# CS 744: BIG DATA SYSTEMS

Shivaram Venkataraman

Fall 2018

# ADMINISTRIVIA

- Assignment 2 grades

- Midterm review session on Nov 2 at 5pm at 1221 CS

- Course Project Proposal feedback

# STREAM PROCESSING

# DASHBOARDS

## Sales Dashboard

| Total Sales | Number of Deals | Avg Deal Size | Rev. per Salesperson |
|---|---|---|---|
| $3,256.8M | 17,164 | $189,545 | $20.5M |

**Week of Date Closed**
December 6, 200   December 25, 20

**Region**
(All)

**Country**
(All)
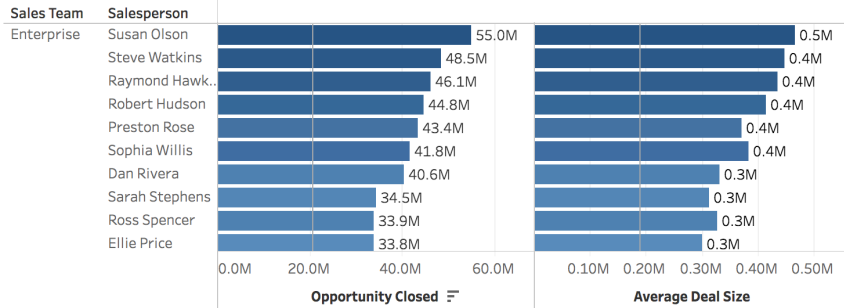
**Sales Team**
- (All)
- Small and Midmarket
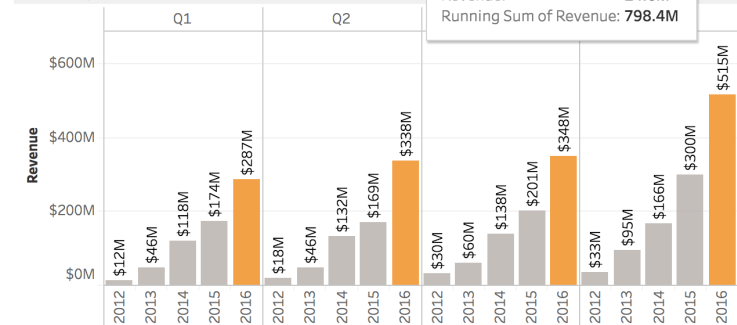- Enterprise

**Avg Deal Size/Salesperson**
$147,043   $336,519

### Revenue Over Time



### Sales Team Performance

| Sales Team | Salesperson | Opportunity Closed | Average Deal Size |
|---|---|---|---|
| Enterprise | Susan Olson | 55.0M | 0.5M |
| | Steve Watkins | 48.5M | 0.4M |
| | Raymond Hawk.. | 46.1M | 0.4M |
| | Robert Hudson | 44.8M | 0.4M |
| | Preston Rose | 43.4M | 0.4M |
| | Sophia Willis | 41.8M | 0.4M |
| | Dan Rivera | 40.6M | 0.3M |
| | Sarah Stephens | 34.5M | 0.3M |
| | Ross Spencer | 33.9M | 0.3M |
| | Ellie Price | 33.8M | 0.3M |

