

Welcome!

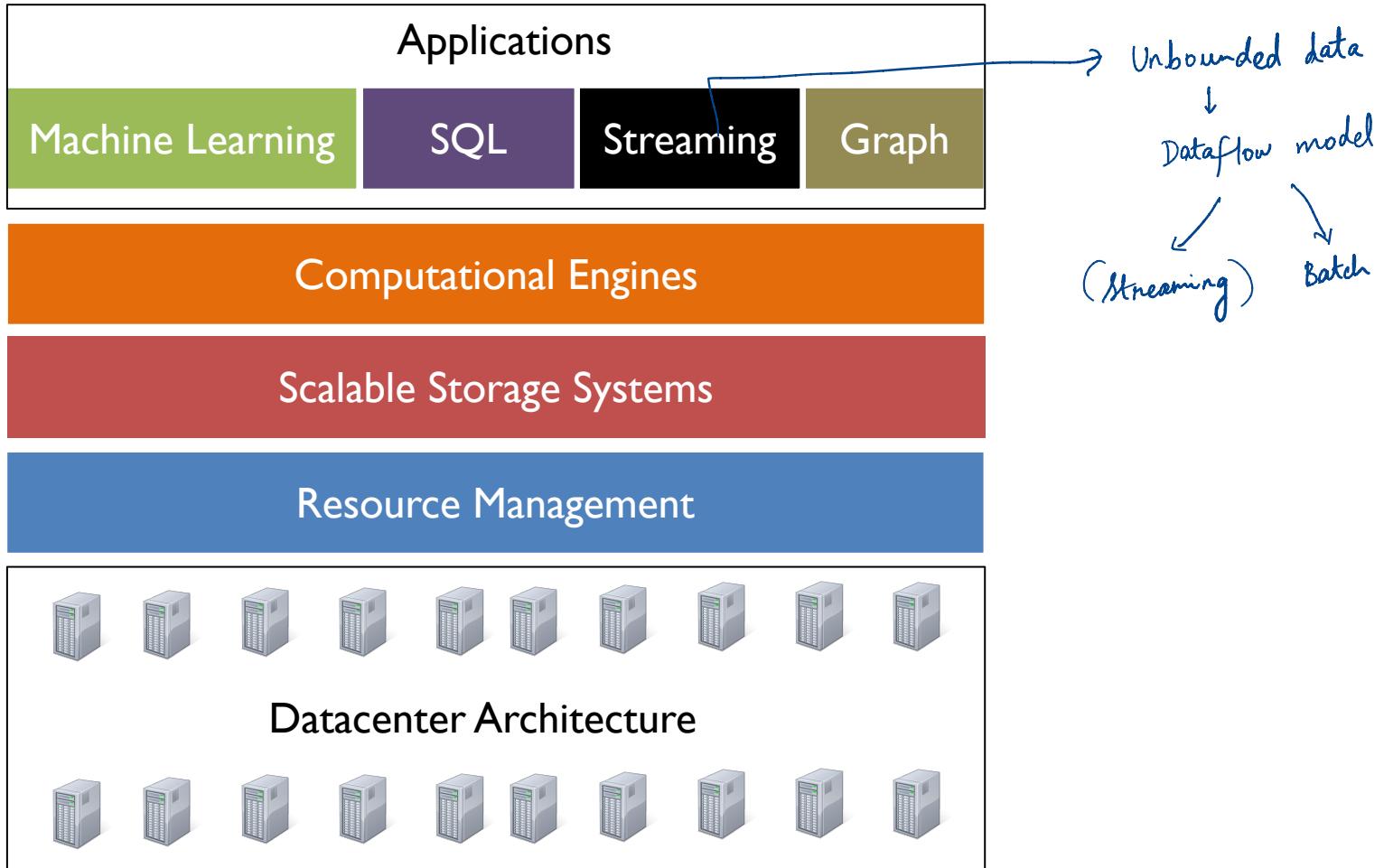
CS 744: NAIAD

Shivaram Venkataraman

Fall 2020

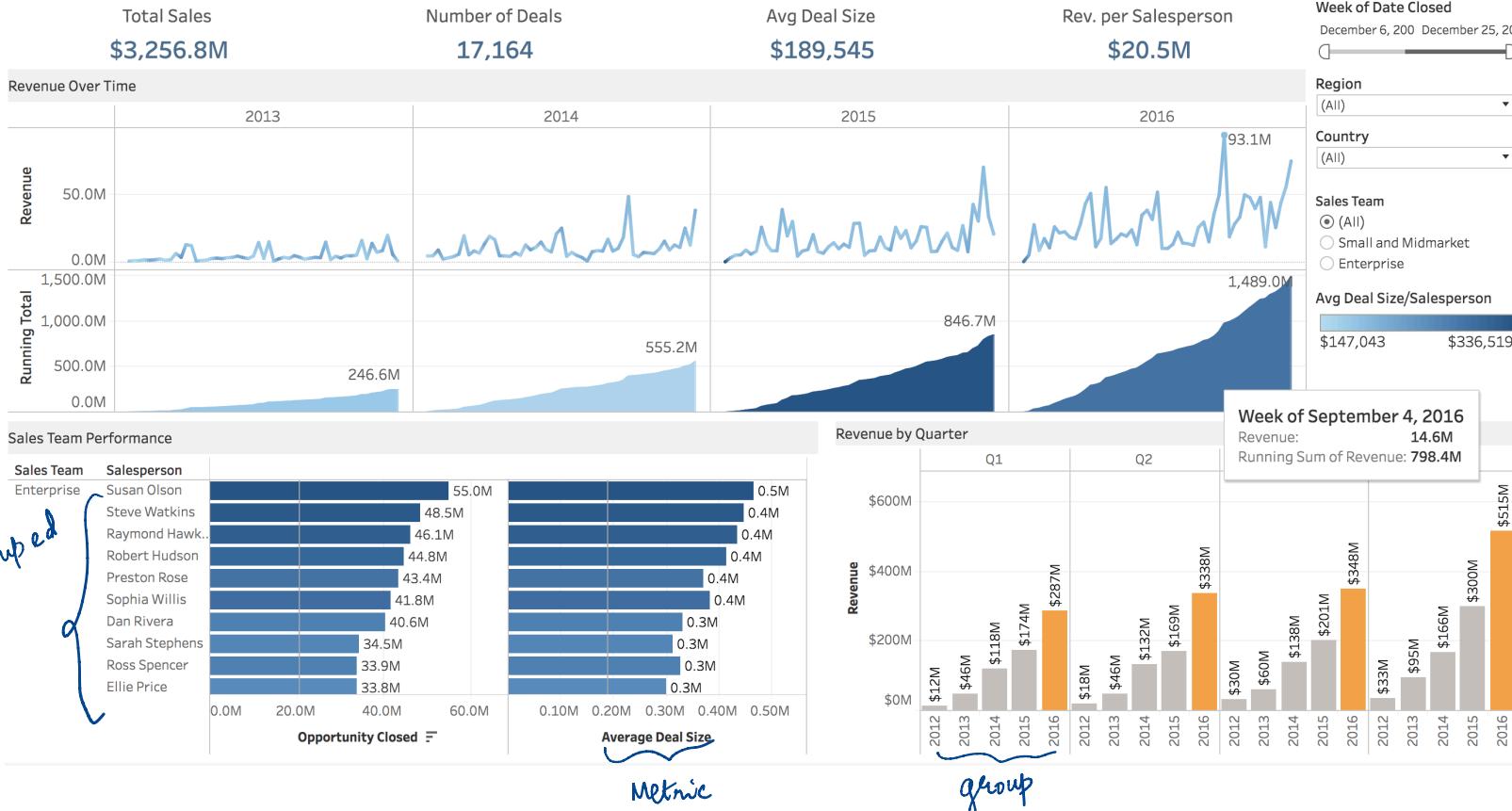
ADMINISTRIVIA

- Course Project Proposal feedback today!
- Midterm grading in progress  Hot CRP system
- Assignment regrades?
- Shivarajan OH poll → Piazza



