

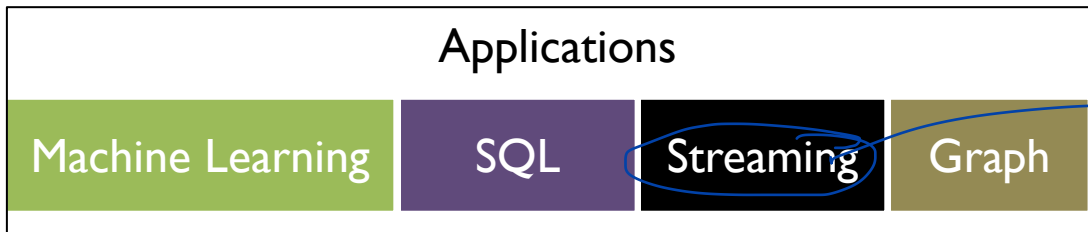
CS 744: NAIAD

Shivaram Venkataraman

Fall 2019

ADMINISTRIVIA

- Course Project Proposal feedback
- Midterm grades
- Checkins?

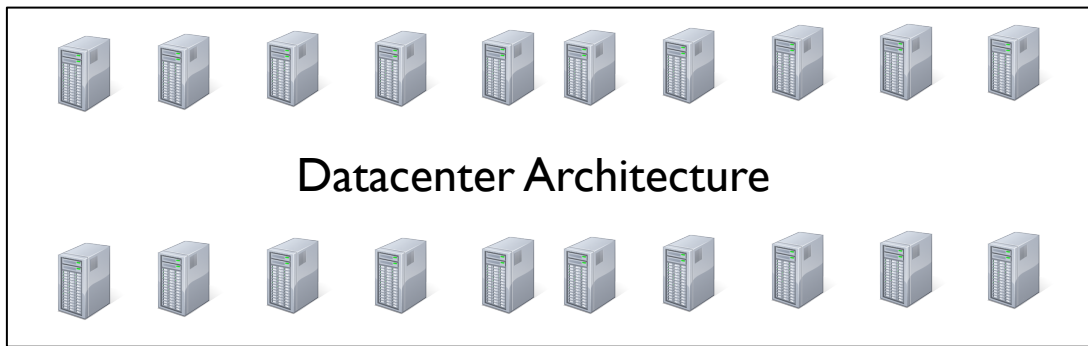


→ Dataflow model

Computational Engines

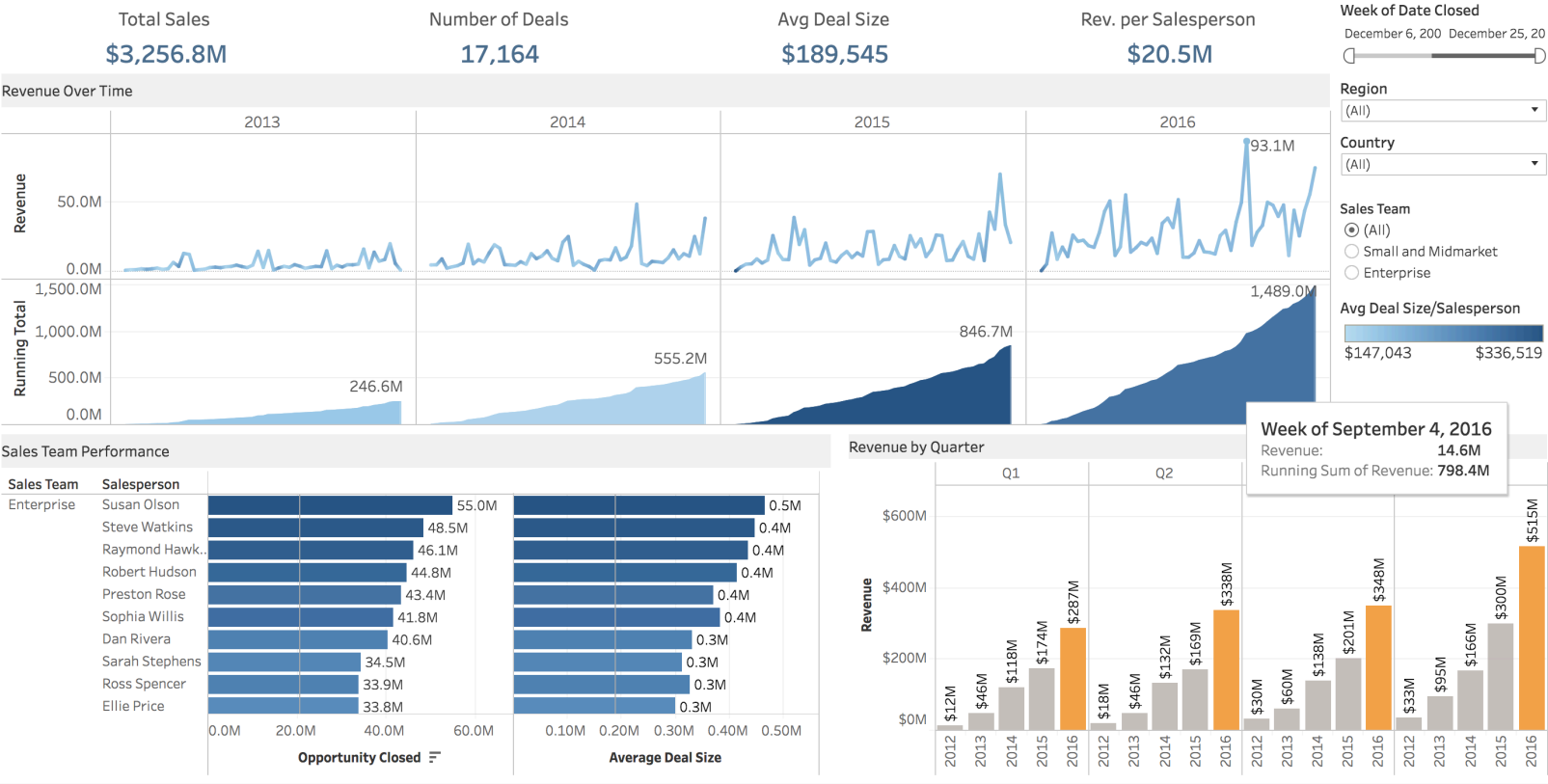
Scalable Storage Systems

Resource Management



DASHBOARDS

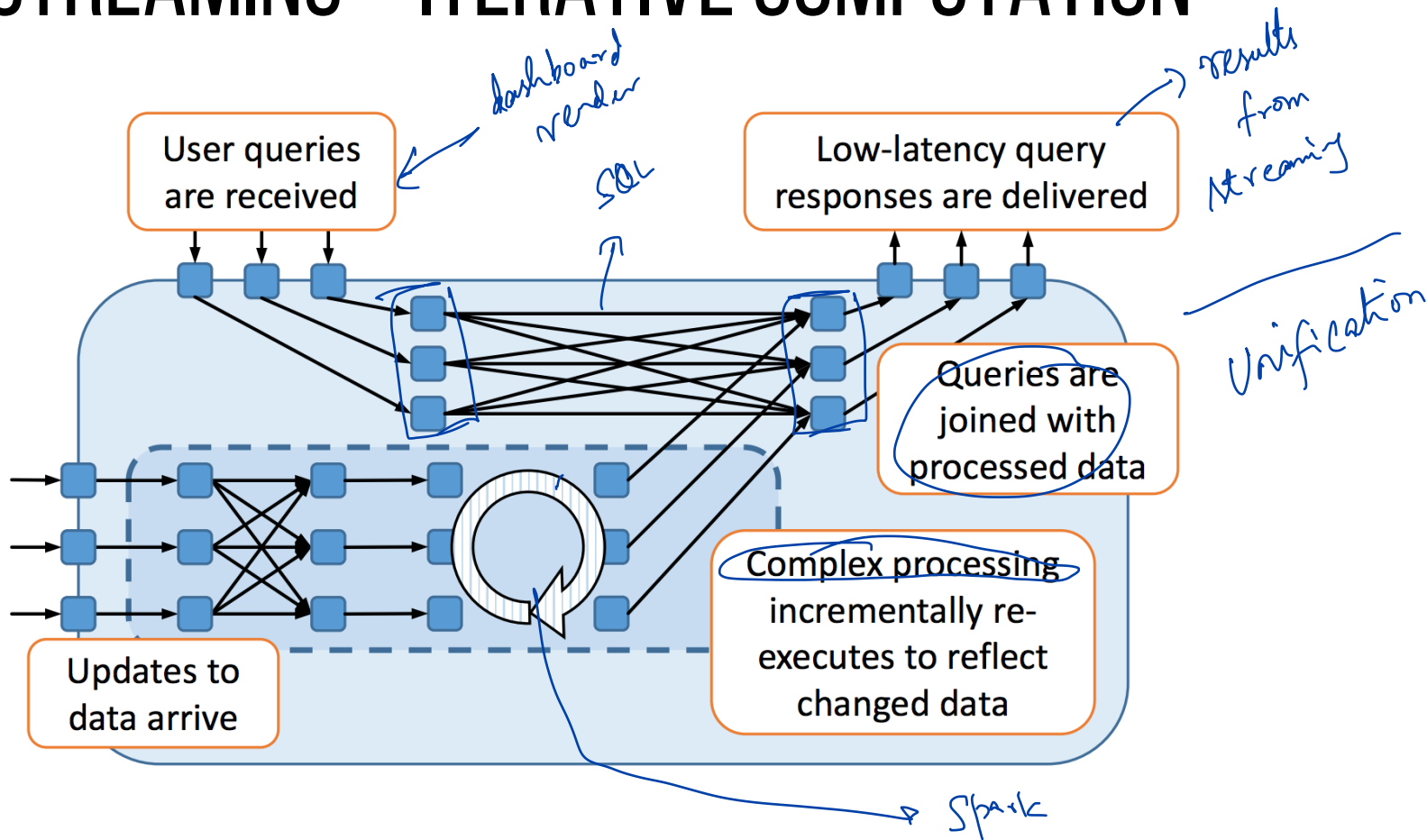
