

Stream Loops on Flink

Reinventing the wheel for the streaming era



Paris Carbone

Systems Researcher@KTH/SICS <parisc@kth.se>
Committer@Flink <senorcarbone@apache.org>

Stream Compute Landmarks

A Timeline

**Sliding
Window
Semantics**

State Management
-Fault Tolerance
-Reconfiguration

**Arbitrary
Computation
Nesting**

**Task Parallel
Execution**

**Out-of-order
Processing
(watermarks)**

Has Streaming Subsumed Batch?

short answer: NO

Flink
DataSet



Flink
DataStream



single-pass

- | | |
|---|--|
| <ul style="list-style-type: none"> • Staged Processing • Allows Bulk/Delta Iterative Computation | <ul style="list-style-type: none"> • Continuous Processing • Allows correct Acyclic Computation |
|---|--|

long answer: not YET

Why Iterations Matter

- Iterations are fundamental building blocks for:
 - Graph Analysis (PageRank, Conn.Comp., SSSP etc.)
 - Machine Learning (Gradient Descent, PCA etc.)
 - Transactions (e.g., optimistic concurrency control)



Iterative Processing Primitives are currently not present in today's production-grade stream processing systems.

Current Challenges

- Data is born and evolves continuously as a data stream (e.g., user interactions, sensor events, server logs)



Current Challenges

- To run iterative analysis users have manually do the...
- bucketing



Day 1



Day 2



Day 3



Day 4



Current Challenges

- To run iterative analysis users have manually do the...
 - bucketing
 - laundry (iterative job scheduling)

Day 1



Day 2



Day 3



Day 4



Current Challenges

- To run iterative analysis users have manually do the...
 - bucketing
 - laundry (iterative job scheduling)
 - sorting (model integration)

