

NAIAD (TIMELY DATAFLOW) & STREAMING SYSTEMS

CS 848: MODELS AND APPLICATIONS OF DISTRIBUTED DATA SYSTEMS

MON, NOV 7TH 2016

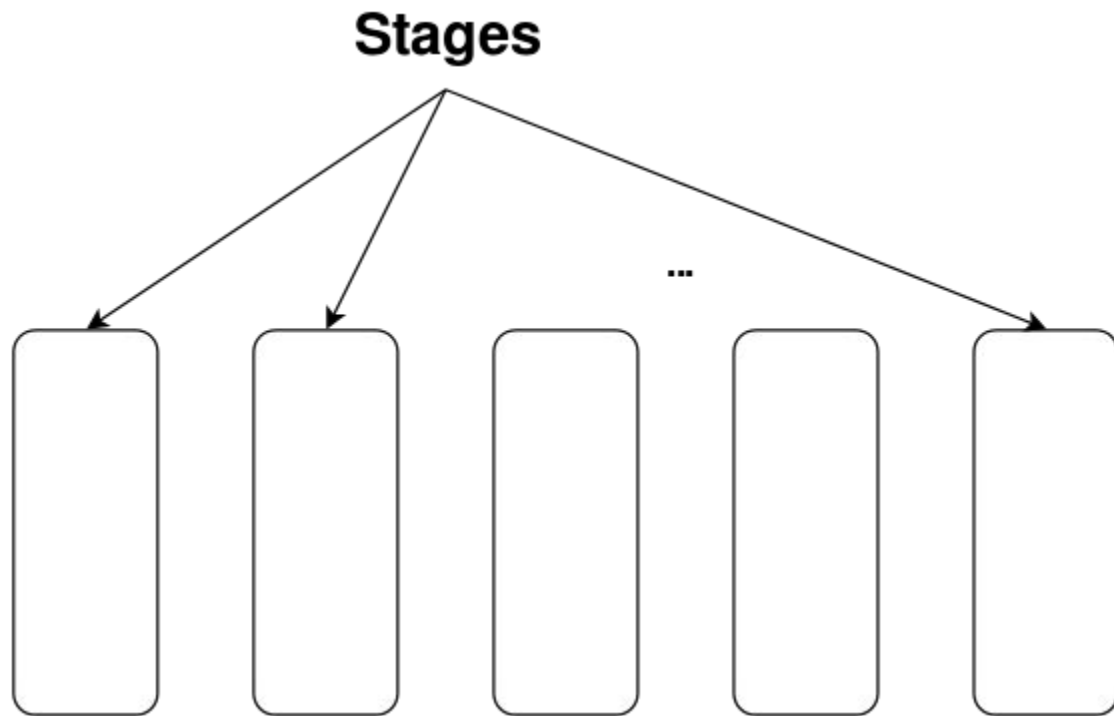
Amine Mhedhbi

WHAT IS 'TIMELY' DATAFLOW ?!

WHAT IS ITS SIGNIFICANCE?

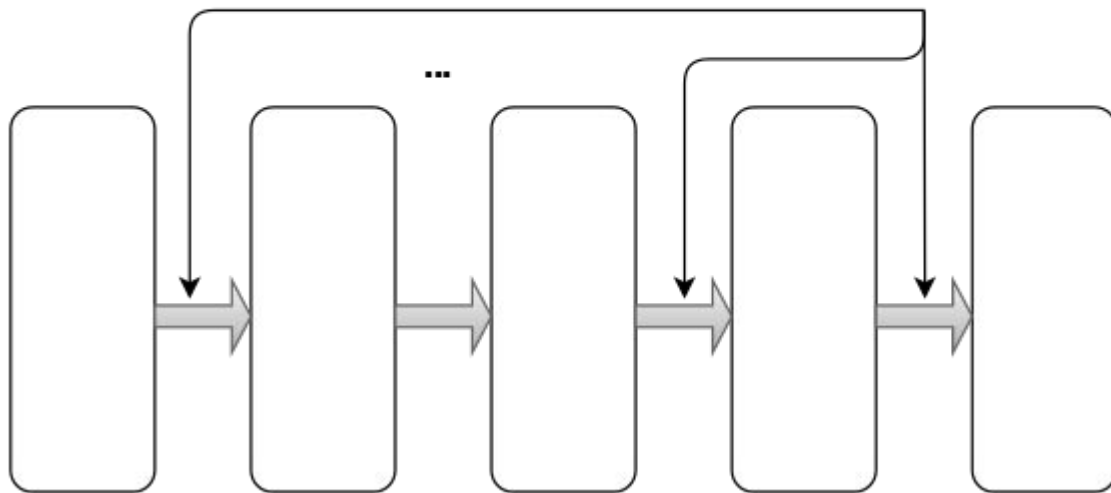
DATAFLOW ?!

DATAFLOW?!



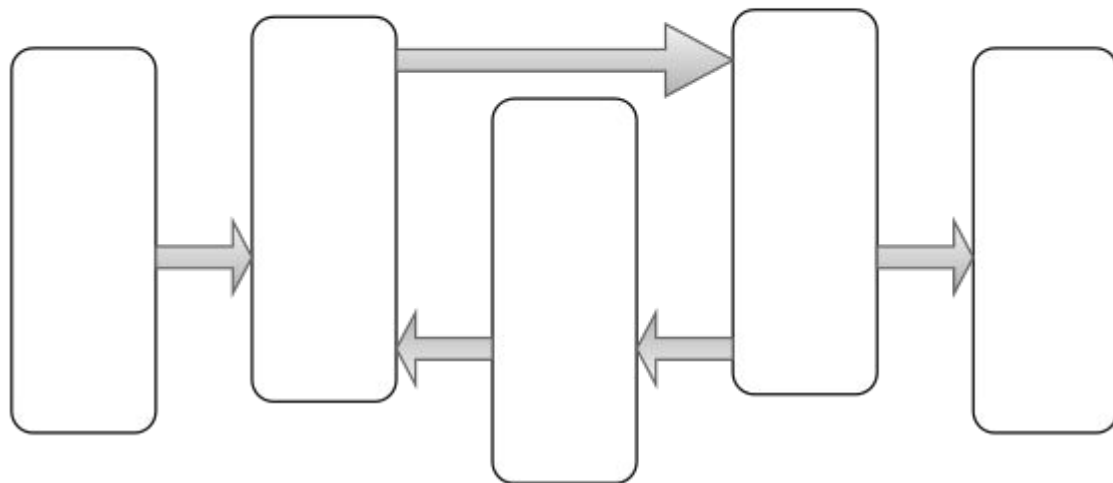
DATAFLOW?!

Stages makeup a directed graph
where data flows along
connectors



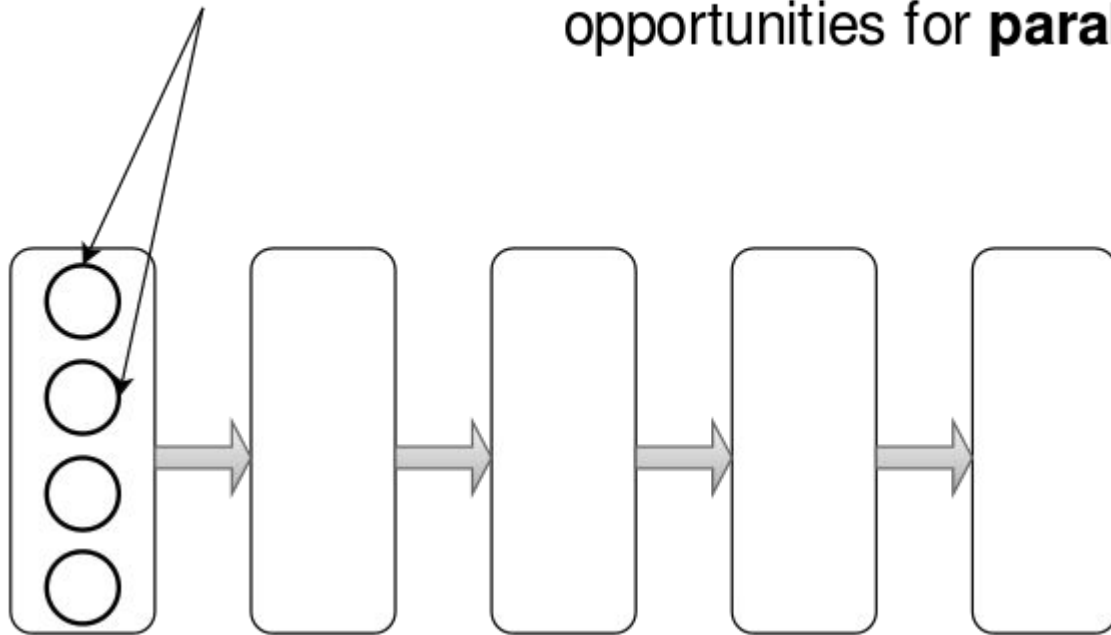
DATAFLOW?!

Stages makeup a directed graph
so it allows **iterations**



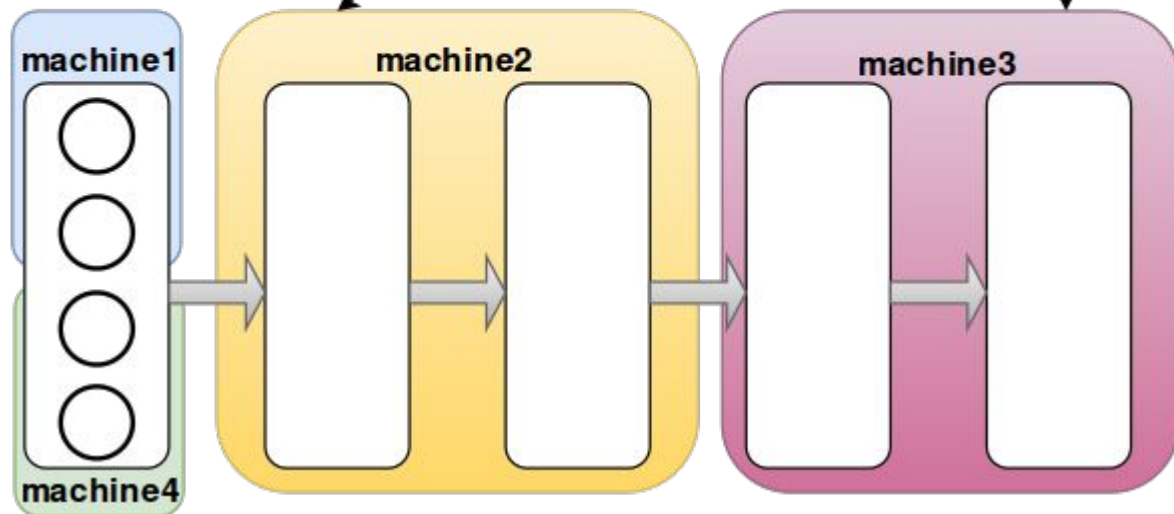
DATAFLOW?!