Sales Dashboard



STREAMING + ITERATIVE COMPUTATION

Computation
is
querying

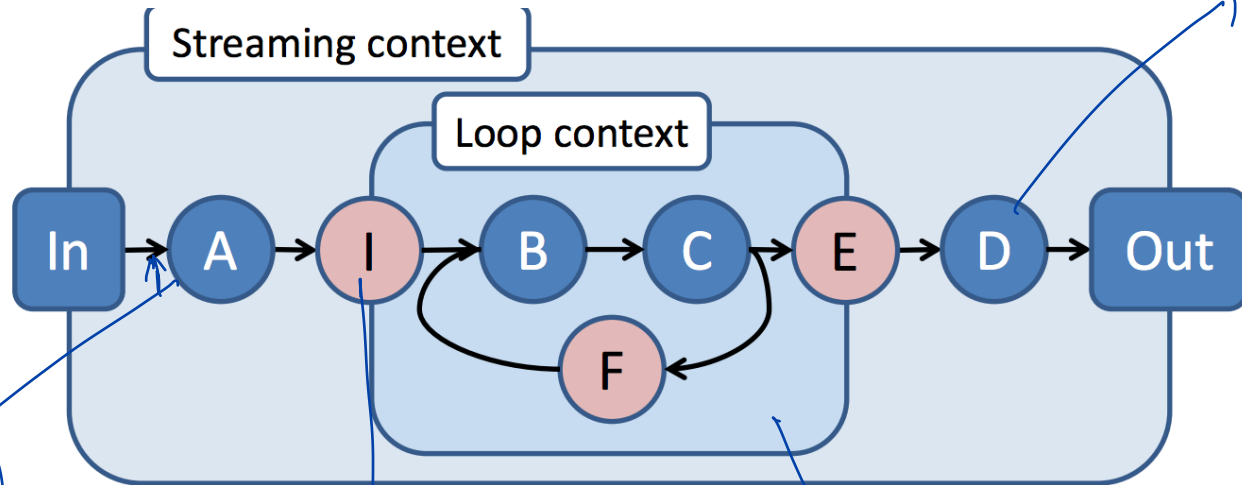
Stream
of
data



TIMELY DATAFLOW

DAGs

Let you reason
when all messages
for 1 are here
"notification"



stateful
vertices
↓
Tensorflow

Timestamp
on these
messages
"track"

Extra
vertices

Loops

TIMELY DATAFLOW

for (i)
 for (j)
 for (p)
 → k:3

