

Good morning

CS 744: SPARK STREAMING

Shivaram Venkataraman

Fall 2020

ADMINISTRIVIA

- Midterm grades this week → ASAP
- Course Projects feedback → Hot CRP

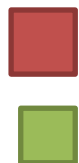
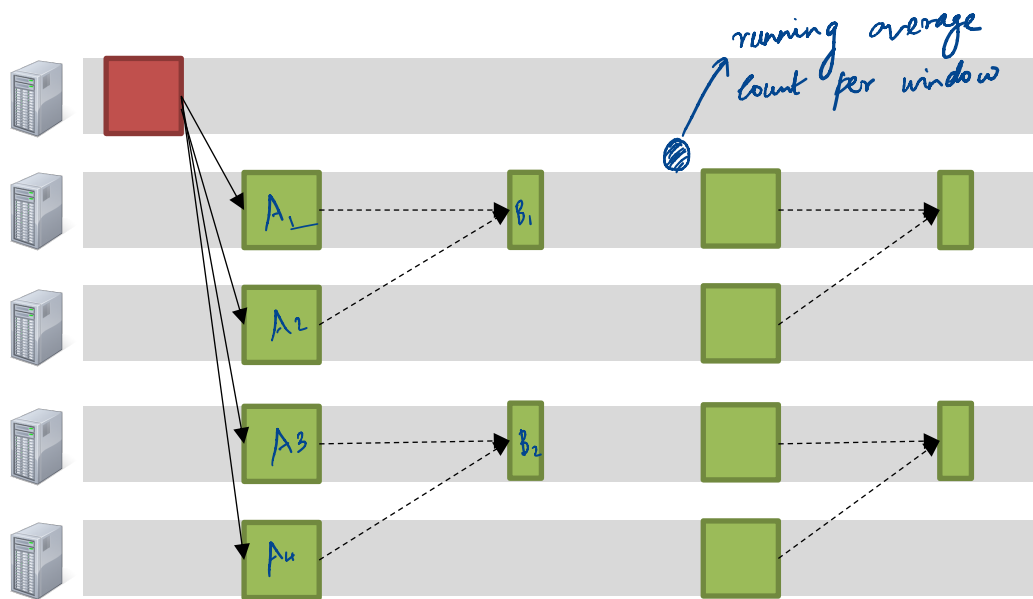


Hopefully you are working on this!

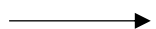
- (1) Assign grades for project proposals
- (2) mid semester update → Nov 20th

A Hg B

CONTINUOUS OPERATOR MODEL



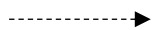
Driver



Control Message



Task



Network Transfer

Long-lived operators

Mutable State

Distributed Checkpoints
for Fault Recovery

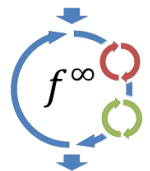
Stragglers ?

Roll back
all
operators
to
checkpoint

Avoid
stragglers

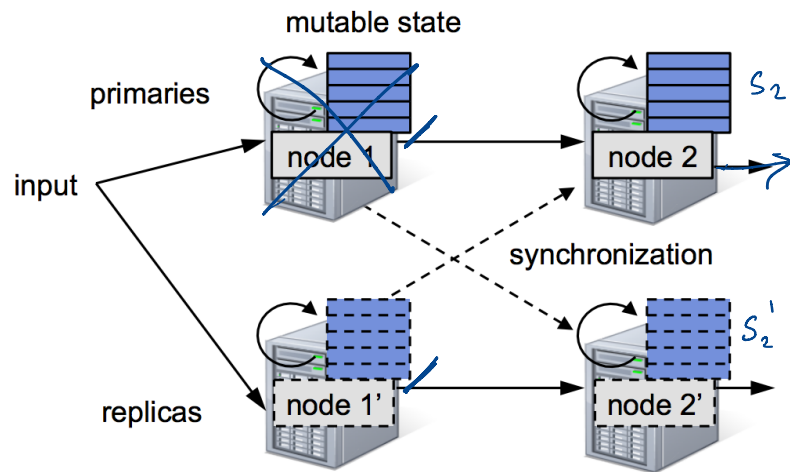


Flink



Naiad

CONTINUOUS OPERATORS



Replication to provide fault tolerance
⇒ Multiple copies (say 2?) of each operator

(1) Overhead of 2x resources required

(2) Replicas need to be in sync
⇒ S_2 and S_2' should be the same

⇒ need to make sure replicas are synchronized during normal computation
↓
overhead!