Physical **stages** is where systems exploit opportunities for **parallelism**



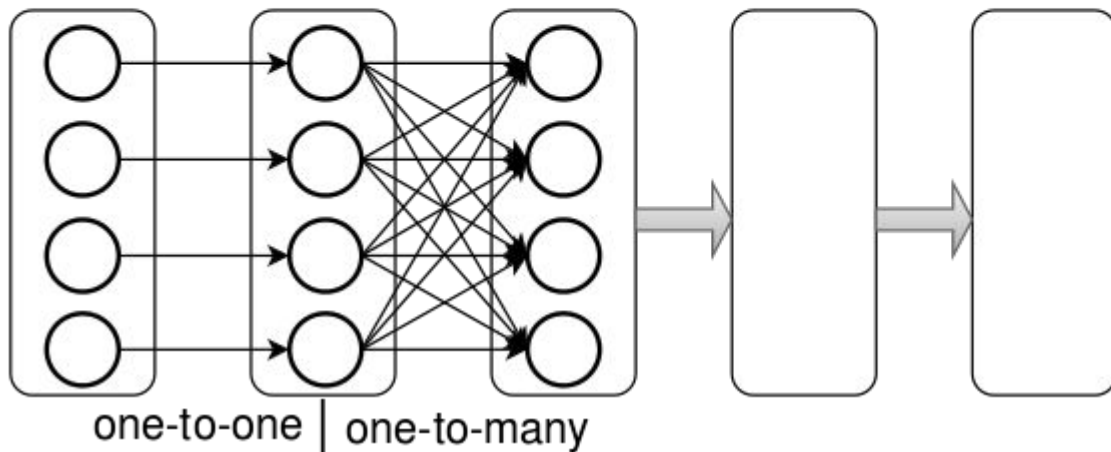
DATAFLOW?!

Also **stages** can run on different machines in the cluster



DATAFLOW?!

With **stages** parallelized as such.
connectors presenting the flow,
have different types of **mapping**

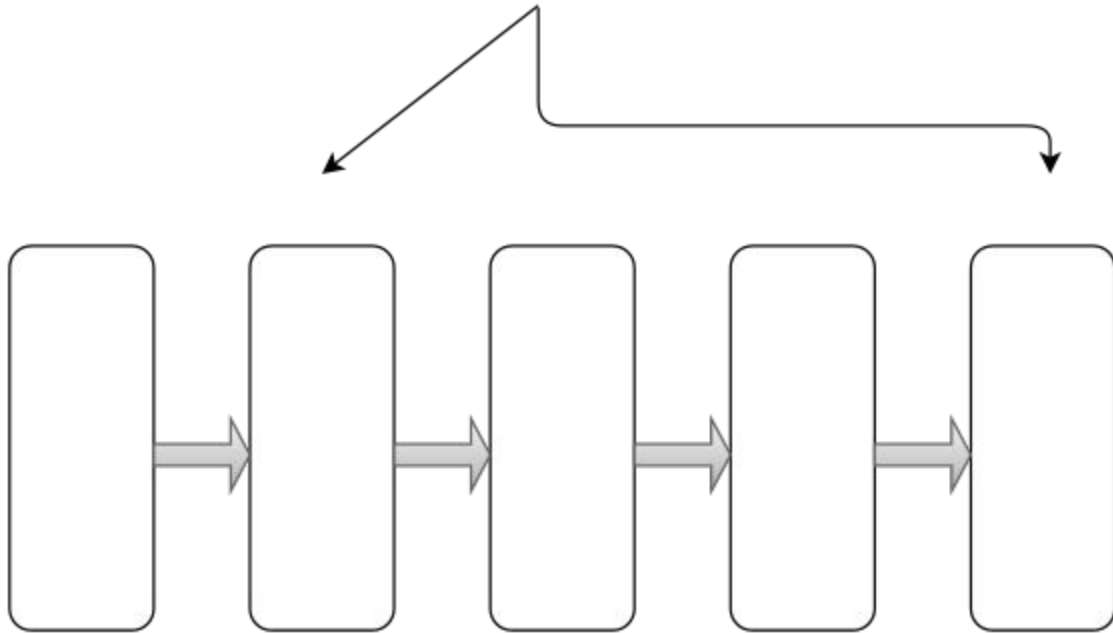


DATAFLOW

- Batch Processing e.g. MapReduce, Spark
- Asynchronous Processing e.g. Storm, MillWheel
- Variations for Graph Processing e.g. Pregel, GraphLab

