

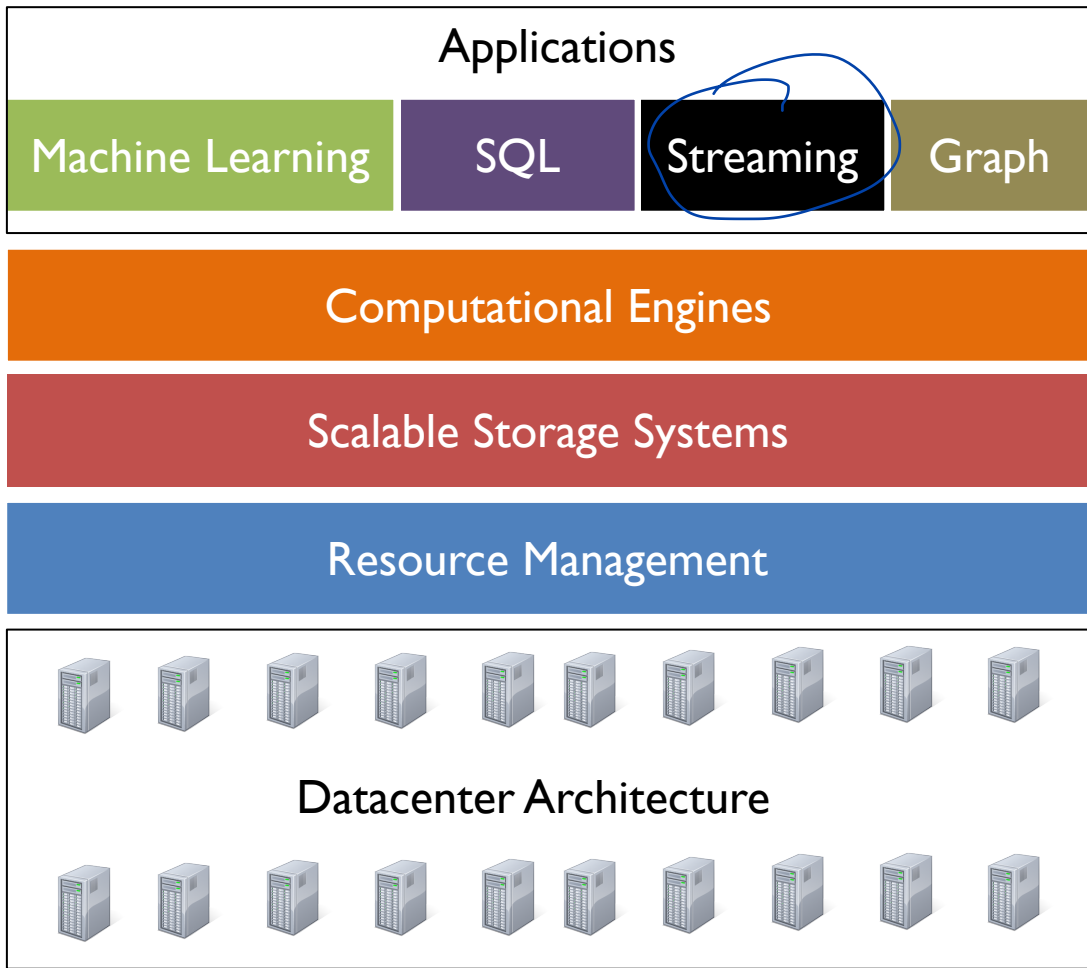
# CS 744: SPARK STREAMING

Shivaram Venkataraman

Fall 2019

# ADMINISTRIVIA

- Midterm grades this week
- Course Projects sign up for meetings → Nov 21, 25



Dataflow  
model

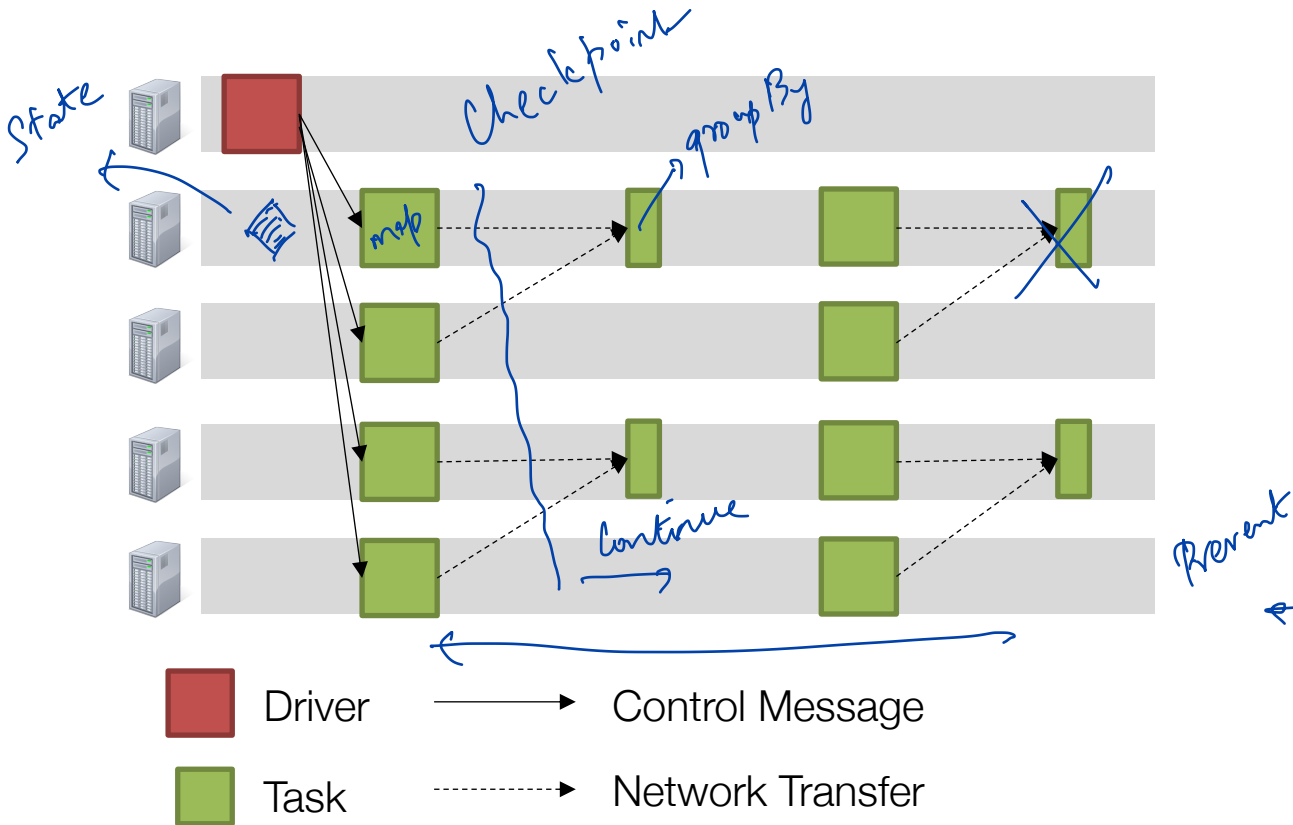
→ Queries

properties