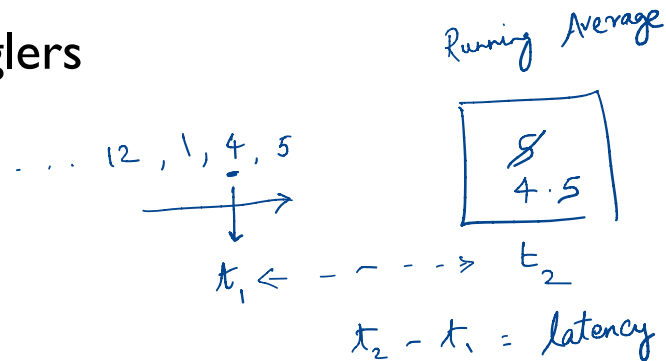
SPARK STREAMING: GOALS

1. Scalability to hundreds of nodes \rightarrow To handle high input streams

2. Minimal cost beyond base processing (no replication)

3. Second-scale latency \rightarrow time between when event arrives to
time when event is reflected in the output

4. Second-scale recovery from faults and stragglers



DISCRETIZED STREAMS (DSTREAMS)

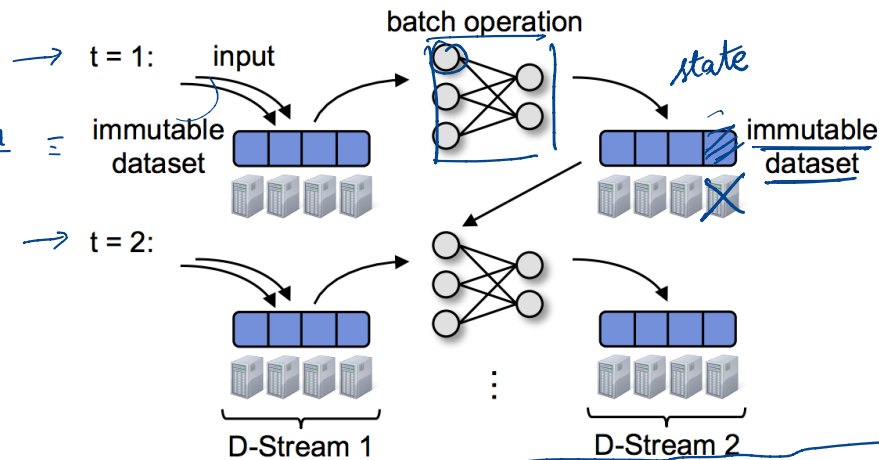
- every micro batch run short, deterministic tasks to compute incremental output

micro batch duration = $\frac{1}{s}$ → overheads in batch computation

- every batch operation is stateless
state is stored as immutable dataset

- each part of the state can be recovered independently
→ Use lineage to do recovery

can be made deterministic
Task
← random < 5 : output 0 else output 1



non-deterministic is opposite

→ if you re-run output might be difficult

EXAMPLE

Handwritten: [google.com, 5
yahoo.com, 12
...]

```
pageViews =  
  readStream(http://...,  
             File system, or from HTTP etc.
```

"1s"
↳ micro batch duration

```
ones = pageViews.map(  
  event => (event.url, 1))  
interval [0, 1)
```

```
counts =  
  ones.runningReduce(  
    (a, b) => a + b)  
interval [1, 2)
```