DASHBOARDS

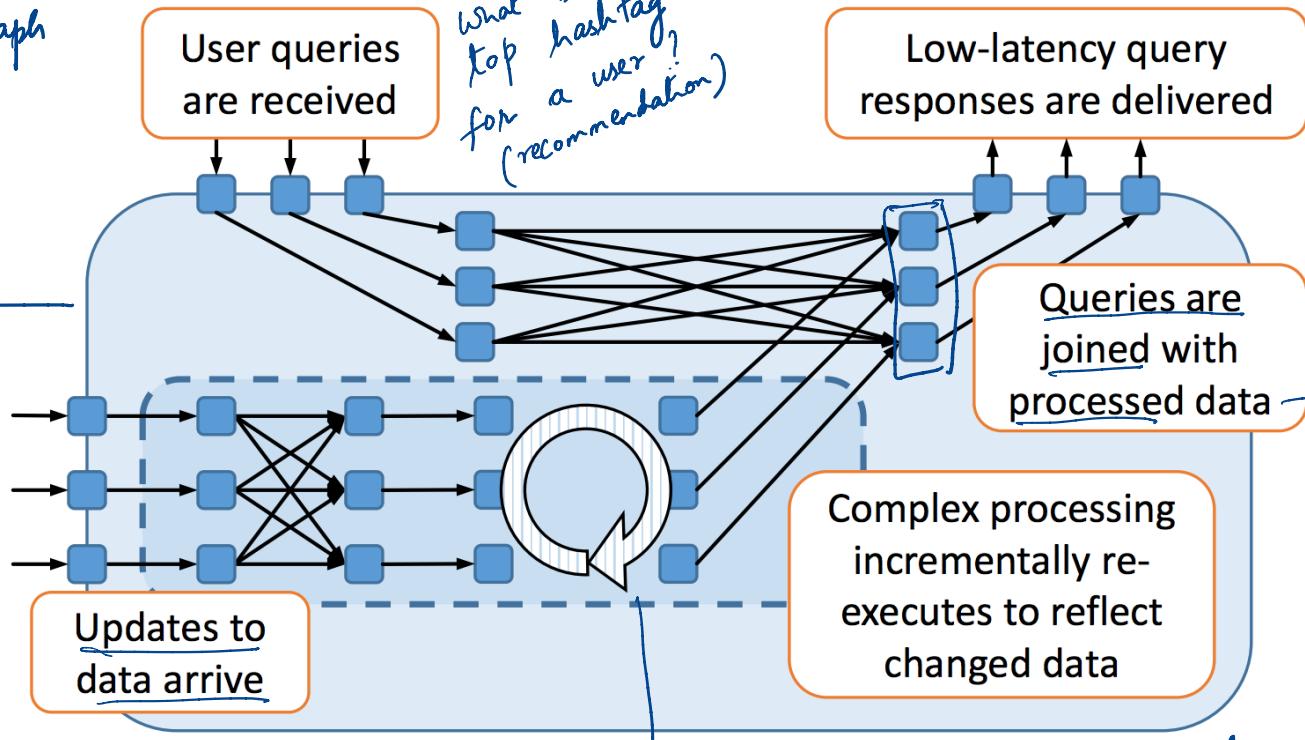
Sales Dashboard



STREAMING + ITERATIVE COMPUTATION



