

Good morning!

CS 744: NAIAD

Shivaram Venkataraman

Fall 2021

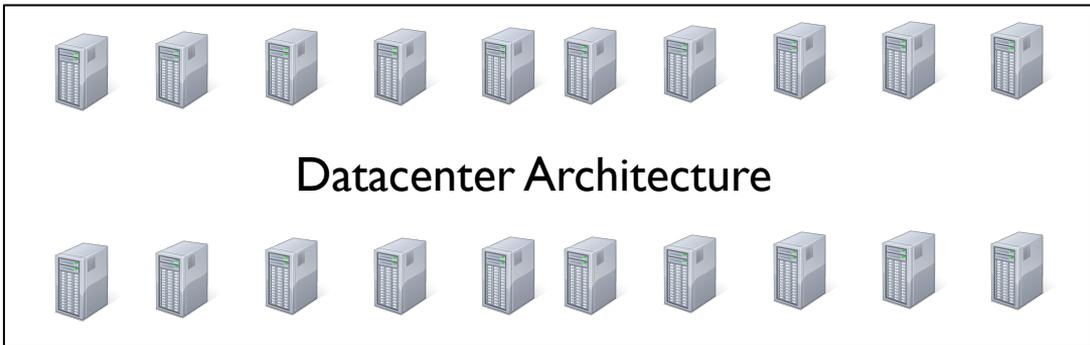
ADMINISTRIVIA

In Progress

- Course Project Proposal feedback
- Midterm grading



Workload
API
Dataflow model



DASHBOARDS

↳ continuously updated as data arrives
 → tput of events, latency to update this

Sales Dashboard

Total Sales
\$3,256.8M

Number of Deals
17,164

Avg Deal Size
\$189,545

Rev. per Salesperson
\$20.5M

Week of Date Closed

December 6, 2016 - December 25, 2016



Region

(All)

Country

(All)

