.  Scalability to hundreds of nodes → *high throughput*

2.  Minimal cost beyond base processing (no replication)

3.  Second-scale latency → *time between event arriving & it being reflected in the output*

4.  Second-scale recovery from faults and stragglers
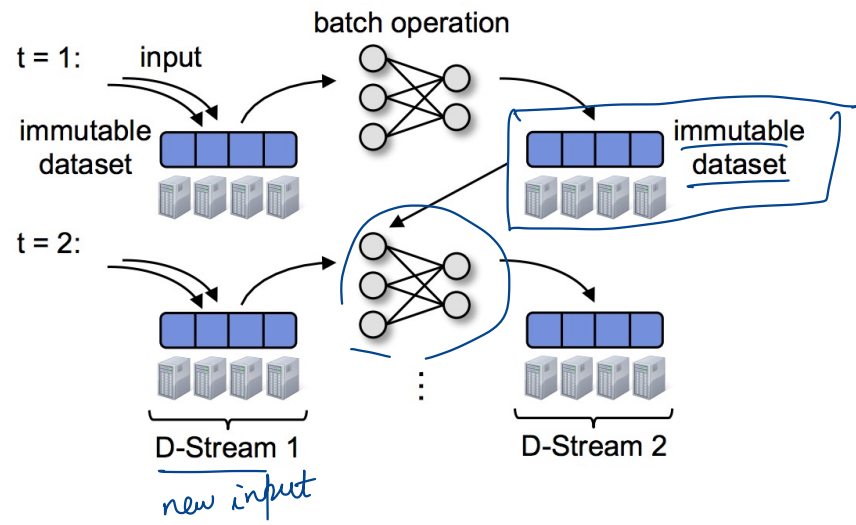
# DISCRETIZED STREAMS (DSTREAMS)

Could contain diff event times

Processing time

0    1    2

→ batch → submit this for computation

↳ re-using batch computation frameworks

→ Save some state at the end of each batch and use that as input for next batch

## batch operation

t = 1:    input

immutable dataset

immutable dataset

t = 2:

D-Stream 1      D-Stream 2

new input

# EXAMPLE

pageViews =
    readStream(http://...,
            "1s")

*read input*

*batch size*

ones = pageViews.map(
    event =>(event.url, 1))

*Key*

counts =
    ones.runningReduce(
        (a, b) => a + b)

*aggregate number of times URL occurs*

interval
[0, 1]

interval
[1, 2]



pageViews DStream    ones DStream    counts DStream

*shuffle*

*map*    *reduce*

# DSTREAM API

Transformations

    Stateless: map, reduce, groupBy, join

*similar to batch API*

*do not have dependencies across time steps*

Stateful:

    Sliding window("5s") → RDDs with data in [0,5), [1,6), [2,7)

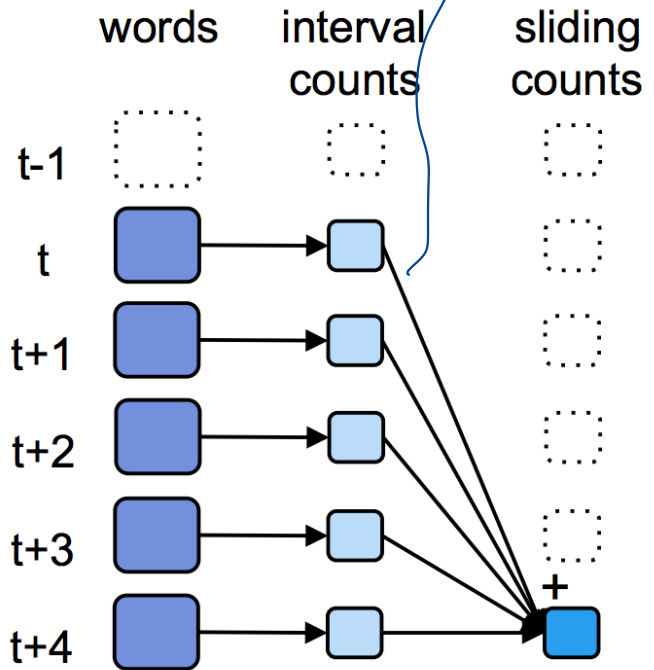    reduceByWindow("5s", (a, b) => a + b)