Timestamp : $(\underbrace{e \in \mathbb{N}}_{\substack{\text{epoch} \\ 1}}, \underbrace{\langle c_1, \dots, c_k \rangle \in \mathbb{N}^k}_{\substack{\text{loop counters} \\ 1}})$

$(0, \langle 1, 1 \rangle)_{t_1}$
 $(0, \langle 2 \rangle)_{t_2}$
 $t_2 \rightarrow t_1$

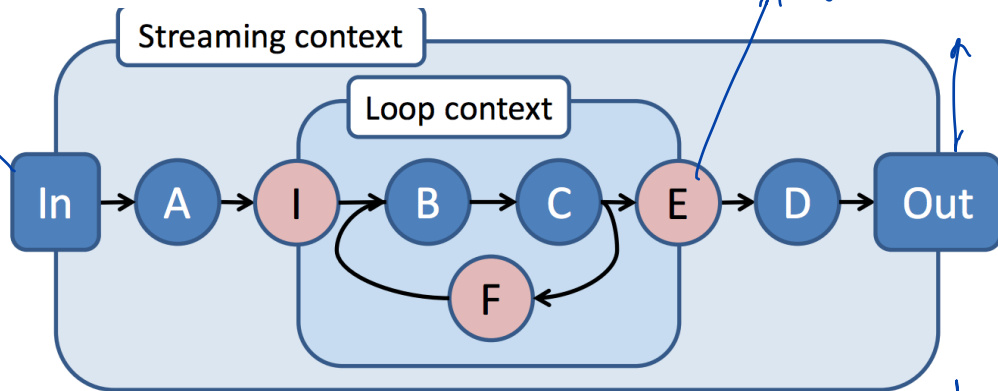
Vertex
 Ingress
 Egress
 Feedback

Input timestamp
 $(e, \langle c_1, \dots, c_k \rangle)$
 Egress
 $(e, \langle c_1, \dots, c_k, c_{k+1} \rangle)$
 Feedback
 $(e, \langle c_1, \dots, c_k \rangle)$

Output timestamp

$(e, \langle c_1, \dots, c_k, 0 \rangle)$
 $(e, \langle c_1, \dots, c_k \rangle)$
 $(e, \langle c_1, \dots, c_k+1 \rangle) \rightarrow \text{Inc loop counter}$

epoch starts



loop starts

New k+1

loop iter

Inc loop counter

VERTEX API

Receiving Messages

v.OnRecv(e : Edge, m : Msg, t : Time)