partitioned from input source
↳ read one chunk from storage

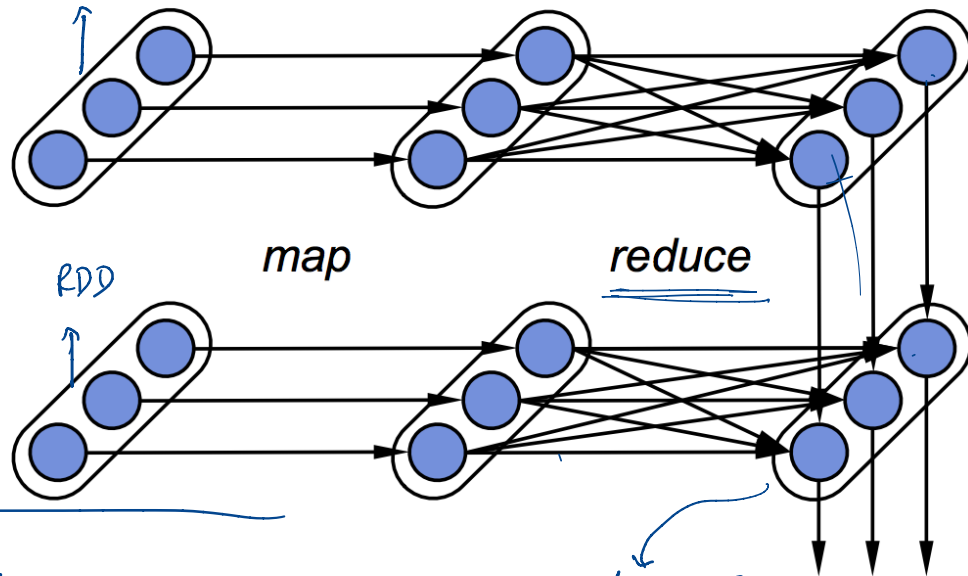
pageViews
DStream

RDD

ones
DStream

counts
DStream

hash, shuffle



Handwritten: [google.com, 7
...]

DSTREAM API

Transformations

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

Stateful:

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

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

↓
creates a sliding window and aggregates RDDs that belong to it.

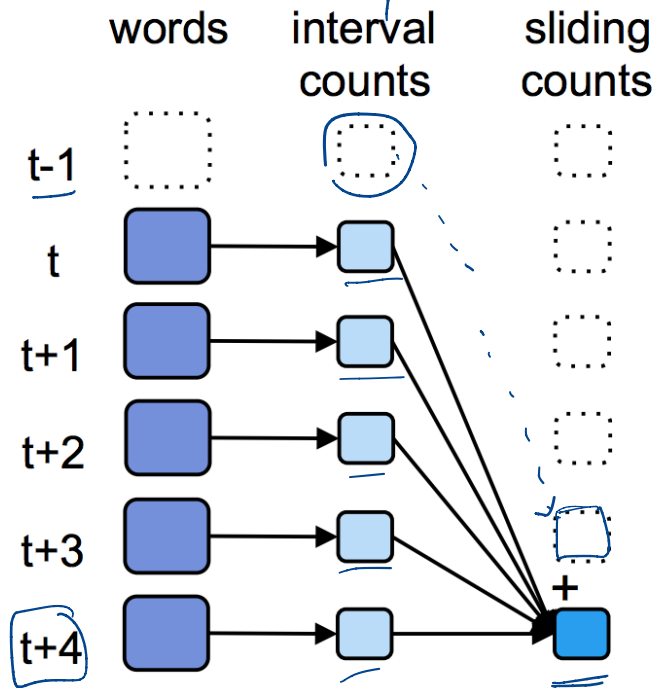
micro batch duration = 1s
 window duration = 5s
 $[t, t+5)$

SLIDING WINDOW

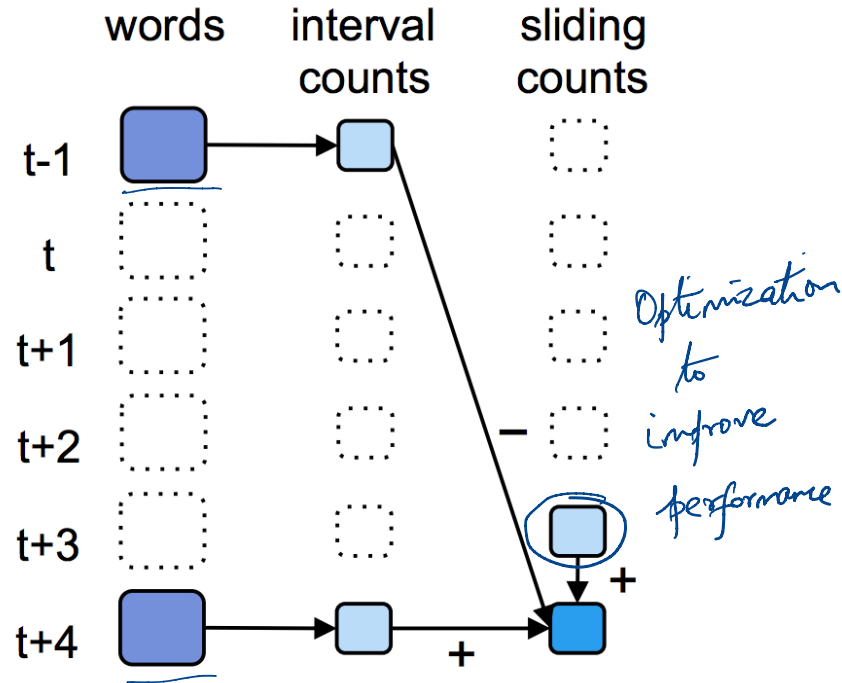
overhead

Add
previous 5
each time

reduce By
Window



(a) Associative only



(b) Associative & invertible

STATE MANAGEMENT

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

Session which has
all events for a
user satisfying some
criteria [login \rightarrow logout]