*creates a window & uses this reduction function*
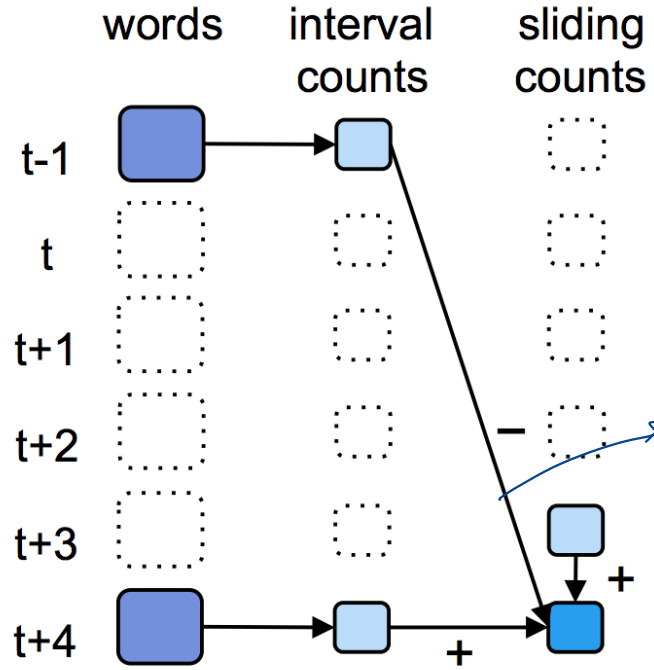
# SLIDING WINDOW



5 seconds    duration

depend on previous four

Add previous 5 each time

(a) Associative only

you have 3 dependencies vs 5

(b) Associative & invertible

# STATE MANAGEMENT

Tracking State: streams of (Key, Event) → (Key, State)

Similar in spirit to mutable state in flink

DStream

operator

```
events.track(
    (key, ev) => 1,

    (key, st, ev) => ev == Exit ? null : 1,

    "30s")
```

→ Initialize state

Given key, prev state, new event
↓
new state

→ Forget

key, event → what is initial value

↳ examples include event time range / session Id etc.

# STATE MANAGEMENT

Tracking State: streams of (Key, Event) → (Key, State)

```
events.track(
    (key, ev) => 1,

    (key, st, ev) => ev == Exit ? null : 1,

    "30s")
```
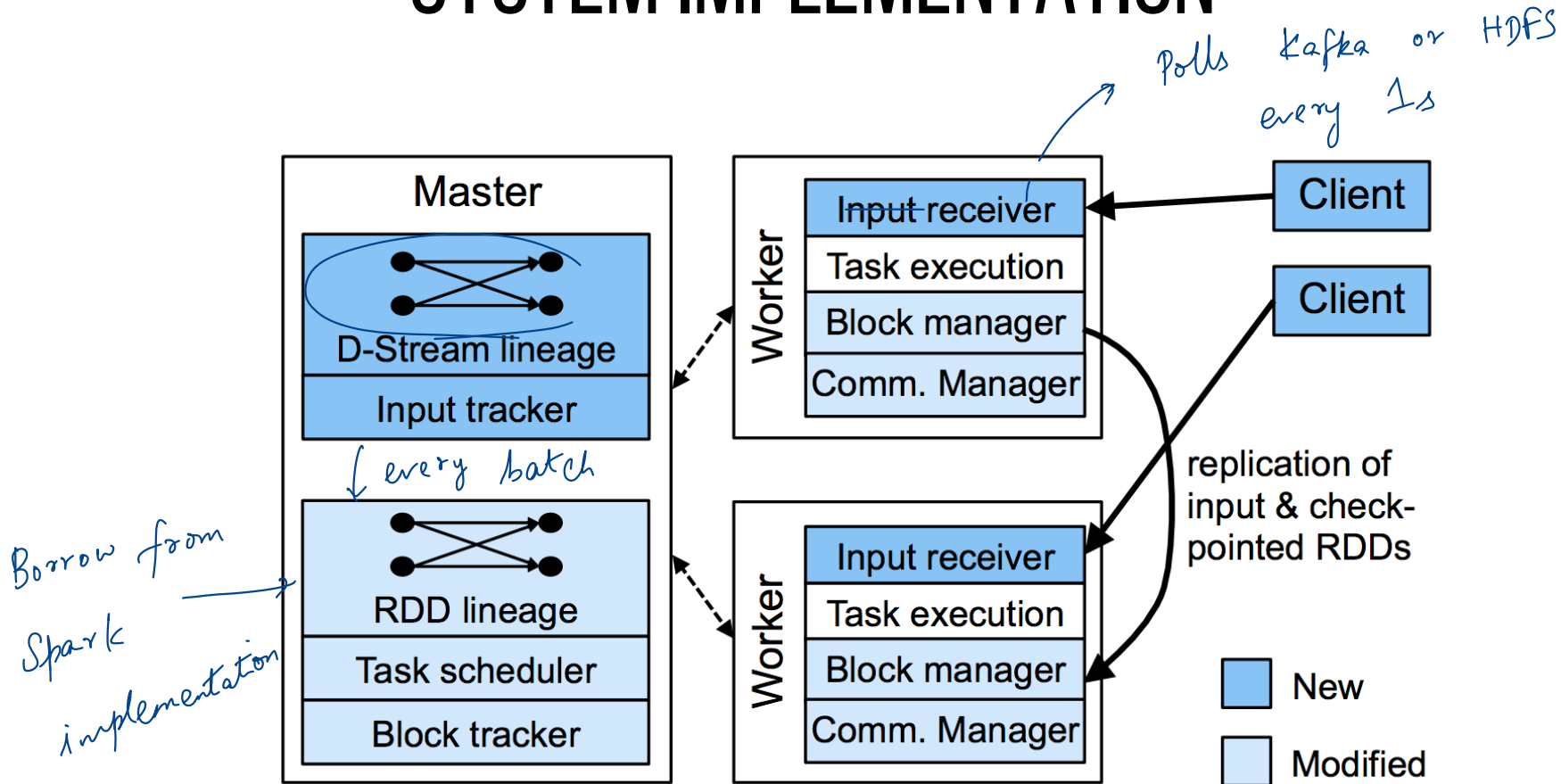
Computation requires
state !!