of

Data, Queries

# CONTINUOUS OPERATOR MODEL

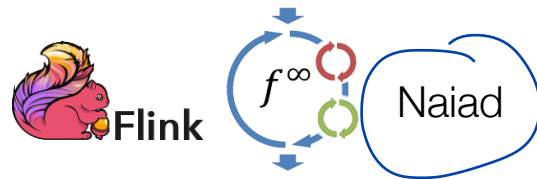


## Long-lived operators

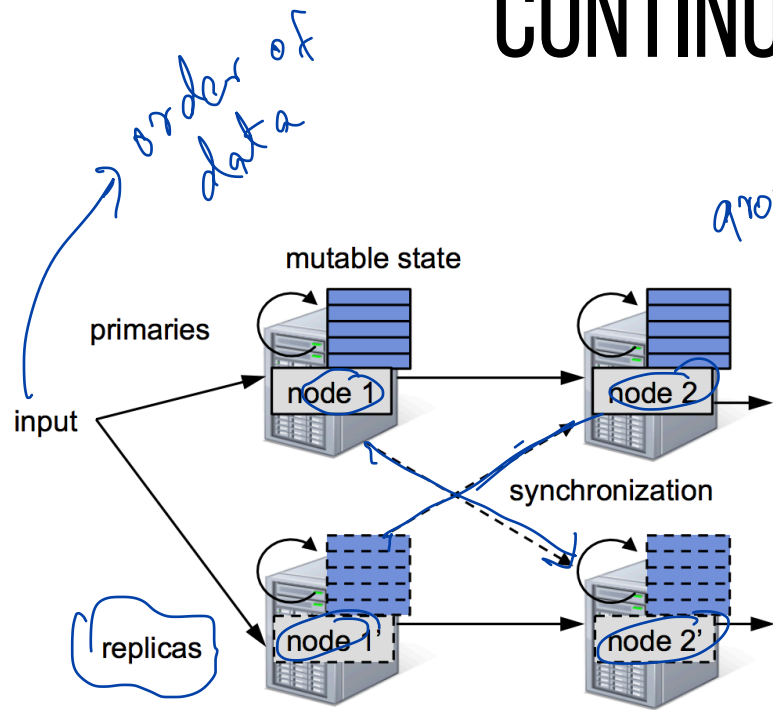
# Mutable State

# Distributed Checkpoints for Fault Recovery

## Stragglers ?



# CONTINUOUS OPERATORS



group by (country)