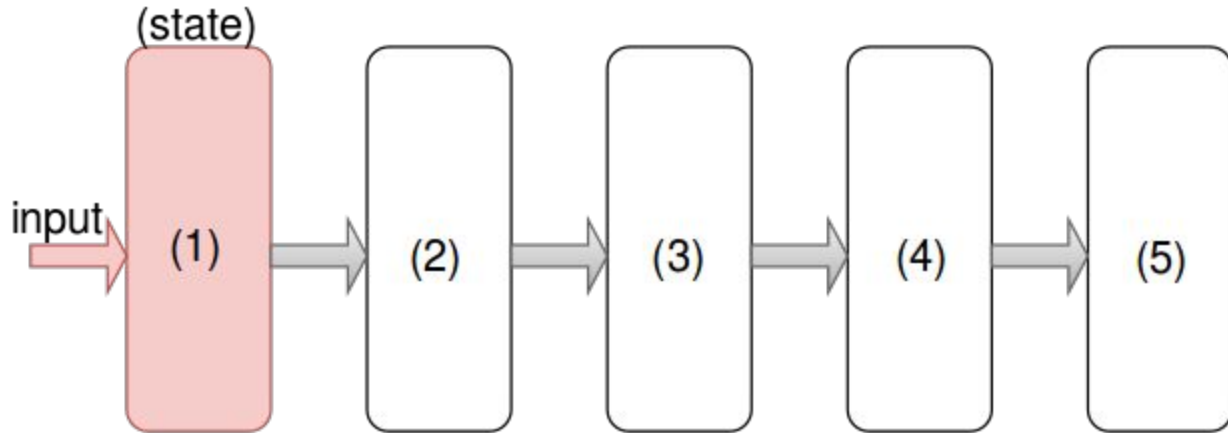
DATAFLOW: BATCH PROCESSING

stages require synchronization



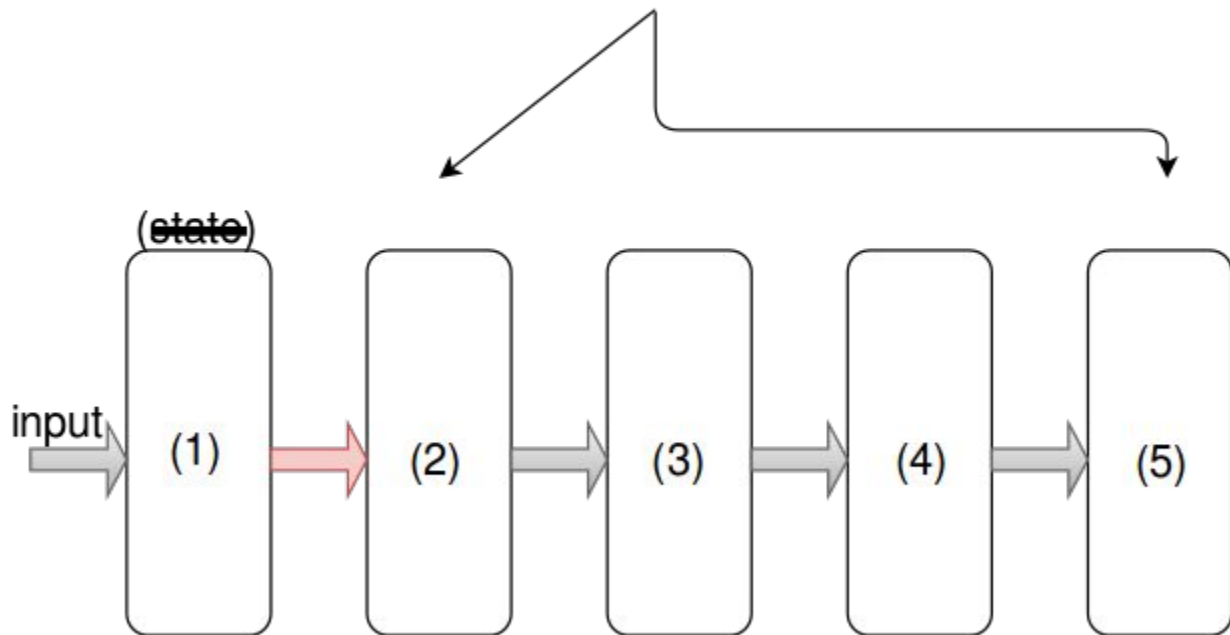
DATAFLOW: BATCH PROCESSING

stages require synchronization

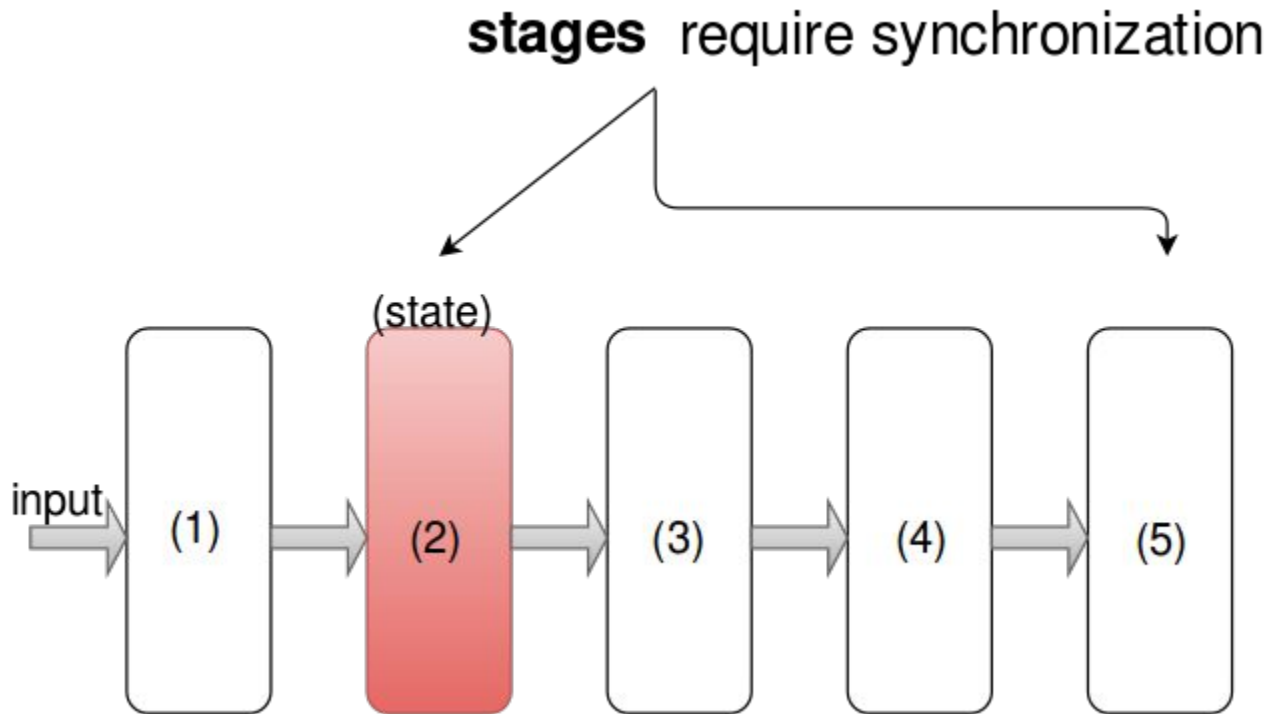


DATAFLOW: BATCH PROCESSING

stages require synchronization

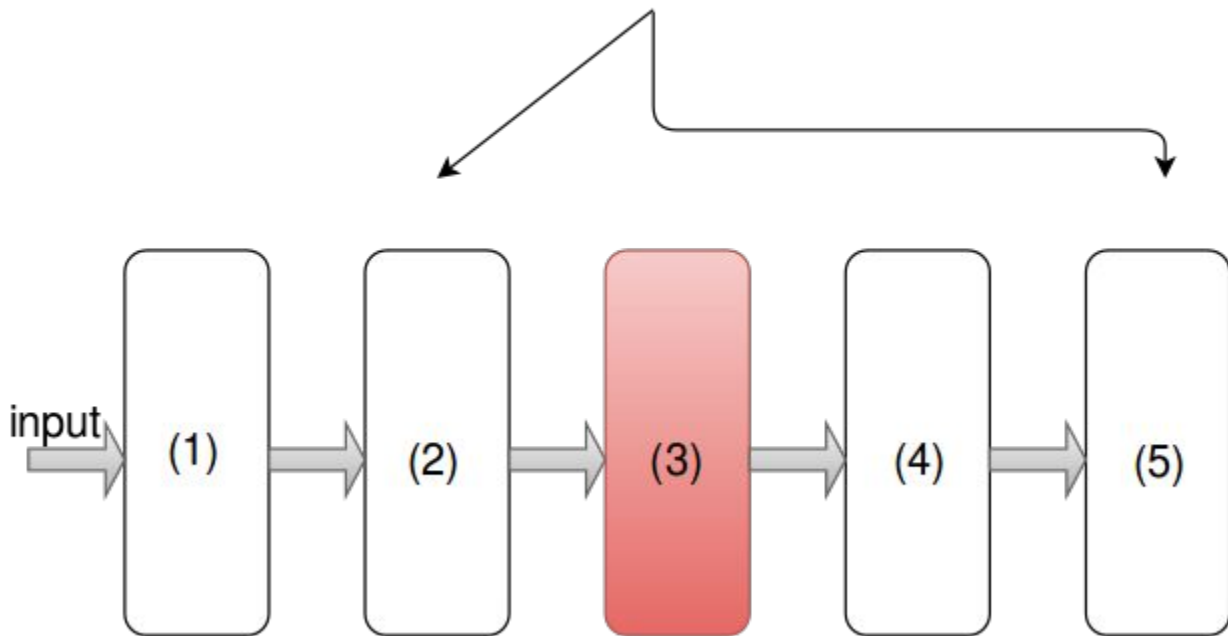


DATAFLOW: BATCH PROCESSING



DATAFLOW: BATCH PROCESSING

stages require synchronization

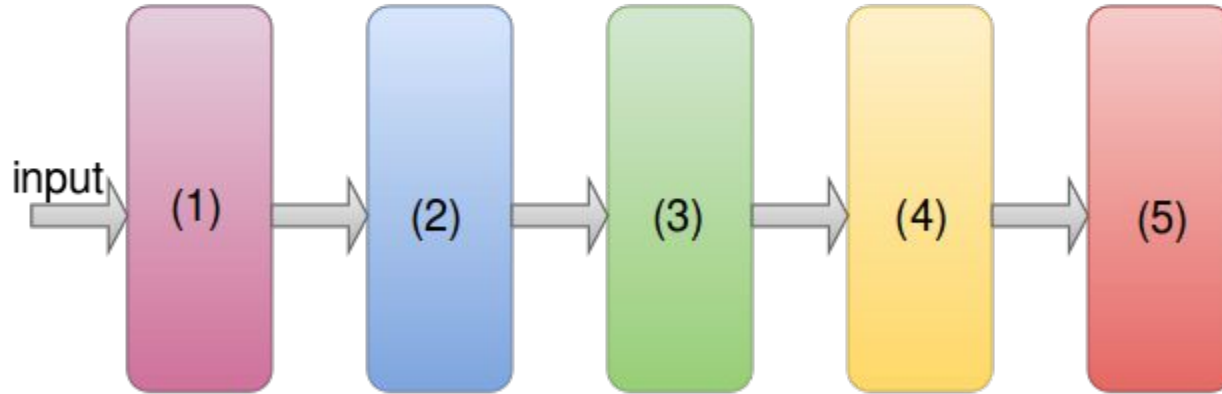


DATAFLOW: BATCH PROCESSING

- Iterations make use of synchronization.
- The cost is latency.

DATAFLOW: ASYNCHRONOUS PROCESSING

stages do NOT require synchronization
all stages are active and output data
after processing input data.



DATAFLOW: ASYNCHRONOUS PROCESSING

- Compared with batch:
 - latency is lower.
 - Aggregations are incremental and data changes over time.
- More efficient for distributed systems.
 - Stages do not need coordination.
- Correspondence between input & output is lost.

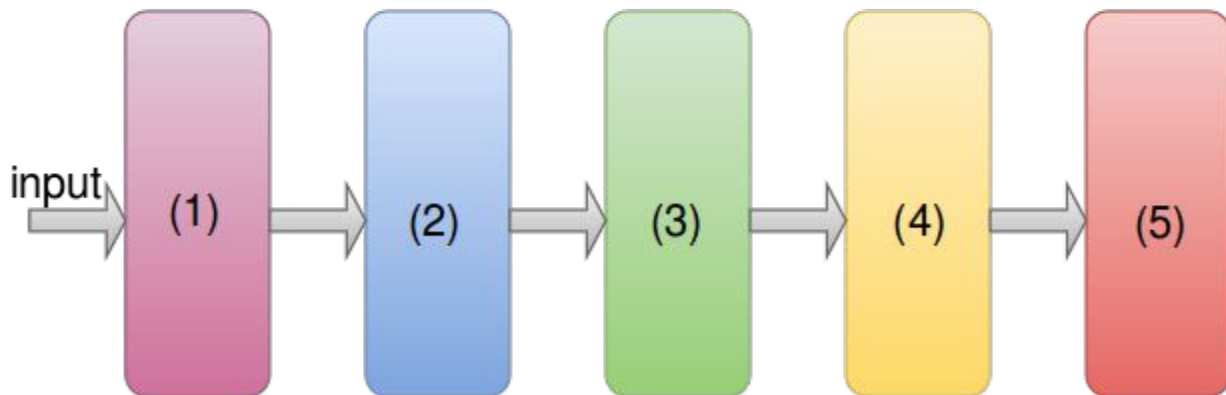
SO, WHAT IS (NAIAD)
TIMELY DATAFLOW ?!

TIMELY DATAFLOW?!

stages are asynchronous.

They can synchronize if needed.

Make use of **logical timestamps**.



TIMELY DATAFLOW

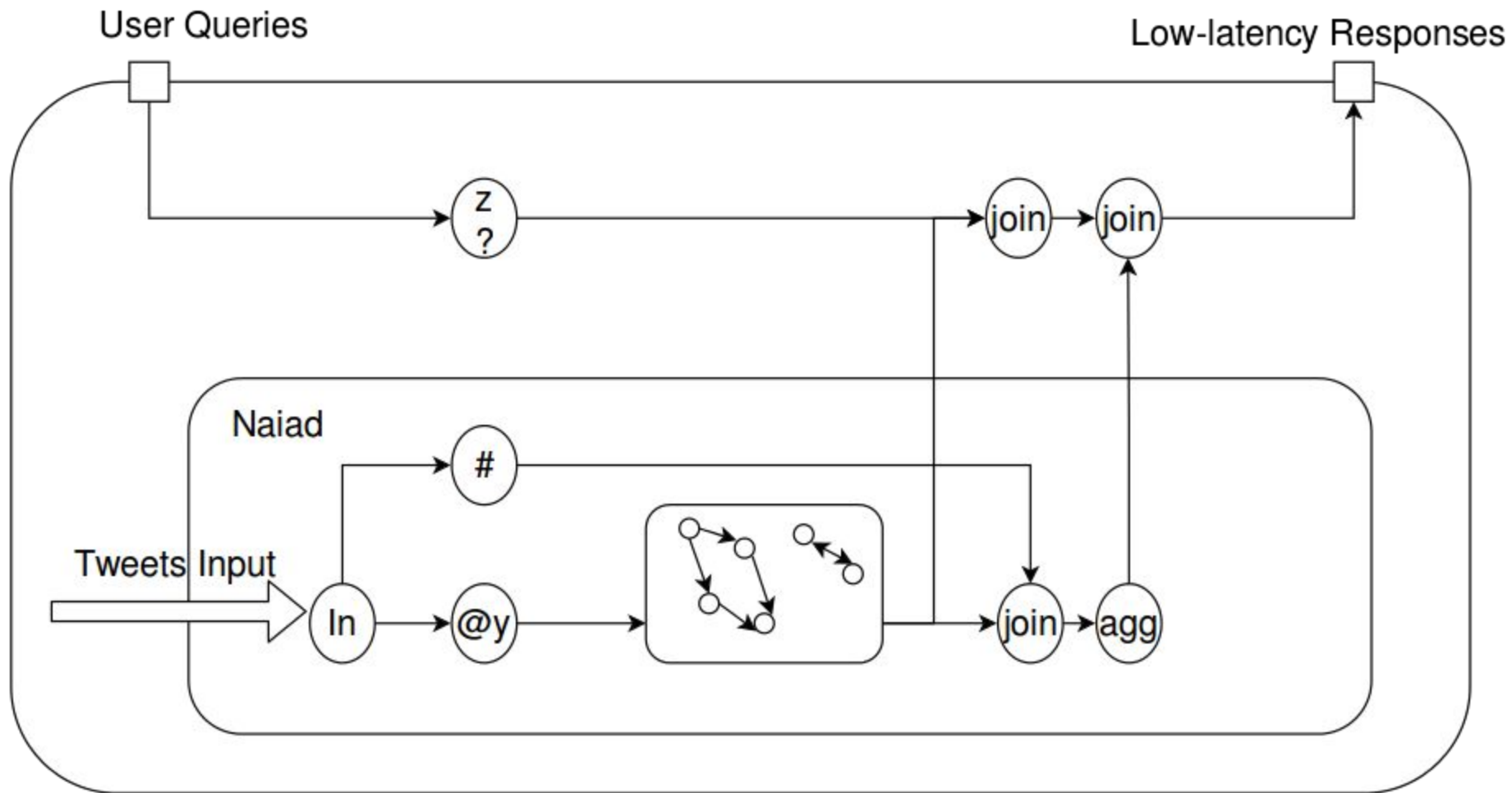
- Reconcile both models batch and async.
- Low-latency and high-throughput.