Stateless
operator

System
manages
state

vs.

Operator
maintains
state

System
helps with
checkpoints
etc.

# SYSTEM IMPLEMENTATION

# OPTIMIZATIONS

Timestep Pipelining

    No barrier across timesteps unless needed

    Tasks from the next timestep scheduled before current finishes

*start computation*

*$t = 2$ when $t = 1$*

*is still running*

Checkpointing

*simple to realize*

    Async I/O, as RDDs are immutable

    Truncate lineage after checkpoint

*lineage can grow infinitely*

*background*
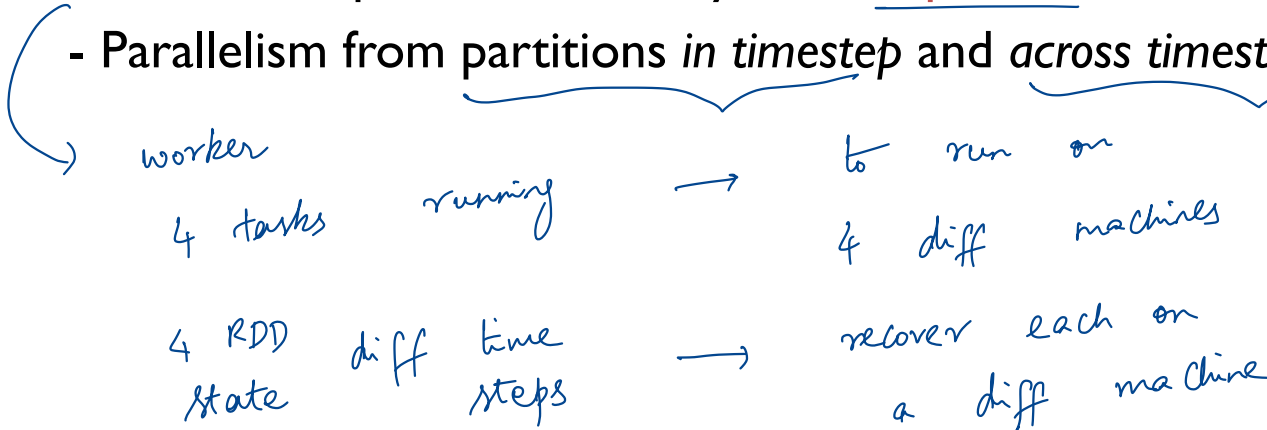
# FAULT TOLERANCE: PARALLEL RECOVERY

↳ *second scale*

Worker failure

- Need to recompute state RDDs stored on worker ⤳ *only need to replay tasks on this worker*
- Re-execute tasks running on the worker

