

Beyond Analytics

The Evolution of Stream Processing Systems

Prospects

Paris Carbone, Marios Fragkoulis, Vasiliki Kalavri, Asterios Katsifodimos

Slides: streaming-research.github.io/Tutorial-SIGMOD-2020



Tutorial overview

- Part I: Introduction & Fundamentals (Vasia)
- Part II: Time, Order, & Progress (Marios)
- Part III: State Management (Paris)
- Part IV: Fault Recovery & High Availability (Marios)
- Part V: Load Management & Elasticity (Vasia)
- Part VI: Prospects (All)

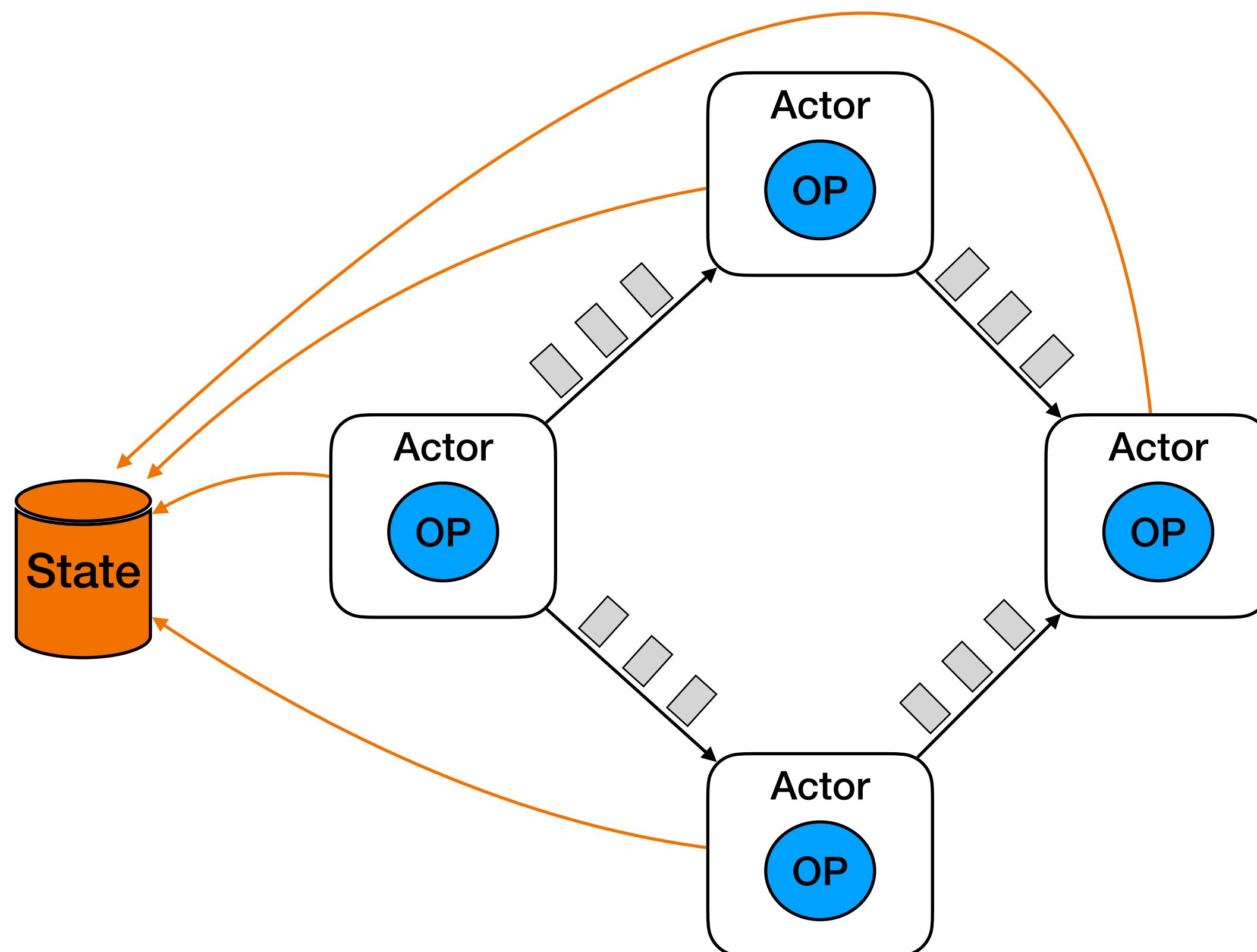
Prospects

Agenda

1. **Emerging applications (Cloud, ML, Streaming Graphs)**
2. **New requirements**
3. Conclusion, references, contacts

Cloud programming frameworks

Actors & Streams



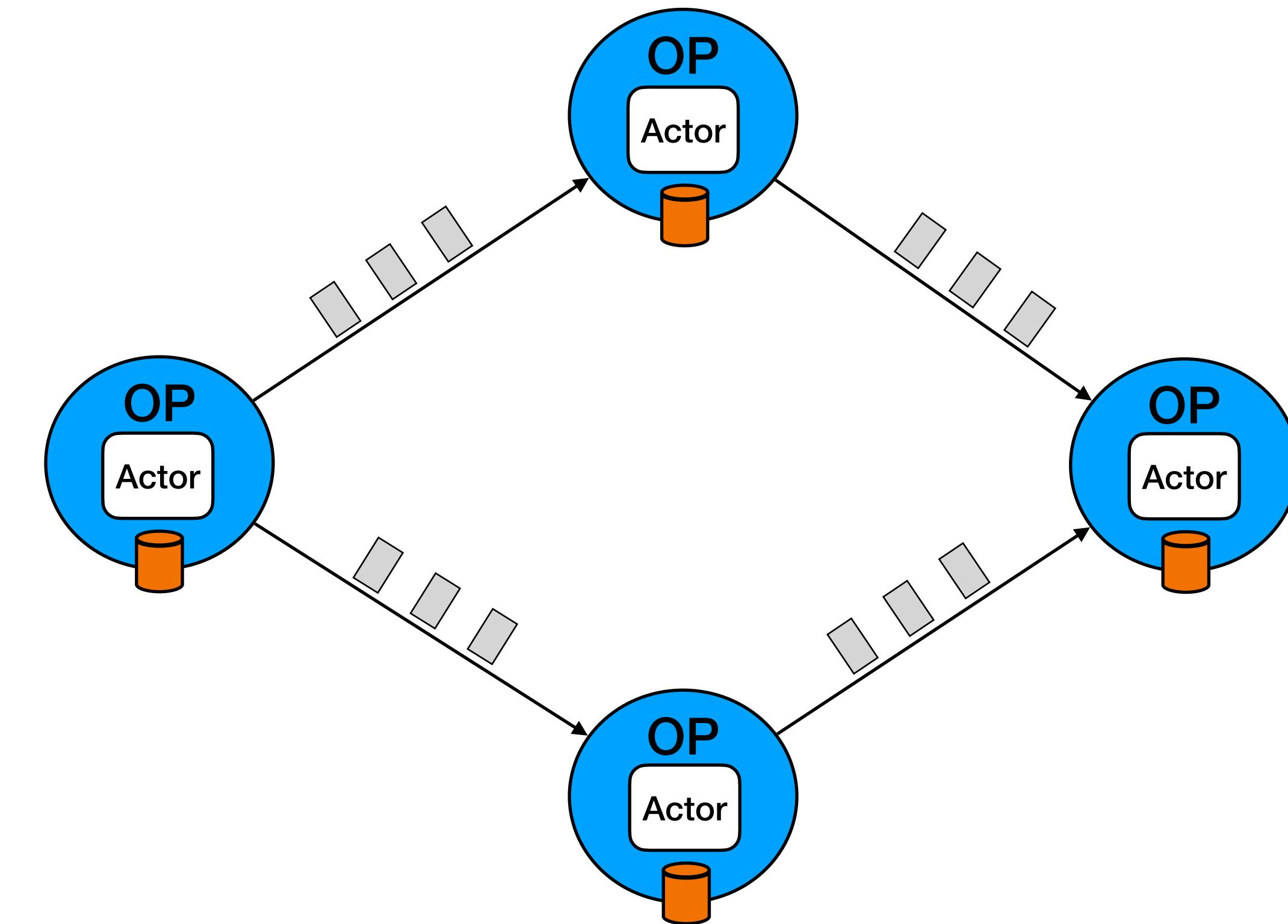
Streams on Actors



Orleans
(CIDR'19)



RAY
(OSDI'18)



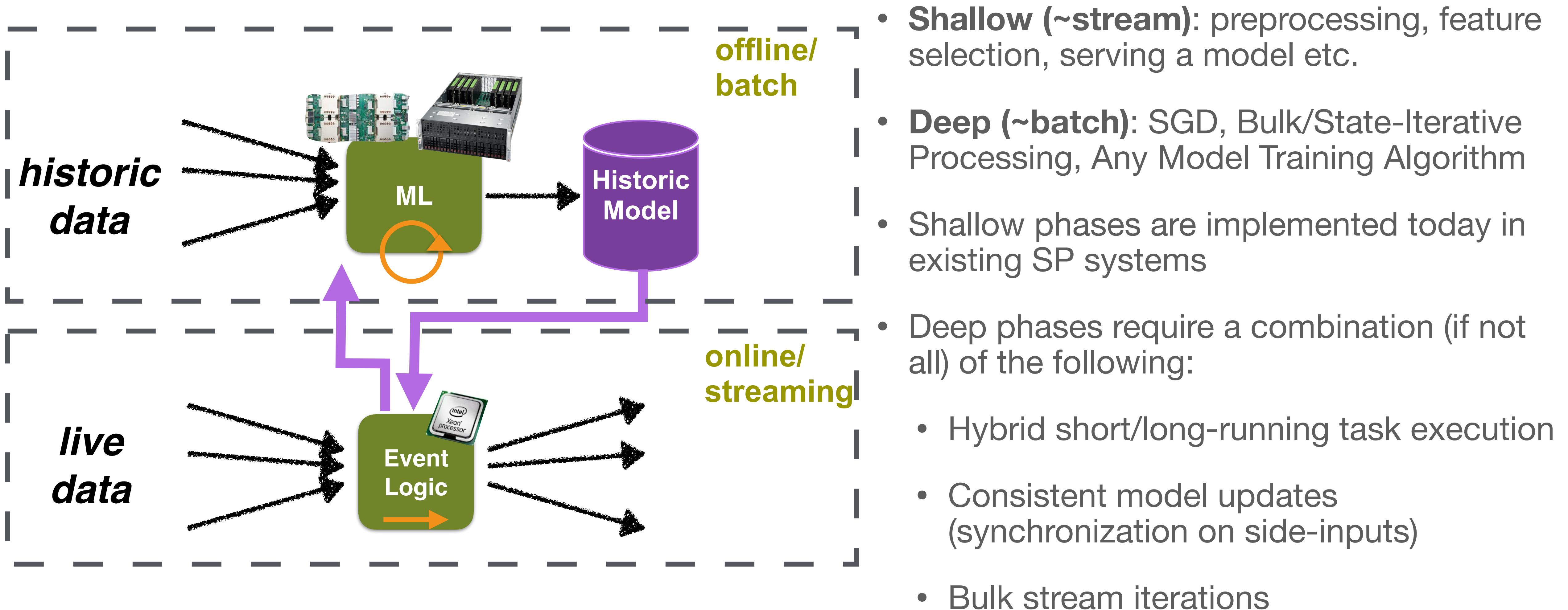
Actors on Streams

sFaaS (VLDB'19)