Country	Pageview

Country	Pageview

mutable state

differences

across  
replicas

resource / memory  
overhead

strict ordering

↳ Synchronization  
might not mitigate stragglers

upstream backup

recovery  
is  
very  
fast

# SPARK STREAMING: GOALS

1. Scalability to hundreds of nodes

$\Leftrightarrow 2\times$  overhead

2. Minimal cost beyond base processing (no replication)

3. Second-scale latency  $\rightarrow$  *revisit*

4. Second-scale recovery from faults and stragglers

# DISCRETIZED STREAMS (DSTREAMS)

Short, deterministic tasks

→ Stateless operation

→ replicate / recover operation

Low overhead

→ run ops in parallel

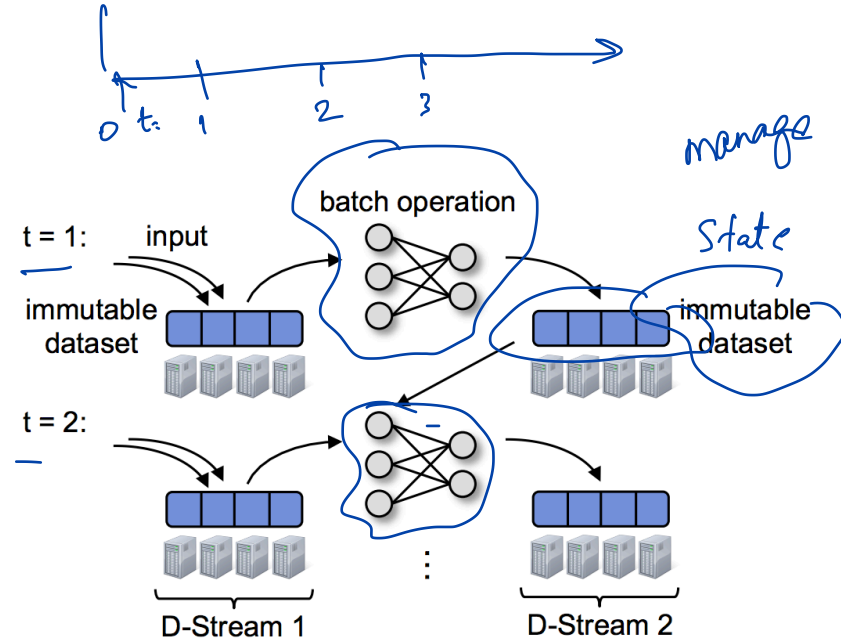
→ each batch can be run in parallel

→ recovery or straggler mitigation in parallel

from  $t=0$

$t=1$

$t=2$

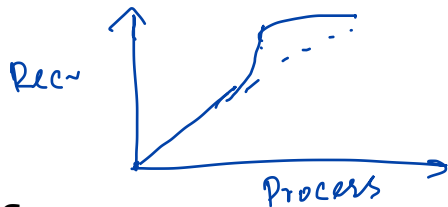


Low latency recovery

Dependencies within / across timestamp

# EXAMPLE

similar to RDD



pageViews =

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

batch  
duration

"1s"

who gives  
TS?  
Processing  
Time

interval  
[0, 1)

pageViews  
DStream

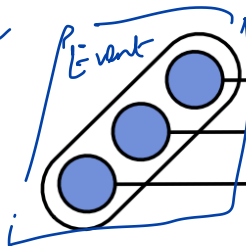
ones

DStream

counts  
DStream

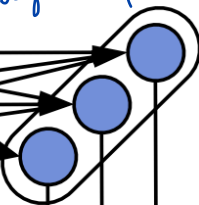
google.com, 5

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



map

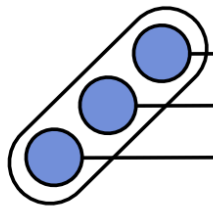
reduce



counts =

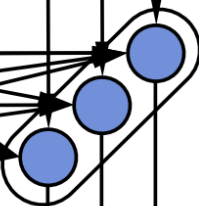
`ones.runningReduce(  
(a, b) => a + b)`

interval  
[1, 2)



map

reduce



(google.com, 100)