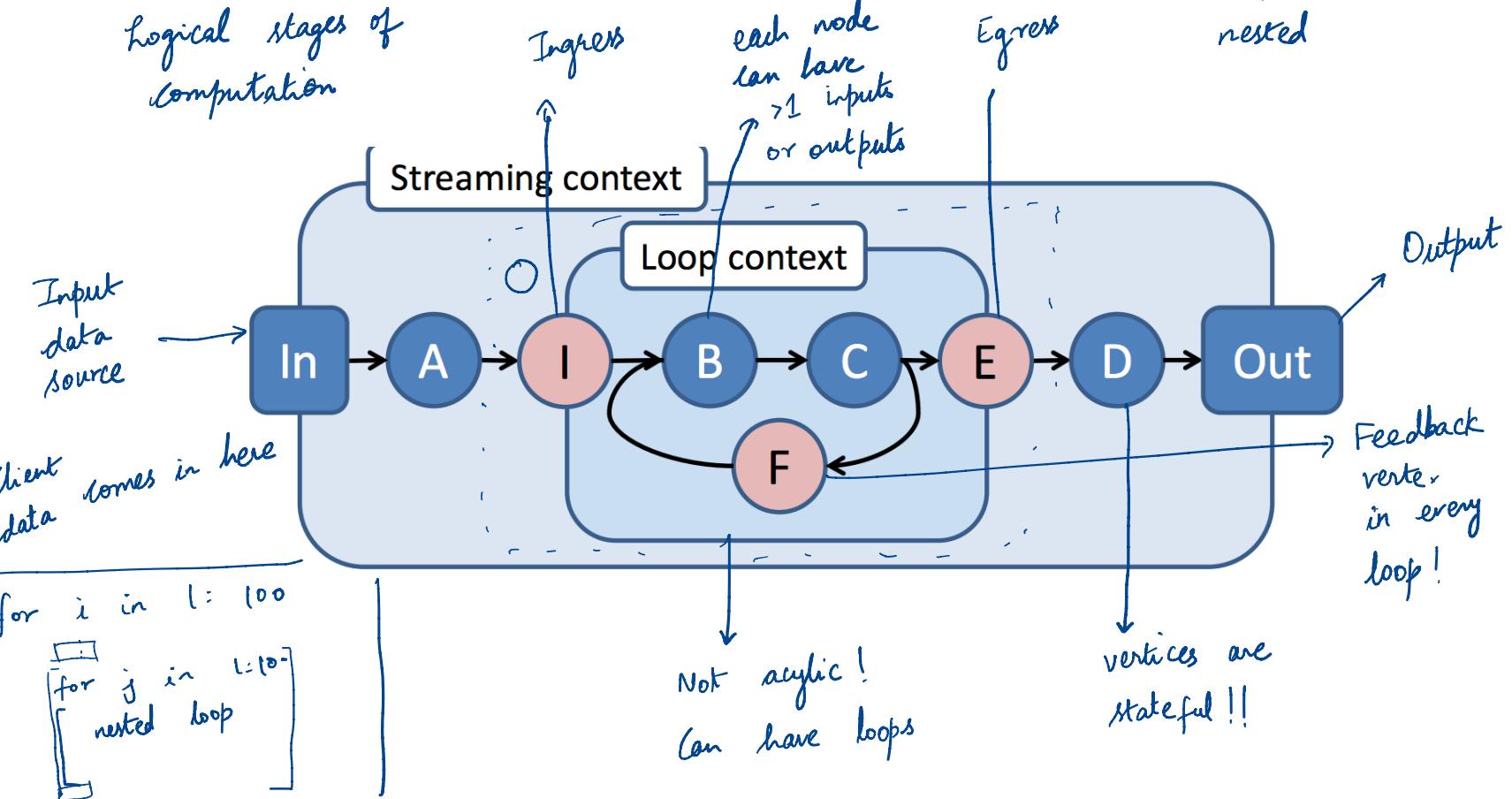
Tweets
Streaming in
to the
system



→ Page Rank
→ ML algorithms
Incremental processing !