WHERE DOES NAIAD FIT?!

NAIAD?!

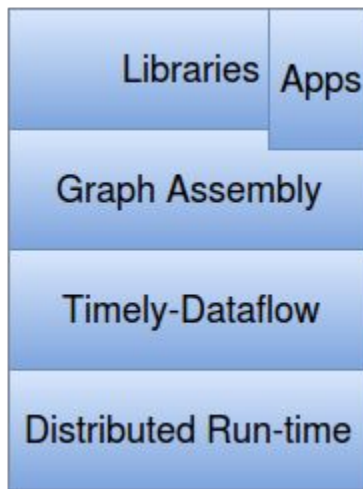
- It is the prototype built by Microsoft Research underlying **Timely dataflow** Computational model.
- Iterative and incremental computations.
- The logical timestamps allow **coordination**.
- Provides efficiency, maintainability and simplicity.

LET'S LOOK AT A COMPUTATIONAL
EXAMPLE



NAIAD?!

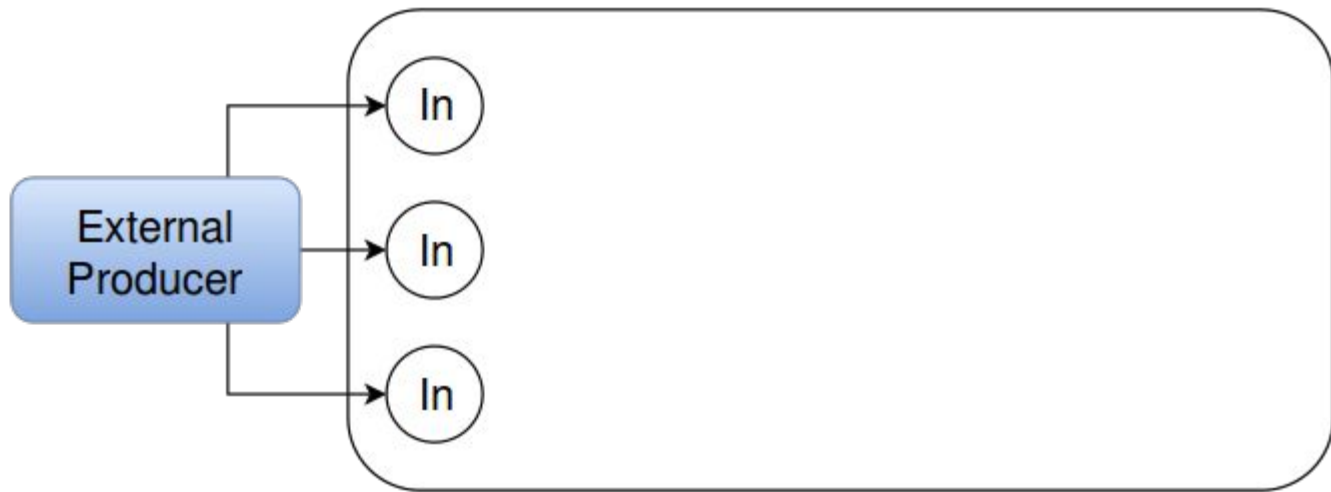
- It is the prototype built by Microsoft Research underlying **Timely dataflow** Computational model.



THE TIMELY DATAFLOW GRAPH STRUCTURE

GRAPH STRUCTURE

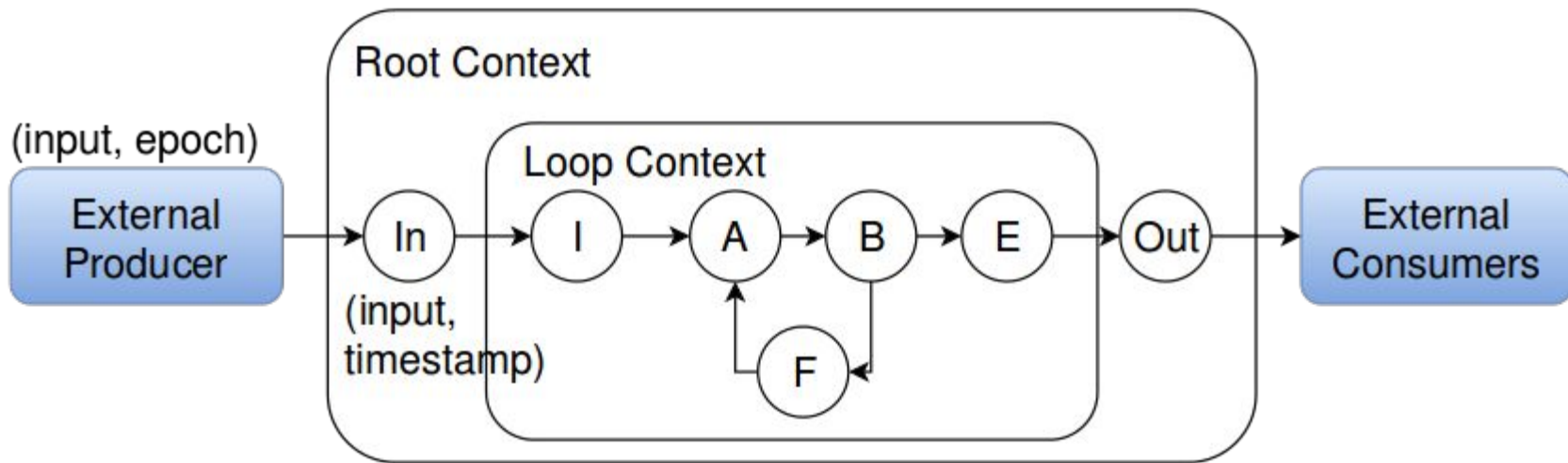
(input, integer epoch)