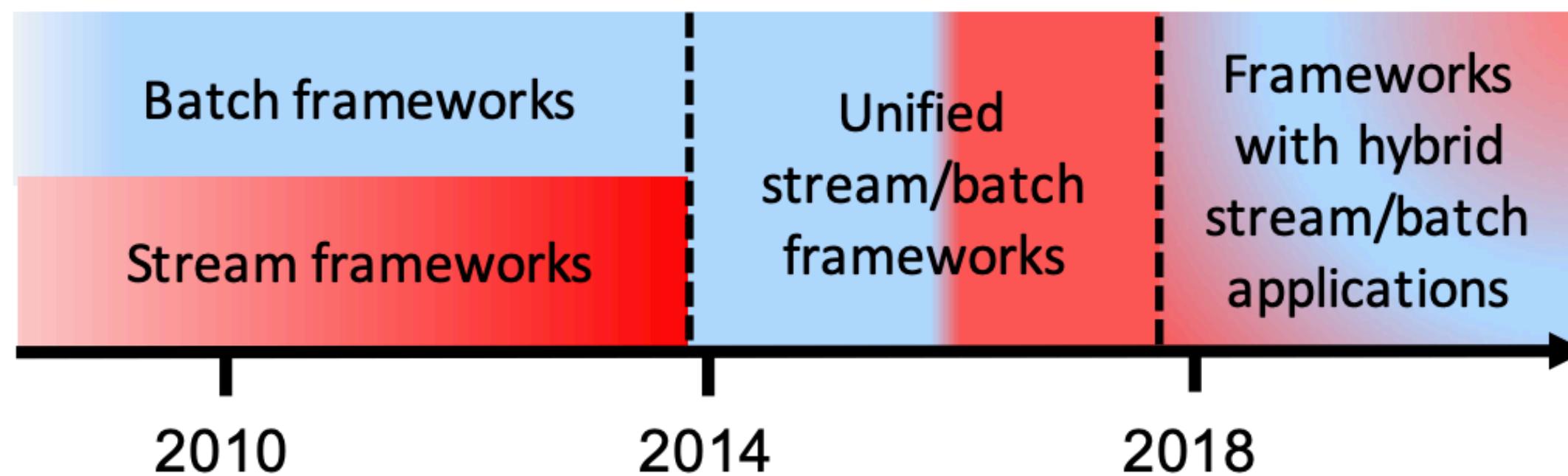
Stateful
Functions

Integrating Shallow+Deep ML Stages



Steps Forward in ML in Data Streaming

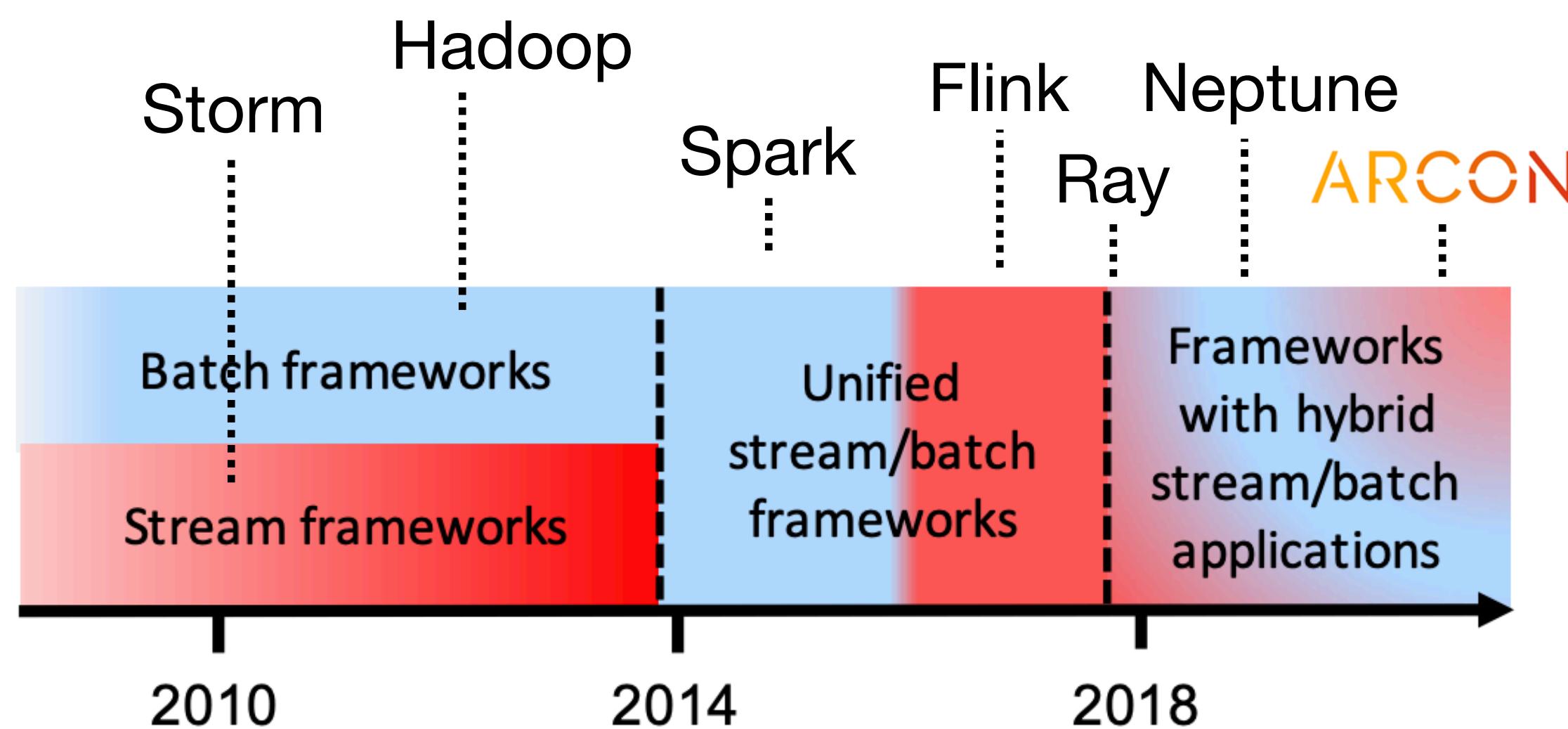
Hybrid Execution



Neptune: Scheduling Suspendable Tasks for Unified Stream/Batch Applications
Garefalakis P, Karanasos K, Pietzuch P in SoCC (2019)

Steps Forward in ML in Data Streaming

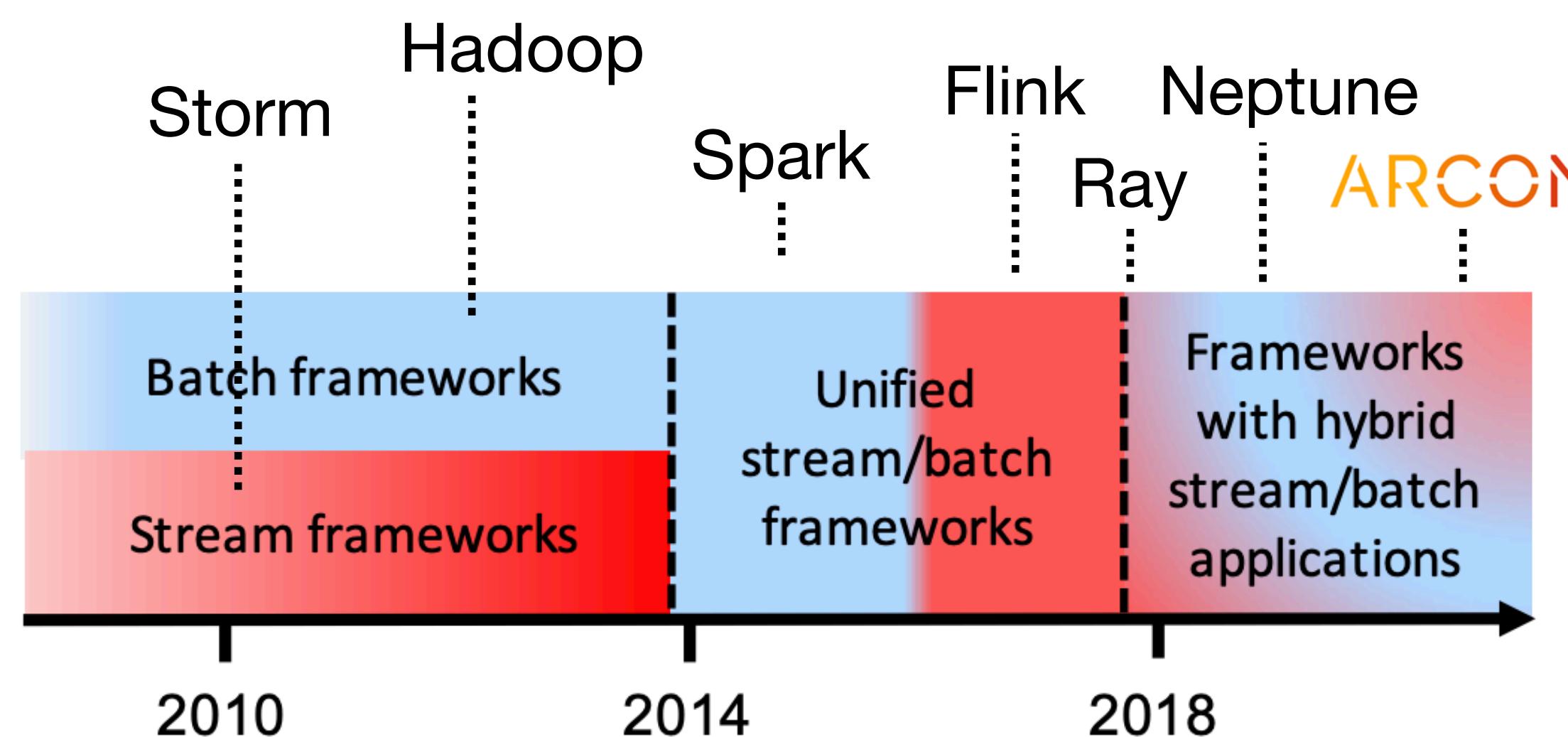
Hybrid Execution



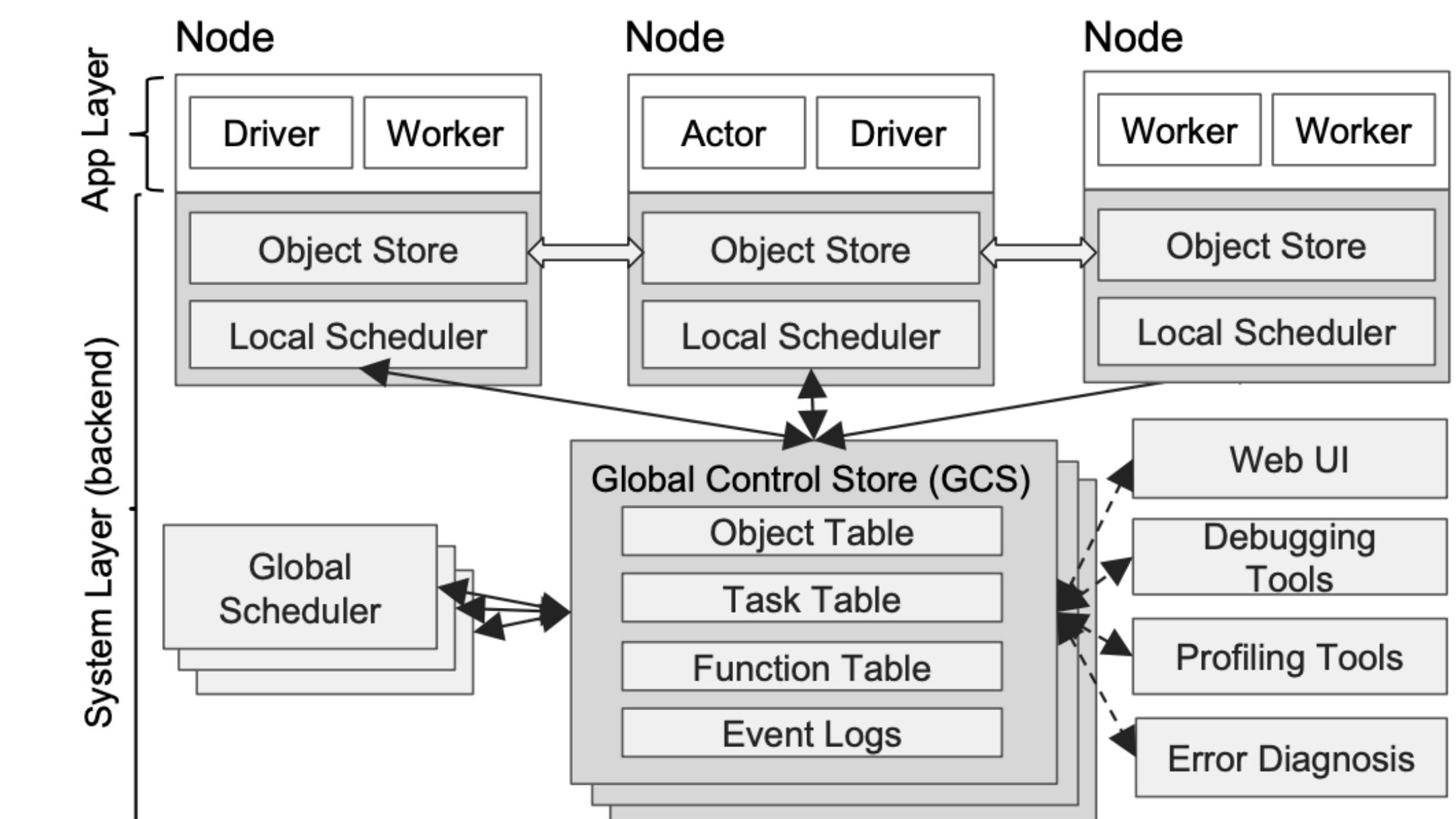
Neptune: Scheduling Suspendable Tasks for Unified Stream/Batch Applications
Garefalakis P, Karanasos K, Pietzuch P in SoCC (2019)

Steps Forward in ML in Data Streaming

Hybrid Execution



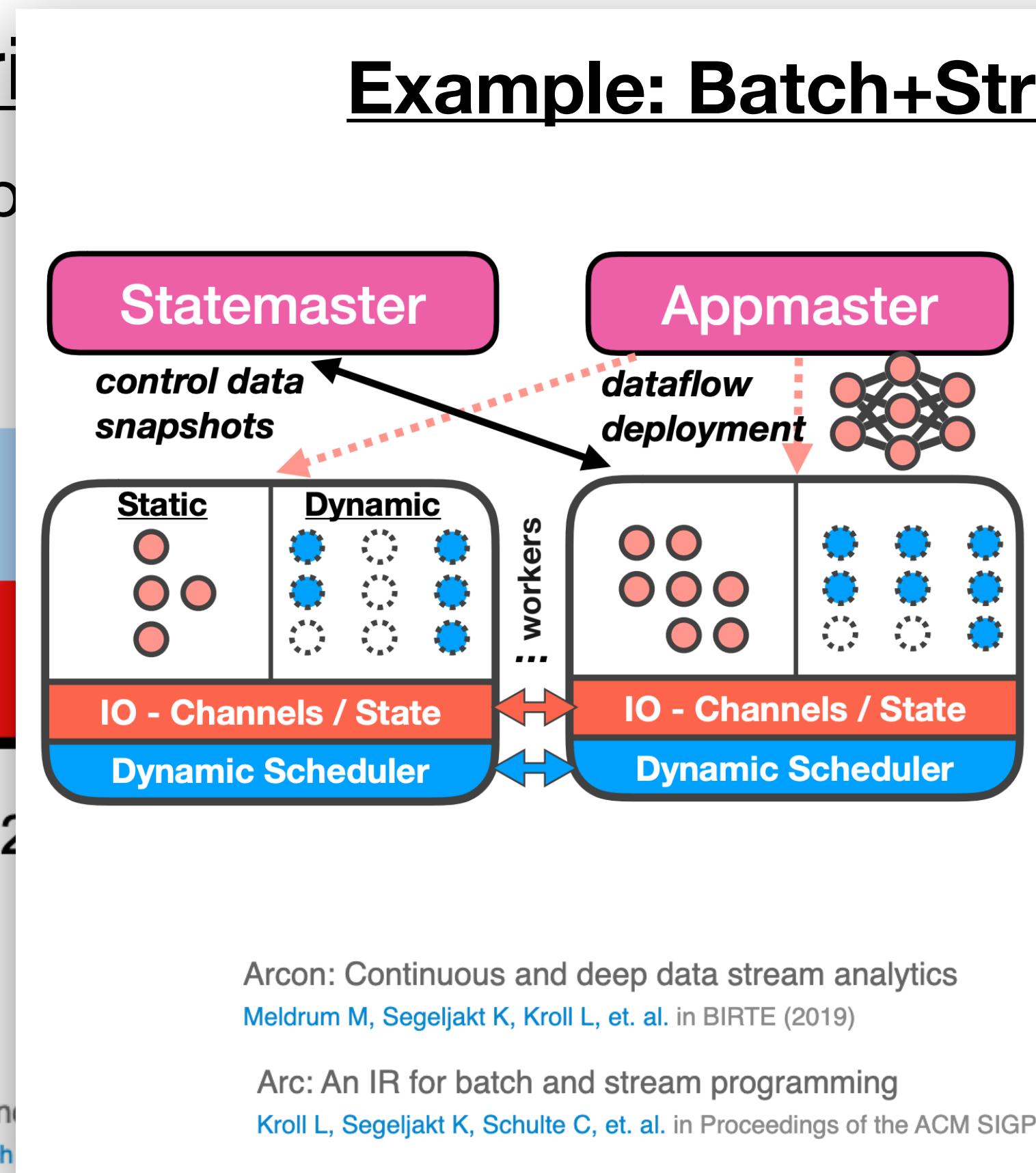
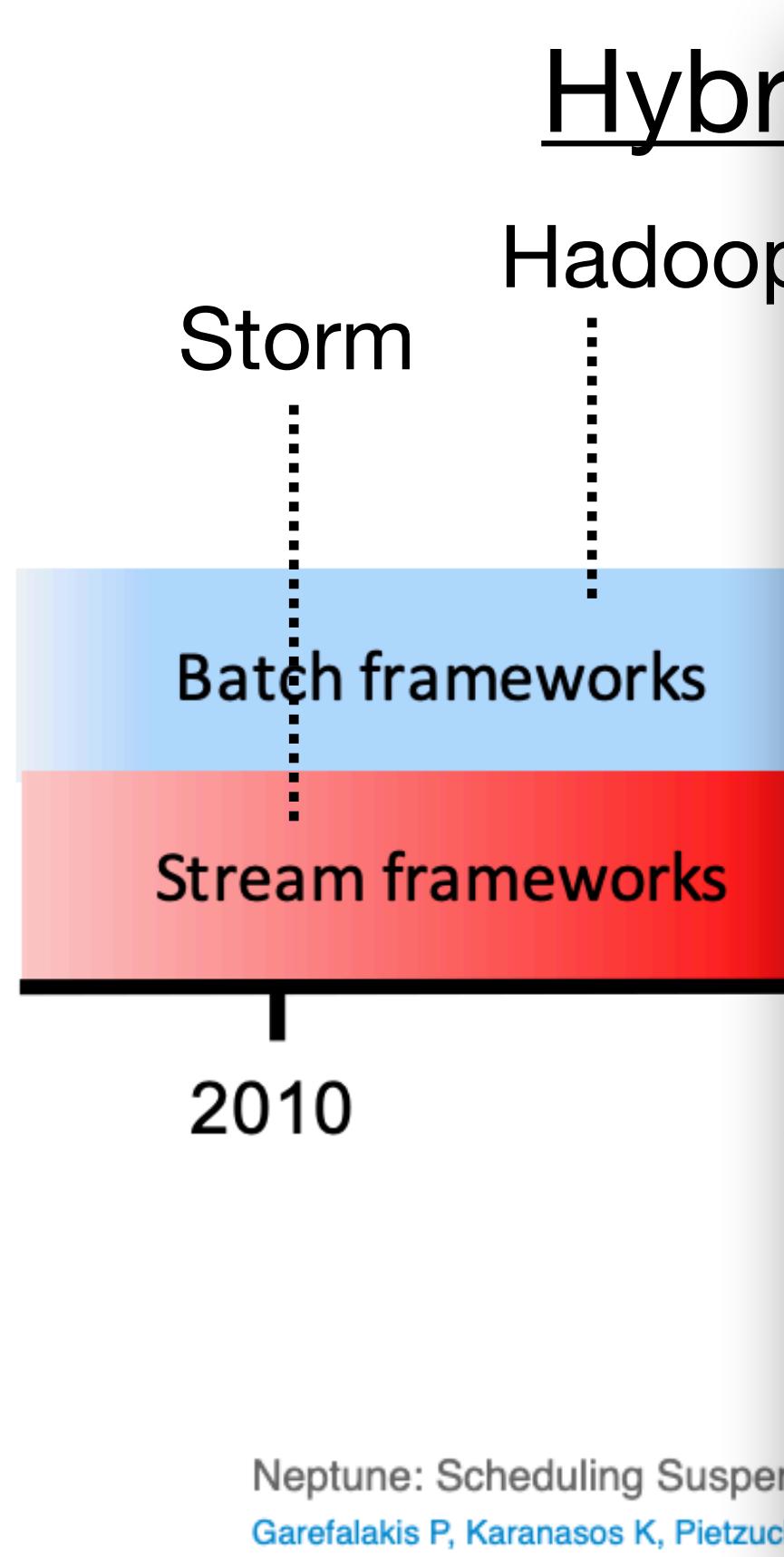
Ray : Actors + Remote Functions for ML