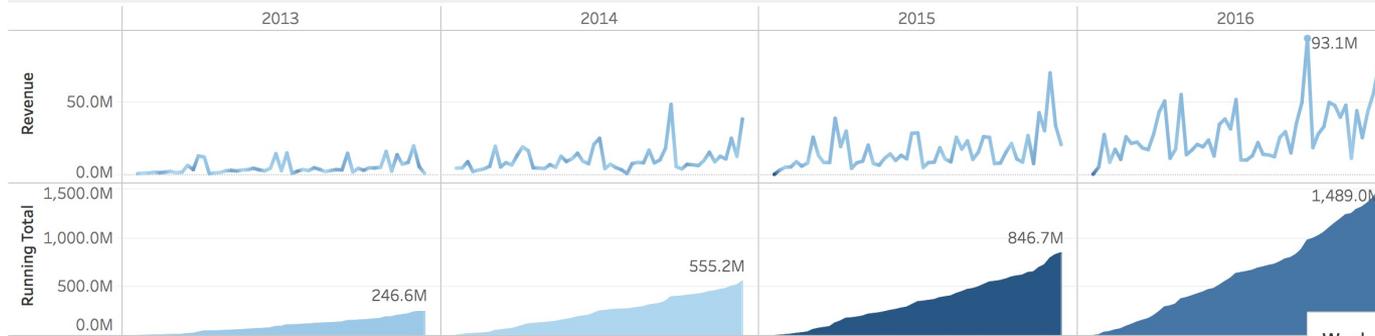
Sales Team

- (All)
- Small and Midmarket
- Enterprise

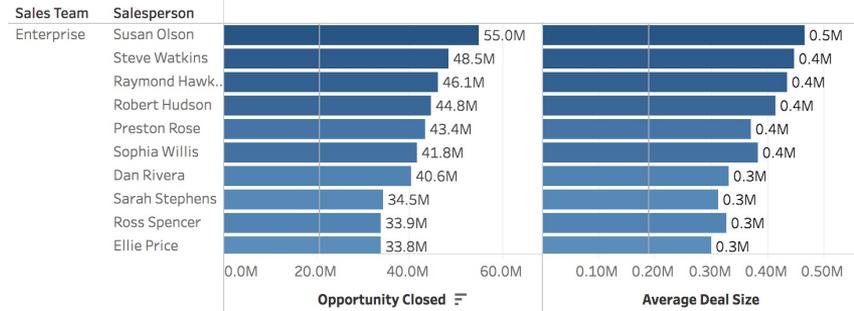
Avg Deal Size/Salesperson



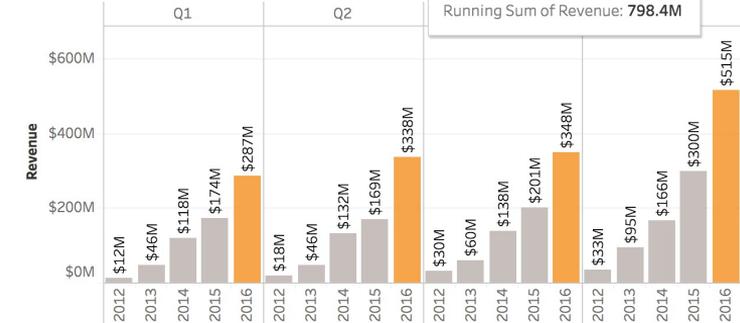
Revenue Over Time



Sales Team Performance

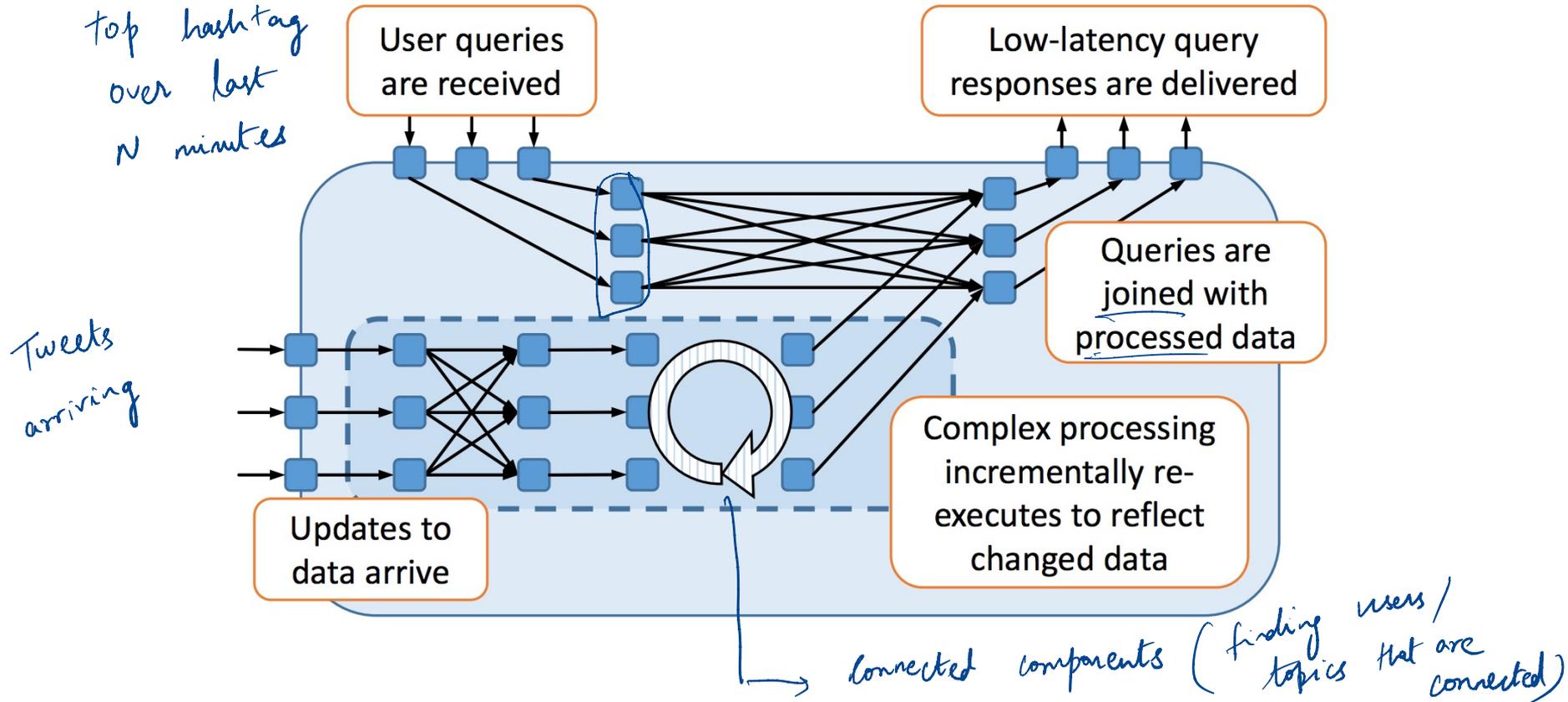


Revenue by Quarter

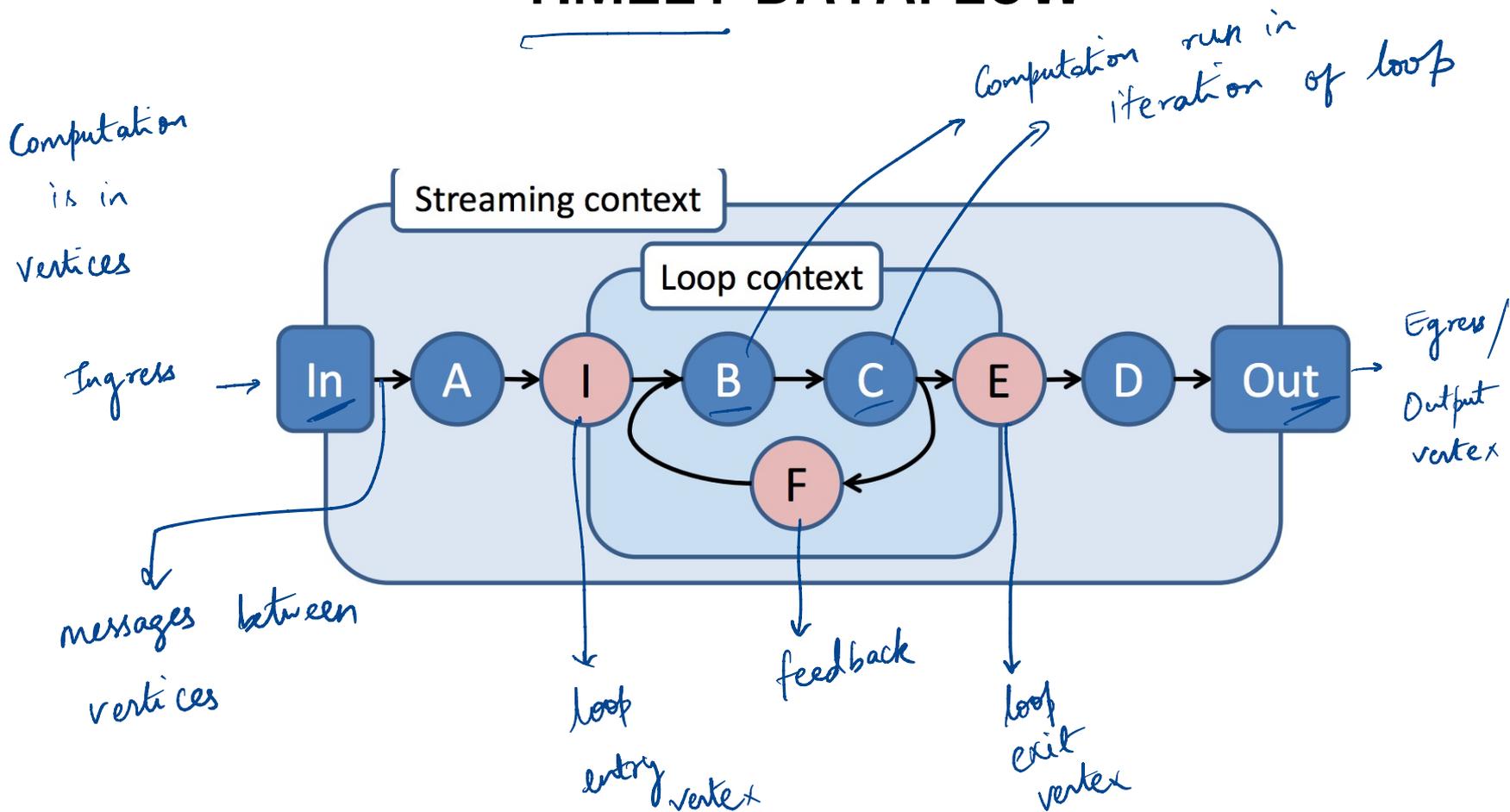


Week of September 4, 2016
 Revenue: 14.6M
 Running Sum of Revenue: 798.4M

STREAMING + ITERATIVE COMPUTATION



TIMELY DATAFLOW



TIMELY DATAFLOW

Epoch
 ↳ input events have an epoch which is given