Day 1



Day 2



Day 3

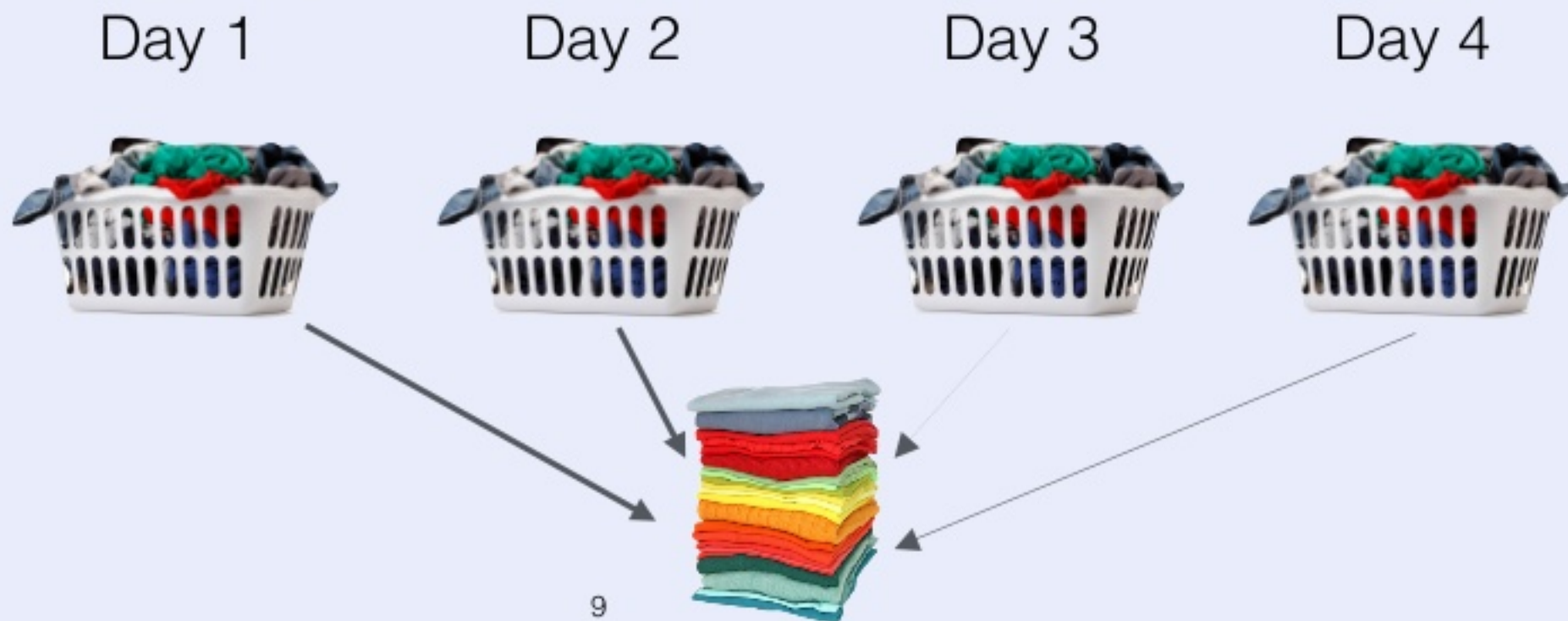


Day 4



Current Challenges

- To run iterative analysis users have manually do the...
 - bucketing
 - laundry (iterative job scheduling)
 - sorting (model integration)



A Partially Solved Problem

- Most challenges mentioned are solved and automated for single-pass aggregation via stream windowing.

window
computation



windows are a **natural fit** for