Neptune: Scheduling Suspendable Tasks for Unified Stream/Batch Applications
 Garefalakis P, Karanasos K, Pietzuch P in SoCC (2019)

Ray: A Distributed Framework for Emerging AI Applications
 Moritz P, Nishihara R, Wang S, et. al. in USENIX OSDI (2018)

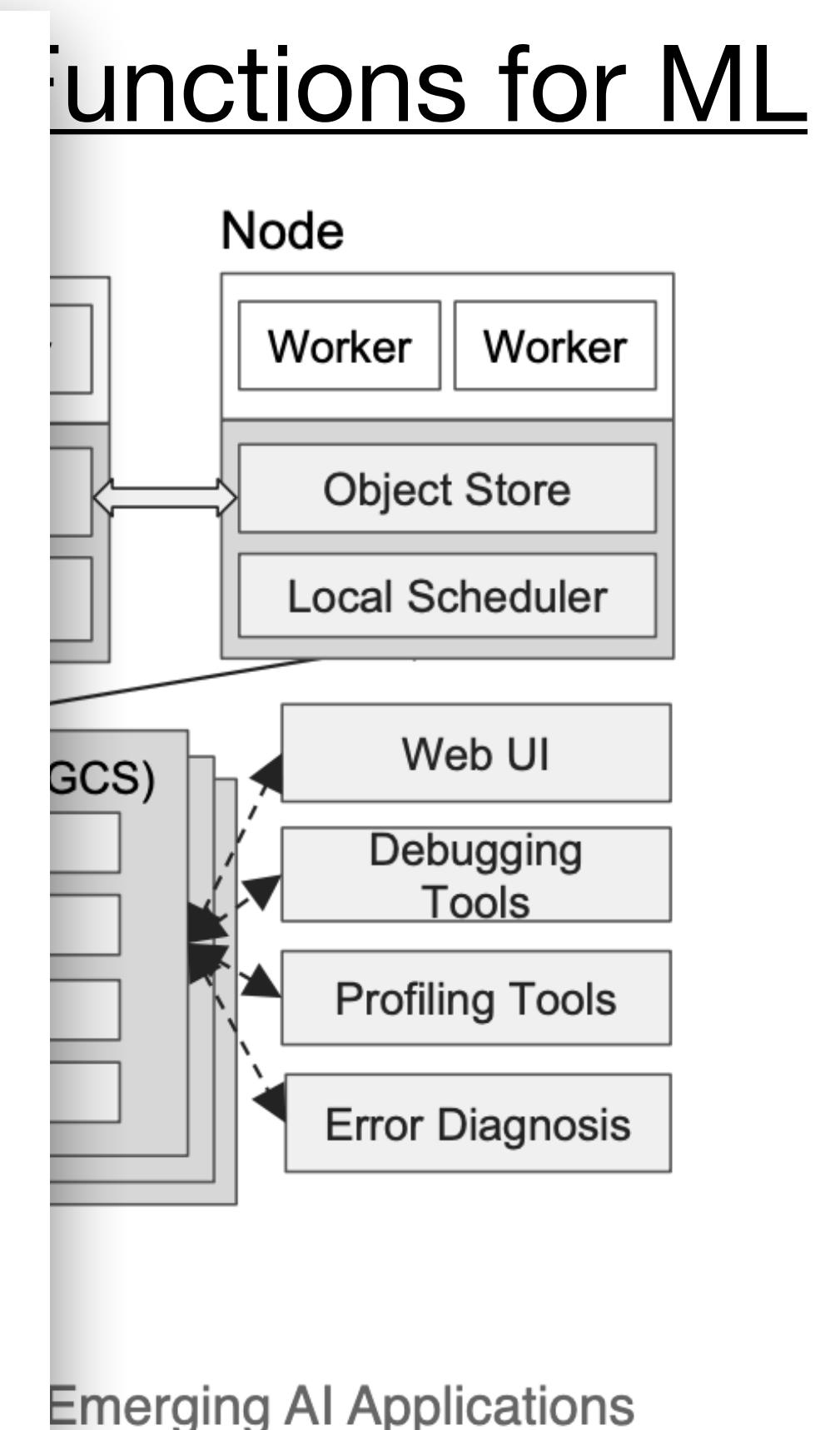
Steps Forward in ML in Data Streaming



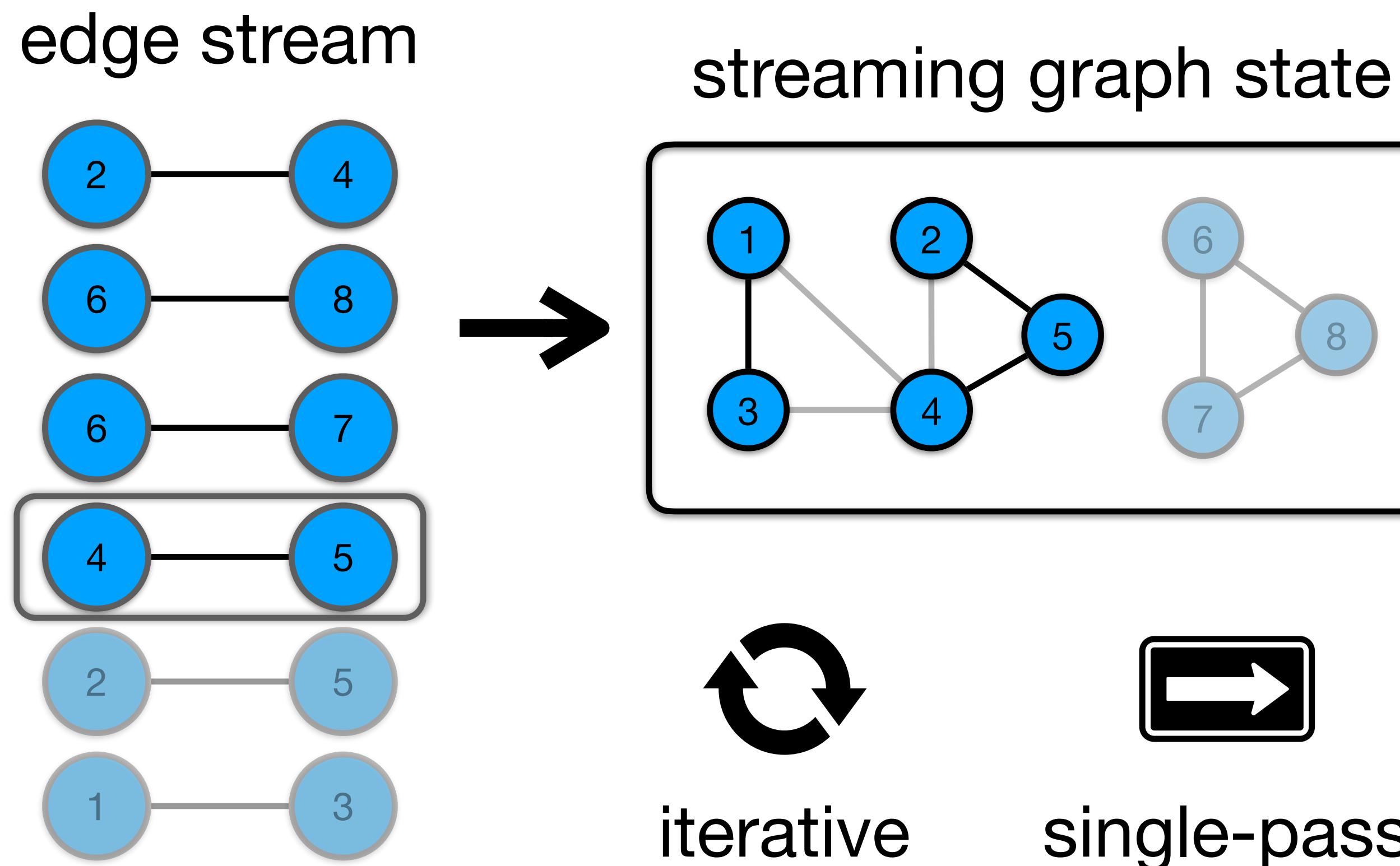
Arc: An IR for batch and stream programming
[Kroll L, Segeljakt K, Schulte C, et. al. in Proceedings of the ACM SIGPLAN Conference on Programming Language Design and Implementation \(2018\)](#)

```
| source: Stream[i32],  
odd_sink: StreamAppender[i32],  
even_sink: StreamAppender[i32] |  
let mapped = result(  
    for(source, StreamAppender[i32], |sink, element|  
        merge(sink, element + 5));  
    for(mapped, odd_sink, |sink, e|  
        if(e % 2 != 0, merge(sink, e), sink));  
    for(mapped, even_sink, |sink, e|  
        if(e % 2 == 0, merge(sink, e), sink))
```

Arc: intermediate representation for batch+stream computation



Graph streams model interactions as events that update an underlying graph structure



Edge events:

A purchase, a movie rating, a like on an online post, a bitcoin transaction, a packet routed from a source to destination

Vertex events:

A new product, movie, user

Open problems in graph streaming

- Streaming graph **state** management
 - See: *LiveGraph: a transactional graph storage system with purely sequential adjacency list scans.* (VLDB'20).
- Streaming graph **partitioning** (of unbounded graphs)
 - See: *Streaming graph partitioning: an experimental study.* (VLDB'18)
- Graph streaming **algorithms**
 - See: *Regular Path Query Evaluation on Streaming Graphs.* (SIGMOD'20 Research 16: Graph and Stream Processing, Thursday 1:30 PM – 3:00 PM).
- **Languages** and operator semantics

For an overview, see: ***Practice of Streaming and Dynamic Graphs: Concepts, Models, Systems, and Parallelism*** (arxiv.org/abs/1912.12740).

Novel APIs and abstractions

Cloud Apps, ML, Graphs

functional APIs

```
val maxTemp = sensorData
  .map(r => Reading(r.id, r.time, r.temp))
  .keyBy(_.id)
  .max("temp")
maxTemp.print()
```

streaming SQL

```
SELECT *
FROM Tumble (
    data => TABLE ( Bid ),
    timecol => DESCRIPTOR ( bidtime ),
    dur = > INTERVAL '10' MINUTES,
    offset = > INTERVAL '0' MINUTES );
```