Strategy

- Run all independent recovery tasks in parallel
- Parallelism from partitions *in timestep* and *across timesteps*

⤳ *worker*

*4 tasks running → to run on 4 diff machines*

*4 RDD state diff time steps → recover each on a diff machine*
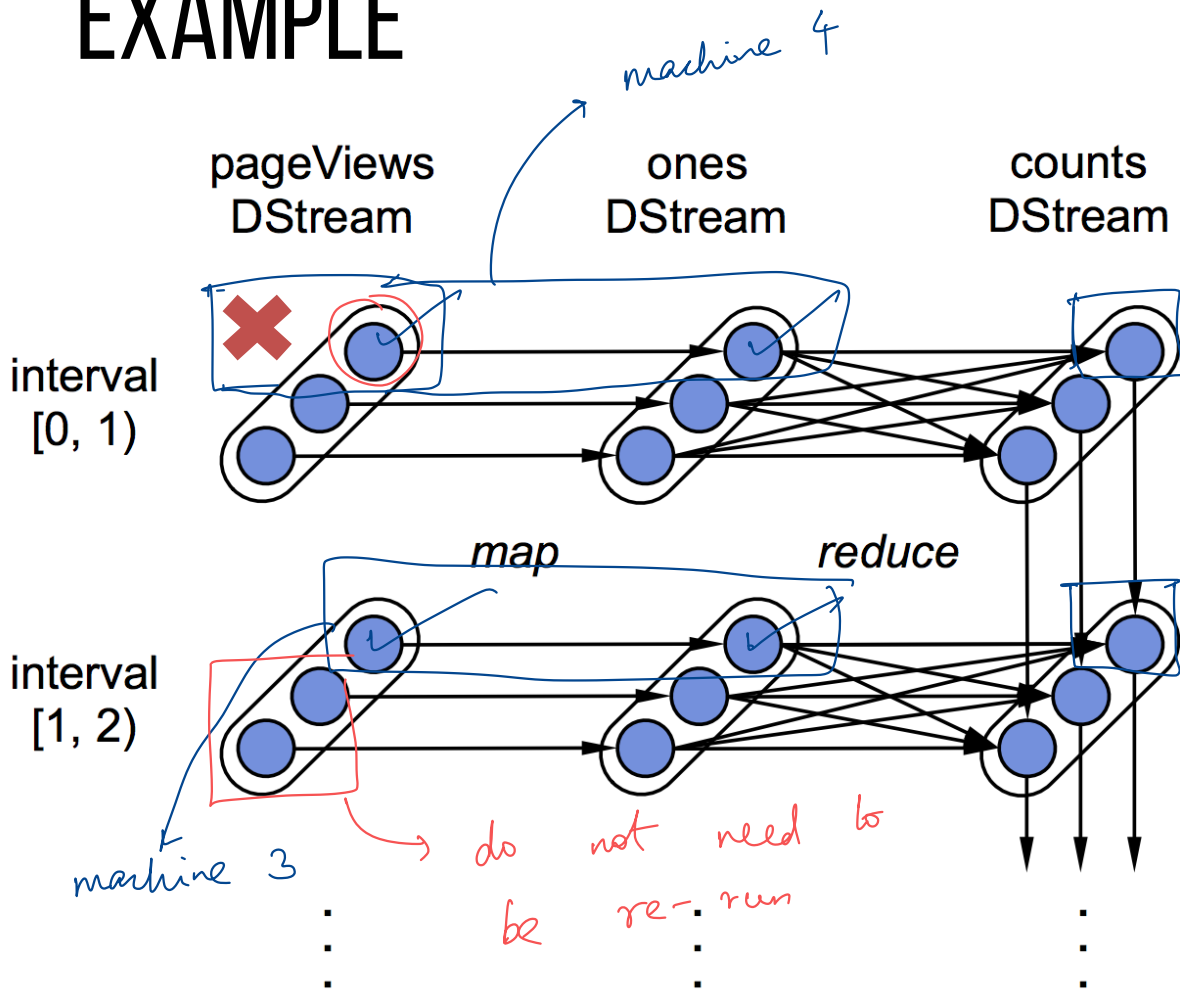
# EXAMPLE

```
pageViews =
    readStream(http://...,
                "1s")


ones = pageViews.map(
    event =>(event.url, 1))


counts =
    ones.runningReduce(
        (a, b) => a + b)
```

# FAULT TOLERANCE

replayable    input
↳ Kafka

Straggler Mitigation: Use speculative execution

→ operators are stateless multiple of them same time

→ the checkpoint driver state

## Driver Recovery
- At each timestep, save graph of DStreams and Scala function objects
- Workers connect to a new driver and report their RDD partitions
- Note: No problem if a given RDD is computed twice (determinism).