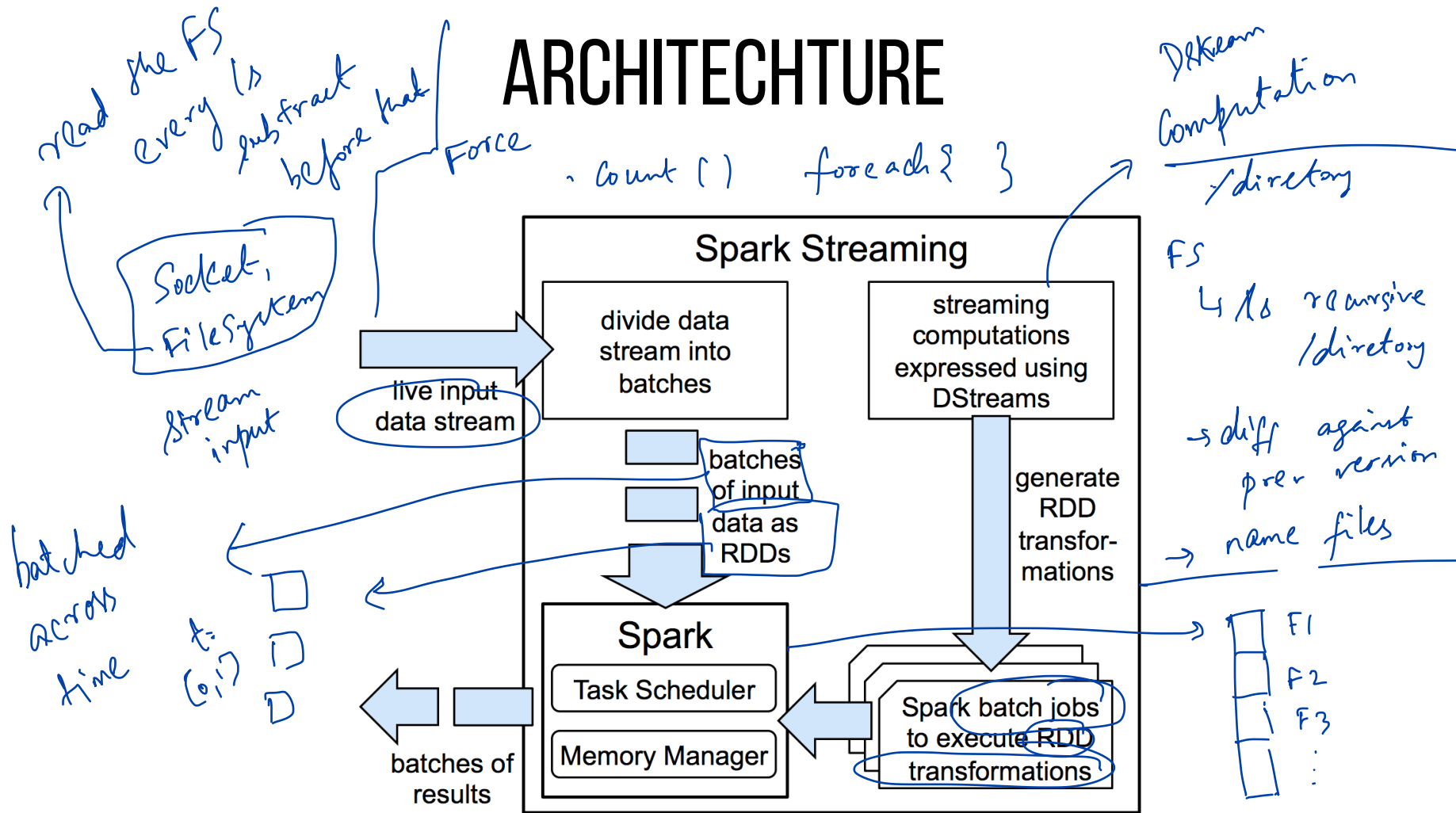
⋮

⋮

⋮



# ARCHITECTURE



# DSTREAM API

Transformations  $\rightarrow$  on DStream

Stateless: map, reduce, groupBy, join

map  
Stream  $\rightarrow$  Stream  
1:1 mapping

Stateful:

window("5s")  $\rightarrow$  RDDs with data in [0,5), [1,6), [2,7)

Sliding window

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

$\downarrow$   
operate across windows

not available on RDDs

reduce by window (5s, a+b)