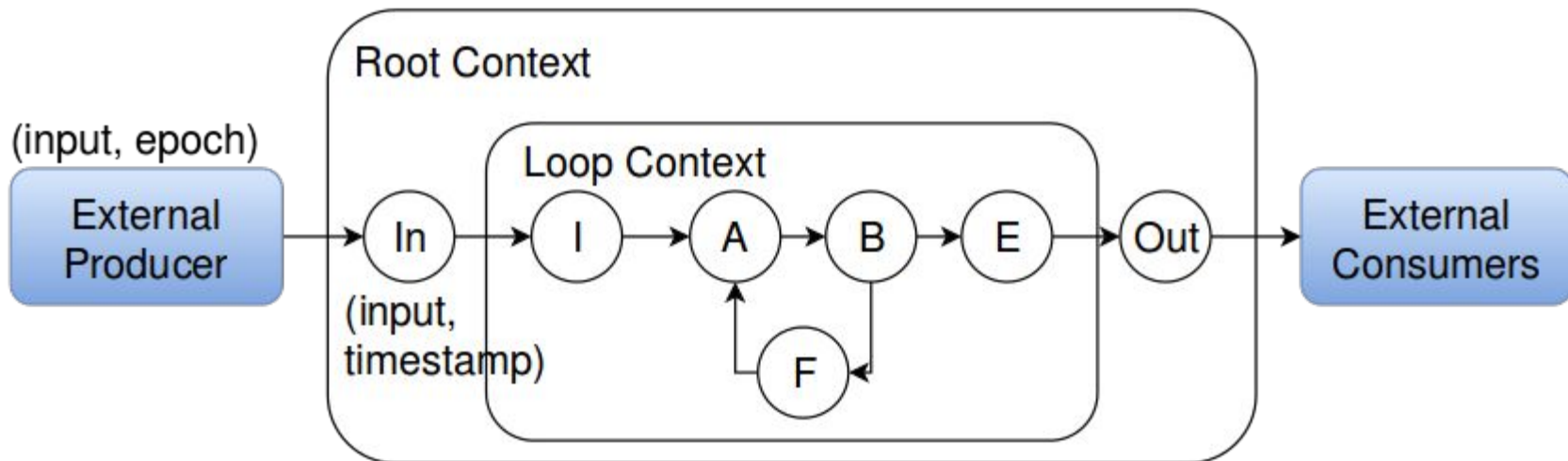
GRAPH STRUCTURE



GRAPH STRUCTURE



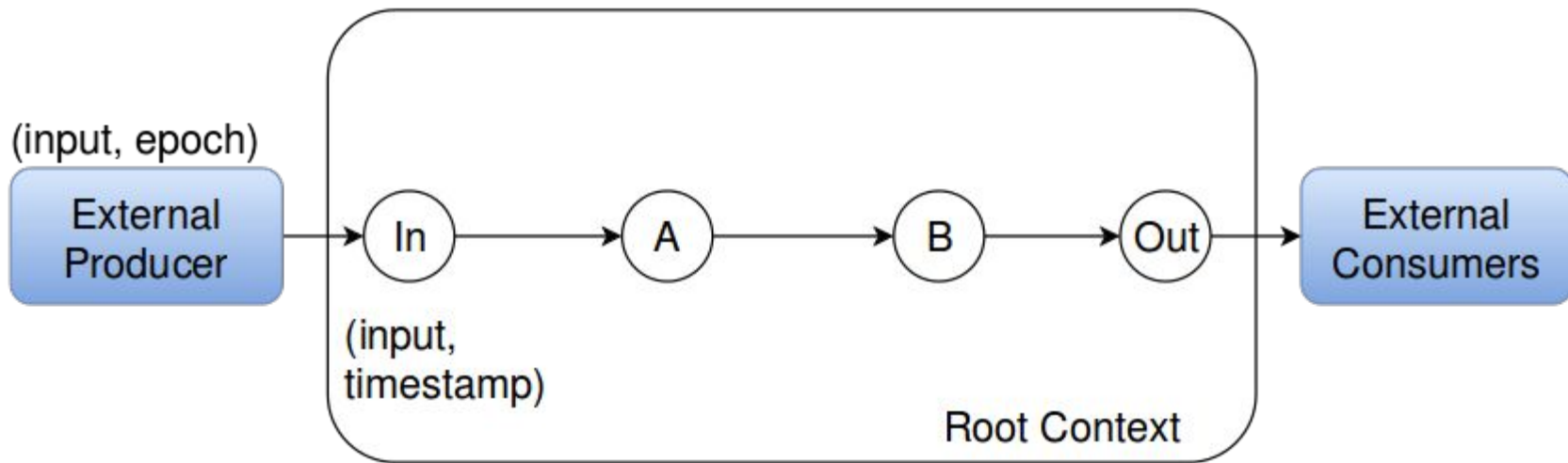
GRAPH STRUCTURE



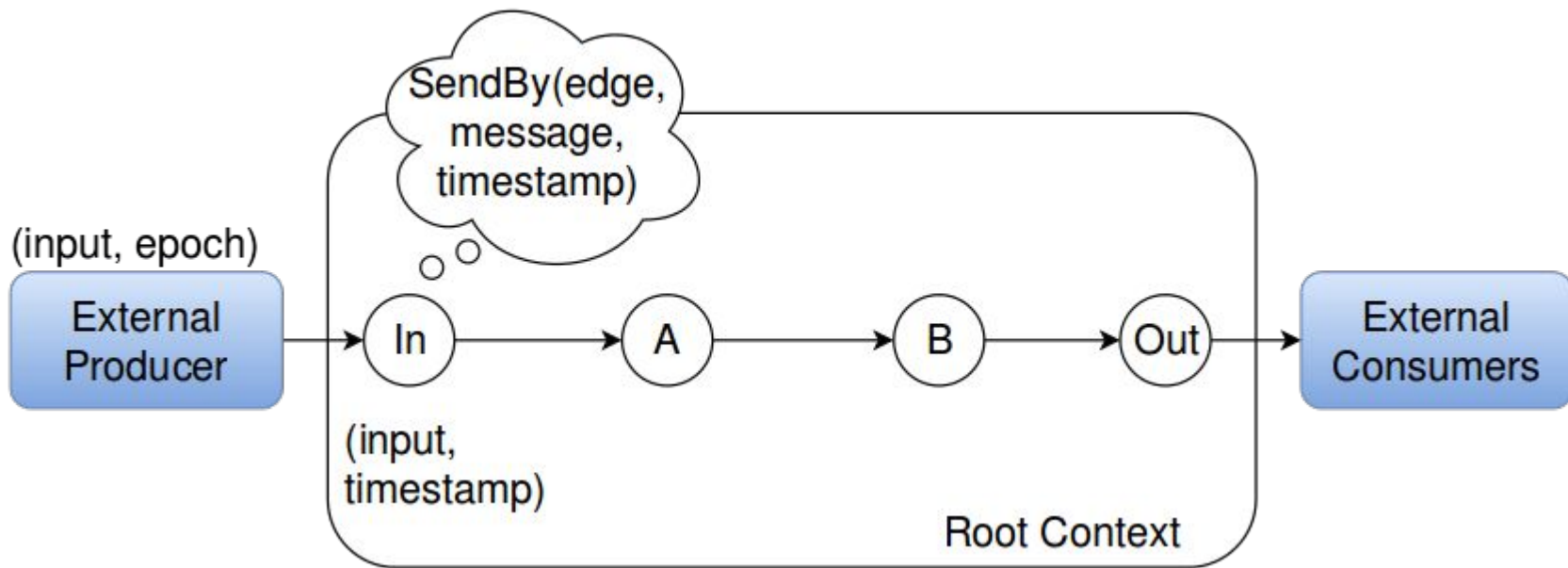
- input comes in as $(data, 0)$, $(data, 1)$, $(data, 2)$
 - Within a loop, I adds a loop counter so it is $(data, epoch, 0)$
F in each iteration increments the loop counter $(data, epoch, 1)$ etc.
E removes the loop counter and it is back to $(data, epoch)$

PROGRAMMING MODEL USING THE TIMESTAMPS

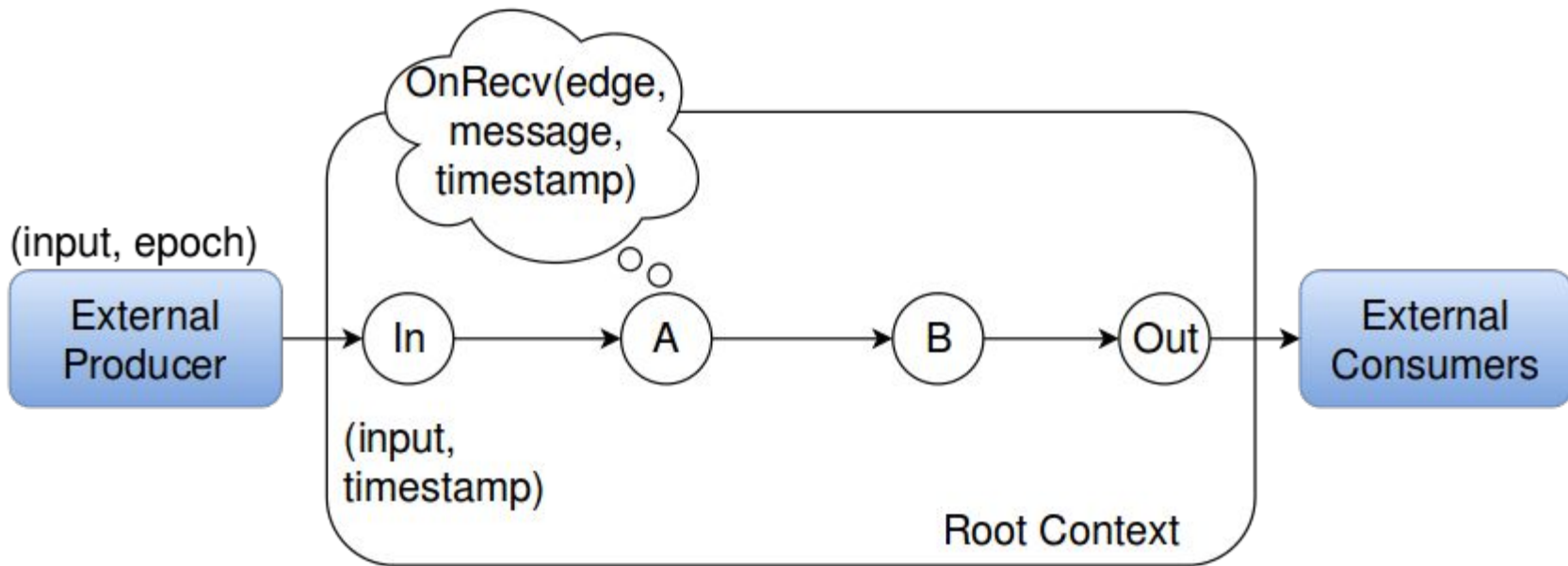
PROGRAMMING MODEL



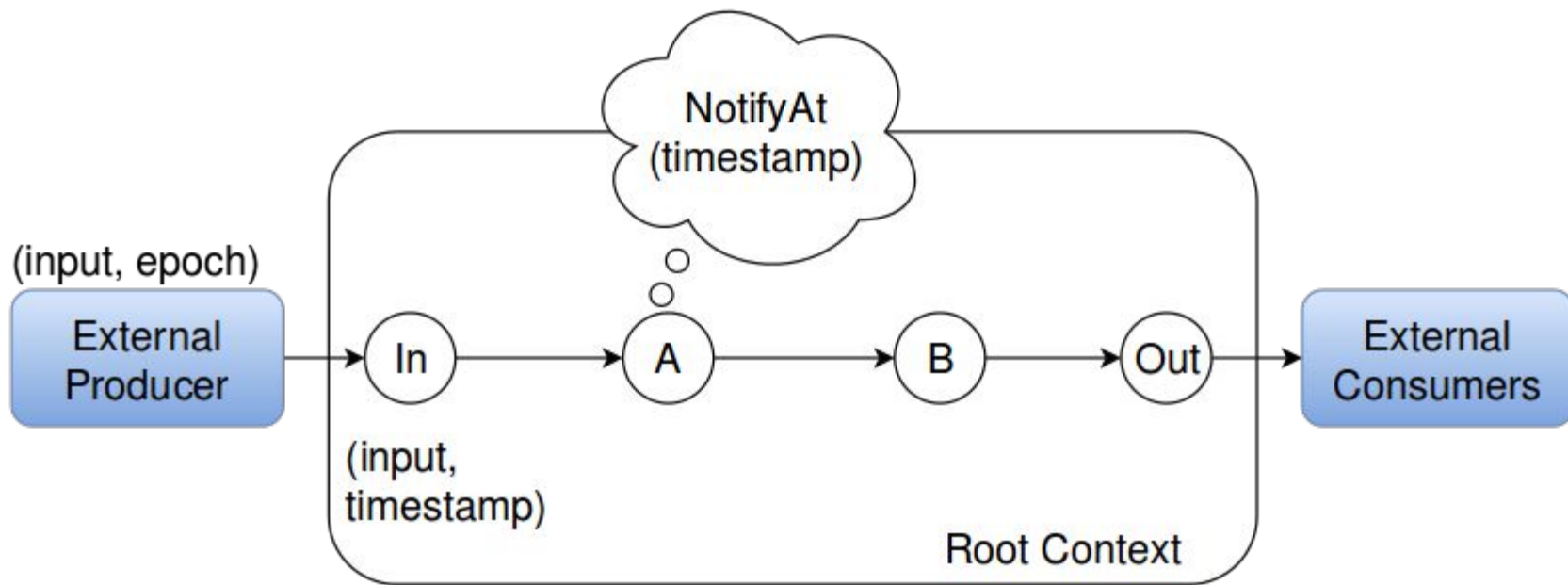
PROGRAMMING MODEL



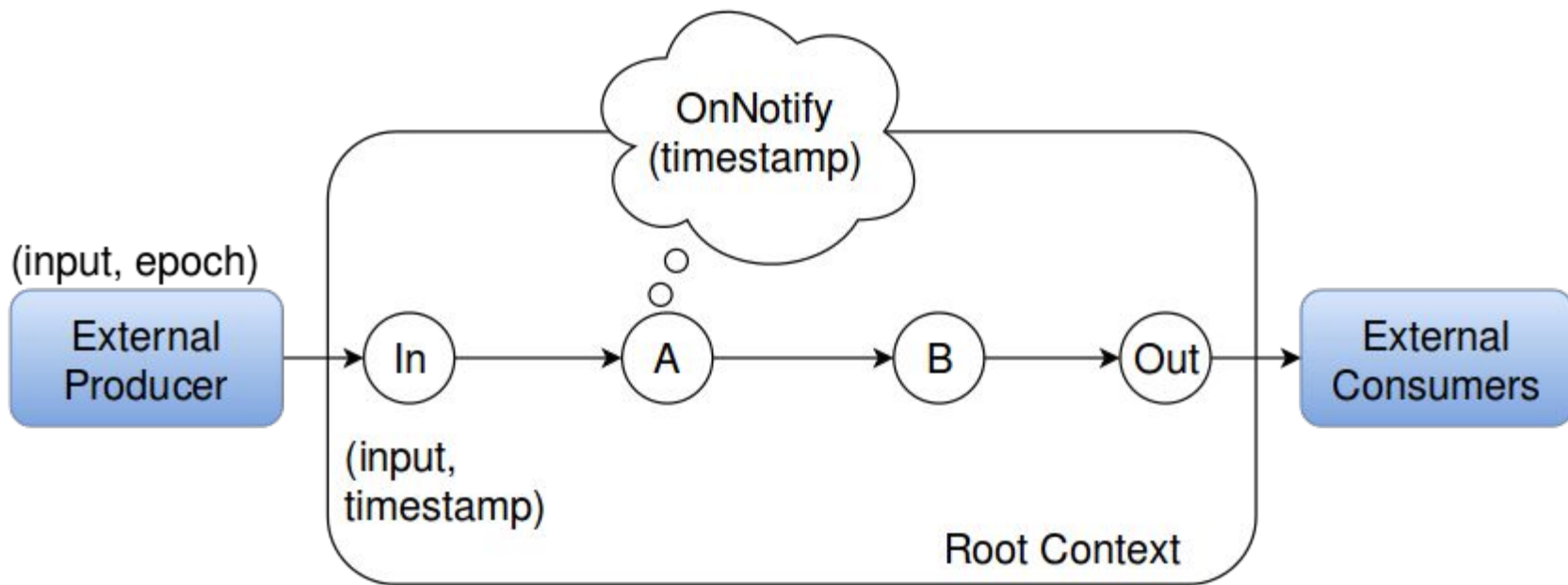
PROGRAMMING MODEL



PROGRAMMING MODEL



PROGRAMMING MODEL



PROGRAMMING MODEL SUMMARY

- `SendBy(edge, message, timestamp)`
- `OnRecv(edge, message, timestamp)`
- `NotifyAt(timestamp)`
- `OnNotify(timestamp)`