declaring iterative computation

How would a **distributed iterative
computation** look like as a **special
type of window aggregation**?

Anatomy of Iterative Computation

Solution (finite state)

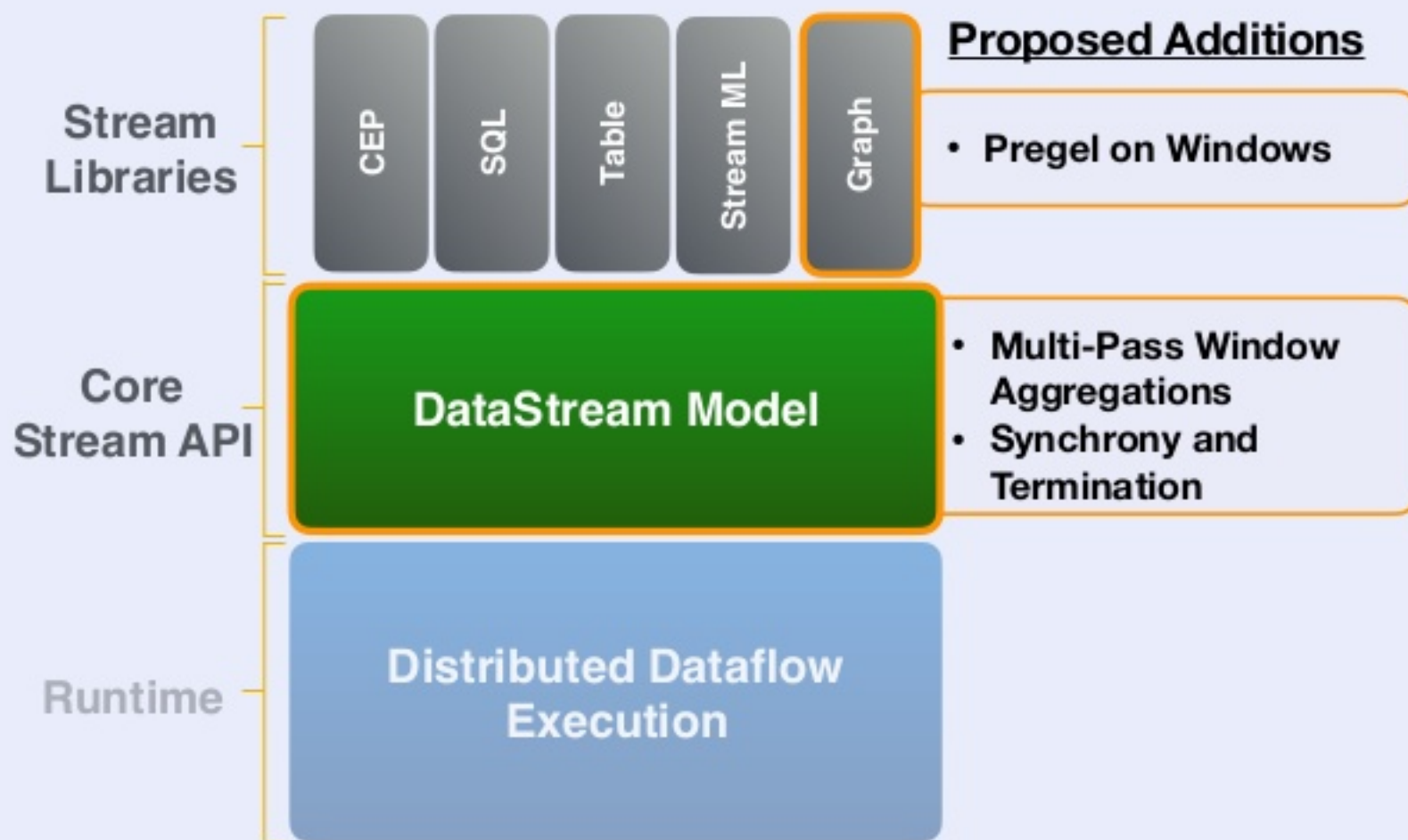
Loop Step Function (computation)

```
1  def iterate(c, L) ( $x_{init}$ )  
2       $x$  =  $x_{init}$   
3      do {  
4           $x$  = L( $x$ )  
5      } while ( $\neg$  c( $x$ , L( $x$ )))  
6       $x$ 
```