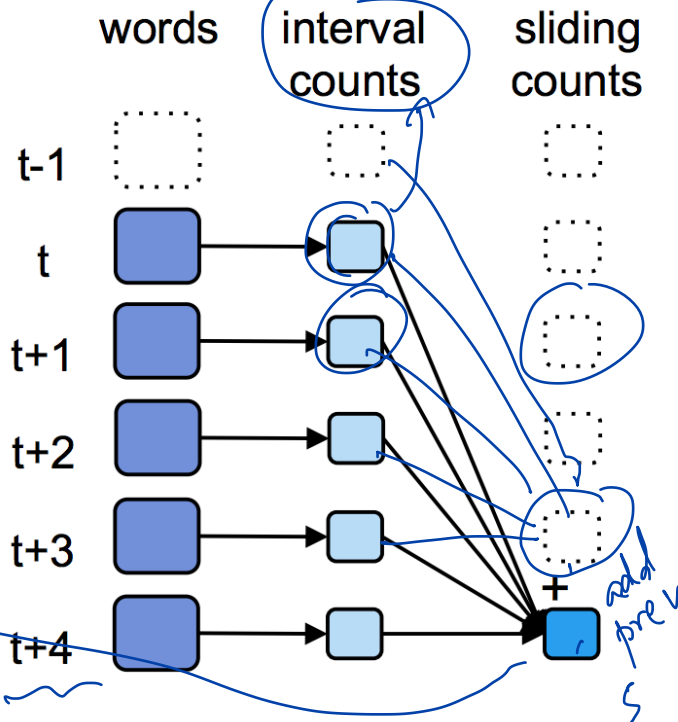
# SLIDING WINDOW

State

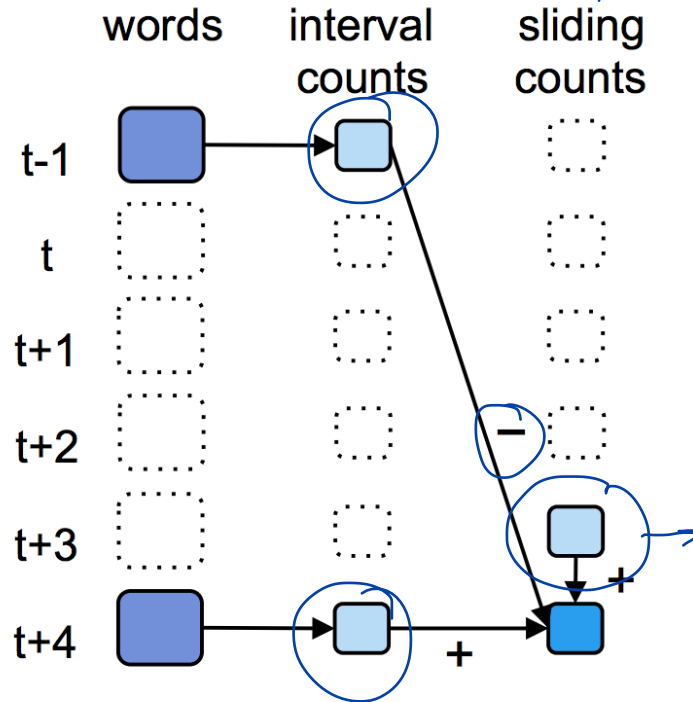
output

2 operation  
5s  
5 operations

Add  
previous 5  
each time



(a) Associative only



(b) Associative & invertible

# STATE MANAGEMENT

Tracking State: streams of (Key, Event)  $\rightarrow$  (Key, State)

events.track(  
 (key, ev) => 1,  $\rightarrow$  Initialize

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

"30s")