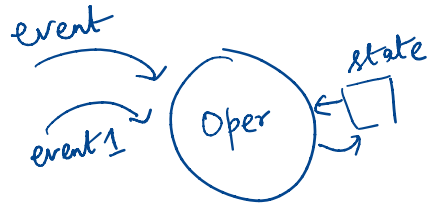
events.track(
user-id (key, ev) \Rightarrow 1, \rightarrow Initialize state

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

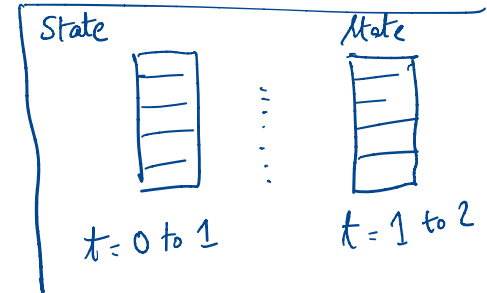
"30s")

update: given prev. state and
return new state

Timeout: forget old states



a new event

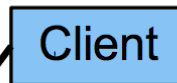
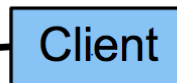
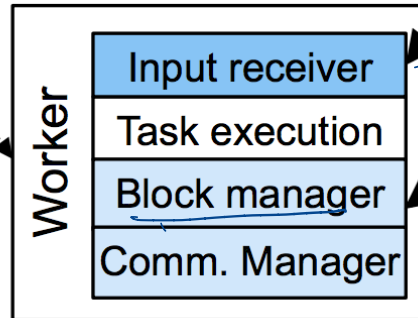
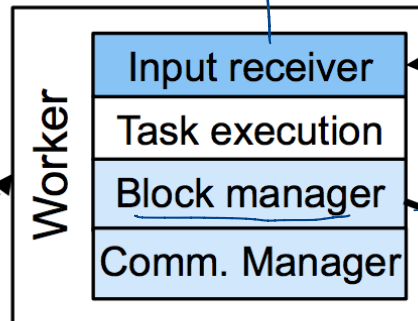
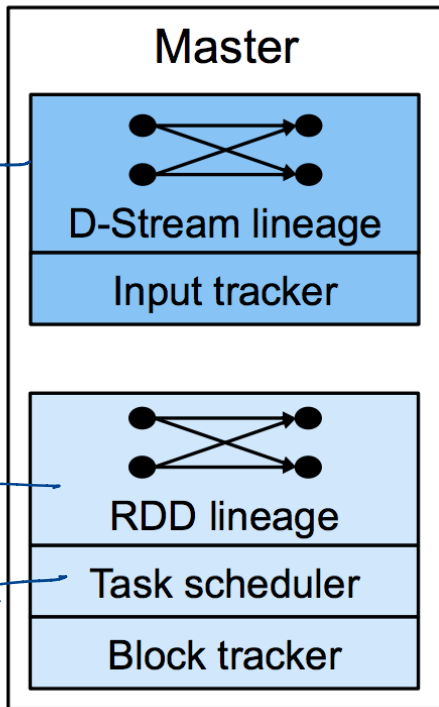