↳ notifications / output when computation for epoch finishes
 epoch loop counters

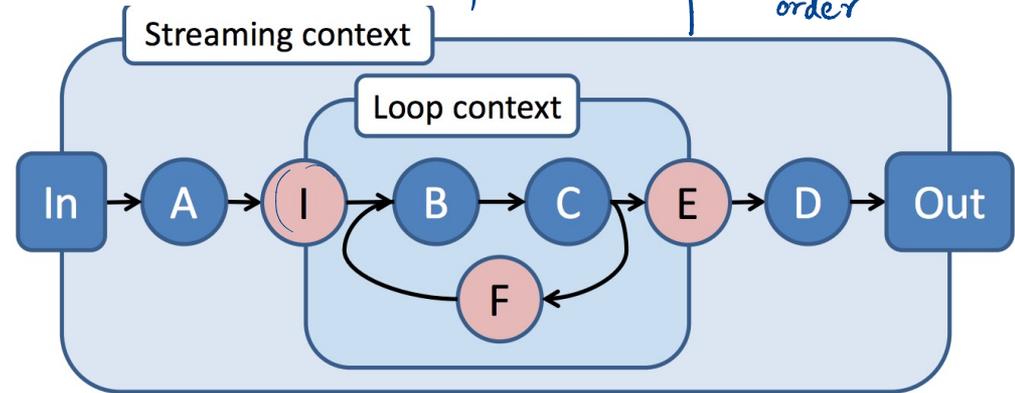
Timestamp : $(e \in \mathbb{N}, \langle c_1, \dots, c_k \rangle \in \mathbb{N}^k)$

for $i = 0 : 10$
 for $j = 0 : 10$

for loop $k+1$

Vertex	Input timestamp
<u>Ingress</u>	$(e, \langle c_1, \dots, c_k \rangle)$
<u>Egress</u>	$(e, \langle c_1, \dots, c_k, c_{k+1} \rangle)$
<u>Feedback</u>	$(e, \langle c_1, \dots, c_k \rangle)$

epoch numbers with timestamp?
 → flexibility → Out or order



Output timestamp	nested loops
$(e, \langle c_1, \dots, c_k, 0 \rangle)$	iter for $(k+1)$ loops
$(e, \langle c_1, \dots, c_k \rangle)$	
$(e, \langle c_1, \dots, c_{k+1} \rangle)$	

VERTEX API

Actor System
(Ray)

Receiving Messages

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

`v.OnNotify(t : Timestamp)`

↳ called when all msgs $\leq t$ have been delivered.

function that a vertex implements

Sending Messages

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

`this.NotifyAt(t : Timestamp)`

functions that vertex can call

↳ send a msg m on this edge e , with t as timestamp

↳ this vertex will not produce any outputs with $< t$.