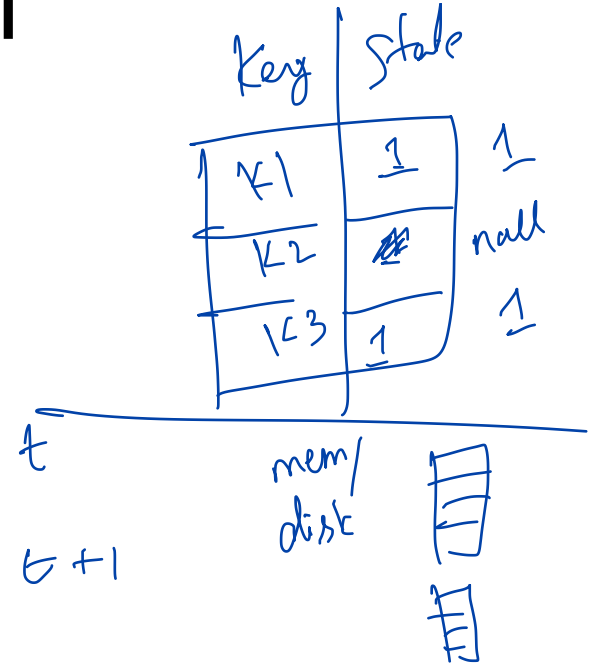
forget  
state

new event-  
old state

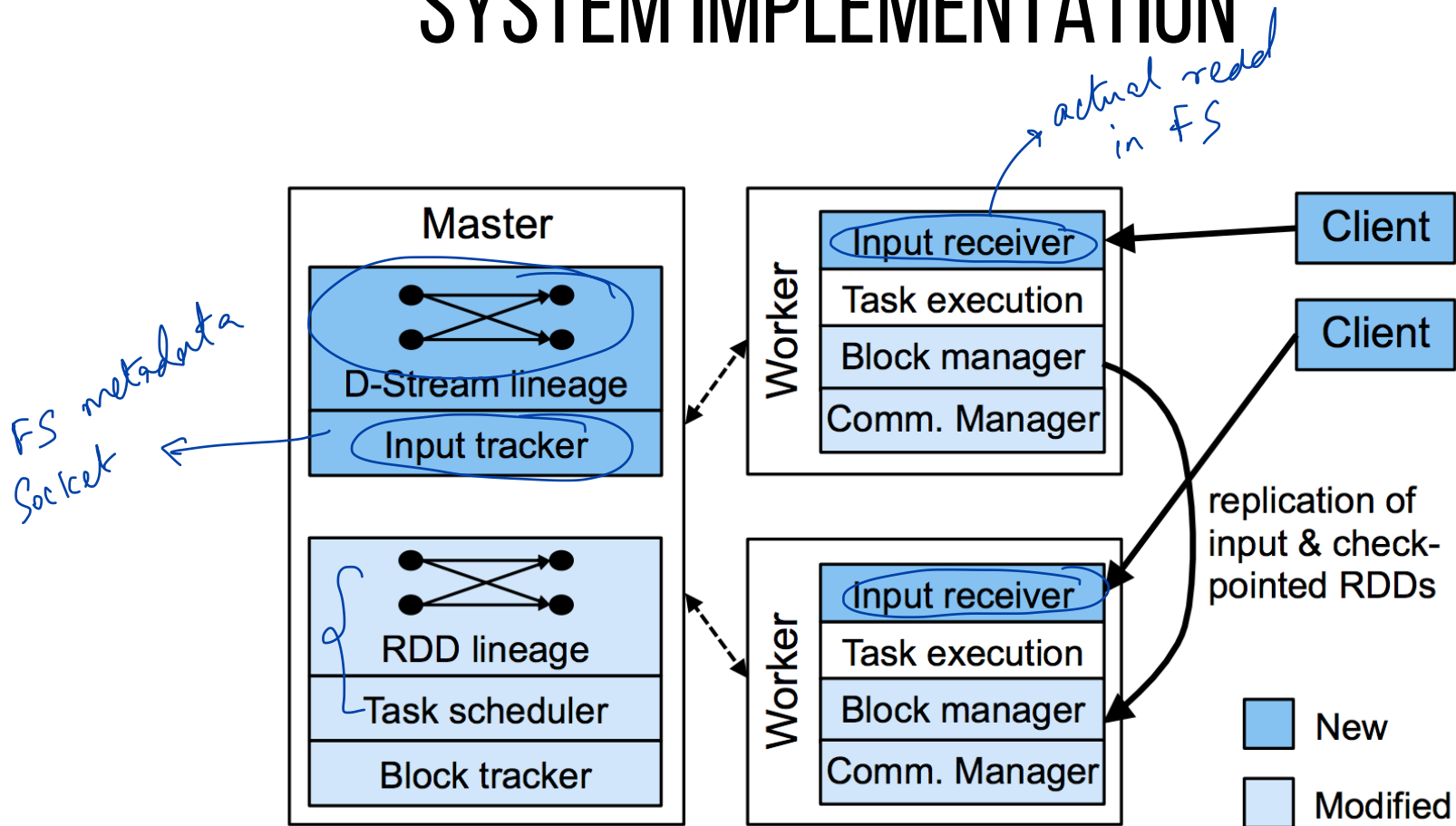
If event is  
exit

state

key, old state, event  
 $\rightarrow$  new state



# SYSTEM IMPLEMENTATION



# OPTIMIZATIONS

## Timestep Pipelining

No barrier across timesteps unless needed

Tasks from the next timestep scheduled before current finishes

running {  $t \in [0, 1]$   $\rightarrow$  Submit Job (rdd1, )

$t \in [1, 2]$   $\rightarrow$  submit Job (

[Checkpointing]

Async I/O, as RDDs are immutable

Forget lineage after checkpoint

State  
window

lineage chain  
very long

$\downarrow$   
map  $t \in [1, 2]$

# FAULT TOLERANCE: PARALLEL RECOVERY

## Worker failure

- Need to recompute state RDDs stored on worker
- Re-execute tasks running on the worker

## Strategy

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

# EXAMPLE

parallelism available across timesteps

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

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

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

Input  
replayed  
FS

interval  
[0, 1)