→ Output from connected components

TIMELY DATAFLOW

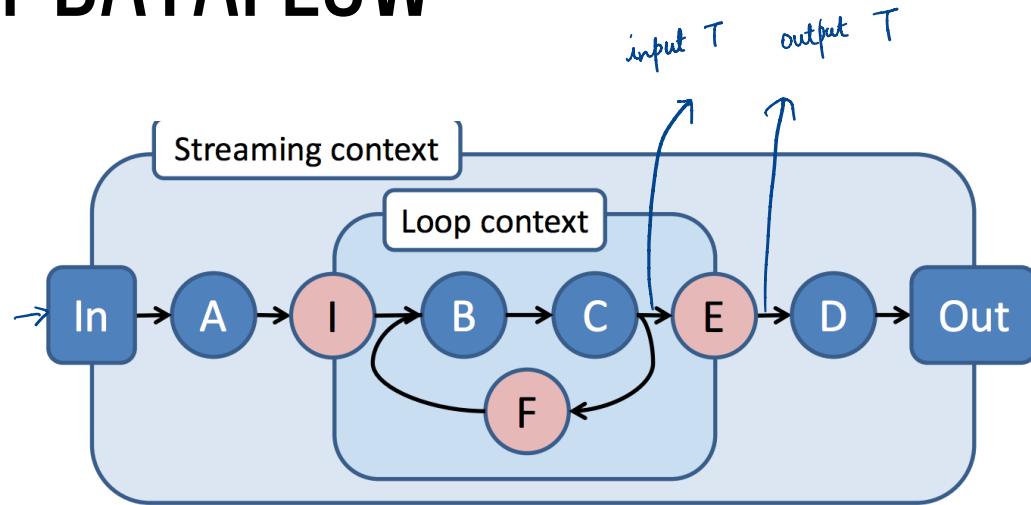


input grows in epochs
 → input events already have an epoch
 → notify when epoch is done

TIMELY DATAFLOW

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

Example : $(0, \langle 0, 0, 17 \rangle)$
 $\langle 0, 0, 27 \rangle$
 $\langle 0, 1, 0 \rangle$
 \vdots



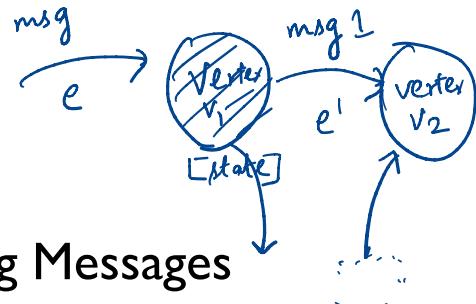
Vertex	Input timestamp
Ingress	$(e, \langle c_1, \dots, c_k \rangle)$
Egress	$(e, \langle c_1, \dots, c_k, c_{k+1} \rangle)$
Feedback	$(e, \langle c_1, \dots, c_k \rangle)$

Output timestamp

$(e, \langle c_1, \dots, c_k, 0 \rangle)$

$(e, \langle c_1, \dots, c_k \rangle)$ for this new loop context

$(e, \langle c_1, \dots, c_{k+1} \rangle)$ done with $k+1^{\text{th}}$ loop increment last loop counter



VERTEX API

Actor Model = Similar

Receiving Messages

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

`v.OnNotify(t : Timestamp)`

↳ called when all messages with
ts $\leq t$ have been delivered

Sending Messages

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

`this.NotifyAt(t : Timestamp)` → will lead to
on Notify(t)

useful to finalize the values for
an epoch

this vertex
will not produce
any output for t
after this.

```
class Vertex {
    running_sum = 0
    onRecv(e: Edge, m: Msg, t: Time) {
        // got msg m from edge E
        // at time T
        running_sum += (int)m;
        sendBy(e', running_sum, t')
    }
}
```

↳ invokes onRecv on v2 with
msg = running_sum, t = t'.