EXAMPLE

```
class DistinctCount<S,T> : Vertex<T> {
    Dictionary<T, Dictionary<S,int>> counts;
    void OnRecv(Edge e, S msg, T time) {
        if (!counts.ContainsKey(time)) {
            counts[time] = new Dictionary<S,int>();
            this.NotifyAt(time);
        }
        if (!counts[time].ContainsKey(msg)) {
            counts[time][msg] = 0;
            this.SendBy(output1, msg, time);
        }
        counts[time][msg]++;
    }
    void OnNotify(T time) {
        foreach (var pair in counts[time])
            this.SendBy(output2, pair, time);
        counts.Remove(time);
    }
}
```

is a vertex

"state" which is inside this vertex

notify for older timestamps

initialize dict for this timestamp

↳ Tell output1 that we have a new msg!

is this a new event (msg) if so counter = 0

(msg, counter) write this as output at the end of epoch

IMPLEMENTING TIMELY DATAFLOW

Need to track when it is safe to notify

Path Summary

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

Scheduler

Occurrence and Precursor count

Precursor count = 0 \rightarrow Frontier



When is it safe to notify a vertex ??

Create possible paths using
could-result in

Identical workers
(no centralized Master)

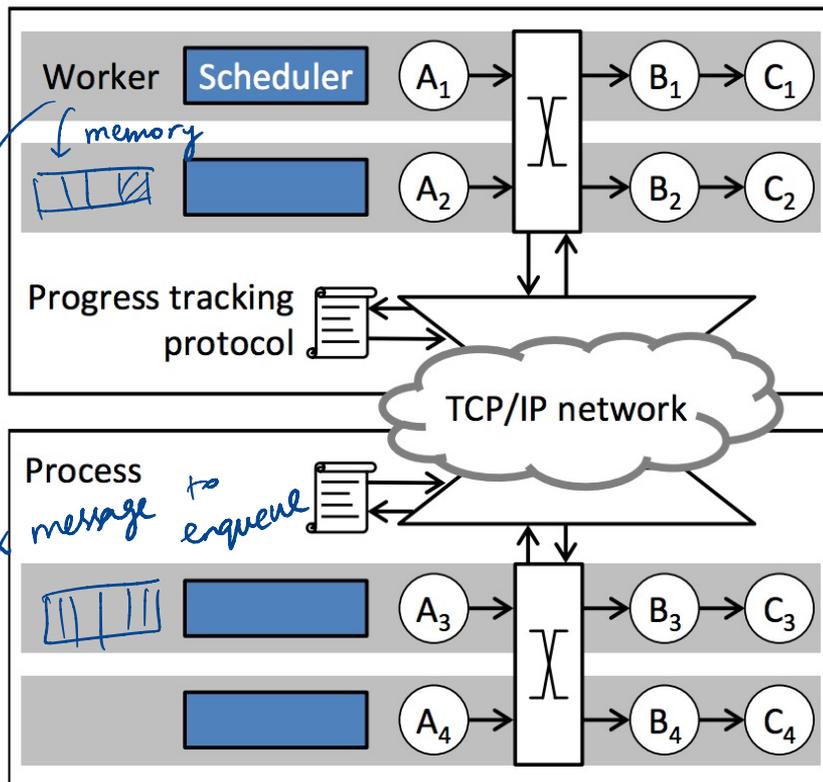
ARCHITECTURE

Workers communicate using
Shared Queue

Batch messages delivered
Account for cycles

Vertex single threaded

Logical graph **A** → **H(m)** → **B** → **C** *Per Core*