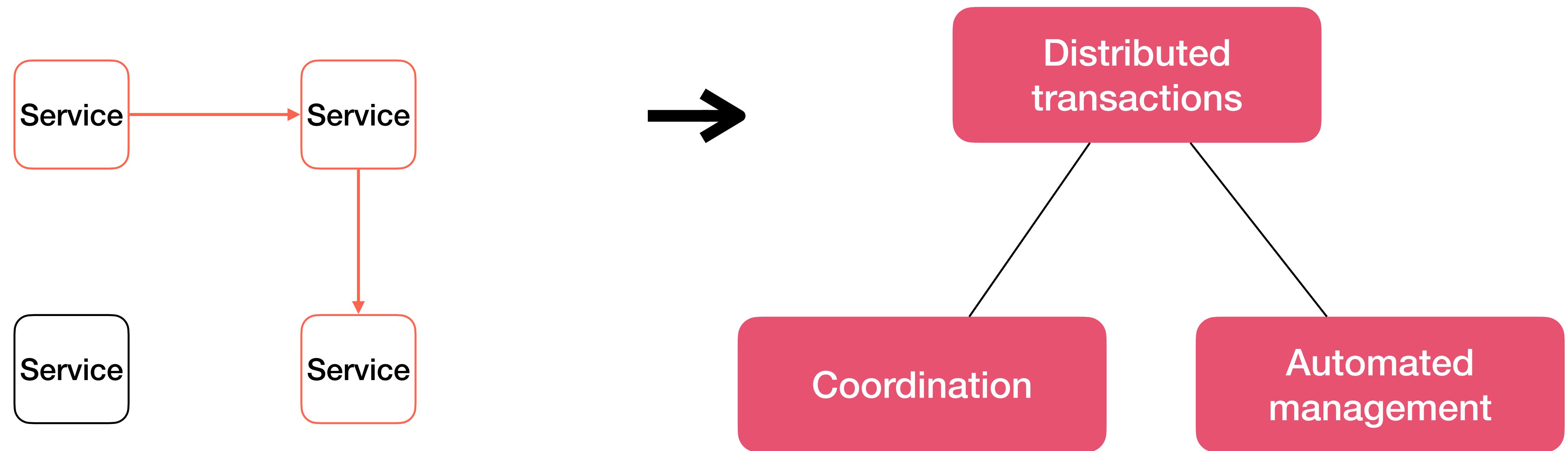
Actor-like APIs

Streaming ML APIs

Streaming graph abstractions

Application transactions

Cloud Apps

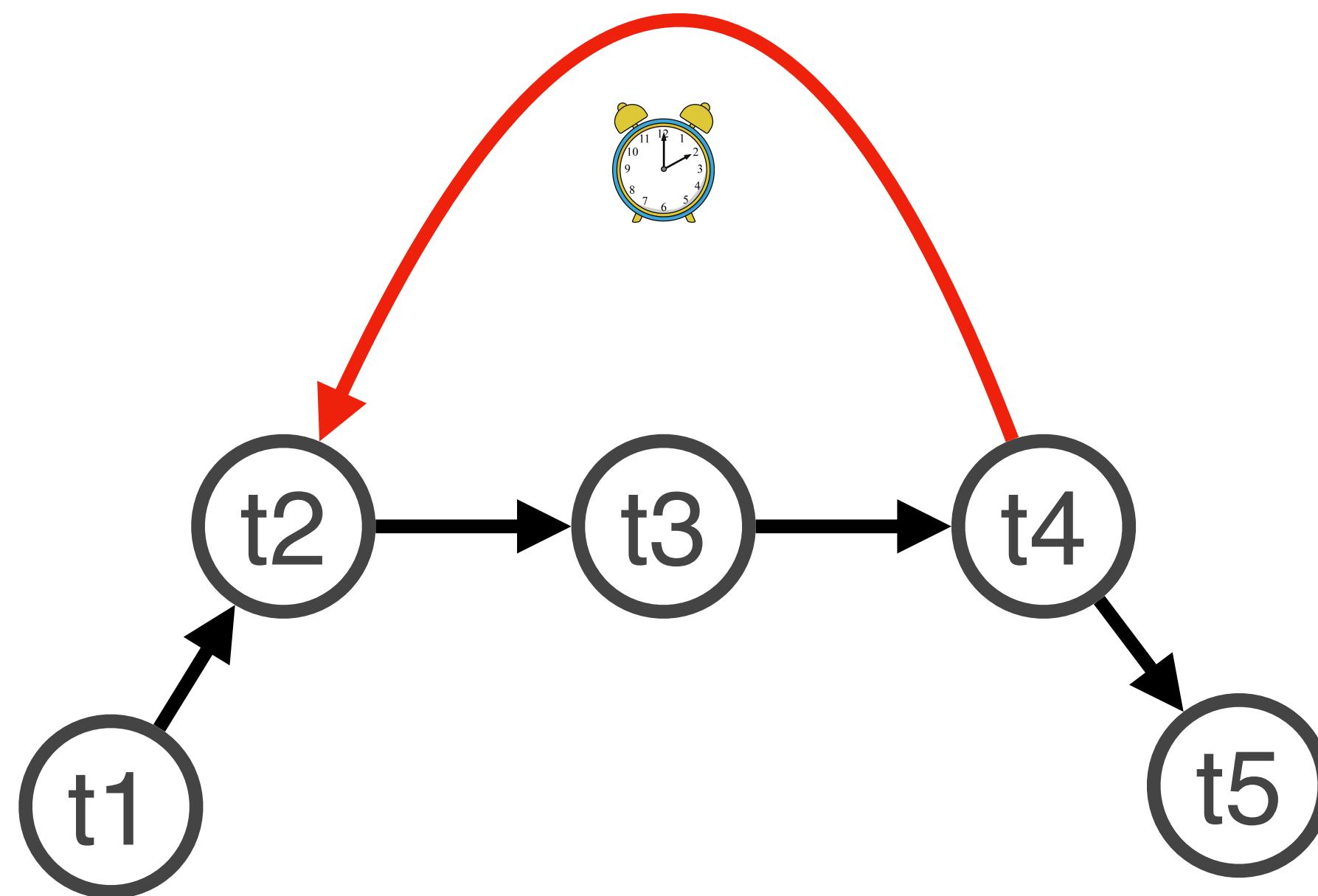


S-Store (VLDB '14)

Iterative computation on streams

Cloud Apps, ML, Graphs

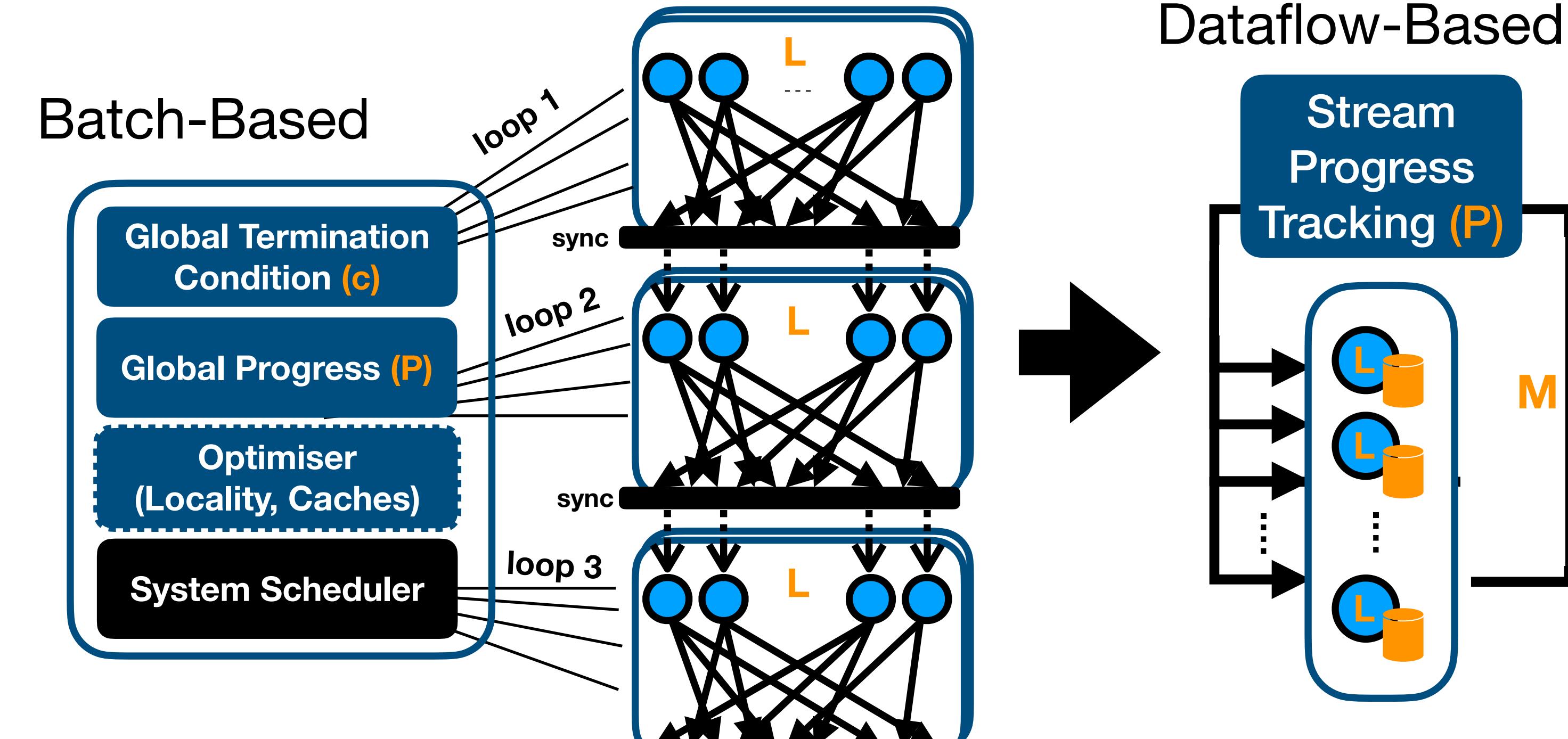
Asynchronous Loops



Problems

- Deadlocks - Cyclic Flow Control
- Time Progress?

Synchronous Loops



Problems

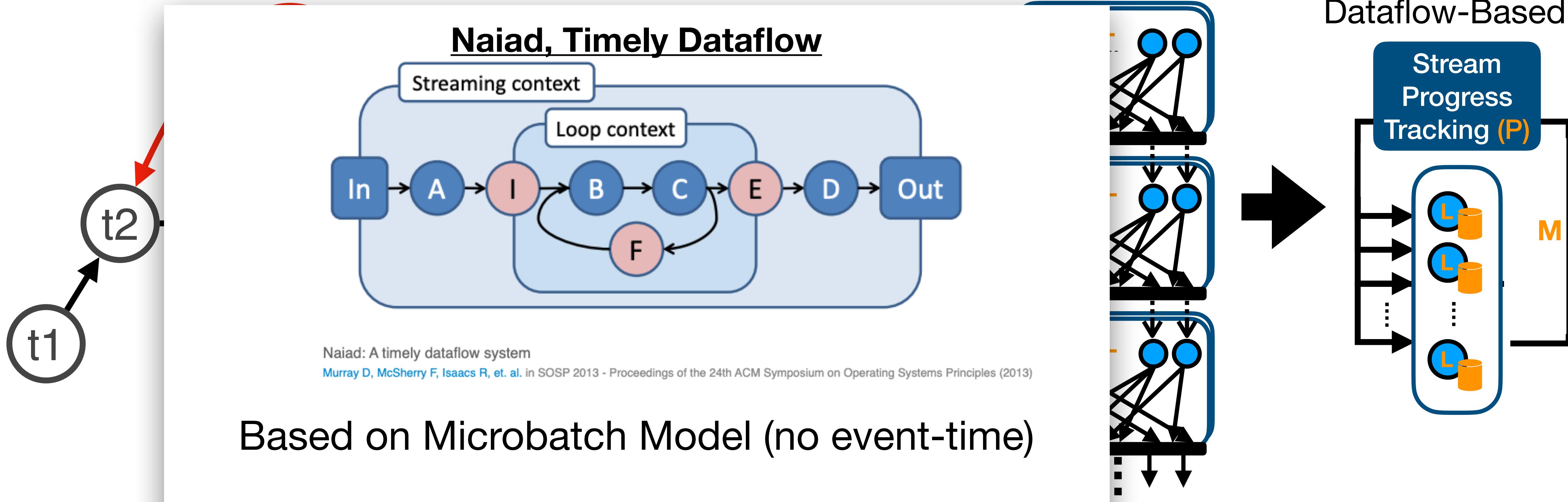
- Decentralised Global Progress Tracking
- Semantics + Arbitrary Nesting

Iterative computation on streams

Cloud Apps, ML, Graphs

Asynchronous Loops

Synchronous Loops

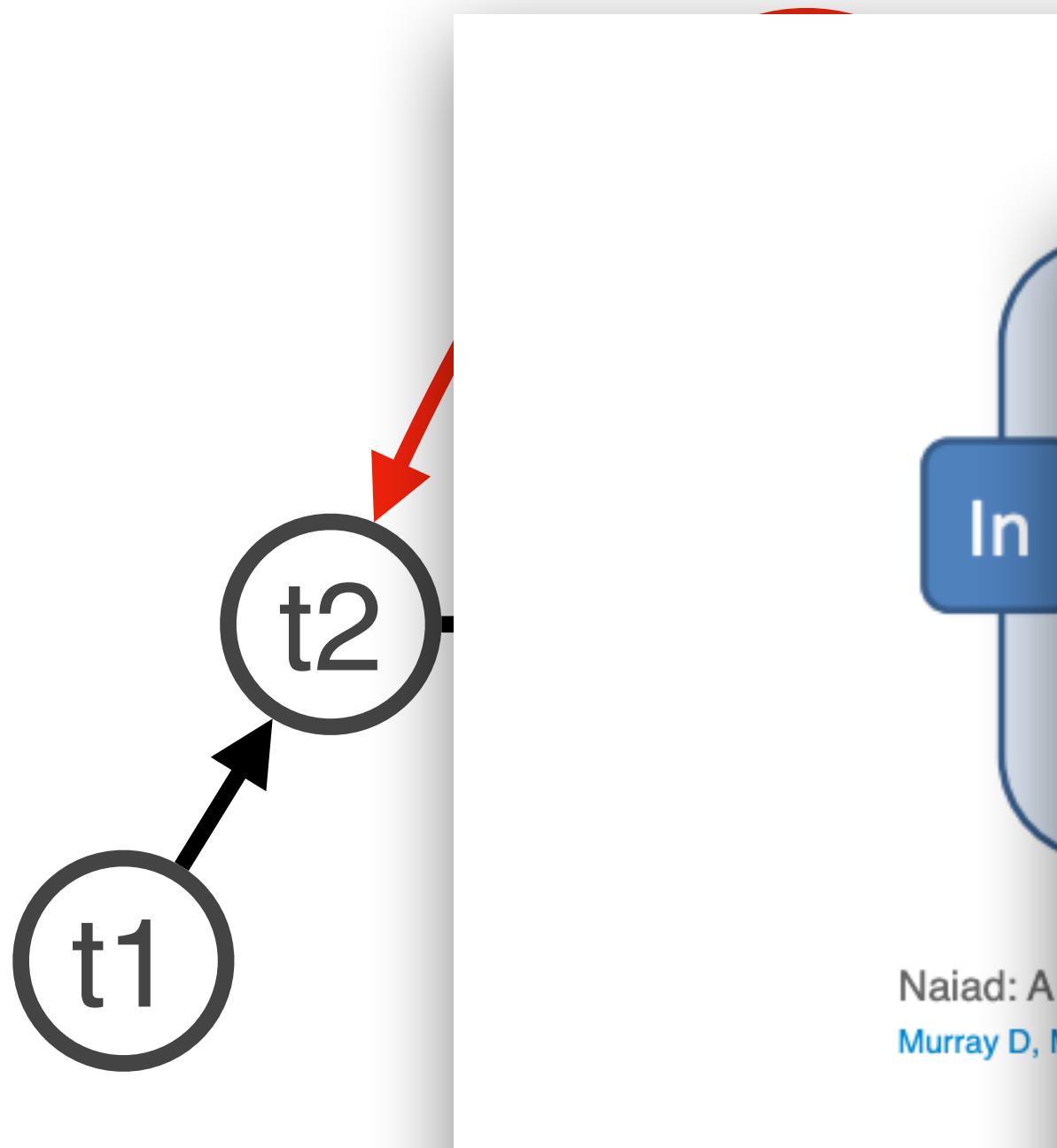


- Deadlocks - Cyclic Flow Control
- Time Progress?