### Revenue by Quarter

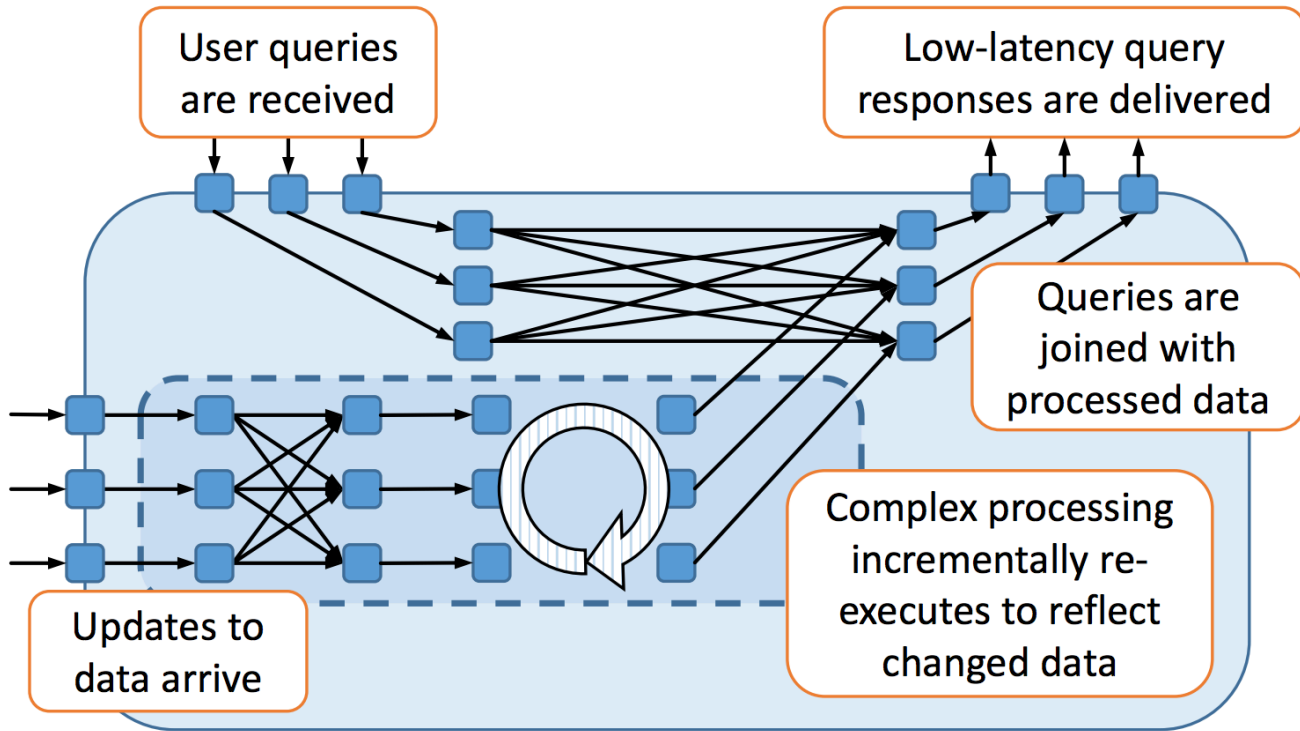**Week of September 4, 2016**
Revenue:                         14.6M
Running Sum of Revenue: 798.4M

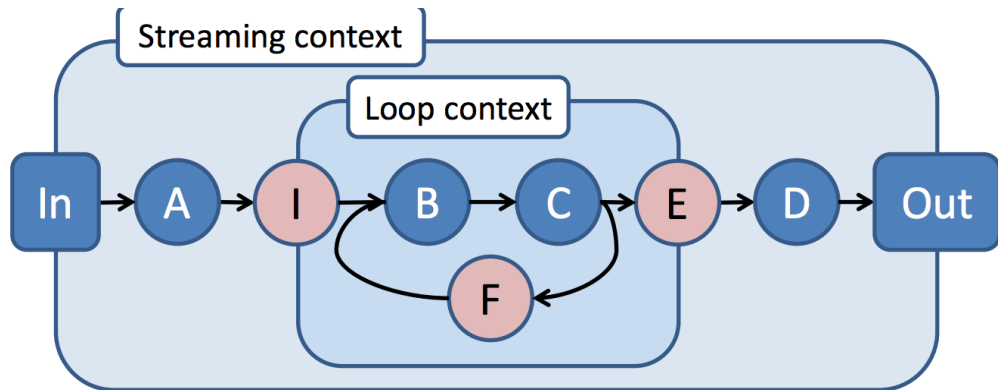# REAL-TIME ANALYSIS

# STREAMING + ITERATIVE COMPUTATION

# TIMELY DATAFLOW

# TIMELY DATAFLOW



epoch        loop counters

$$\text{Timestamp} : (\overbrace{e \in \mathbb{N}}^{\text{epoch}}, \overbrace{\langle c_1, \ldots, c_k \rangle \in \mathbb{N}^k}^{\text{loop counters}})$$