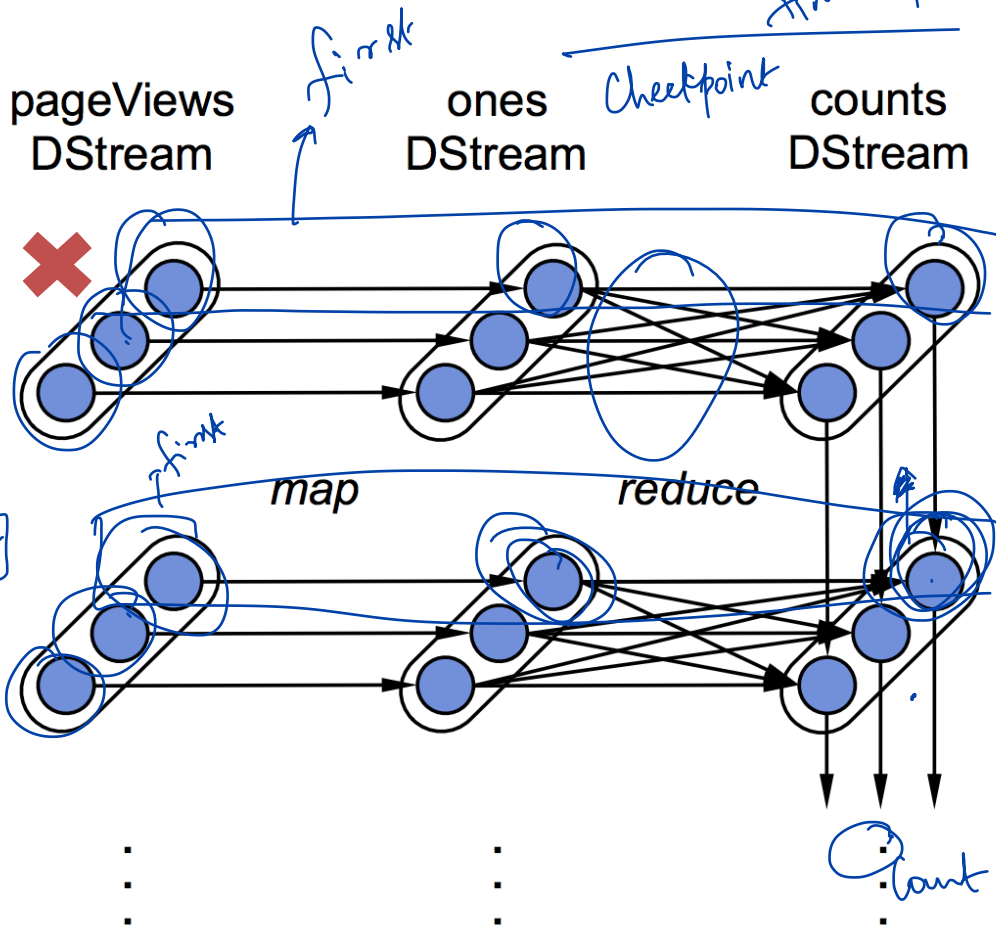
[0, 2]

interval  
[1, 2)

Associative



[2, 3)