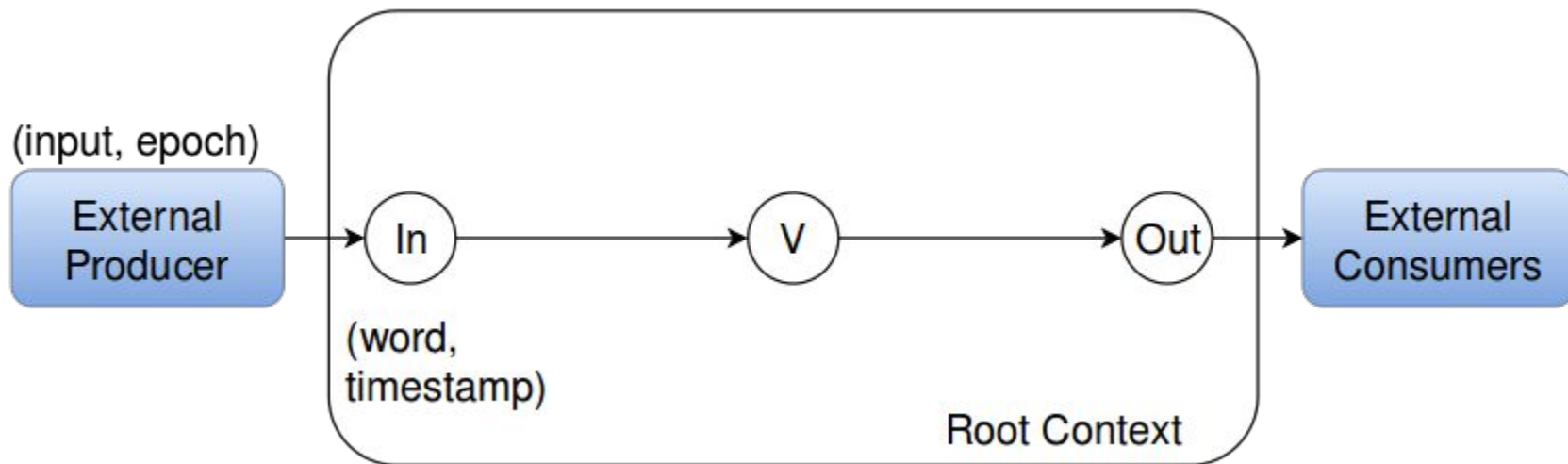
PROGRAMMING MODEL IN PRACTICE

NOTICE

- Project was discontinued in 2014.
Silicon Valley lab closed.
- The paper uses C#.
The latest one is open sourced and is in Rust.

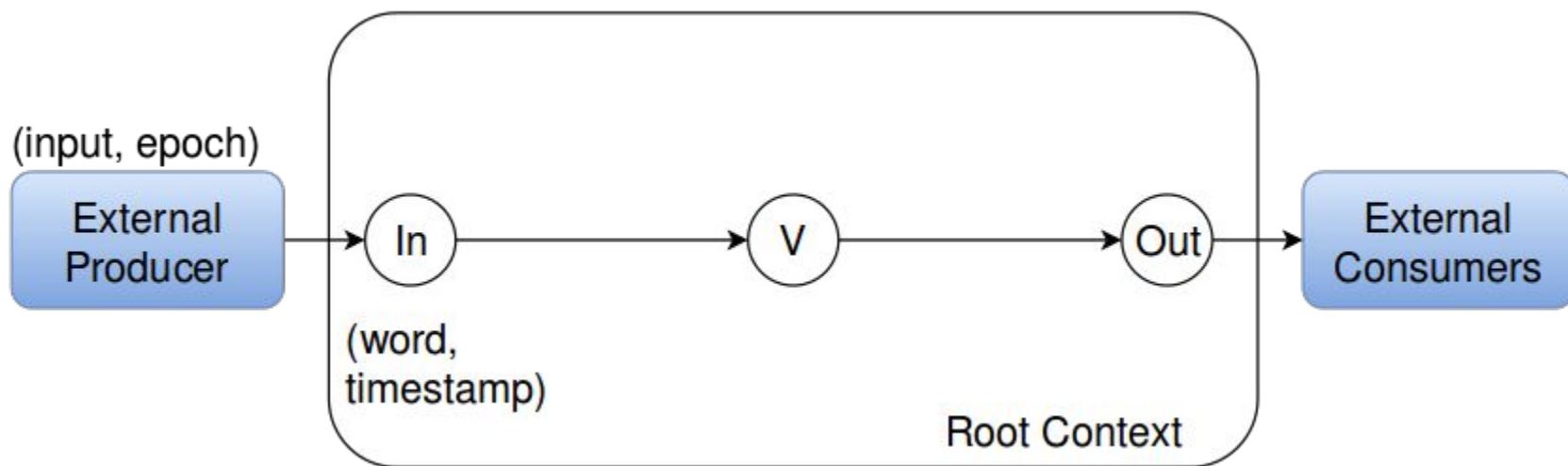
WORD COUNT EXAMPLE

```
Class V<Msg, Time>: Vertex<Time> { ... }
```



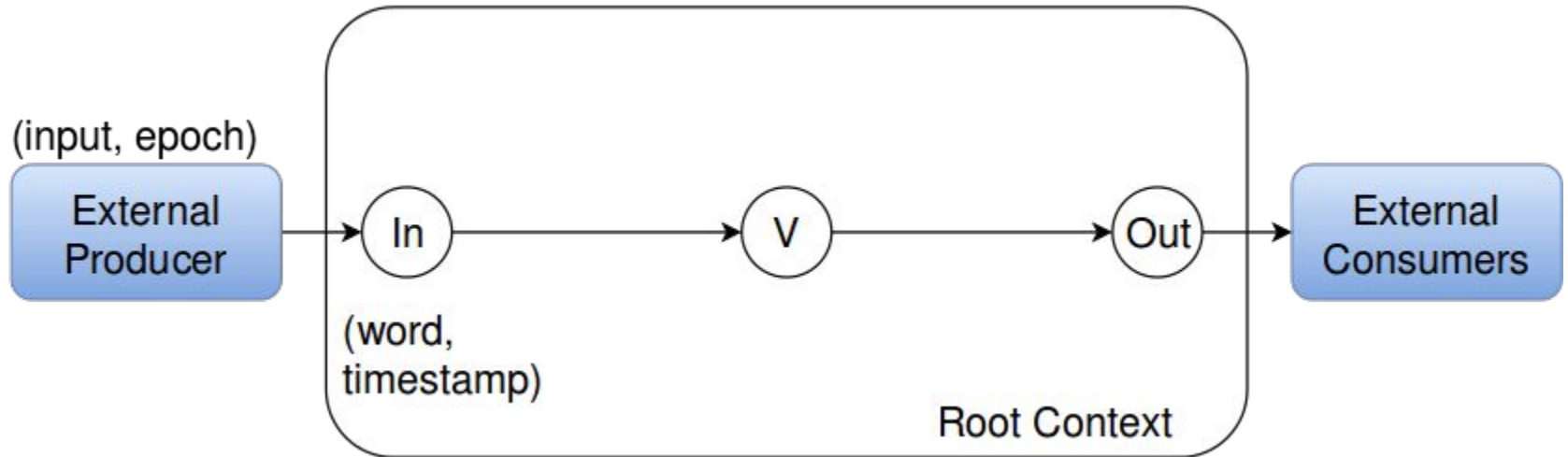
WORD COUNT EXAMPLE

```
{ Dict<Time, Dict<Msg, int> > counts; ... }
```



WORD COUNT EXAMPLE (2 DIFFERENT IMPLEMENTATIONS)

```
{ void OnRecv (Edge e, Msg m, Time t) { ... }  
  void OnNotify (Time t) { ... } }
```

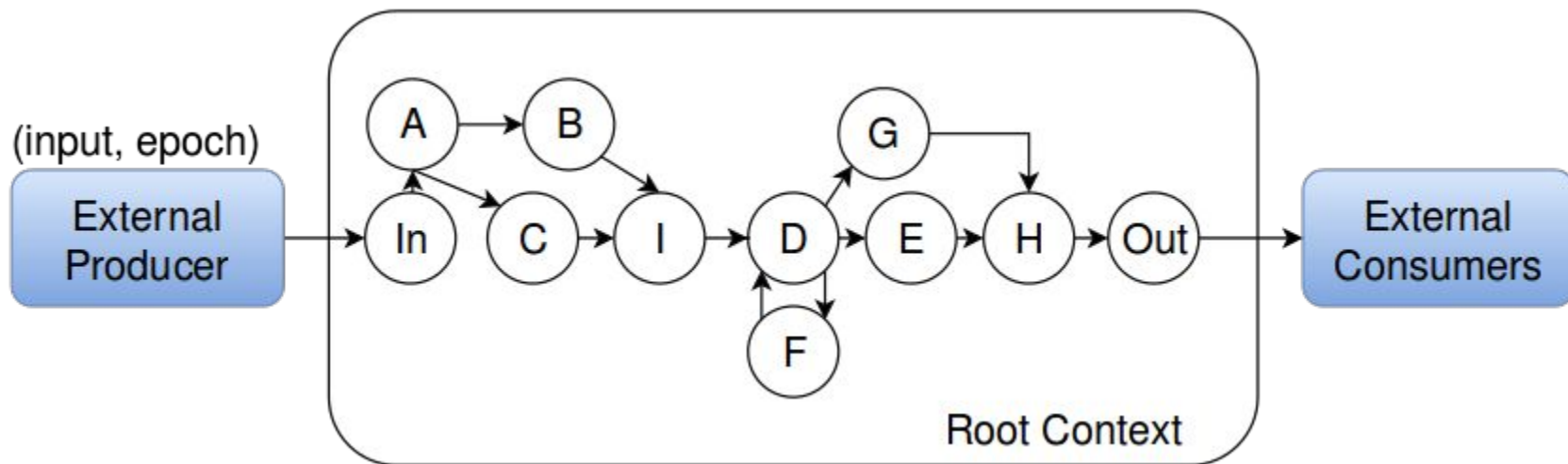


WRITING PROGRAMS IN GENERAL

- It is possible to write programs against the Timely Dataflow abstraction.
- It is possible to use libraries (MapReduce, Pregel, PowerGraph, LINQ etc.)
- In General:
 - Define Input, computational & Output vertices.
 - Create a timely dataflow graph using the appropriate interface.
 - Supply labeled data to input stages.
 - Stages follow a push-based model.