→ similar to GFS recovery

# SUMMARY

Micro-batches: New approach to stream processing

Simplifies fault tolerance, straggler mitigation

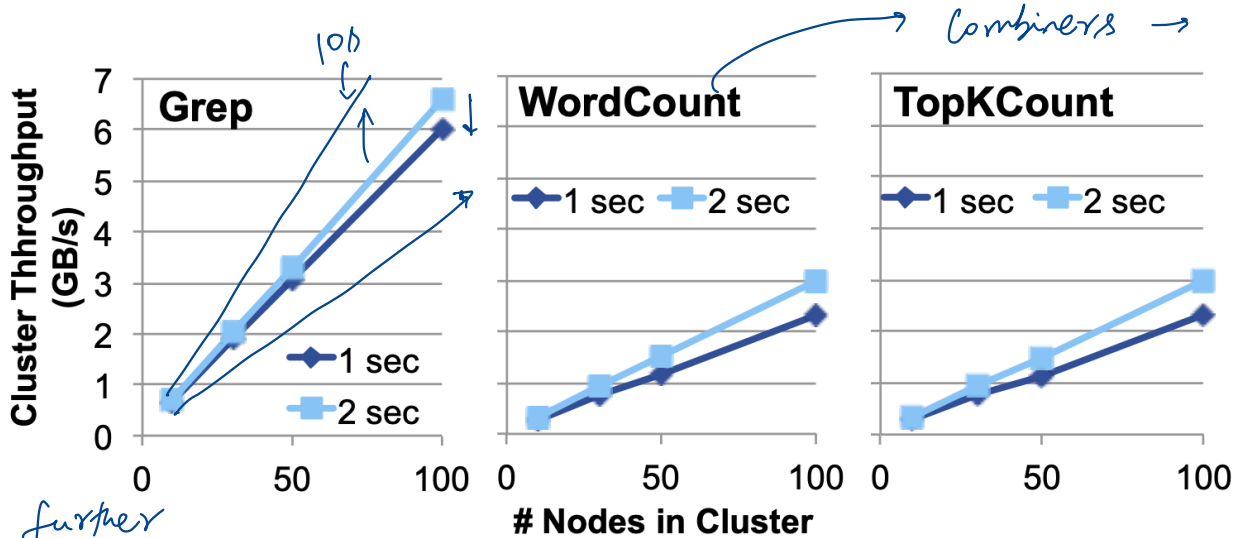Unifying batch, streaming analytics

       ↳ Share Code

# DISCUSSION

https://forms.gle/RVtChgDQzbX16tqT7

If the latency bound was made to 100ms, how do you think the above figure would change? What could be the reasons for it?

What about 10s?

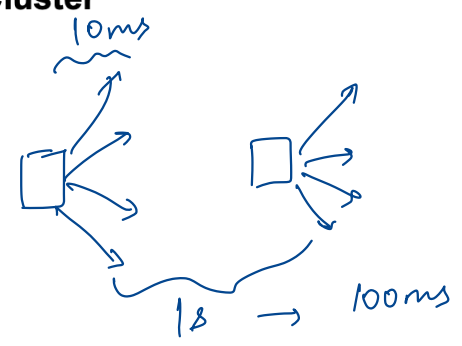Combiners → not as effective at 100 ms



10s
↳ higher tput

Capped at some limit resource limits of cluster.

① 100ms
tput will further go down?
↳ "overheads" will be higher
↳ scheduling, task launch
→ input reading

Consider the pros and cons of approaches in Flink vs Spark Streaming. What application properties would you use to decide which system to choose?

Flink
→ checkpoints are more expensive
    ↳ unreliable hardware
        ↳ don't want to use Flink
    ↳ cluster is small
      → FT is less of concern?

Spark Streaming
    ↳ low latency
      < second scale
  then not Spark Streaming

→ streaming join with historical data spark has advantages

# NEXT STEPS

Next class: Graph processing!

Midterm grades soon!