SYSTEM IMPLEMENTATION

- Persist data safely
- (1) Disk locally
 - (2) Memory remote
 - (3) Disk remote

Windowing
State

Inherit
from
Spark



replication of
input & check-
pointed RDDs

persist
(2 machines)
memory

New

Modified

OPTIMIZATIONS

$(0, 1)$ map
○
○

Timestep Pipelining → fuse together map operations

$(1, 2)$ ○
○
○ schedule before
prev finishes

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

can be done by storing to remote memory

FAULT TOLERANCE: PARALLEL RECOVERY

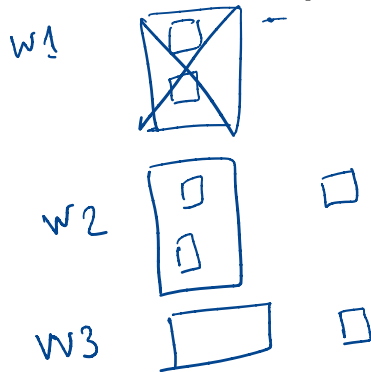
Worker failure

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

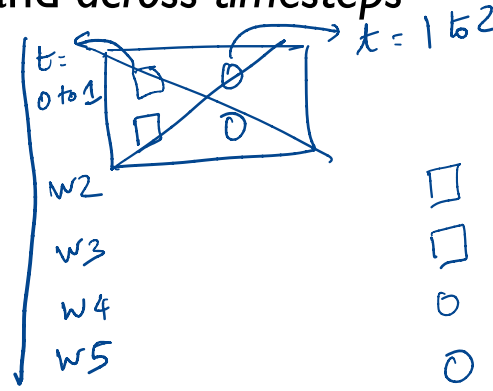
these might be
used for future
outputs

Strategy

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



parallelism
across
partitions!

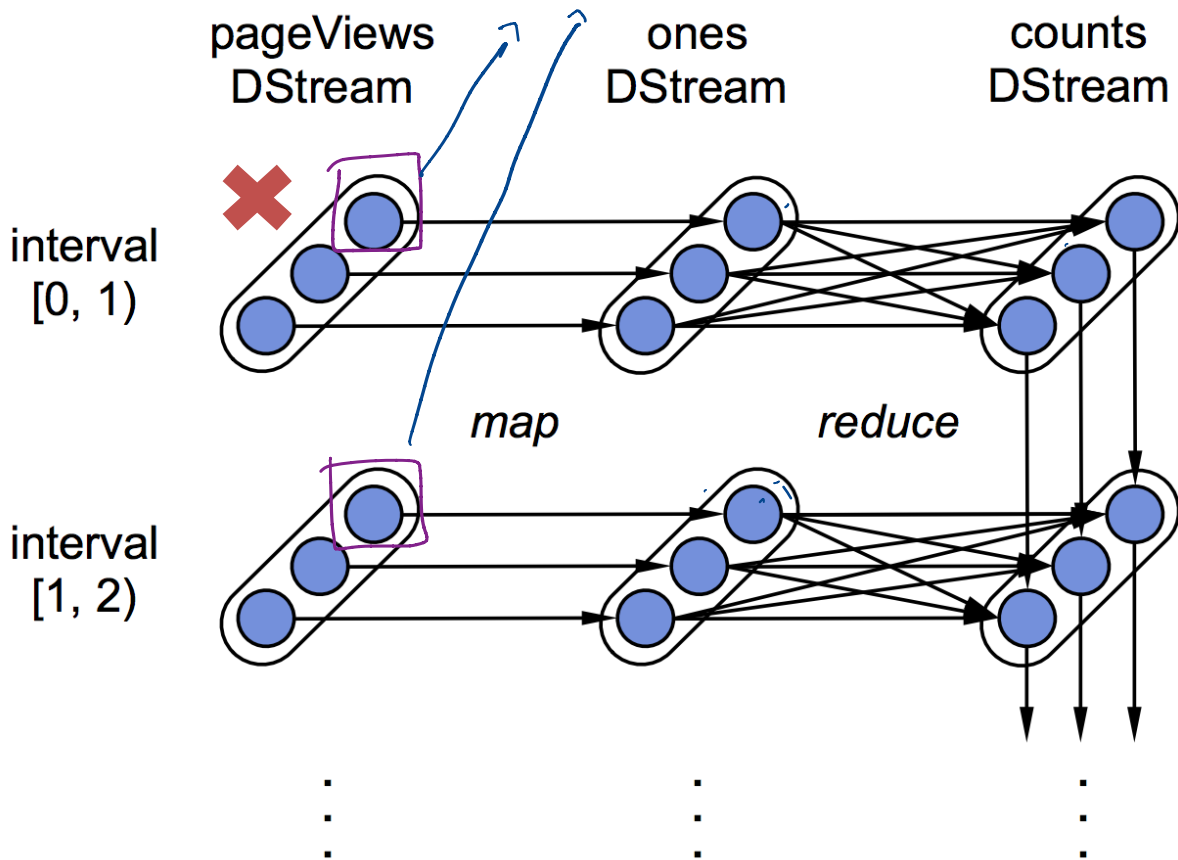


pending
any
dependencies!