# FAULT TOLERANCE

## Straggler Mitigation

Use speculative execution

Task runs more than 1.4x longer than median task → straggler

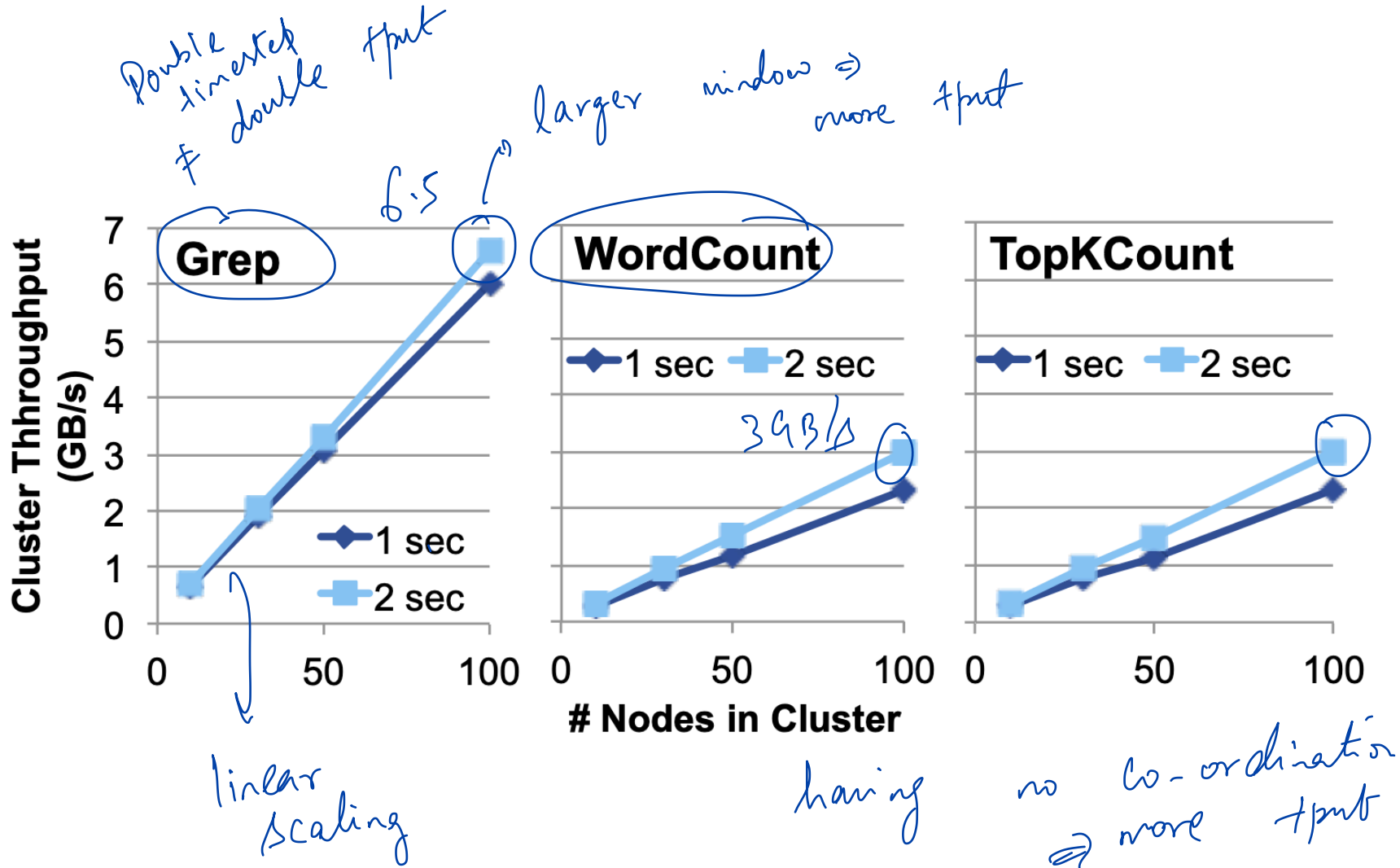
## Master Recovery

→ NEW!

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

# DISCUSSION

<https://forms.gle/xUvzCIbdV7H48mTM8>



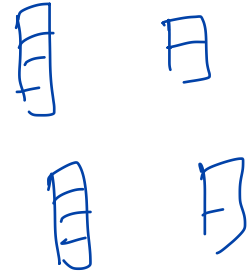
If the latency bound was made to 100ms, how do you think the above figure would change? What could be the reasons for it?

→ Throughput might decrease

Combiner  
Example

↳ Overhead of creating RDDs

↳ Take more time to process  
than 100ms?



→ linear scaling might not work → Lots of tasks  
→ Network



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

Naiad - low latency  
- incremental output

\* Applications that might  
have fewer stragglers

Spark Streaming

lower latency  
recovery

API simple

# NEXT STEPS

Next class: Graph processing

Sign up for project check-ins!

# SHORTCOMINGS?

## Expressiveness

- Current API requires users to “think” in micro-batches

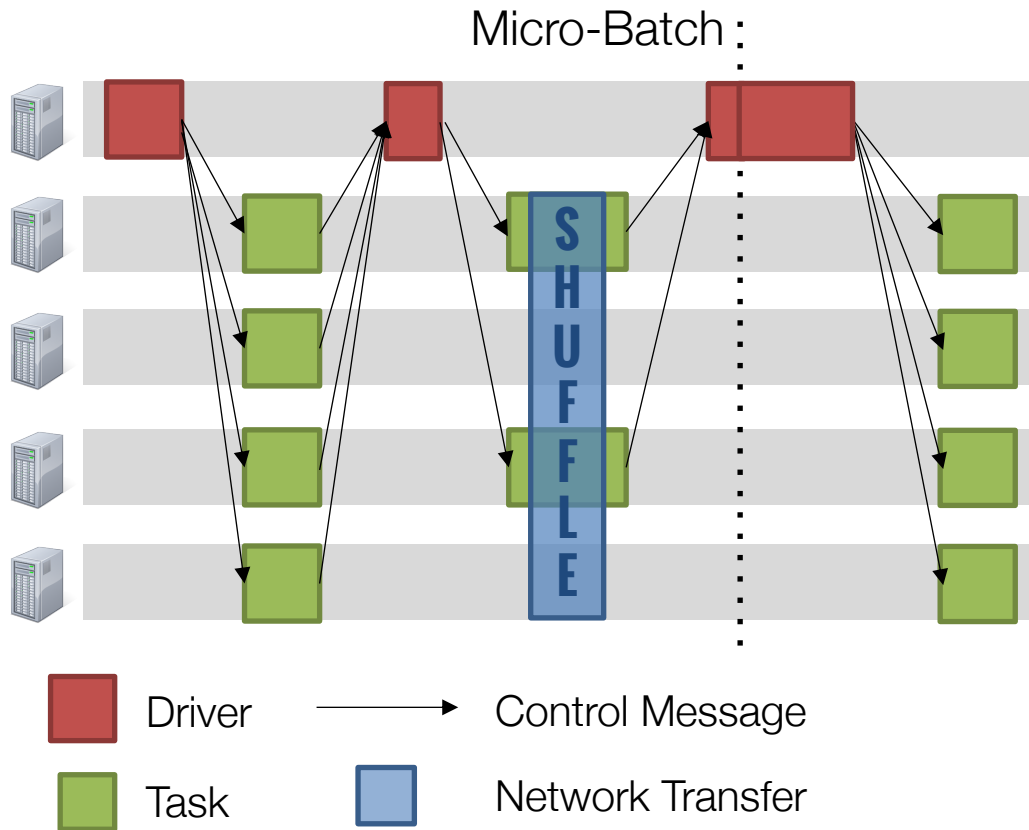
## Setting batch interval

- Manual tuning. Higher batch → better throughput but worse latency

## Memory usage

- LRU cache stores state RDDs in memory

# COMPUTATION MODEL: MICRO-BATCHES





# SUMMARY

Micro-batches: New approach to stream processing

Higher latency for fault tolerance, straggler mitigation

Unifying batch, streaming analytics