Final Result-Solution

Termination (fixpoint/fixed)

Programming Model Extensions



based on Flink 1.6

Window Iterations

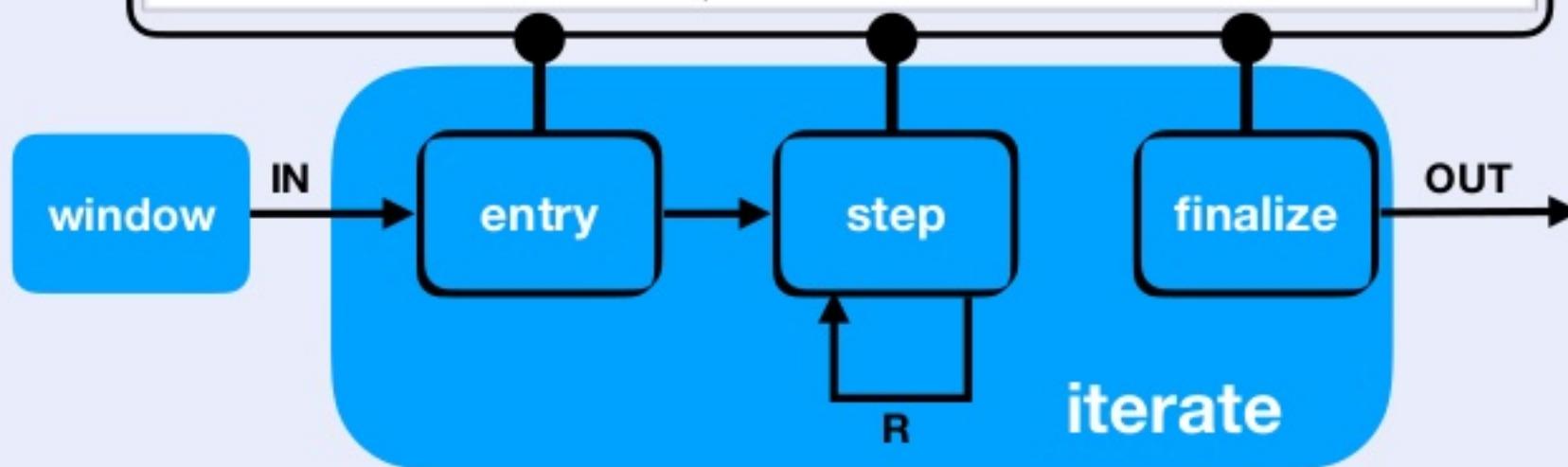
An extension of existing window aggregation to multi-pass aggregates

```
1 | val stream: DataStream[IN] = ...
2 | stream.keyBy(...).window(...)
3 |   .iterate(<Termination>, <Synchrony>, <Key>, <Entry>, <Step>, <Finalize>)
```