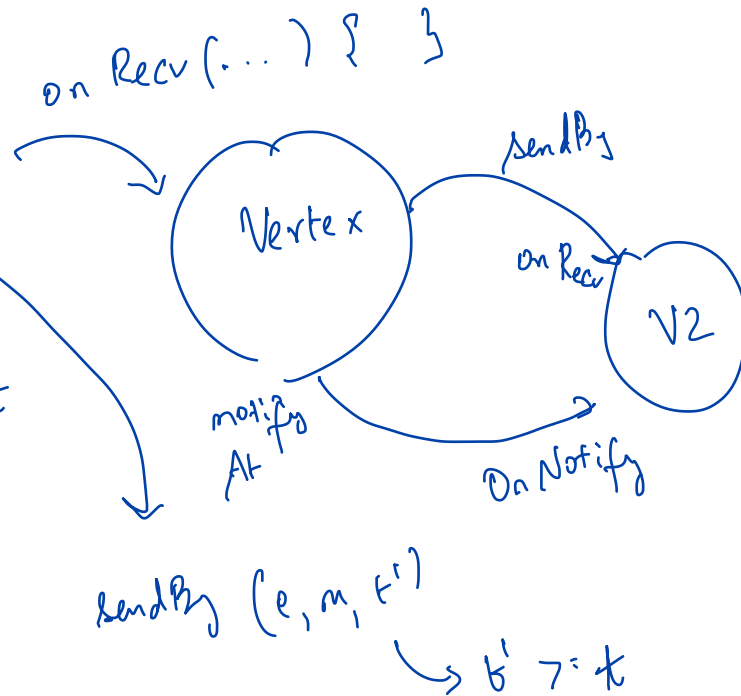
v.OnNotify(t : Timestamp) → "on Recv
if $t' \leq t$

Sending Messages

this.SendBy(e : Edge, m : Msg, t : Time)

this.NotifyAt(t : Timestamp)

not specifying
any vertex



IMPLEMENTING TIMELY DATAFLOW

Need to track when it is safe to notify

→ know that all messages with $t' \leq t$ have been delivered before

Path Summary

Check if (t_1, l_1) could-result-in (t_2, l_2)

whenever new message queue

oc

pc



v_1

v_2

v_3 note

Scheduler

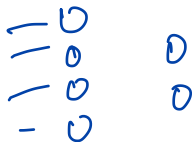
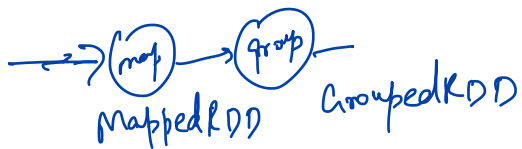
Occurrence and Precursor count

Precursor count = 0 → Frontier

sendBy (t, e)
 $oc[t, e] \leftarrow 1$
 Scheduler
 $rec(t, e) \quad oc[t, e] \leftarrow 1$

t_1

ARCHITECTURE



Workers communicate using
Shared Queue

Centralized master?

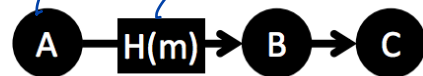
Batch messages delivered
Account for cycles