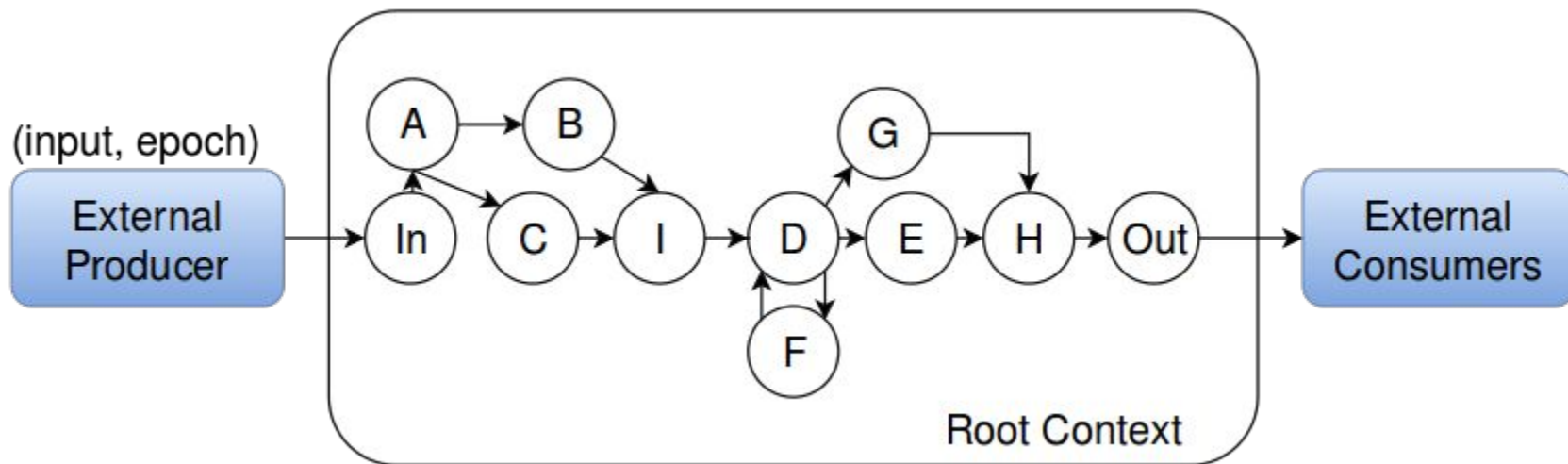
TIMELY GUARANTEES

HOW IS TIMELY DATAFLOW ACHIEVED



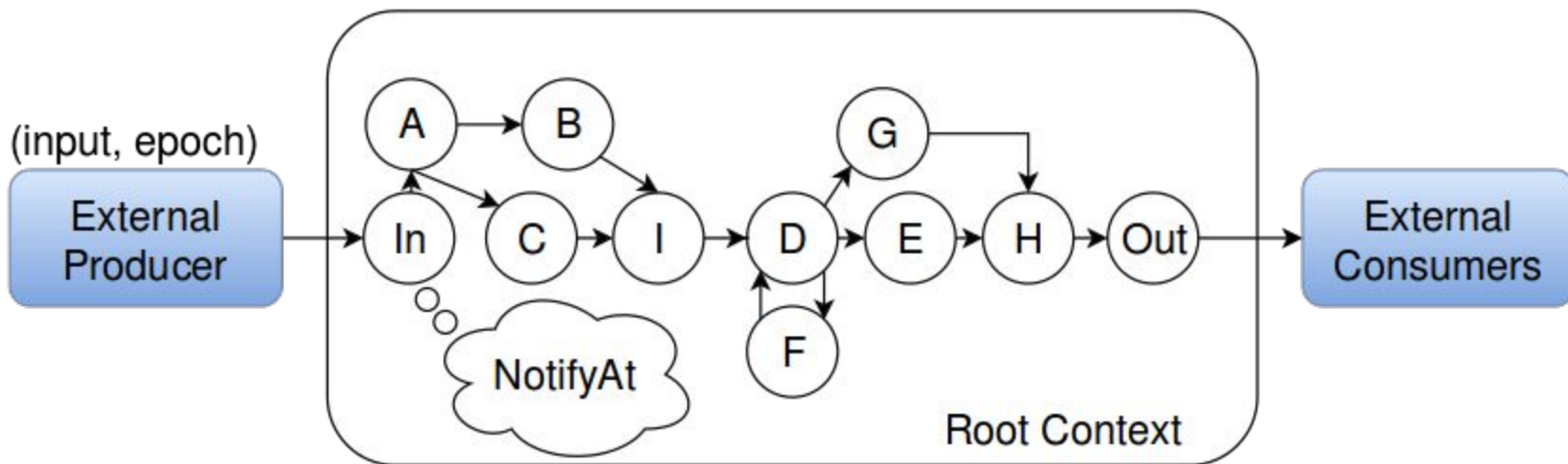
HOW IS TIMELY DATAFLOW ACHIEVED

- **Key point:** timestamps at which future message can occur depends on: 1. Unprocessed events & 2. Graph Structure.



HOW IS TIMELY DATAFLOW ACHIEVED

- **Pointstamp** of an event (timestamp, location: E or V)
 - SendBy \rightarrow Msg event of pointstamp (t, e)
 - NotifyAt \rightarrow Notif event of pointstamp (t, v)

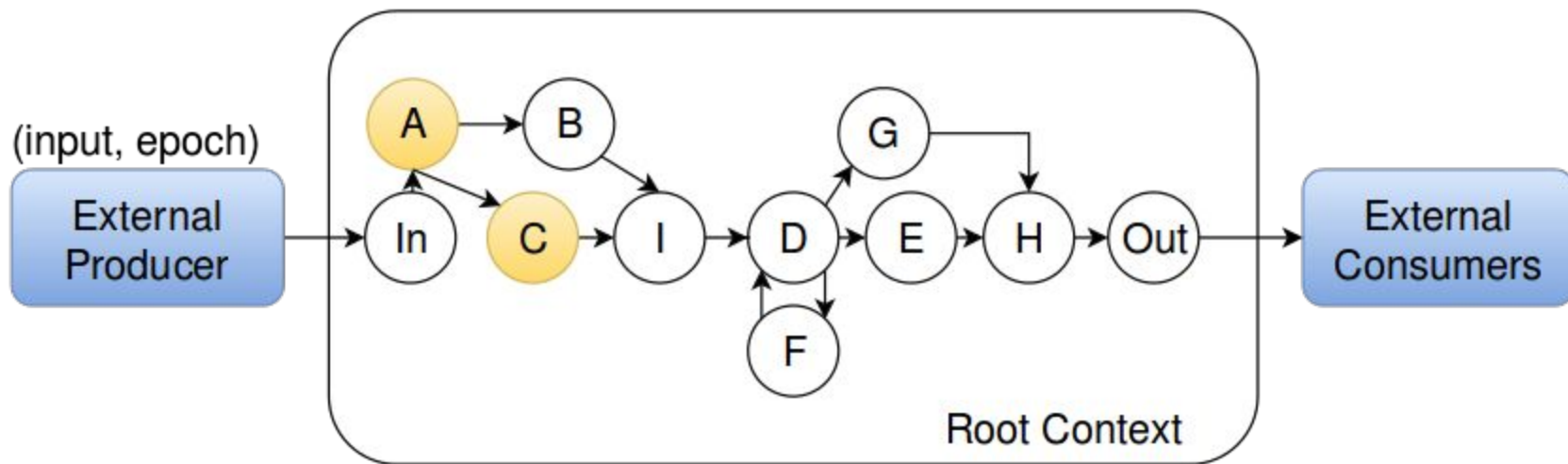


HOW IS TIMELY DATAFLOW ACHIEVED

- Pointstamp(t_1 , l_1) **could-result-in** Pointstamp(t_2 , l_2)
If there is a path between l_1 and l_2 presented by $f()$
i.e. $f(t_1) \leq t_2$

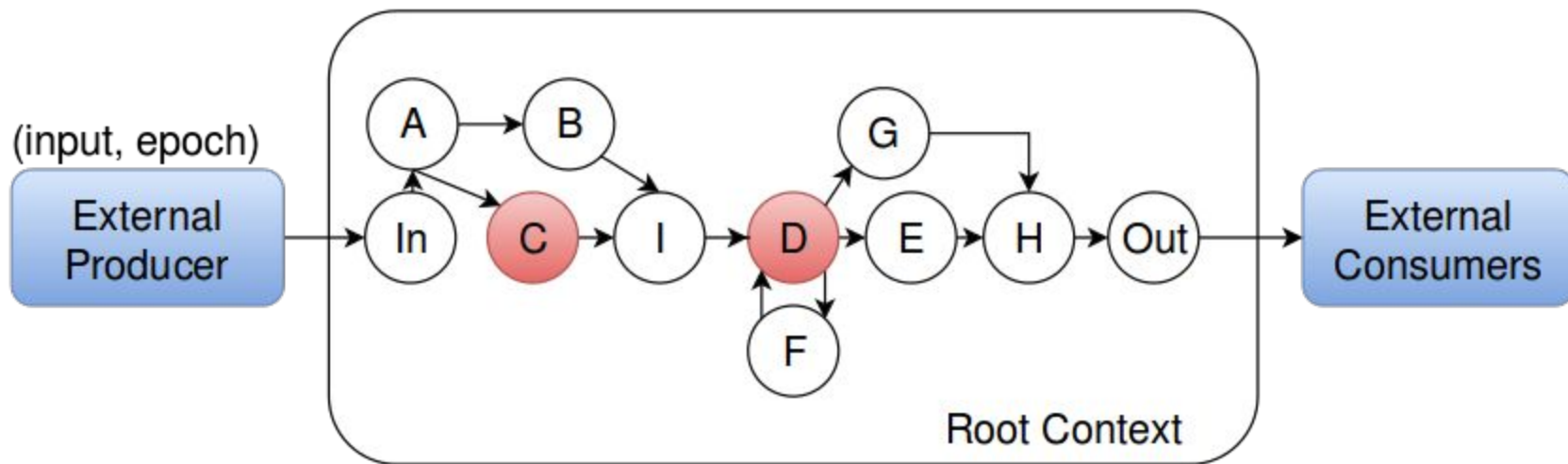
HOW IS TIMELY DATAFLOW ACHIEVED (CORRECTNESS GUARANTEES)

- Path Summary between A and C: “”