| Vertex | Input timestamp | Output timestamp |
|---|---|---|
| Ingress | $(e, \langle c_1, \ldots, c_k \rangle)$ | $(e, \langle c_1, \ldots, c_k, 0 \rangle)$ |
| Egress | $(e, \langle c_1, \ldots, c_k, c_{k+1} \rangle)$ | $(e, \langle c_1, \ldots, c_k \rangle)$ |
| Feedback | $(e, \langle c_1, \ldots, c_k \rangle)$ | $(e, \langle c_1, \ldots, c_k + 1 \rangle)$ |

# VERTEX API

Receiving Messages

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

    v.OnNotify(t : Timestamp)

Sending Messages

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

    this.NotifyAt(t : Timestamp)

Conditions

    OnNotify(t) invoked after
    all OnRecv(e, r, t') for all $t' \leq t$

    SendBy or NotifyAt only called
    with $t' >= t$

# IMPLEMENTING TIMELY DATAFLOW

Need to track when it is safe to notify

Path Summary

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

| Operation | Update |
|---|---|
| $v.\text{SENDBY}(e, m, t)$ | $\text{OC}[(t, e)] \leftarrow \text{OC}[(t, e)] + 1$ |
| $v.\text{ONRECV}(e, m, t)$ | $\text{OC}[(t, e)] \leftarrow \text{OC}[(t, e)] - 1$ |
| $v.\text{NOTIFYAT}(t)$ | $\text{OC}[(t, v)] \leftarrow \text{OC}[(t, v)] + 1$ |
| $v.\text{ONNOTIFY}(t)$ | $\text{OC}[(t, v)] \leftarrow \text{OC}[(t, v)] - 1$ |

