

Hi

CS 744: SPARK STREAMING

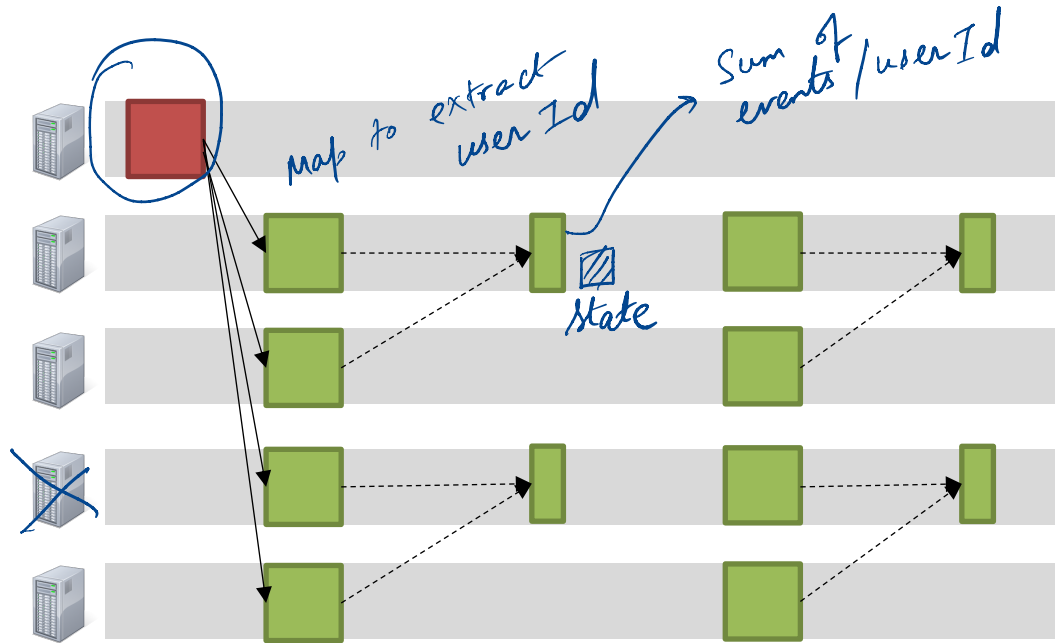
Shivaram Venkataraman

Fall 2021

ADMINISTRIVIA

- Midterm grades this week → Friday
 - Course Projects feedback → TODAY
 - Google Cloud Credits
 - ↳ \$50 / student email address
- Private Piazza / e-mail

CONTINUOUS OPERATOR MODEL



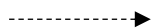
Driver



Control Message



Task



Network Transfer

✓ Long-lived operators

✓ Mutable State

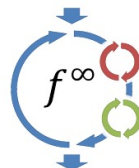
Distributed Checkpoints
for Fault Recovery

Stragglers ?