multiple stage same time

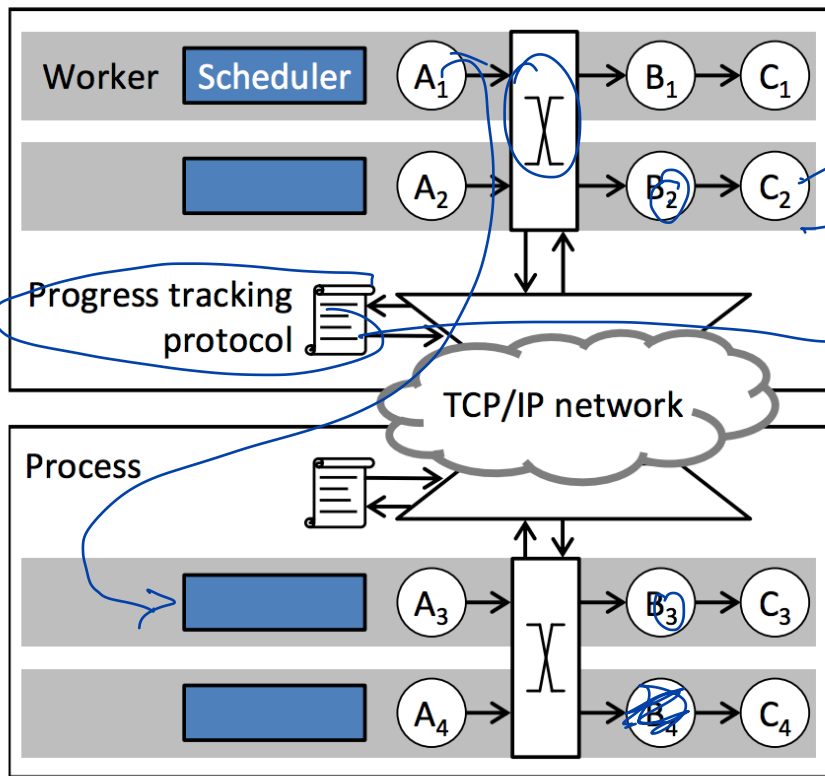
Vertex single threaded

no queue per vertex

Logical graph



vertices → Objects
Partitioner → B₃



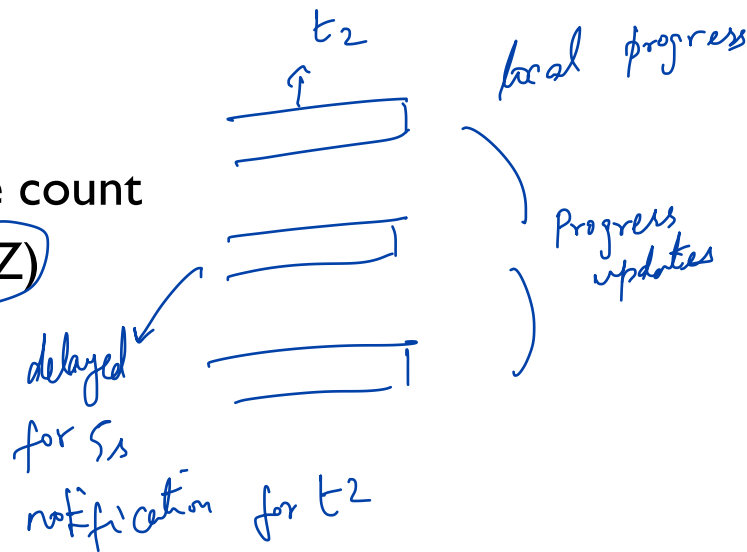
DISTRIBUTED PROGRESS TRACKING

Broadcast-based approach

Maintain local precursor count, occurrence count

Send progress update ($p \in \text{Pointstamp}, \delta \in \mathbb{Z}$)

Local frontier tracks global frontier



Optimizations

Batch updates and broadcast

Use projected timestamps from logical graph

vertices
Stateful

FAULT TOLERANCE

Checkpoint

Log data as computation goes on
Write a full checkpoint on demand

Restore

Reset all workers to checkpoint
Reconstruct state

Pause worker threads
Flush message queues OnRecv

Recovery time could be large
failures are rare?

Resume execution

Driver → centralized point of contention
→ Simplifies

MICRO STRAGGLERS

What is different from stragglers in MapReduce?