Scheduler

    Occurrence and Precursor count

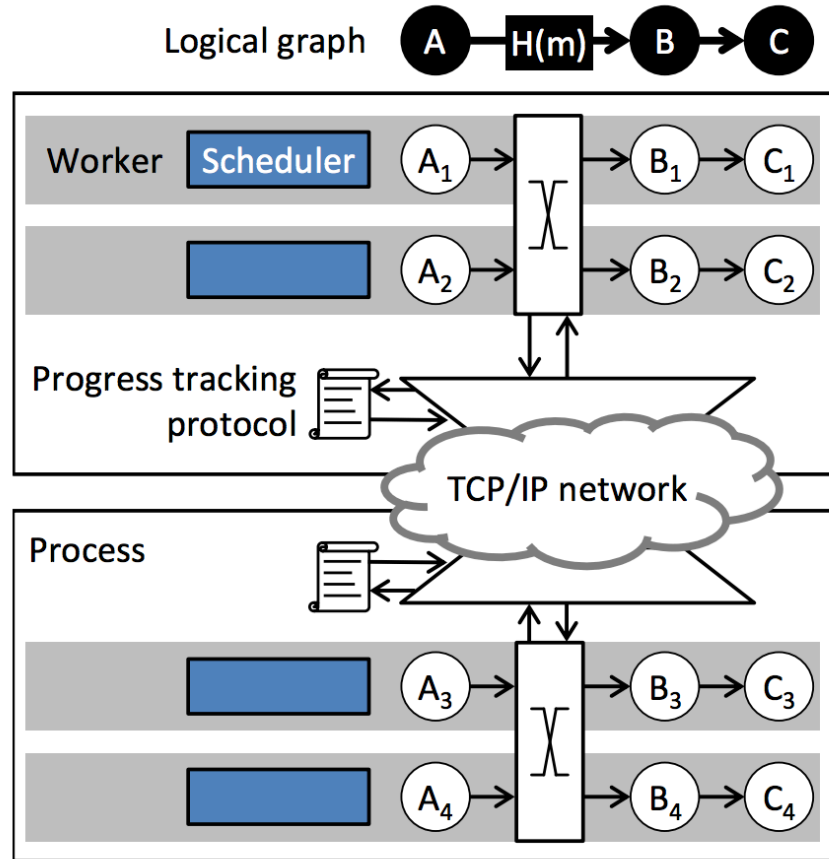    Precursor count = 0 → Frontier

# ARCHITECHTURE

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 $\in$ Pointstamp, $\delta \in Z$)

    Local frontier tracks global frontier

Optimizations

    Batch updates and broadcast

    Use projected timestamps from logical graph

# FAULT TOLERANCE

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

Trade-off between mutable updates and recovery time!

# MICRO STRAGGLERS

Networking

    - Disable Nagle's algorithm

    - Reduce TCP retransmission window


Concurrency

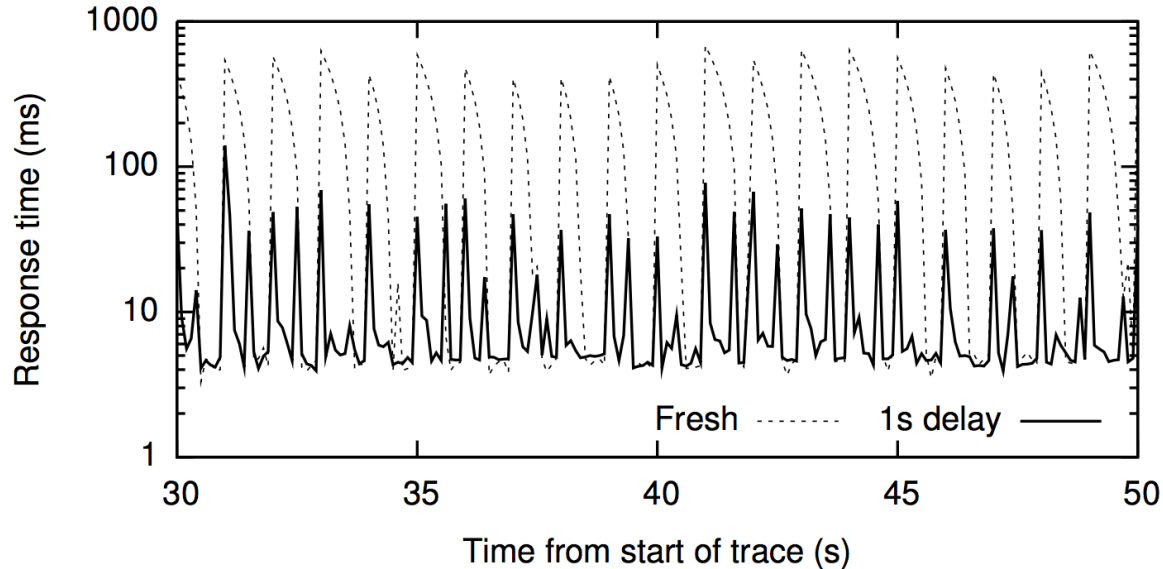    - Reduce clock granularity to avoid spin lock delay


Garbage Collection

    - Arrays of value types (similar to Plain-Old Java Objects)

    - Buffer pool

# DIFFERENTIAL DATAFLOW

```csharp
// 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();
```

# END-TO-END EXAMPLE



Stream of incoming tweets: username, raw tweet
Incremental connected components, top hashtag computation
Queries: username → return top hashtag from their component

# SUMMARY

Stream processing → Increasingly important workload trend

Timely dataflow: Principled approach to model batch, streaming together

Vertex message model
- Compute frontier
- Distributed progress tracking