- Decentralised Global Progress Tracking
- Semantics + Arbitrary Nesting

Iterative computation on streams

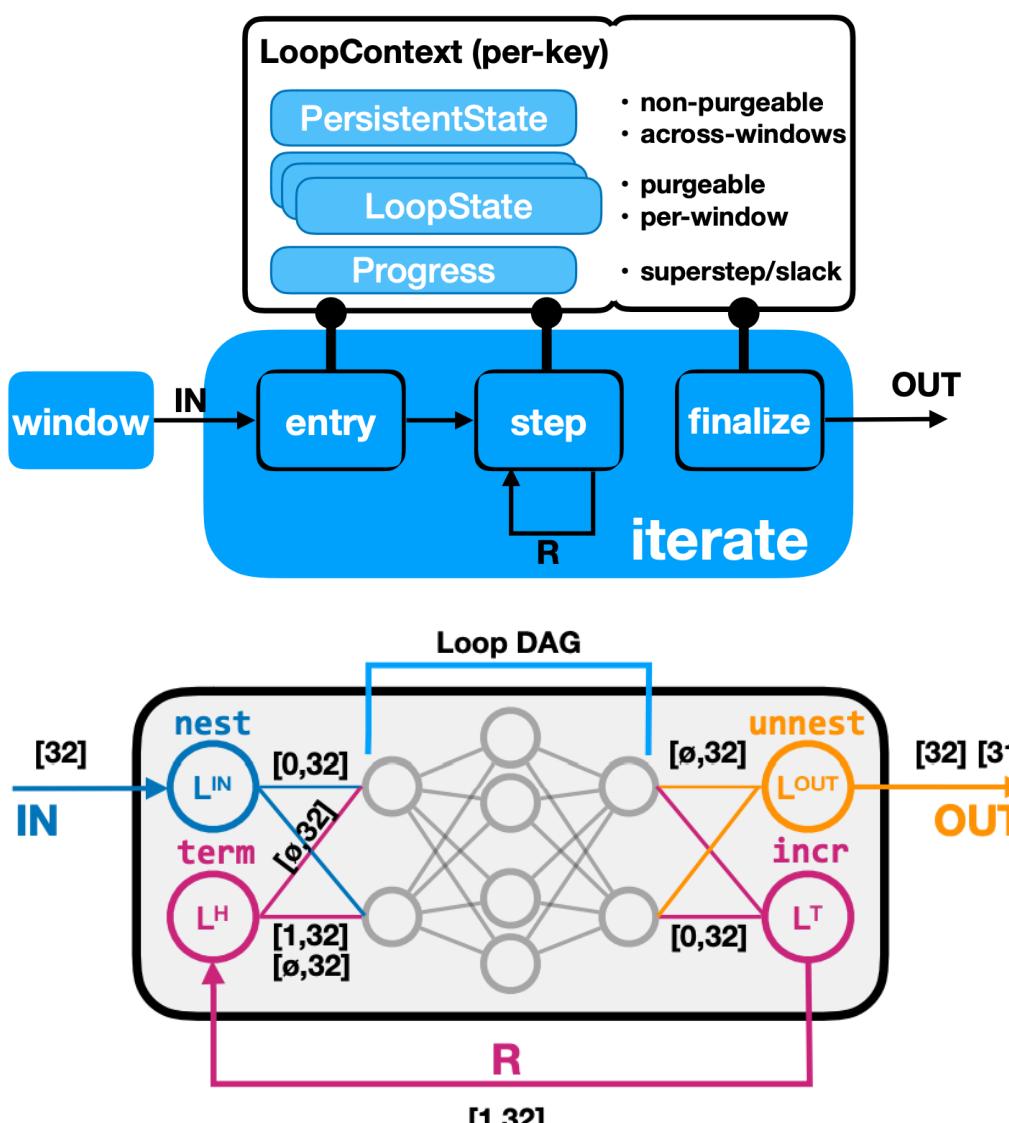
Asynchronous Loops



Based on

Naiad, Timely Dataflow

Flink Window Loops



```

1 case class GraphEdge(from:Long, to:Long)
2 case class VComponent(id:Long, component:Long)
3
4 val input: DataStream[GraphEdge] = getEdgeStream()
5 input
6   .flatMap(e => List(e,GraphEdge(e.to,e.from))) //add inverse
7   .keyBy(vertex => vertex.id) //Scope Computation by Key
8   .timeWindow(1 Min)
9   .iterate(Termination.Fixpoint, Synchrony.Strict,
10   (vComponent => vComponent.id), //Loop Key
11   (ctx:LoopContext,in:Iterable[GraphEdge],out:Collector[VComponent]) =>
12   { //ENTRY FUNCTION
13     ctx.loopState("neighbors").setList(in.map(e => e.to));
14     ctx.loopState("cc").setValue(ctx.key)
15     ctx.loopState("neighbors").foreach(n => out.collect(VComponent(n,ctx.key)))
16   },
17   (ctx:LoopContext,in:Iterable[VComponent],out:Collector[VComponent]) =>
18   { //STEP FUNCTION
19     val newcc = in.minBy(c => c.component).component
20     var cc = ctx.loopState("cc")
21     if(cc.value != newcc){
22       cc.setValue(newcc)
23       ctx.loopState("neighbors").foreach(n => out.collect(VComponent(n,newcc)))
24     }
25   },
26   (ctx:LoopContext,out:Collector[VComponent]) =>
27   { //FINALIZE FUNCTION
28     out.collect(VComponent(ctx.key, ctx.loopState("cc").value))
29   });

```

Based on out-of-order event time model

Scalable and Reliable Data Stream Processing

P Carbone
KTH Royal Institute of Technology

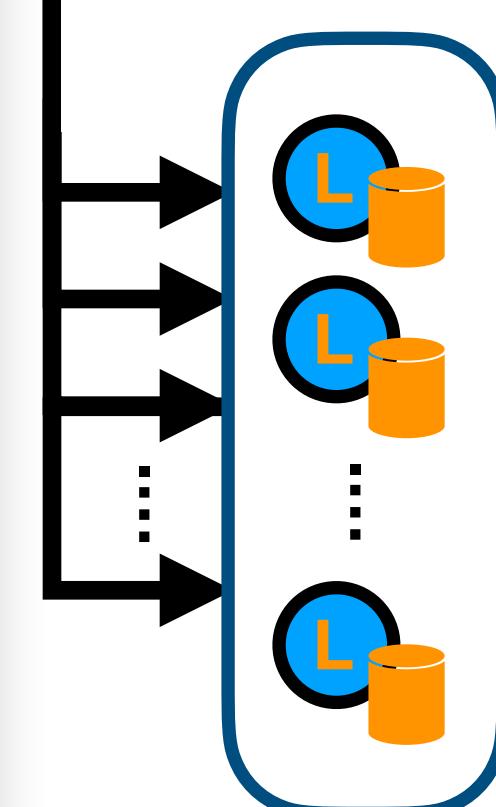
- Deadlocks - Cyclic FIFO
- Time Progress?

Cloud Apps, ML, Graphs

Synchronous Loops

Dataflow-Based

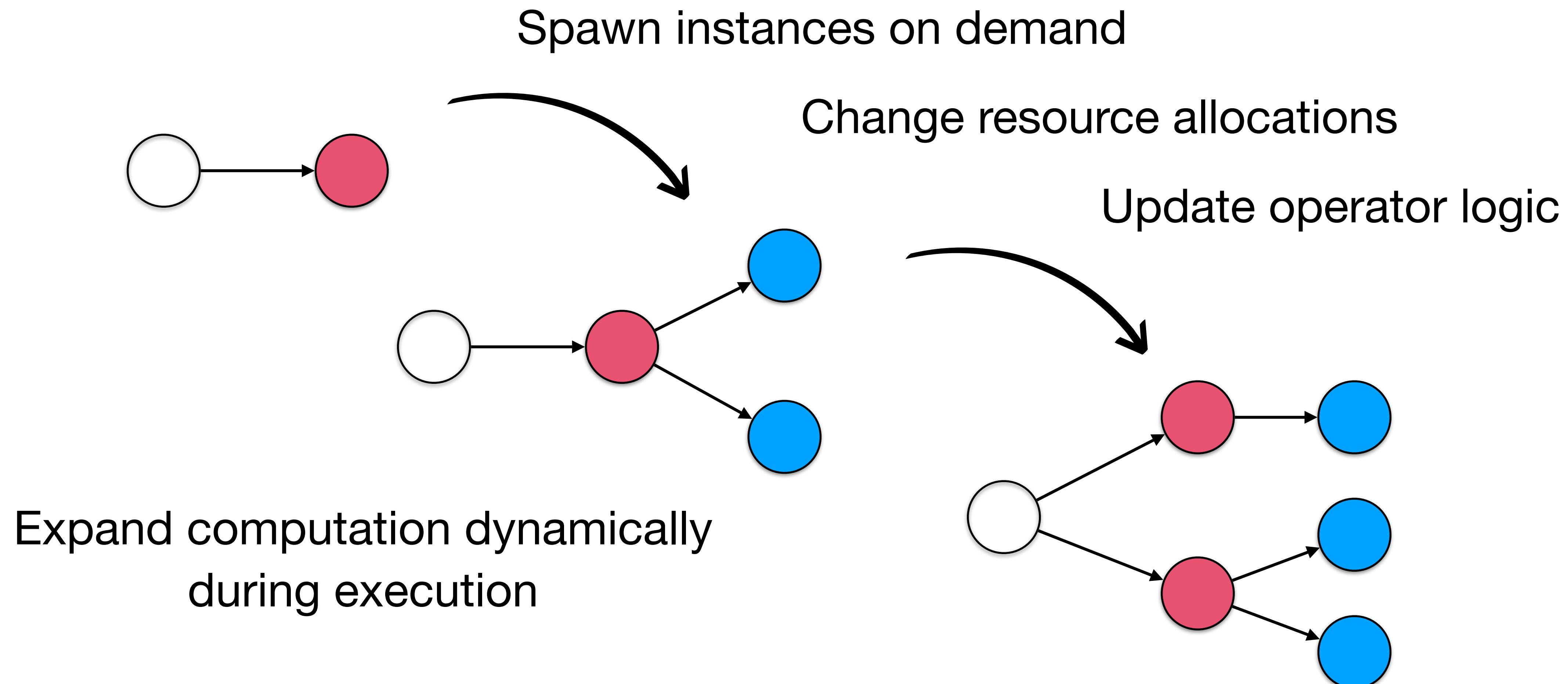
Stream
Progress
Tracking (P)



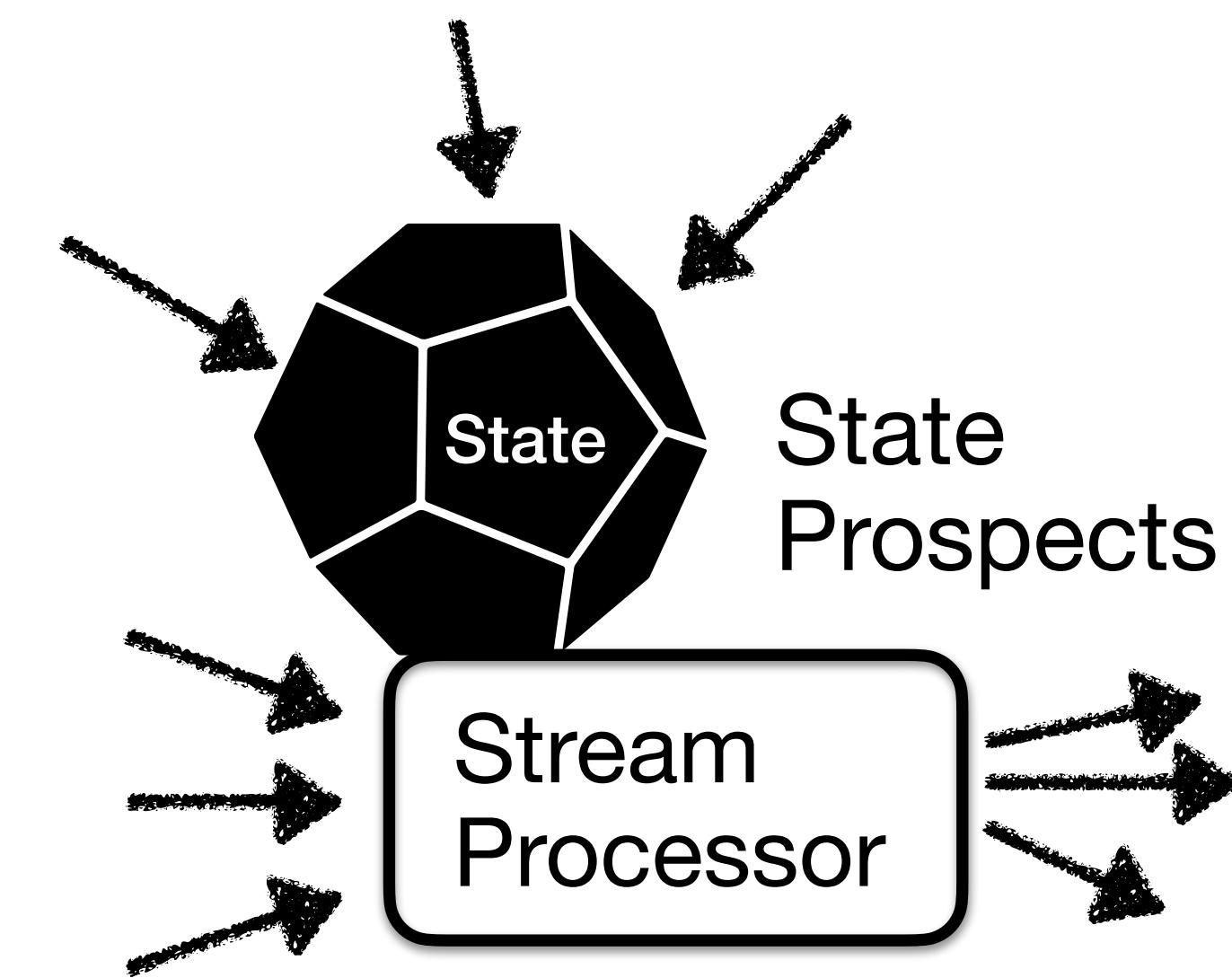
- Semantics + Arbitrary Nesting

Elasticity, adaptivity, dynamic plans

Cloud Apps, ML



Taking Stateful Processing Further

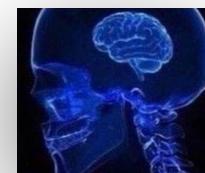


Taking Stateful Processing Further

Workload-Aware State Backends



Databases : Optimal plans using different indexes (e.g., B+Trees)

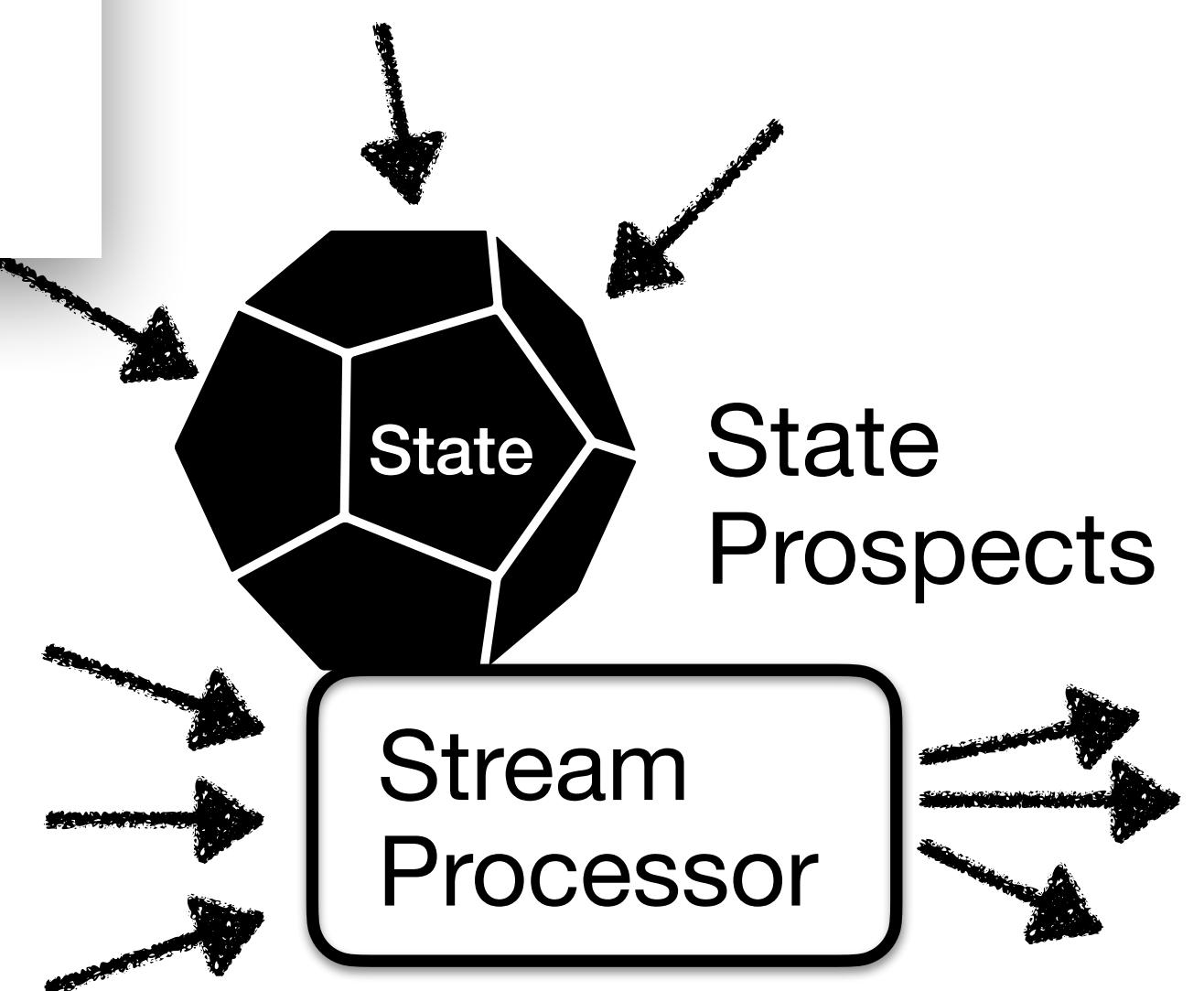


Stream Processors : Use LSM for every operator.