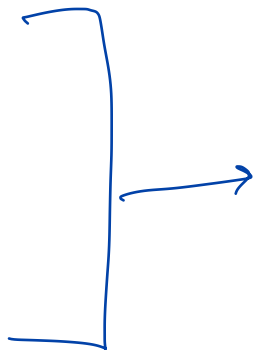
Stateful

Sources of stragglers

Network

Concurrency

Garbage Collection



Systems tricks

Preventive

DIFFERENTIAL DATAFLOW

```
// 1a. Define input stages for the dataflow.  
var input = controller.NewInput<string>();  
// 1b. Define the timely dataflow graph.  
// Here, we use LINQ to implement MapReduce.  
var result = input.SelectMany(y => map(y))  
                   .GroupBy(y => key(y),  
                             (k, vs) => reduce(k, vs));  
// 1c. Define output callbacks for each epoch  
result.Subscribe(result => { ... });  
// 2. Supply input data to the query.  
input.OnNext(/* 1st epoch data */);  
input.OnCompleted();
```

SUMMARY

Stream processing → Increasingly important workload trend

Timely dataflow: Principled approach to model batch, streaming together

Vertex message model

- Compute frontier
- Distributed progress tracking

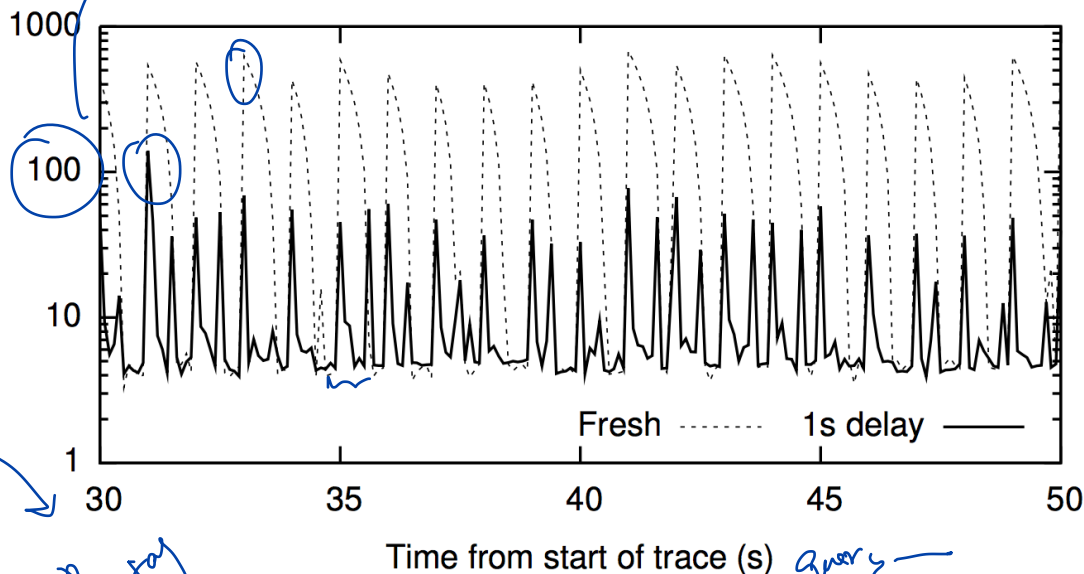
DISCUSSION

<https://forms.gle/v3YsWlHvnqsxCuPu5>



Query when
Contention: compute

Response time (ms)



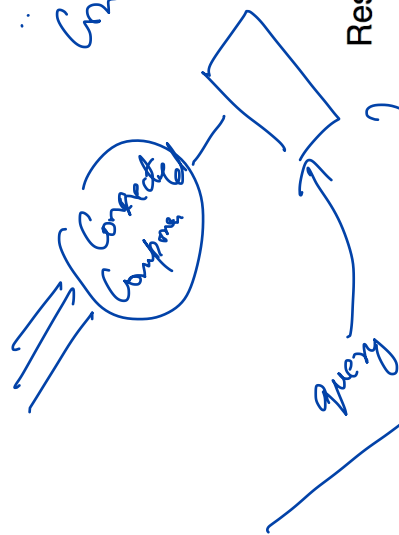
Top track toy

Incremental
Computation

queries —

Connected
Component

t_1



What are some example scenarios discussed in the dataflow paper that are NOT a good fit for implementation using Naiad?

Stale updates \rightarrow watermark

Triggering \rightarrow

Good but not perfect is fine!

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?