Iteration Primitive	Definition	Description
<i>Termination</i>	Termination.Fixpoint Termination.Fixed($n \in \mathbb{N}$)	Loop Termination Criterion
<i>Synchronization</i>	Synchrony.Stale($n \in \mathbb{N}$) Synchrony.Strict	Loop Synchronization Type
<i>Loop Key</i>	$R \rightarrow K$	Key Extractor for Structured Loop
<i>Entry Function</i>	$(ctx, [IN]) \rightarrow (ctx, [R])$	Initialization Logic
<i>Step Function</i>	$(ctx, [R]) \rightarrow (ctx, [R])$	Iteration Step Logic
<i>Finalize Function</i>	$(ctx) \rightarrow (ctx, [OUT])$	Finalization Logic after Termination

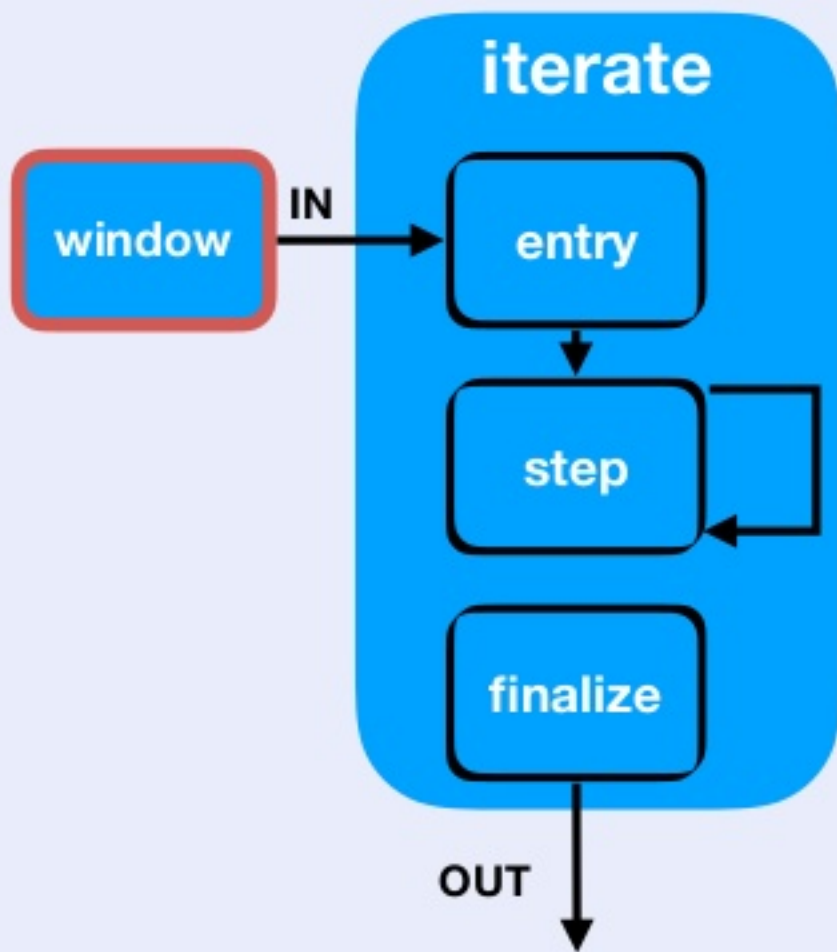
Loop Context

LoopContext Properties	Definition
<i>key</i> : K	The current key of the computation
<i>localProgress</i> : \mathbb{N}	The current superstep of the local computation
<i>globalProgress</i> : \mathbb{N}	The last parallel complete superstep
<i>eventTime</i> : \mathbb{N}	The event time (ctx) of local computation
<i>loopState</i> : <i>ManagedState</i>	keyed state, purgeable per iterative process
<i>persistentState</i> : <i>ManagedState</i>	keyed state, persistent across iterative processes