In support of workload-aware streaming state management

V Kalavri, J Liagouris

12th {USENIX} Workshop on Hot Topics in Storage and File Systems (HotStorage 20)

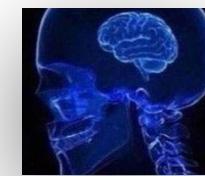


Taking Stateful Processing Further

Workload-Aware State Backends



Databases : Optimal plans using different indexes (e.g., B+Trees)



Stream Processors : Use LSM for every operator.

In support of workload-aware streaming state management

V Kalavri, J Liagouris

12th {USENIX} Workshop on Hot Topics in Storage and File Systems (HotStorage 20)

Queryable State

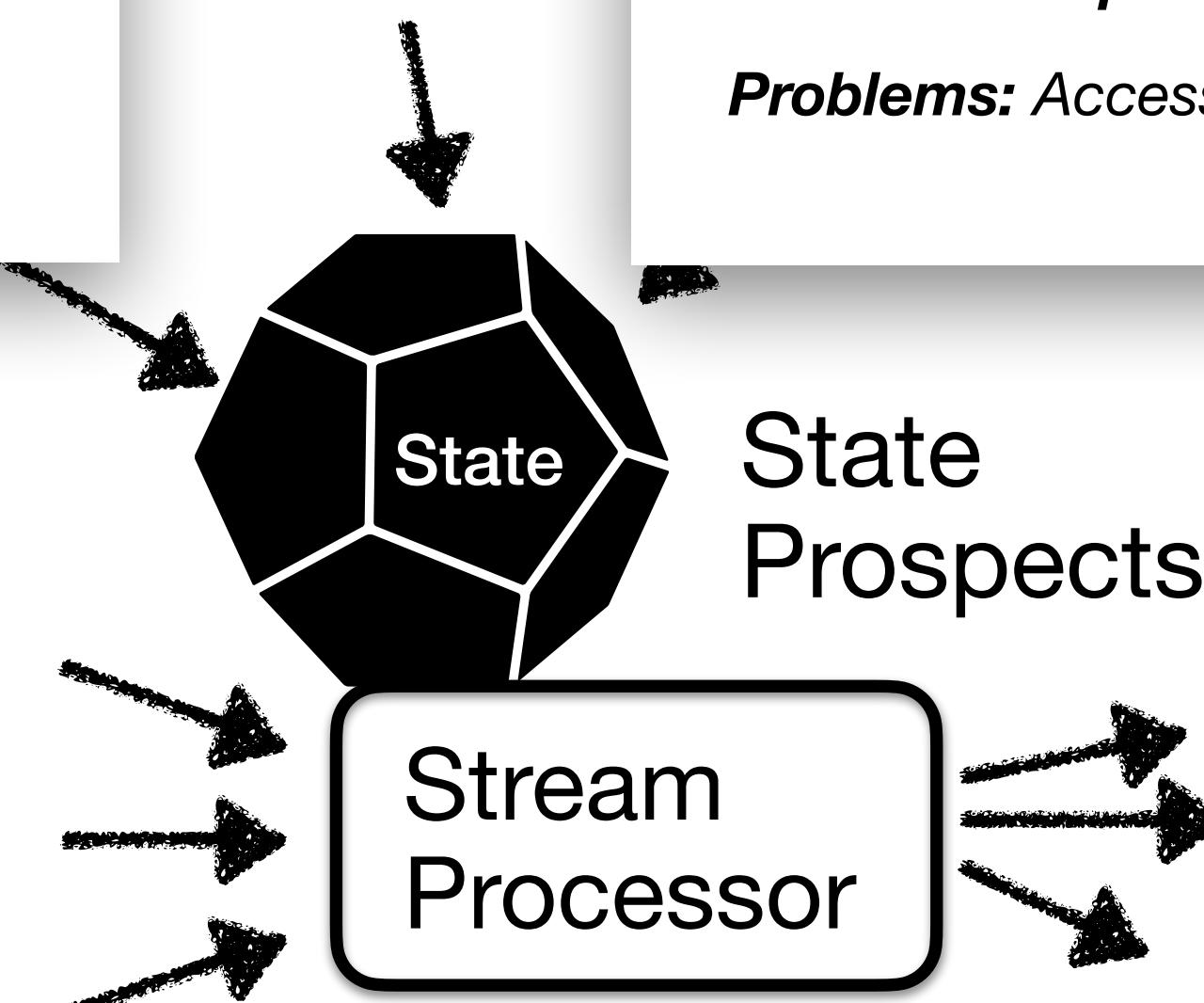
Currently : State is Black Box to the User

Queryable:

- **Materialized Views**
- **Exploratory Data Analysis**



Problems: Access Isolation? Indexes?



State
Prospects

Stream
Processor

Taking Stateful Processing Further

Workload-Aware State Backends



Databases : Optimal plans using different indexes (e.g., B+Trees)



Stream Processors : Use LSM for every operator.

In support of workload-aware streaming state management

V Kalavri, J Liagouris

12th {USENIX} Workshop on Hot Topics in Storage and File Systems (HotStorage 20)

Queryable State

Currently : State is Black Box to the User

Queryable:

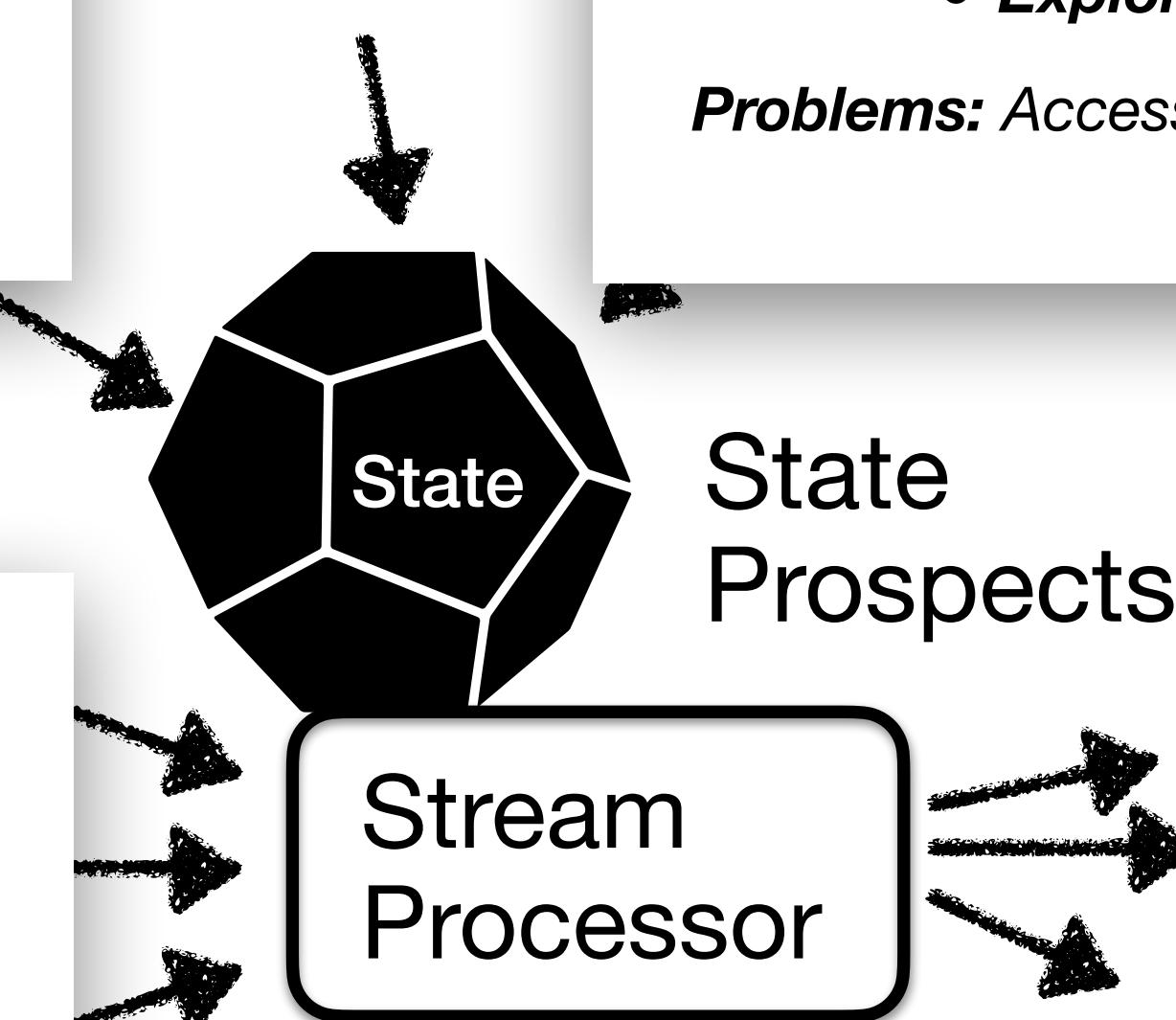
- *Materialized Views*
- *Exploratory Data Analysis*



Problems: Access Isolation? Indexes?

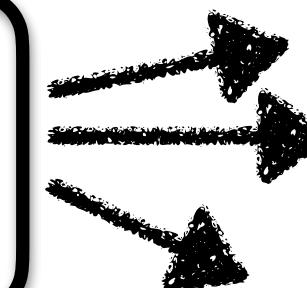
State Versioning

- **Schema Evolution**
- **Time-Travel Stream State Access**
- **MVCC on Stream State**



State
Prospects

Stream
Processor

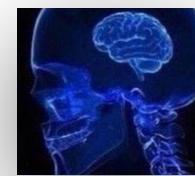


Taking Stateful Processing Further

Workload-Aware State Backends



Databases : Optimal plans using different indexes (e.g., B+Trees)

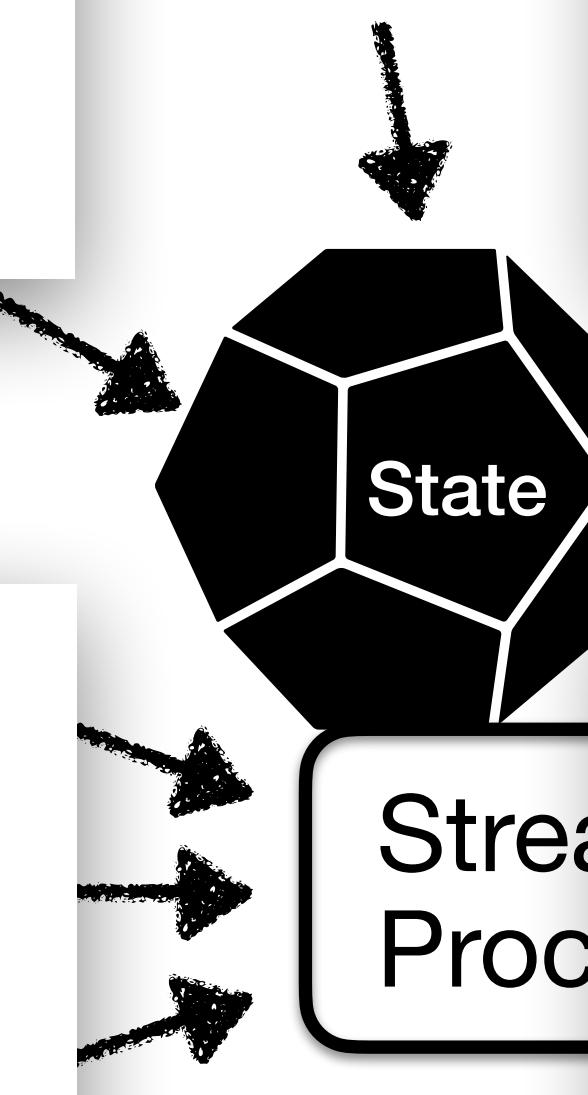


Stream Processors : Use LSM for every operator.

In support of workload-aware streaming state management

V Kalavri, J Liagouris

12th {USENIX} Workshop on Hot Topics in Storage and File Systems (HotStorage 20)



State Versioning

- Schema Evolution
- Time-Travel Stream State Access
- MVCC on Stream State



Queryable State

Currently : State is Black Box to the User

Queryable:

- *Materialized Views*
- *Exploratory Data Analysis*



Problems: Access Isolation? Indexes?