Avoid
stragglers
with good
engineering

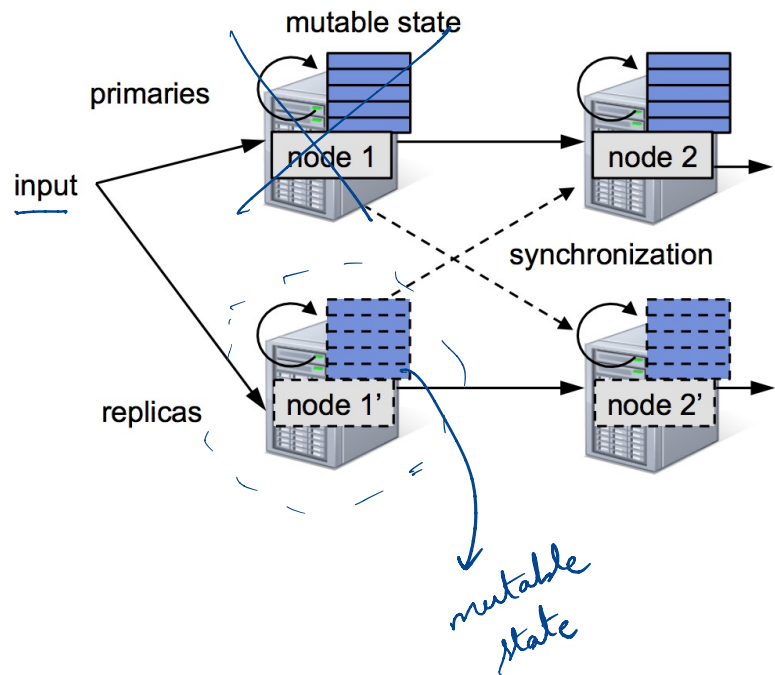


Flink



Naiad

CONTINUOUS OPERATORS



- Replicate every operation to another machine

- Minimizes recovery time

- 2x the resources

- Replicas remain in-sync

- replication protocol which also adds overhead

SPARK STREAMING: GOALS

1. Scalability to hundreds of nodes → high throughput
2. Minimal cost beyond base processing (no replication) → resource efficiency
3. Second-scale latency = time between input arriving to when it is part of the output
4. Second-scale recovery from faults and stragglers

DISCRETIZED STREAMS (DSTREAMS)

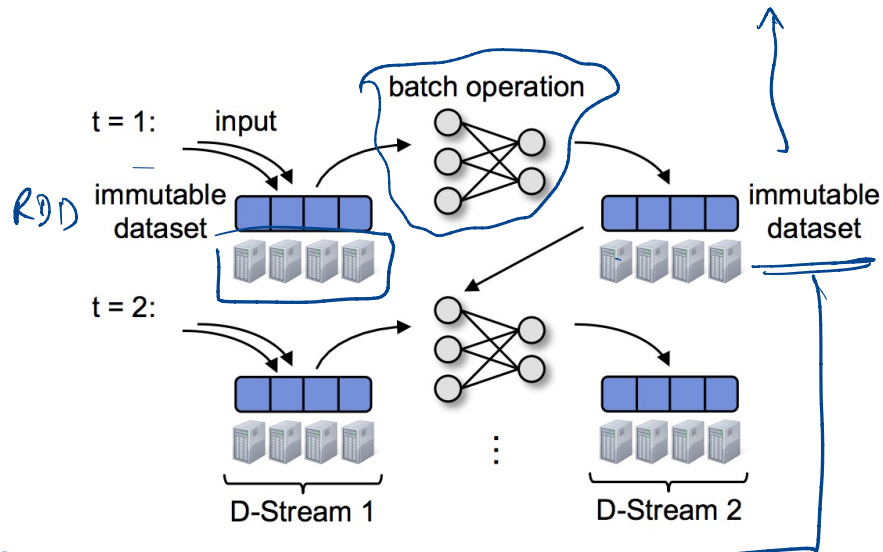
- Divide time into a number of micro-batches

Memory / Disk

- 1s as micro-batch

$t = 1s$ to $2s \leftarrow$ all events

- Run the batch operation on the input events, and save output (operator state) also as an RDD
- Next micro-batch, use the output from previous and compute