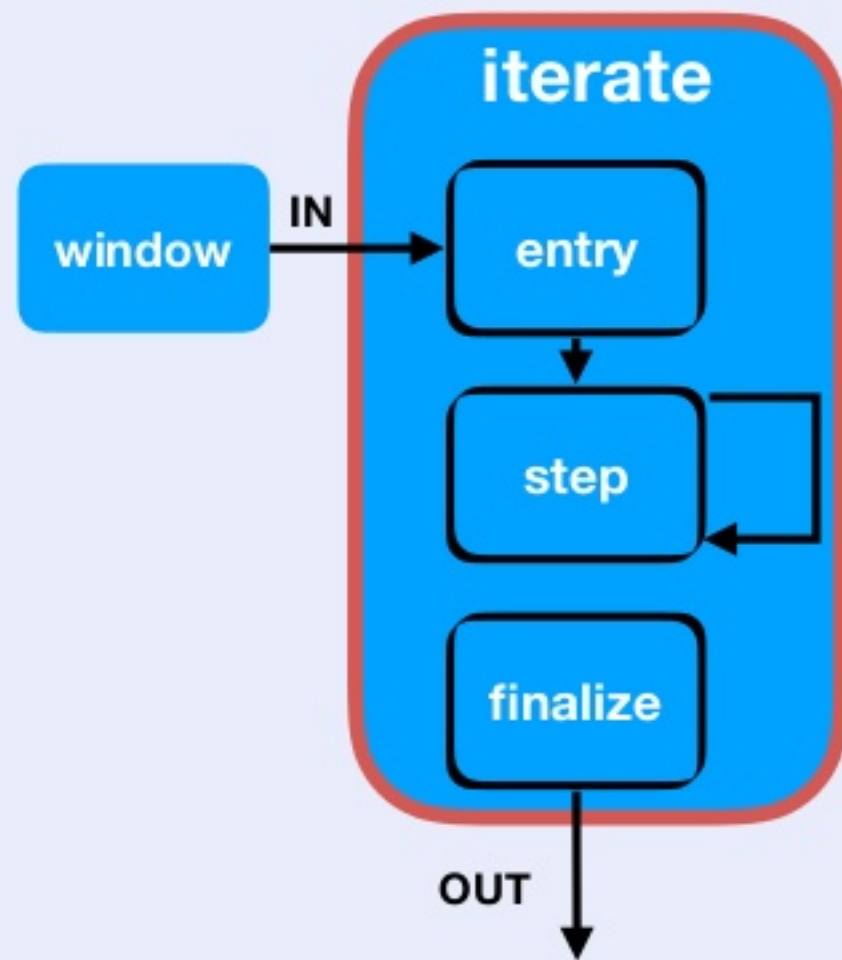


PageRank Example

```
1 case class GraphEdge(from:Long, to:Long)
2 case class VRank(id:Long, rank:Double)
3
4 def computeRank(w:Iterable[Double]) = 0.15/w.size + 0.85 * w.sum
5
6 val input: DataStream[GraphEdge] = getEdgeStream()
7 input
8   .flatMap(e => List(e, GraphEdge(e.to, e.from)))
9   .keyBy(edge => edge.from)
10  .timeWindow(30 Sec)
```