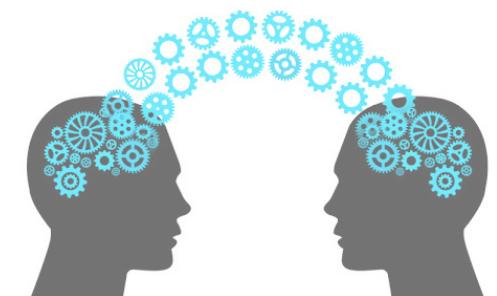
State Prospects

Inter-Query/Task State Sharing

Currently : State restricted to dedicated task access

Shared Access:

- Mutable shared state for ML applications (Ray)
- Inter-Query Optimisation (Aurora, Shared Arrangements)

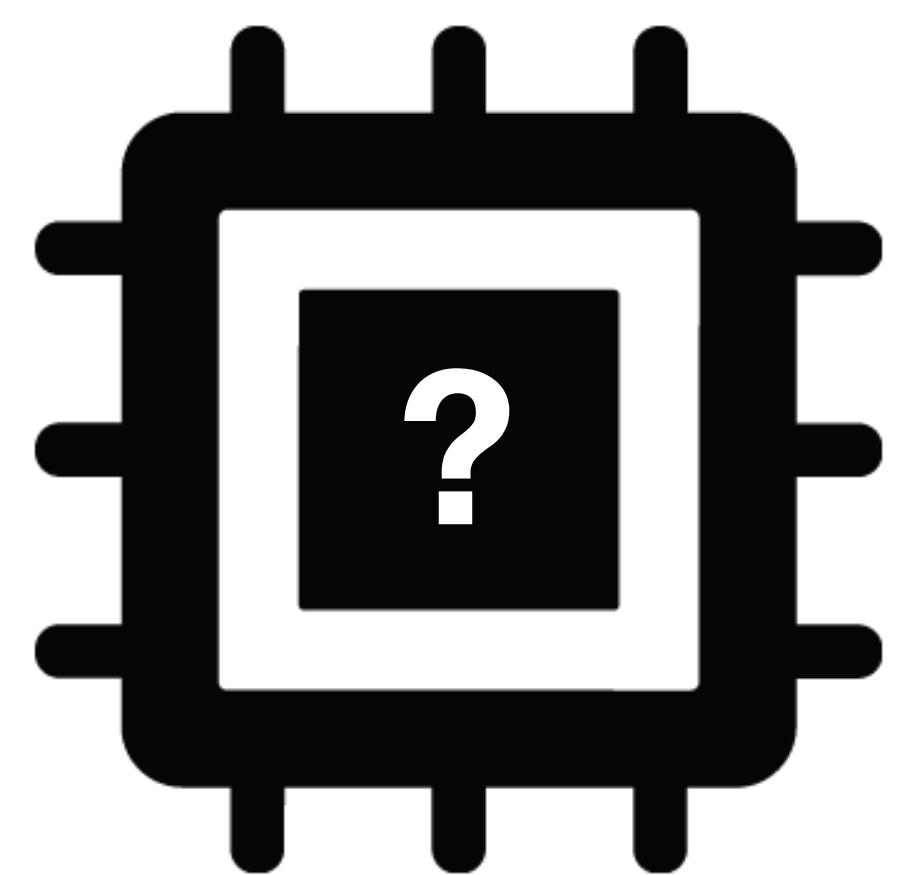


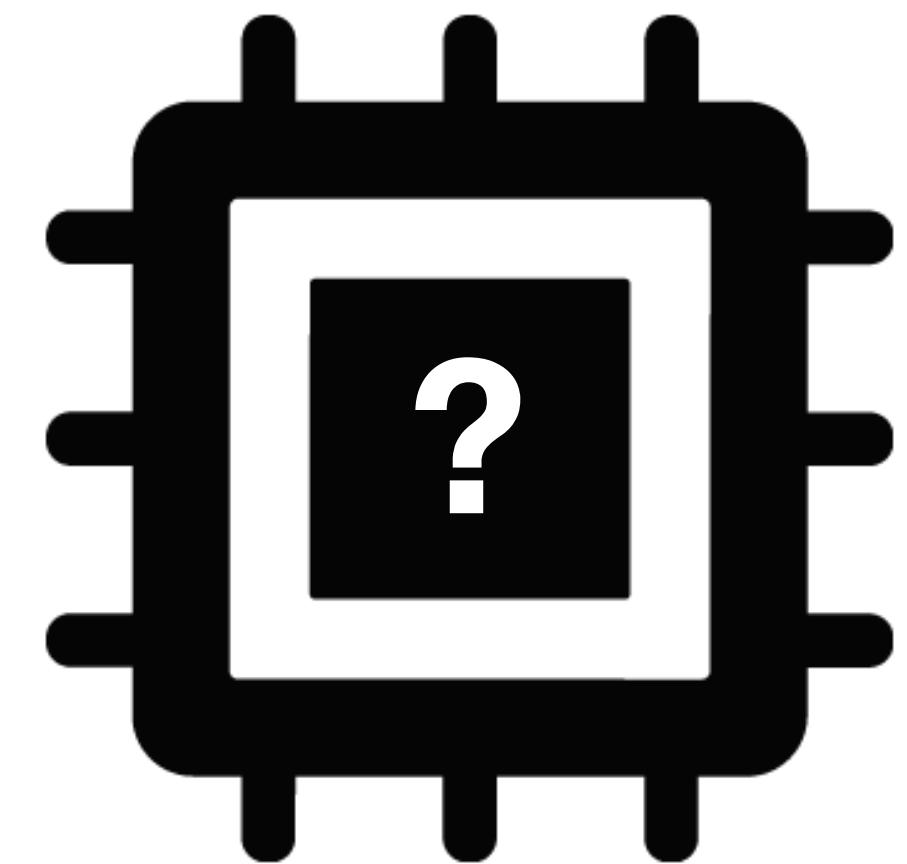
Shared Arrangements: practical inter-query sharing for streaming dataflows

McSherry F, Lattuada M, Schwarzkopf A, Roscoe T, et. al.

II. New requirements

Hardware acceleration





Hardware-conscious stream processing: A survey

[S Zhang, F Zhang, Y Wu, B He, P Johns - ACM SIGMOD Record, 2020 - dl.acm.org](https://dl.acm.org/)

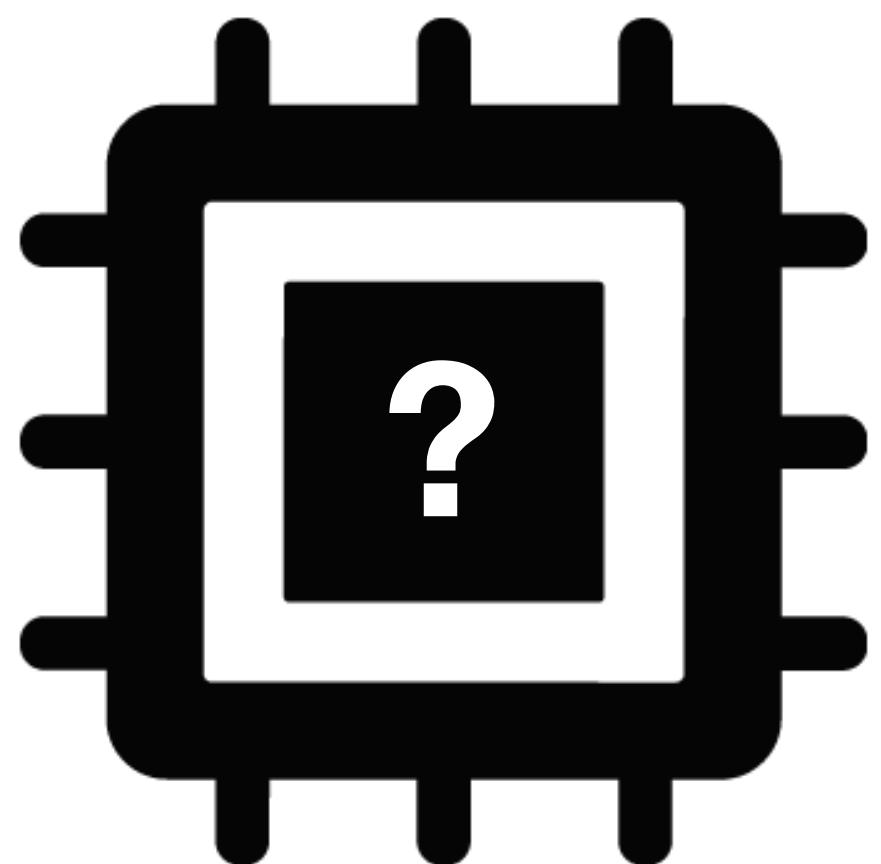
II. New requirements

Hardware acceleration

GPUs



Saber: Window-based hybrid stream processing for heterogeneous architectures
A Koliousis, M Weidlich, R Castro Fernandez, AL Wolf, P Costa, ...
Proceedings of the 2016 International Conference on Management of Data, 555-569



Hardware-conscious stream processing: A survey
[S Zhang, F Zhang, Y Wu, B He, P Johns - ACM SIGMOD Record, 2020 - dl.acm.org](#)

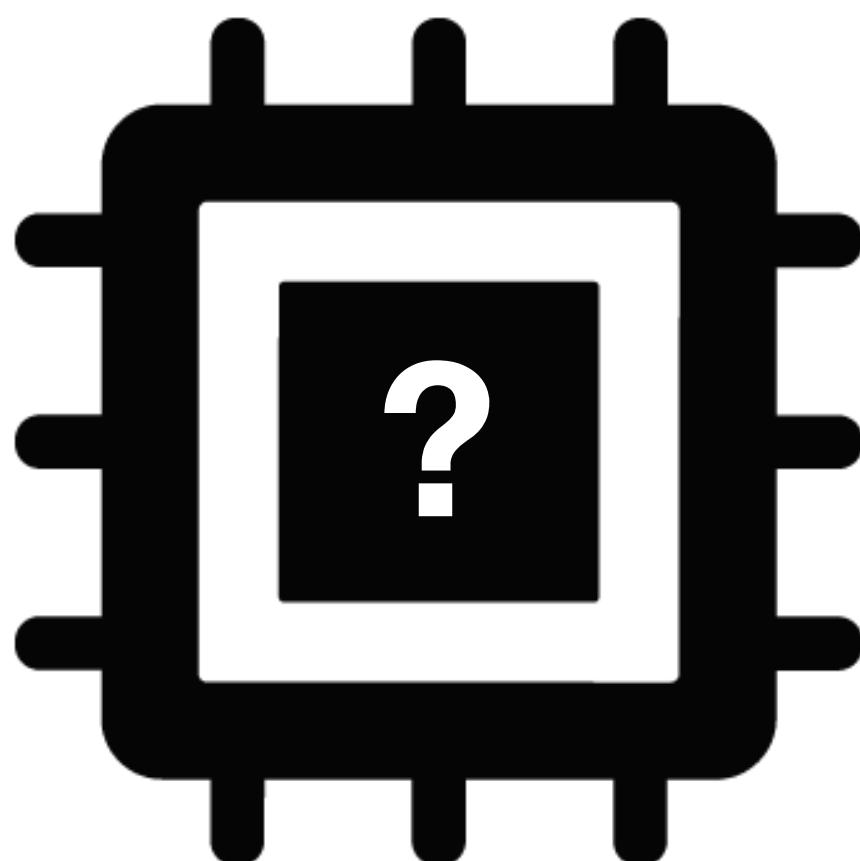
II. New requirements

Hardware acceleration

GPUs

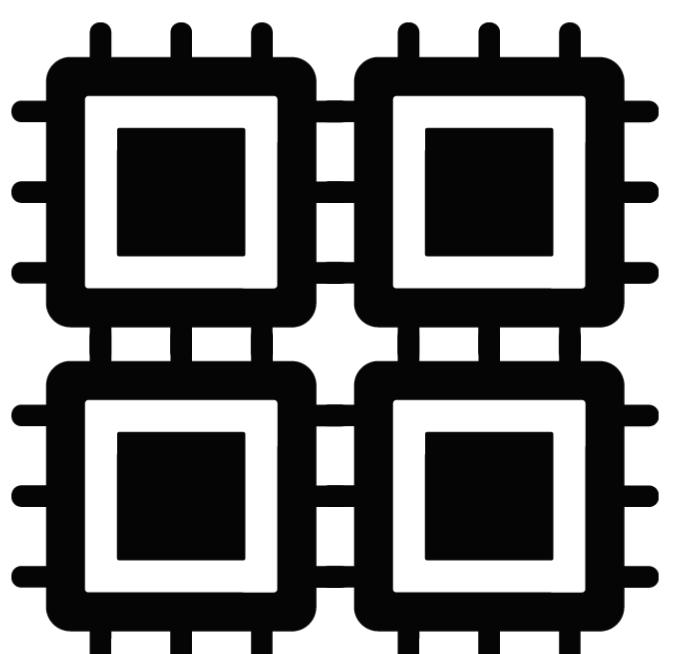


Saber: Window-based hybrid stream processing for heterogeneous architectures
A Koliousis, M Weidlich, R Castro Fernandez, AL Wolf, P Costa, ...
Proceedings of the 2016 International Conference on Management of Data, 555-569



Analyzing Efficient Stream Processing on Modern Hardware
S Zeuch, B Del Monte, J Karimov, C Lutz, M Renz, J Traub, S Breß, T Rabl, ...
Proceedings of the VLDB Endowment (VLDB)

LightSaber: Efficient Window Aggregation on Multi-core Processors
G Theodorakis, A Koliousis, P Pietzuch, H Pirk
Proceedings of the 2020 ACM SIGMOD International Conference on Management of ...



Multicore

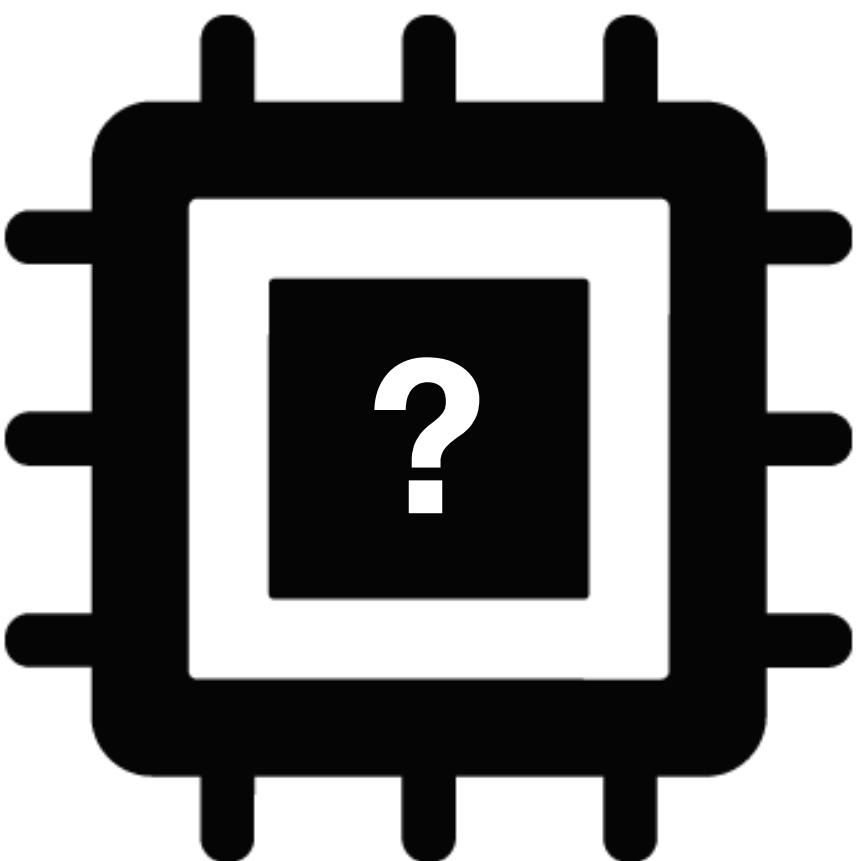
II. New requirements

Hardware acceleration

GPUs

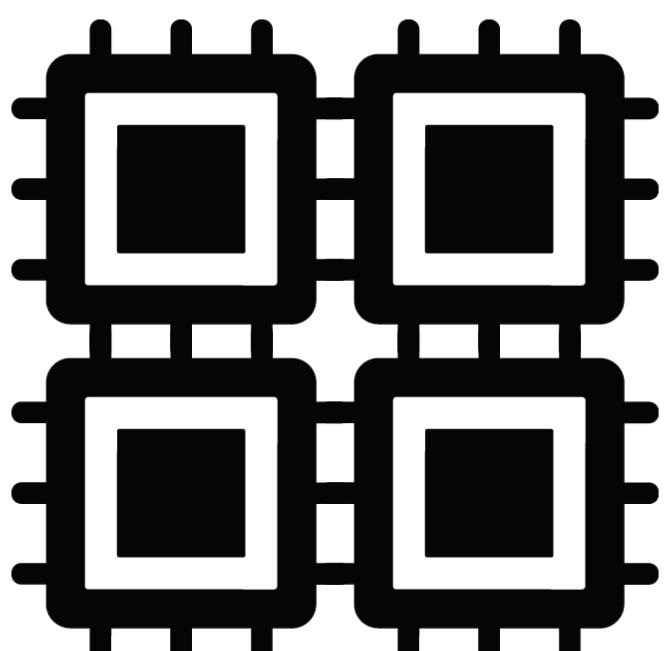


Saber: Window-based hybrid stream processing for heterogeneous architectures
A Koliousis, M Weidlich, R Castro Fernandez, AL Wolf, P Costa, ...
Proceedings of the 2016 International Conference on Management of Data, 555-569