deterministic

→ recompute it, lineage

with this be the same

t = 1-2s



Compute

~~RDD~~

RDD

t = 2-3s



~~RDD~~



Processing time stamps / not event timestamps

EXAMPLE

DStream

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

read a stream
and not a file

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

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

running sum of
url events

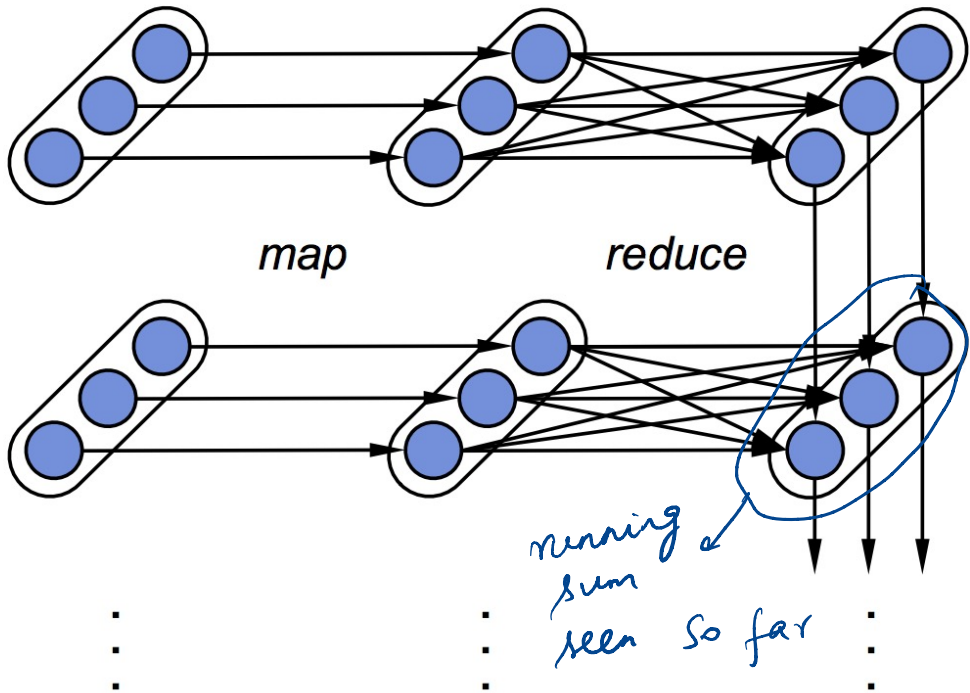
pageViews
DStream

ones
DStream

counts
DStream

interval
[0, 1)

interval
[1, 2)



DSTREAM API

Transformations

Stateless: map, reduce, groupBy, join → very similar to RDD API

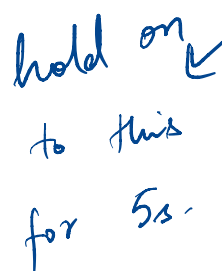
Stateful:

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

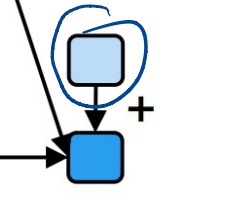
reduceByWindow("5s", (a, b) => a + b)

↳ form a window of 5s
and reduce using sum

Computation is deeper



Computation is deeper



(b) Associative & invertible

STATE MANAGEMENT

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

— User defined state object for every key

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

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

"30s")

old state

new event

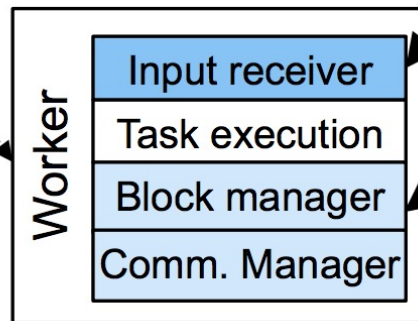
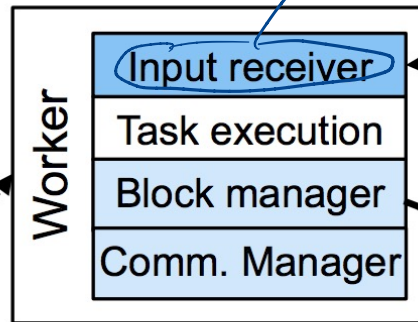
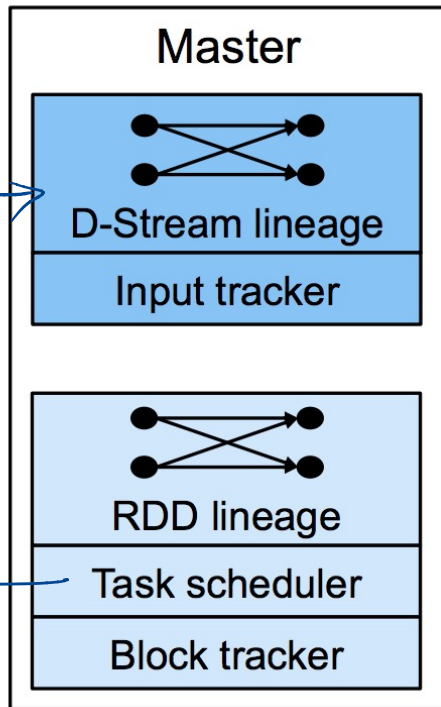
timeout to delete old states