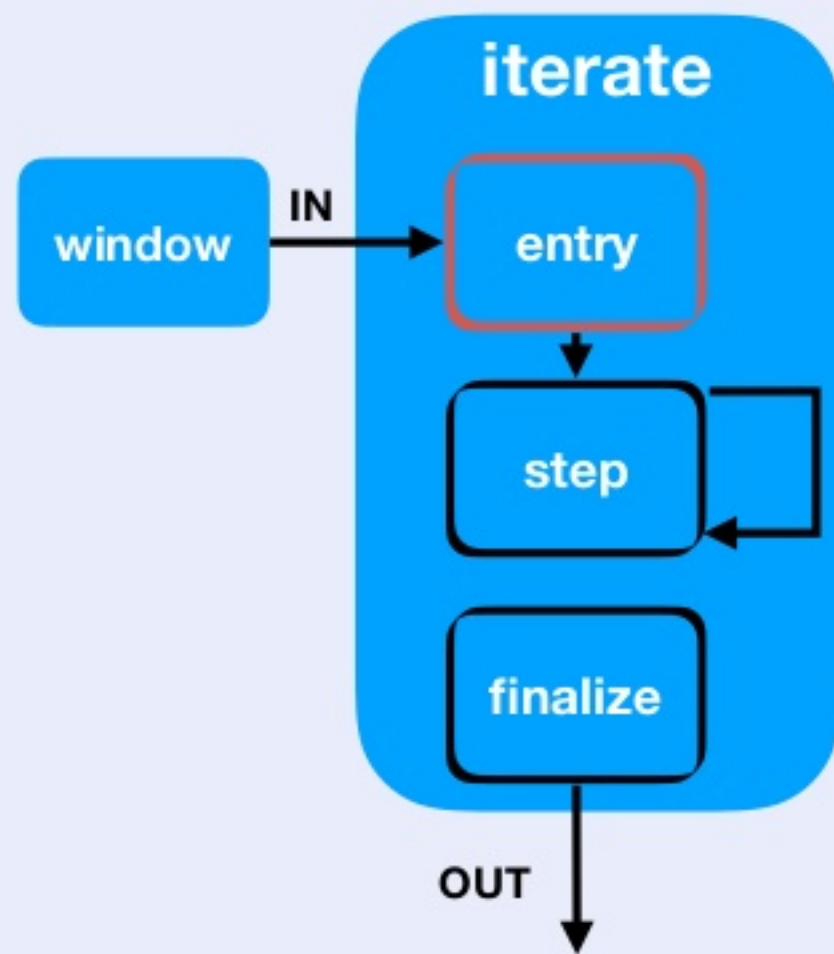
PageRank Example



```

1 case class GraphEdge(from:Long, to:Long)
2 case class VRank(id:Long, rank:Double)
3
4 def computeRank(w:Iterable[Double]) = 0.15/w.size + 0.85 * w.sum
5
6 val input: DataStream[GraphEdge] = getEdgeStream()
7 input
8   .flatMap(e => List(e, GraphEdge(e.to, e.from)))
9   .keyBy(edge => edge.from)
10  .timeWindow(30 Sec)
11  .iterate(Termination.Fixed(40), Synchrony.Strict,
12          (vRank => vRank.id),
  
```

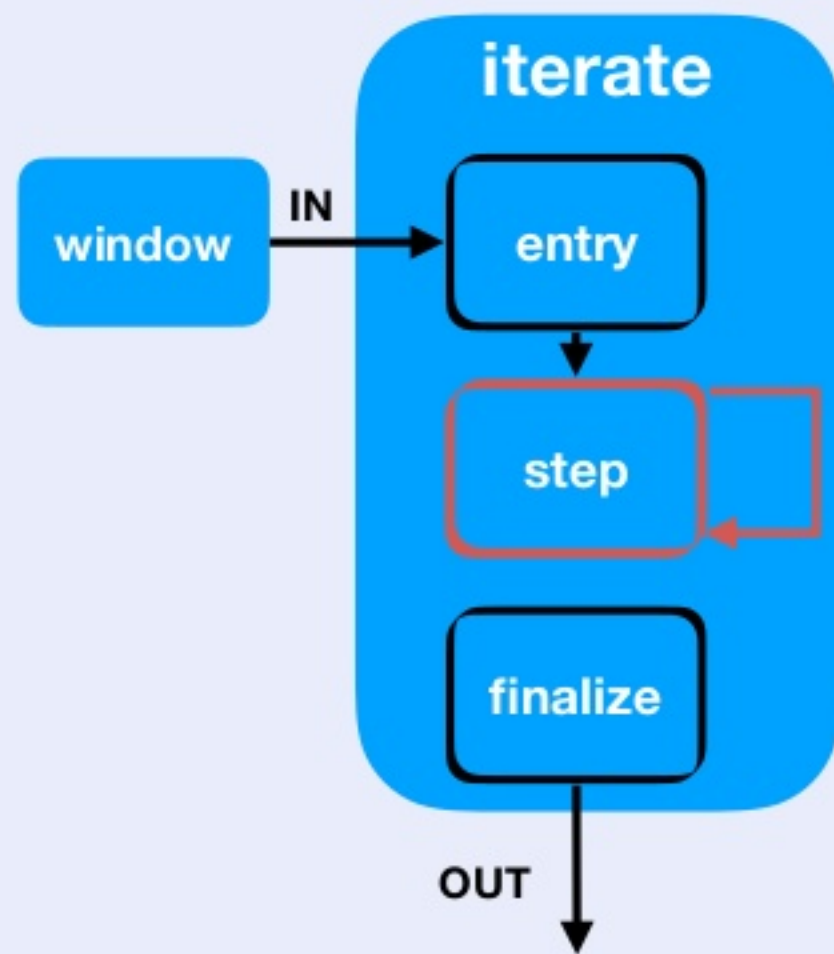
PageRank Example



```

1 case class GraphEdge(from:Long, to:Long)
2 case class VRank(id:Long, rank:Double)
3
4 def computeRank(w:Iterable[Double]) = 0.15/w.size + 0.85 * w.sum
5
6 val input: DataStream[GraphEdge] = getEdgeStream()
7 input
8   .flatMap(e => List(e, GraphEdge(e.to, e.from)))
9   .keyBy(edge => edge.from)
10  .timeWindow(30 Sec)
11  .iterate(Termination.Fixed(40), Synchrony.Strict,
12    (vRank => vRank.id),
13    (ctx:LoopContext, input:Iterable[GraphEdge], out:Collector[VRank]) =>
14      { //ENTRY FUNCTION
15        ctx.loopState("neighbors").setList(input.map(e => e.to))
16        //Start from stored rank
17        val rank = ctx.loopState("rank")
18        rank.setValue(ctx.persistentState("rank").value)
19        ctx.loopState("neighbors").foreach(n => out.collect(VRank(n,
20          rank.value)))
  