EXAMPLE

*These are independent
can be recovered in parallel.*

```
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 → *fall back*

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

Driver

~~Master~~ Recovery → *Runs forever*

MR Master → retry the job on failure

- 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).

→ *AFS master recovery is similar!*

SUMMARY

Micro-batches: New approach to stream processing

Simplifies fault tolerance, straggler mitigation

Unifying batch, streaming analytics

DISCUSSION

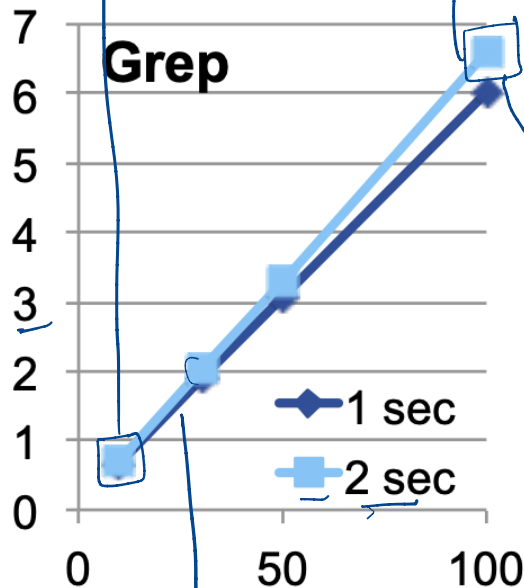
<https://forms.gle/eiqbjjTU95bMQLtm9>

$$y = x \text{ vs. } y = 0.5x$$

slope indicates overhead \Rightarrow more machines

Cluster Throughput (GB/s)

(GB/s)

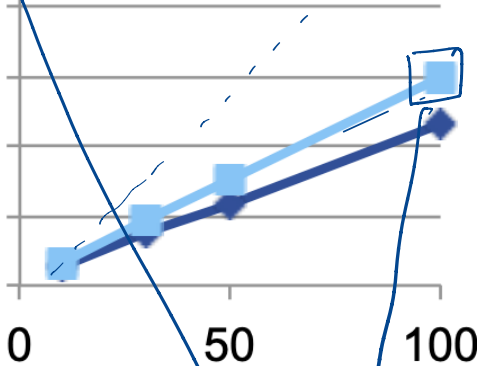


linear growth with cluster size

higher throughput for larger mini batch size overhead per mini-batch

WordCount

1 sec 2 sec

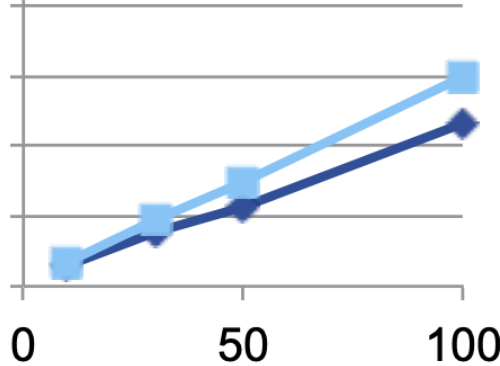


Nodes in Cluster

more tput for grep!

TopKCount

1 sec 2 sec



no co-ordination for grep!

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

too low latency \rightarrow low tput

overheads in task scheduling
tracking RDDs etc.

if we go to 1000 machines \Rightarrow overheads could be large!

\downarrow
linear scaling might not last?

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

latency sensitive

iterative + streaming
workflows

Spark Streaming

failures

stragglers

NEXT STEPS

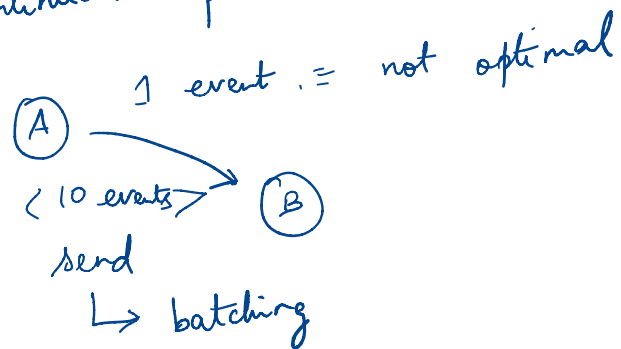
Next class: Graph processing!
Midterm grades ASAP!

lazy execution

↳ 1 sec (micro batch time)
trigger "forces computation"

Batching ?!

↳ Continuous operator



Very low latency

→ MPI - based

→ C++ Actor model

↳ Erlang → Telephone companies