update my state based on event that arrived

SYSTEM IMPLEMENTATION

Tracks D-Streams relate to each other

Exist in spark



poll http server read from FS/storage system or

Client

Client

replication of input & checkpointed RDDs

RDD in memory of 2 machines

New

Modified

OPTIMIZATIONS

Timestep Pipelining

No barrier across timesteps unless needed

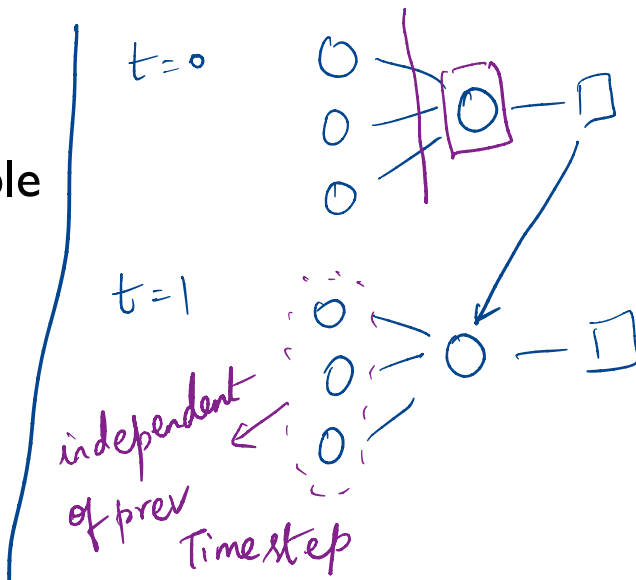
Tasks from the next timestep scheduled before current finishes