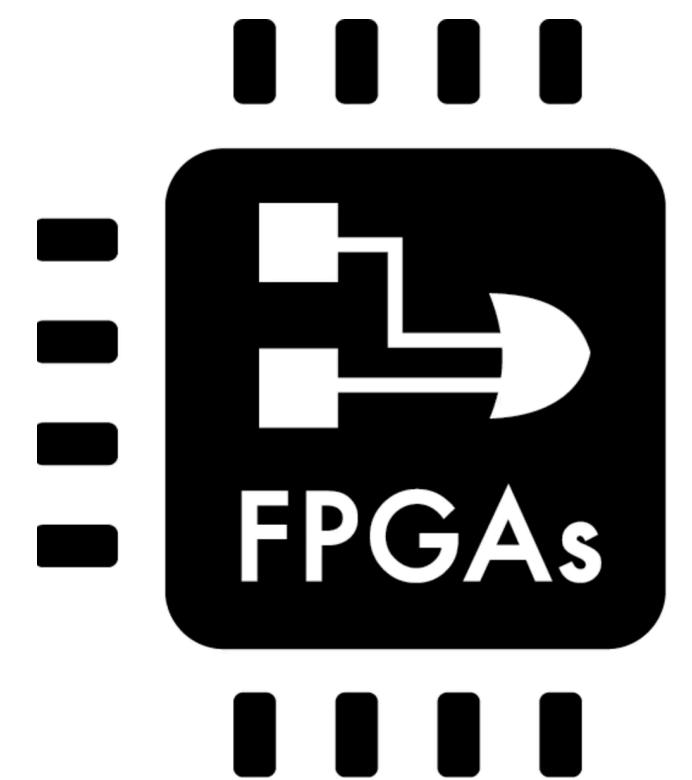


Analyzing Efficient Stream Processing on Modern Hardware
S Zeuch, B Del Monte, J Karimov, C Lutz, M Renz, J Traub, S Breß, T Rabl, ...
Proceedings of the VLDB Endowment (VLDB)

LightSaber: Efficient Window Aggregation on Multi-core Processors
G Theodorakis, A Koliousis, P Pietzuch, H Pirk
Proceedings of the 2020 ACM SIGMOD International Conference on Management of ...



Multicore



Fleet: A framework for massively parallel streaming on FPGAs
J Thomas, P Hanrahan, M Zaharia
Proceedings of the Twenty-Fifth International Conference on Architectural ...

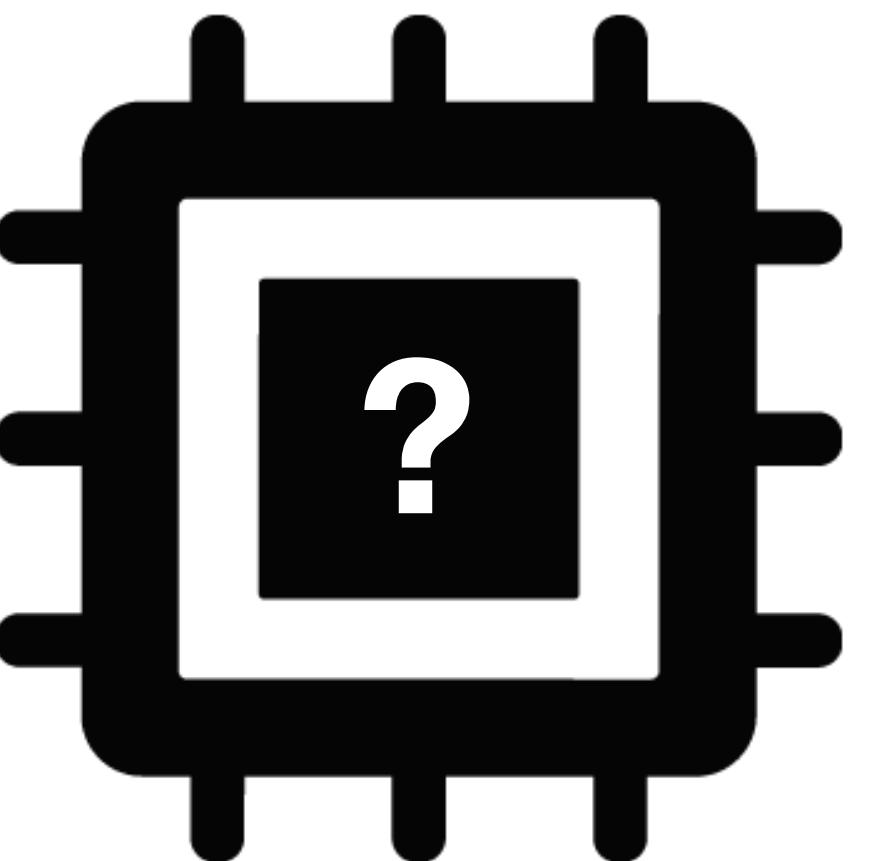
II. New requirements

Hardware acceleration

GPUs

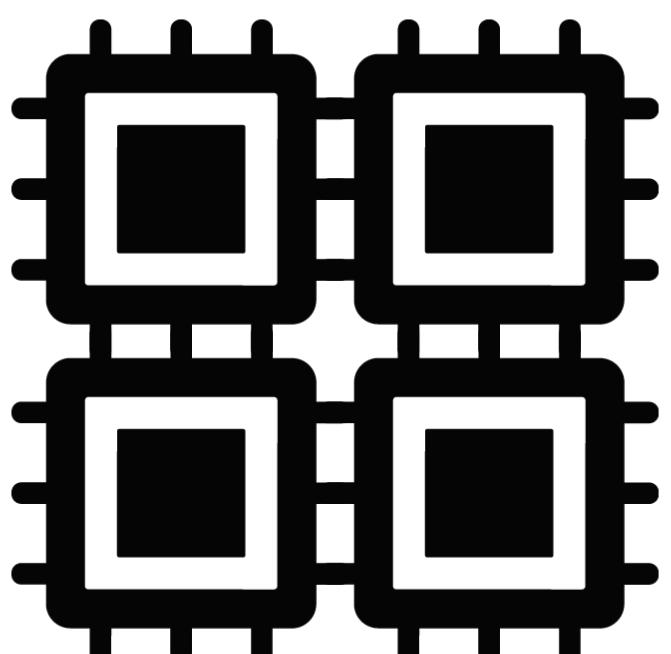


Saber: Window-based hybrid stream processing for heterogeneous architectures
A Koliousis, M Weidlich, R Castro Fernandez, AL Wolf, P Costa, ...
Proceedings of the 2016 International Conference on Management of Data, 555-569

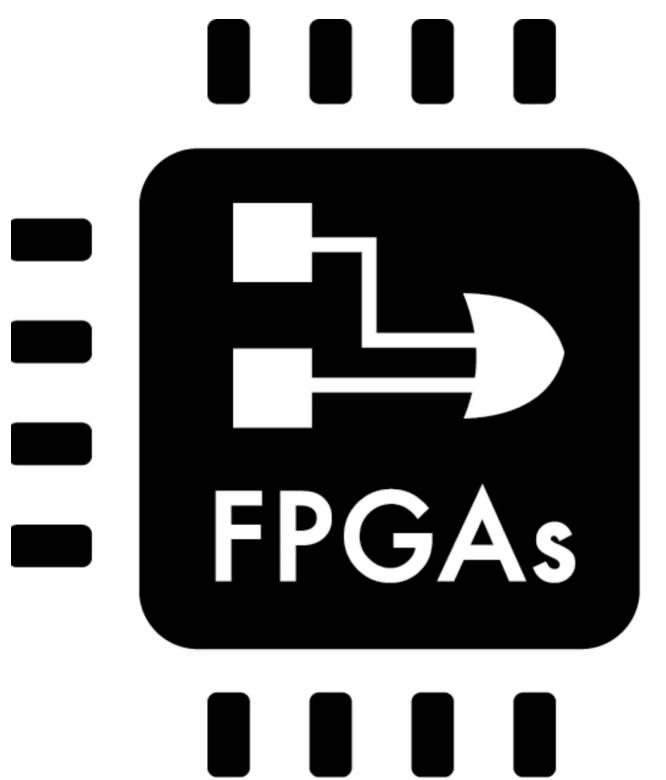


Analyzing Efficient Stream Processing on Modern Hardware
S Zeuch, B Del Monte, J Karimov, C Lutz, M Renz, J Traub, S Breß, T Rabl, ...
Proceedings of the VLDB Endowment (VLDB)

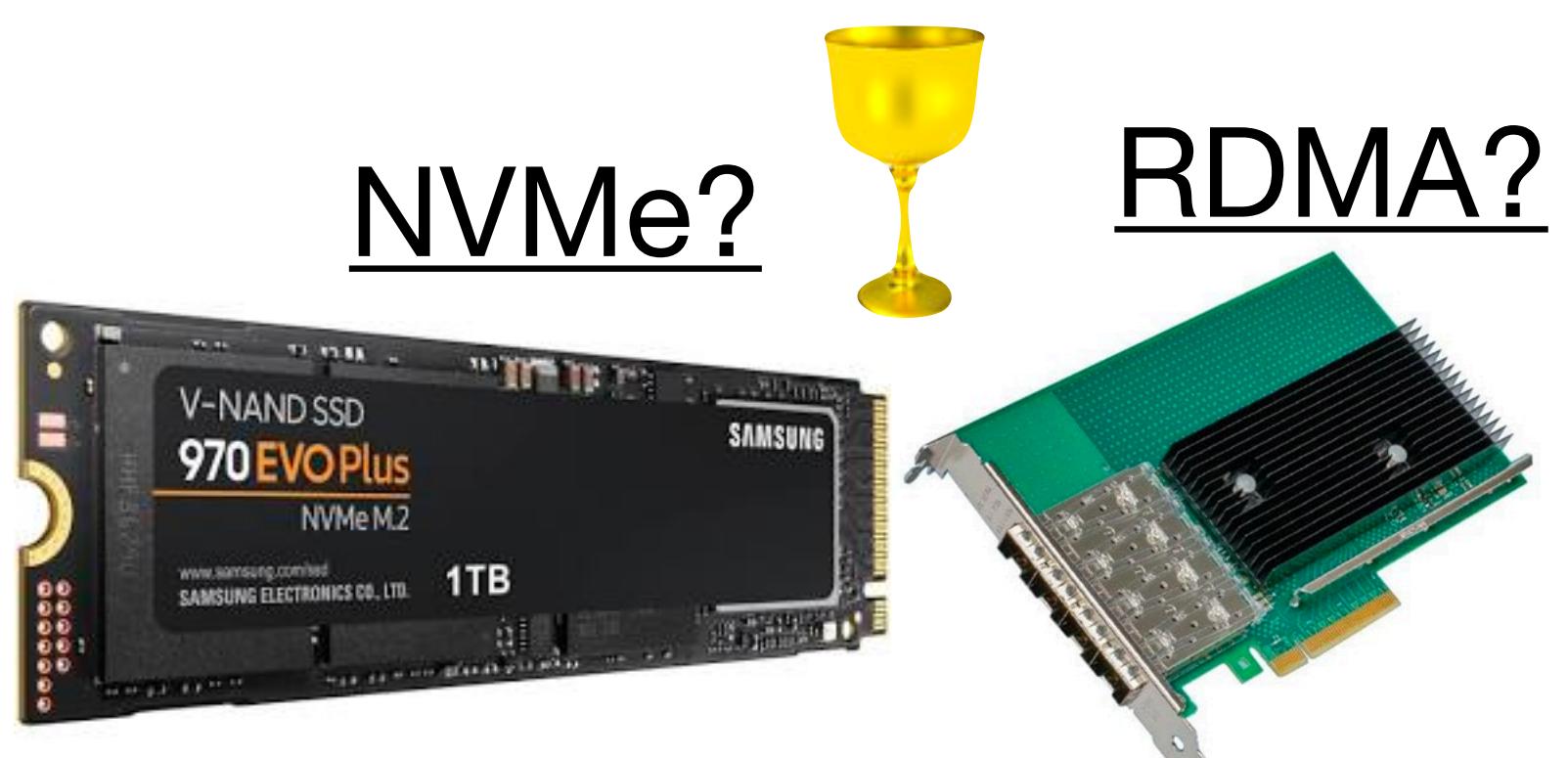
LightSaber: Efficient Window Aggregation on Multi-core Processors
G Theodorakis, A Koliousis, P Pietzuch, H Pirk
Proceedings of the 2020 ACM SIGMOD International Conference on Management of ...



Multicore



Fleet: A framework for massively parallel streaming on FPGAs
J Thomas, P Hanrahan, M Zaharia
Proceedings of the Twenty-Fifth International Conference on Architectural ...



III. Conclusion

Requirements overview

| | Cloud Apps | ML | Graphs |
|------------------------------|------------|----|--------|
| Programming models | ✓ | ✓ | ✓ |
| Transactions | ✓ | | |
| Advanced state backends | ✓ | ✓ | ✓ |
| Loops and cycles | ✓ | ✓ | ✓ |
| Elasticity & reconfiguration | ✓ | | |
| Dynamic topologies | ✓ | ✓ | |
| Shared mutable state | | ✓ | ✓ |
| Queryable state | ✓ | ✓ | |
| State versioning | ✓ | ✓ | |
| Hardware acceleration | | ✓ | ✓ |

Stream processing @SIGMOD'20

Research 28: Stream Processing, Thursday 10:30 AM – 12:00 PM

- Ahmed S. Abdelhamid, Ahmed R. Mahmood, Anas Daghstani, and Walid G. Aref. **Prompt: Dynamic Data-Partitioning for Distributed Micro-batch Stream Processing Systems.**
- Bonaventura Del Monte, Steffen Zeuch, Tilmann Rabl, and Volker Markl. **Rhino: Efficient Management of Very Large Distributed State for Stream Processing Engines.**
- Philipp M. Grulich, Breß Sebastian, Steffen Zeuch, Jonas Traub, Janis von Bleichert, Zongxiong Chen, Tilmann Rabl, and Volker Markl. **Grizzly: Efficient Stream Processing Through Adaptive Query Compilation.**
- Georgios Theodorakis, Alexandros Koliousis, Peter Pietzuch, and Holger Pirk. **LightSaber: Efficient Window Aggregation on Multi-core Processors.**
- Amirhesam Shahvarani and Hans-Arno Jacobsen. **Parallel Index-based Stream Join on a Multicore CPU.**

Beyond Analytics

The Evolution of Stream Processing Systems

Prospects

Paris Carbone, Marios Fragkoulis, Vasiliki Kalavri, Asterios Katsifodimos

Slides: streaming-research.github.io/Tutorial-SIGMOD-2020