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

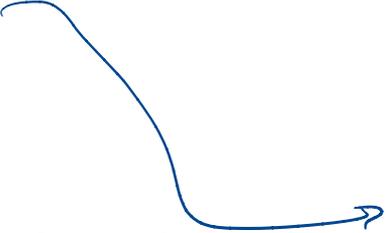
Optimizations

Batch updates and broadcast

Use projected timestamps from logical graph



helps you derive
a local frontier
that tracks
global



When is it safe to
notify at local
machine given
global progress

FAULT TOLERANCE

mutable data in vertex

Checkpoint()

Log data as computation goes on

Write a full checkpoint on demand

Pause worker threads

Flush message queues OnRecv

it is safe to
checkpoint / consistent

Restore

Reset all workers to checkpoint

Reconstruct state

not just the failed
worker!!

Resume execution

↳ redo work from
the checkpoint

When failure is infrequent
get better perf when no
failures!

MICRO STRAGGLERS

What is different from stragglers in MapReduce?

↳ duplicate or back up
tasks

Sources of stragglers

Network

Concurrency

Garbage Collection

minimize
prob. of
straggler!

HIGH LEVEL API

```
// 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

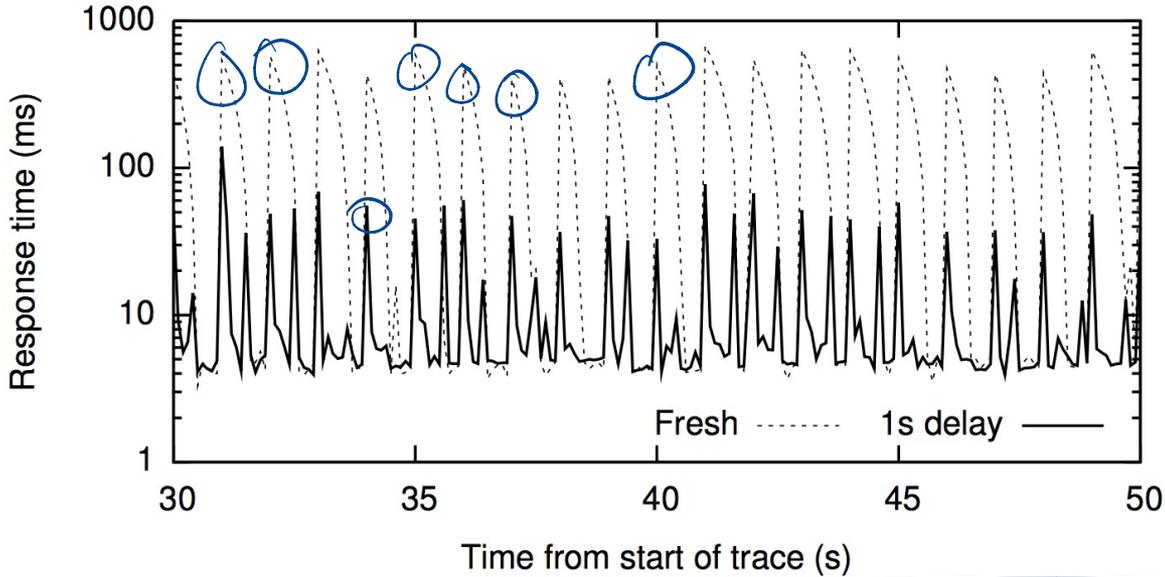
DISCUSSION

<https://forms.gle/cYYyI4JfX7IqBzY66>

Input:

32,000 tweets/second, Queries every 100ms

① $t = \text{Fresh}$
Suffers from bigger slowdown than $t = 1s$



state data already been processed

② Query is waiting for incremented connected components

materialize - 10

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?

SUMMARY

Next class: Spark Streaming

~~Course project peer feedback due tonight!~~