Checkpointing

Async I/O, as RDDs are immutable

Forget lineage after checkpoint

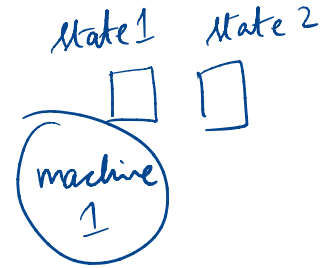
↓
stored in HDFS



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*

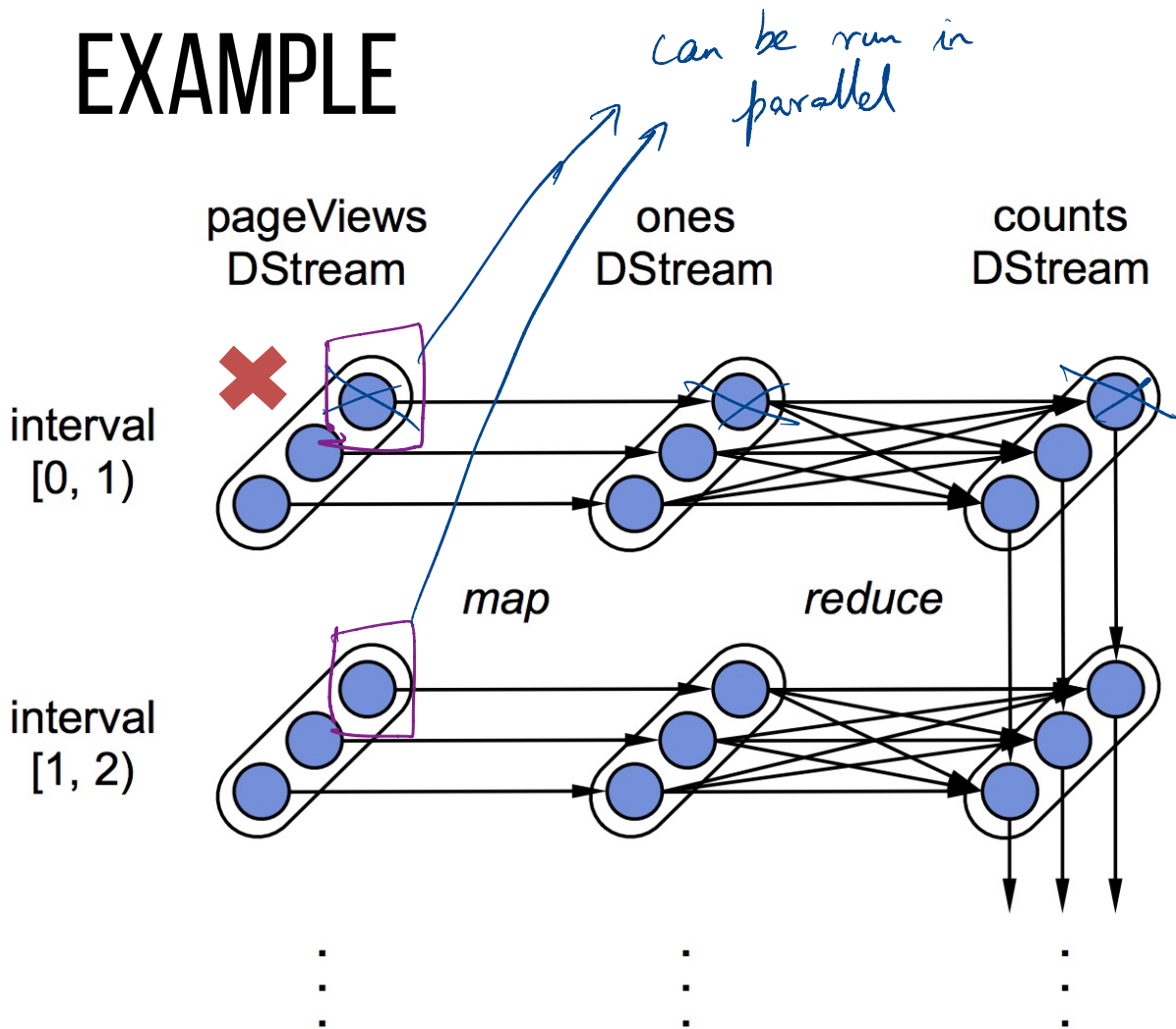
Fault recovery
can run
fast

Note

- Only recompute state / tasks which were lost

EXAMPLE

```
pageViews =  
  readStream(http://...,  
    "1s")  
  
ones = pageViews.map(  
  event =>(event.url, 1))  
  
counts =  
  ones.runningReduce(  
    (a, b) => a + b)
```



FAULT TOLERANCE

Straggler Mitigation

Use speculative execution

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

Master Recovery → *Computation runs 24x7*

- 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).

→ *similar in spirit to HFS recovery*

SUMMARY

Micro-batches: New approach to stream processing

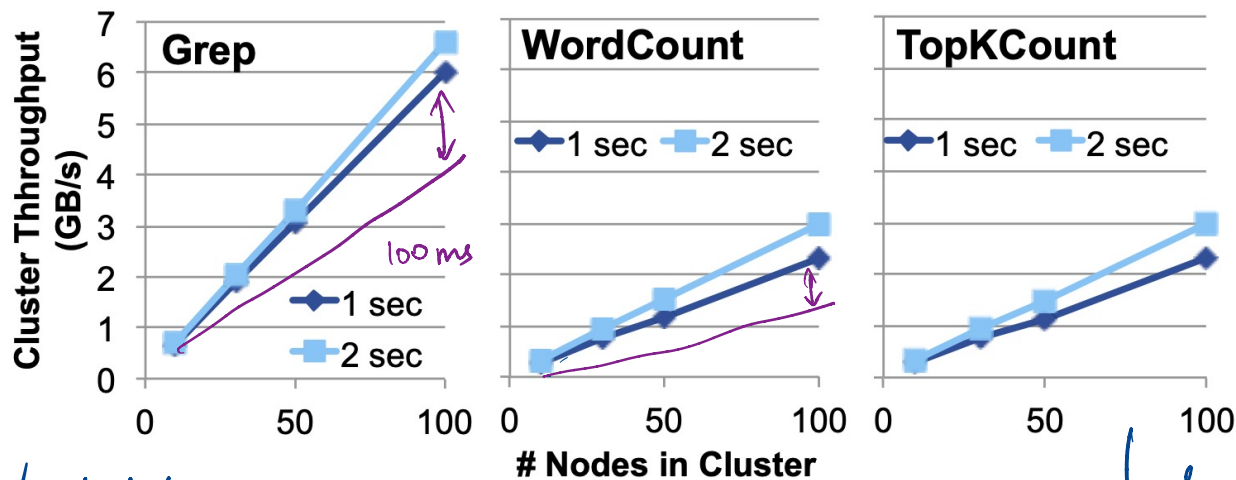
Simplifies fault tolerance, straggler mitigation

Unifying batch, streaming analytics

DISCUSSION

<https://forms.gle/4Xbu9y9KTW5qph8H8>

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



Overhead increases with cluster size

Every microbatch

= Fixed overhead + Time to do Grep / Word Count

~

Bigger factor

Small microbatch
 ⇒ lots of RDDs
 ↳ more metadata
 ↳ tracking / scheduling tasks

2s
 100ms

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 / Flink

- low latency / quick processing

Spark Streaming

- fault recovery
time

NEXT STEPS

Next class: Graph processing!

Midterm grades ASAP!