$t' \geq t$

IMPLEMENTING TIMELY DATAFLOW

one process

Need to track when it is safe to notify

Path Summary

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

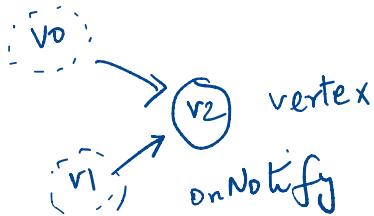
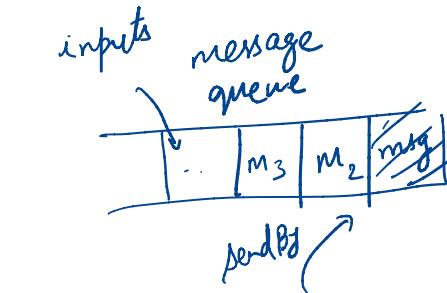
Scheduler

Occurrence and Precursor count

Precursor count = 0 \rightarrow Frontier

when occurrence = 0
update precursor
for downstream
vertices

sendBy \rightarrow increments occurrence count
recv \rightarrow decrements count

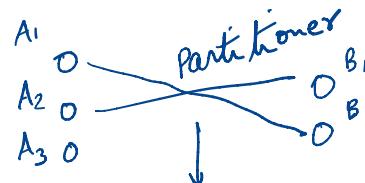


All precursors must be done with time stamp T before onNotify can be called.

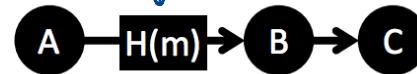
Scheduler
pops a msg
calls v.onRecv(),
:

long running
mutable
inside
vertices
state them

ARCHITECHTURE



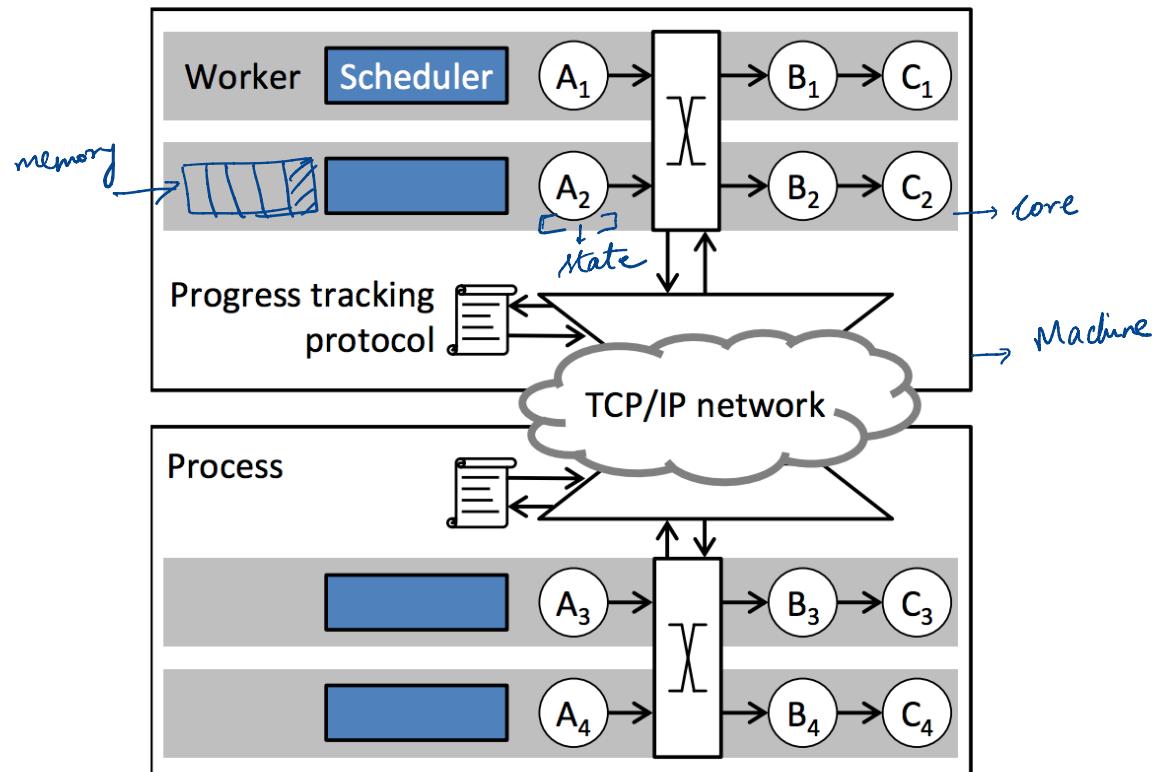
Logical graph



Workers communicate using
Shared Queue

Batch messages delivered
Account for cycles

Vertex single threaded



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

} when is it safe
to notify
send updates to
all workers before
updating locally

Optimizations

Batch updates and broadcast

Use projected timestamps from logical graph

A₁
A₂
A₃} Physical

logical → A

FAULT TOLERANCE

mutable data in vertices

not just failed ones!

Checkpoint

Log data as computation goes on

Write a full checkpoint on demand

Pause worker threads

Flush message queues OnRecv

Restore

Reset all workers to checkpoint

Reconstruct state

Resume execution

→ re-doing
computation)

Inputs can be
replayed !!

Actors in Ray / Pytorch

How frequently
checkpoint



Overhead vs.

Time to
recovery

MICRO STRAGGLERS

What is different from stragglers in MapReduce?

stateful vertices ← IMPORTANT
speculative / retry is difficult

Sources of stragglers

Network

Concurrency

Garbage Collection



Preventive measures
to avoid stragglers
rather than mitigate them.

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 */); → feed inputs  
input.OnCompleted();
```

Streaming

Scope / Spark ?

Timely Dataflow

SUMMARY

Stream processing → Increasingly important workload trend

Timely dataflow: Principled approach to model batch, streaming together

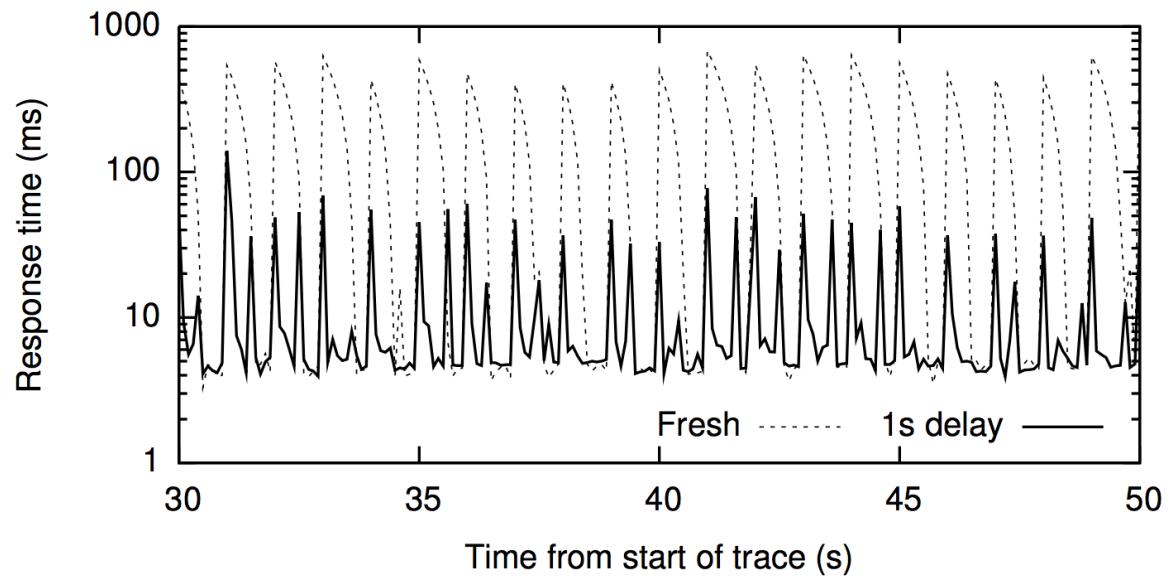
Vertex message model

- Compute frontier
- Distributed progress tracking

long running stateful
vertices

DISCUSSION

<https://forms.gle/qn8AzxHMT9VtyEc37>



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!