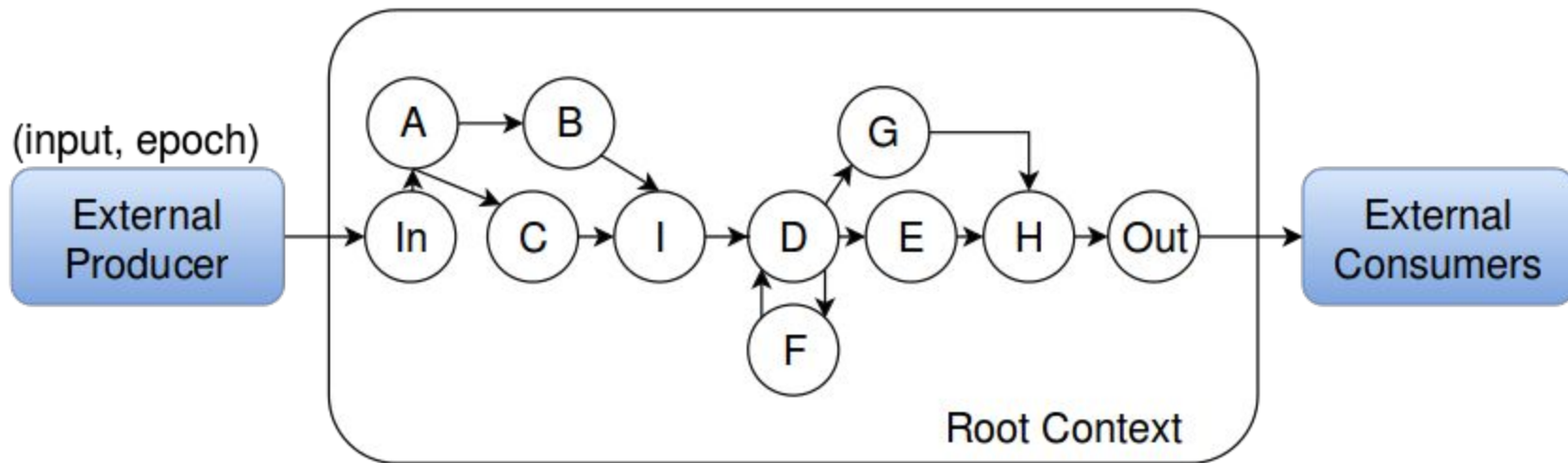
HOW IS TIMELY DATAFLOW ACHIEVED (CORRECTNESS GUARANTEES)

- Path Summary between A and C: “add” or “add-increment(n)”



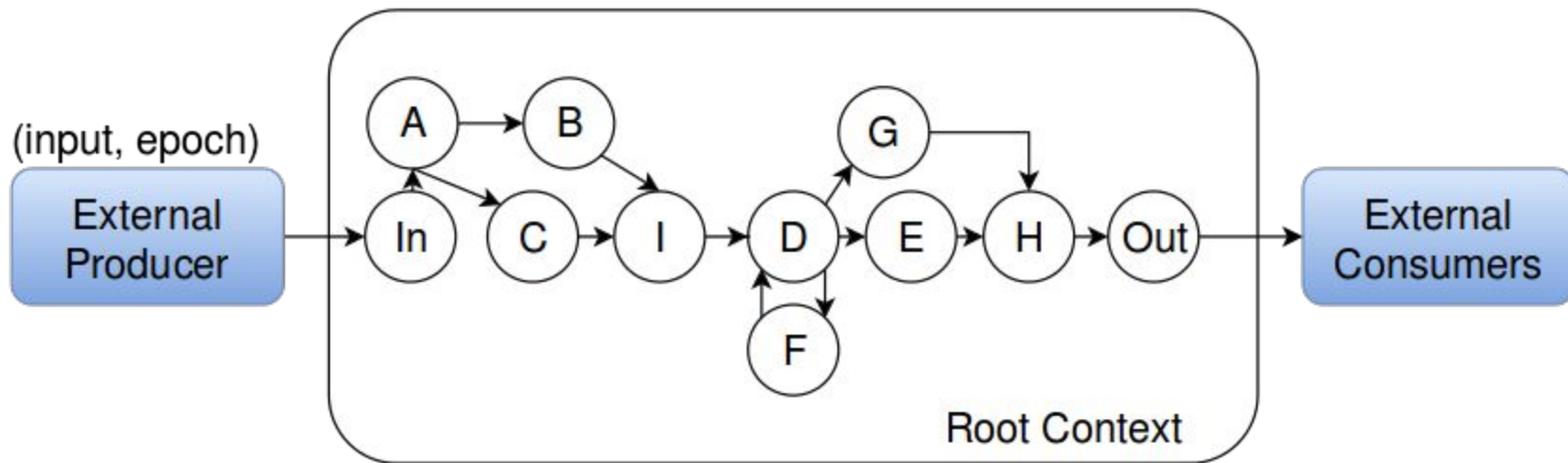
SINGLE-THREADED IMPLEMENTATION

- Scheduler that needs to deliver events.



SINGLE-THREADED IMPLEMENTATION

- Scheduler has active pointstamps \leftrightarrow unprocessed events.



SINGLE-THREADED IMPLEMENTATION

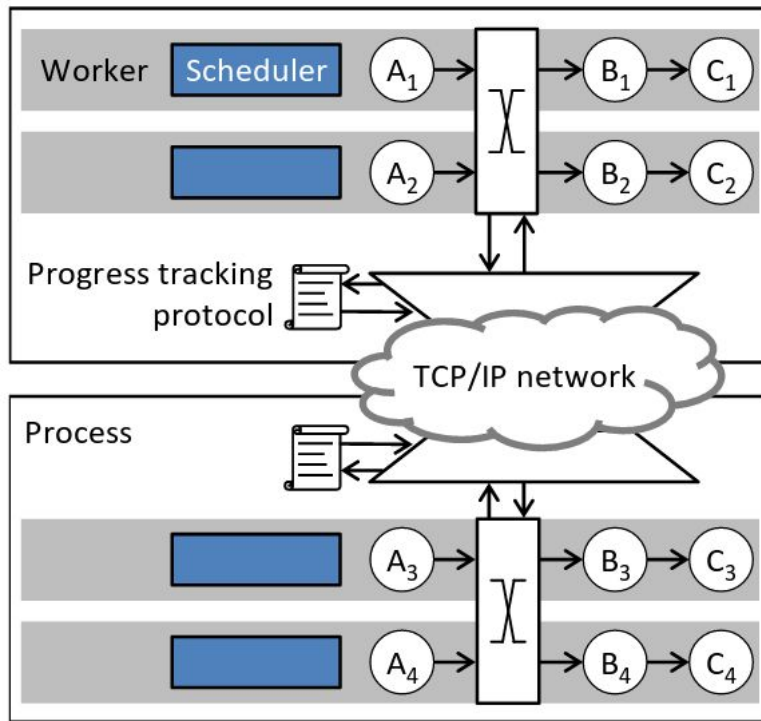
- Scheduler has active pointstamps \leftrightarrow unprocessed events.
- Scheduler has two counts:
 - Occurrence count of not resolved event.
 - Precursor count of how many active pointstamps precede it in the could-result-in order.

SINGLE-THREADED IMPLEMENTATION

- Pointstamp(t , l) becomes **active**.
Precursor count to number of existing active pointstamps that could result in it.
Increment precursor count of any pointstamp it could-result-in.
Becomes not active when occurrence is zero.
When not active, decrement the precursor count for any pointstamp that it could-result-in.

THE DISTRIBUTED ENVIRONMENT

DISTRIBUTED IMPLEMENTATION



DISTRIBUTED PROGRESS TRACKING

- Initial protocol: same as single multi-threaded.
 - Broadcast occurrence count updates.

- Do not immediately update local occurrence count.
 - Broadcast progress updates to all workers including myself.
 - Broadcast from a worker to another delivered in a FIFO manner.

- Use of a projected timestamp.
- A technique to buffer and accumulate updates.

MICRO-STRAGGLERS

- Have a big effect on overall performance.
 - Packet Loss (Networking)
 - Contention on concurrent data
 - Garbage collection

PERFORMANCE EVALUATION

PERFORMANCE EVALUATION

- I invite you to read: “Scalability! BUT at what Cost”

PERFORMANCE EVALUATION

- Comparison with:
 - SQL Server Parallel Data Warehouse (RDBMS)
 - Scalable HyperLink Store (distributed in-memory DB for storing large portions of the web graph)
 - DryadLINQ (data parallel computing using a declarative / high level programming language)
- Algos i.e. PageRank, SCC etc.

CONCLUSION: “OUR PROTOTYPE OUTPERFORMS GENERAL-PURPOSE BATCH PROCESSORS AND OFTEN OUTPERFORMS STATE-OF-THE-ART ASYNC SYSTEMS WHICH PROVIDE FEW SEMANTIC GUARANTEES.”

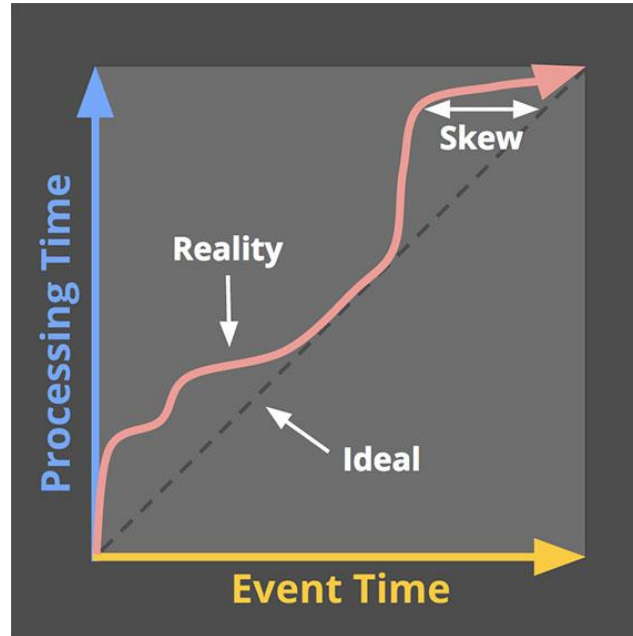
CONCLUSION: “OUR PROTOTYPE OUTPERFORMS GENERAL-PURPOSE BATCH PROCESSORS AND OFTEN OUTPERFORMS STATE-OF-THE-ART ASYNC SYSTEMS WHICH PROVIDE FEW SEMANTIC GUARANTEES.”

STREAMING SYSTEMS AS OF TODAY

STREAMING SYSTEMS

- Systems that have unbounded data in mind.
- They are a superset of batch processing systems.

STREAMING SYSTEMS



Reference: Fig-1: Example of time domain mapping. [Streaming 101](#)

STREAMING SYSTEMS

Design Questions:

- **What** results are calculated?
The types of *transformations* within the pipeline.
- **Where** in event time are results calculated?
The use of *event-time windowing* within the pipeline.
- **When** in processing time are results materialized? The use of watermarks and triggers.
- **How** do refinements of results relate?
Discard or accumulate or accumulate and retract.

FIN.

Thank you!

RESOURCES

- Link to transcribed talk in [pdf format](#).
- Timely Dataflow ([Rust Implementation](#))
- Frank [blog posts](#):
 - Timely dataflow
 - Differential dataflow
- [The world beyond batch: Streaming 101](#)
- [The world beyond batch: Streaming 102](#)