```

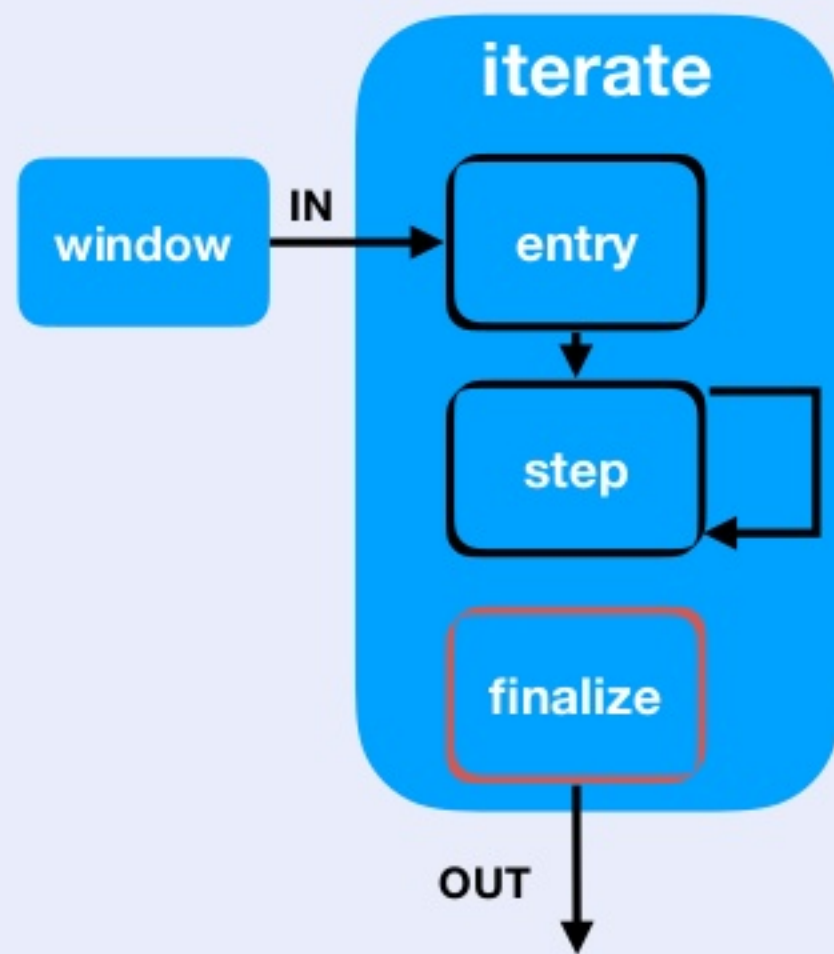

PageRank Example



```

1 case class GraphEdge(from:Long, to:Long)
2 case class VRank(id:Long, rank:Double)
3
4 def computeRank(w:Iterable[Double]) = 0.15/w.size + 0.85 * w.sum
5
6 val input: DataStream[GraphEdge] = getEdgeStream()
7 input
8   .flatMap(e => List(e, GraphEdge(e.to, e.from)))
9   .keyBy(edge => edge.from)
10  .timeWindow(30 Sec)
11  .iterate(Termination.Fixed(40), Synchrony.Strict,
12    (vRank => vRank.id),
13    (ctx:LoopContext, input:Iterable[GraphEdge], out:Collector[VRank]) =>
14      { //ENTRY FUNCTION
15        ctx.loopState("neighbors").setList(input.map(e => e.to))
16        //Start from stored rank
17        val rank = ctx.loopState("rank")
18        rank.setValue(ctx.persistentState("rank").value)
19        ctx.loopState("neighbors").foreach(n => out.collect(VRank(n,
20          rank.value)))
21      },
22      (ctx:LoopContext, input:Iterable[VRank], out:Collector[VRank]) =>
23        { //STEP FUNCTION
24          val newRank = computeRank(input.map(c => c.rank))
25          ctx.loopState("rank").setValue(newRank)
26          ctx.loopState("neighbors").foreach(n => out.collect(VRank(n,
27            newRank)))
28        }
29    )
  
```

PageRank Example



```

1 case class GraphEdge(from:Long, to:Long)
2 case class VRank(id:Long, rank:Double)
3
4 def computeRank(w:Iterable[Double]) = 0.15/w.size + 0.85 * w.sum
5
6 val input: DataStream[GraphEdge] = getEdgeStream()
7 input
8   .flatMap(e => List(e, GraphEdge(e.to, e.from)))
9   .keyBy(edge => edge.from)
10  .timeWindow(30 Sec)
11  .iterate(Termination.Fixed(40), Synchrony.Strict,
12    (vRank => vRank.id),
13    (ctx:LoopContext, input:Iterable[GraphEdge], out:Collector[VRank]) =>
14      { //ENTRY FUNCTION
15        ctx.loopState("neighbors").setList(input.map(e => e.to))
16        //Start from stored rank
17        val rank = ctx.loopState("rank")
18        rank.setValue(ctx.persistentState("rank").value)
19        ctx.loopState("neighbors").foreach(n => out.collect(VRank(n,
20          rank.value)))
21      },
22    (ctx:LoopContext, input:Iterable[VRank], out:Collector[VRank]) =>
23      { //STEP FUNCTION
24        val newRank = computeRank(input.map(c => c.rank))
25        ctx.loopState("rank").setValue(newRank)
26        ctx.loopState("neighbors").foreach(n => out.collect(VRank(n,
27          newRank)))
28      },
29    (ctx:LoopContext, out:Collector[VRank]) =>
30      { //FINALIZE FUNCTION
31        //Log new rank
32        val finalRank = ctx.loopState("rank").value
33        ctx.persistentState("rank").setValue(finalRank)
34        out.collect(VRank(ctx.key, finalRank))
35      }
36    ))
  
```

Higher-Level Abstractions

- Vertex-Centric Model (Part of experimental Gelly-Streams lib)

Example:
single source
shortest paths

```
1 //Stream of edges with nil (empty) state and vertices with 'Long' state
2 val input: EdgeStream<Long,Nil> edges = getEdgeStream()
3 //iterations on directed graph snapshots defined over a 5min tumbling
  window
4 input
5   .snapshot(Minutes(5), Direction.DIRECTED)
6   .runFixpoint(
7     (ctx:VertexContext) => //VERTEX COMPUTATION
8       { //vertex computation
9         if(ctx.superstep == 0 && ctx.isSource()){
10           ctx.setSnapshotState(0L)
11           ctx.neighbors.foreach(vertex => vertex.send(Message(1L)))
12         }
13         else{
14           dist = ctx.getVertexState().getOrElse(Long.max)
15           newDist = (dist :: ctx.getMessages().values).min
16           if(dist != newDist){
17             ctx.setSnapshotState(distance);
18             ctx.neighbors.foreach(vertex =>
19               vertex.send(Message(newDist+1)))
20           }
21         }
22       });
23   .toVertexStream(); // creates a DataStream<GraphVertex(distance)>
```


Enough about what

- We need to examine “how” to implement iterations

window
computation



executed in **parallel** and **out-of-order**
using event-time progress

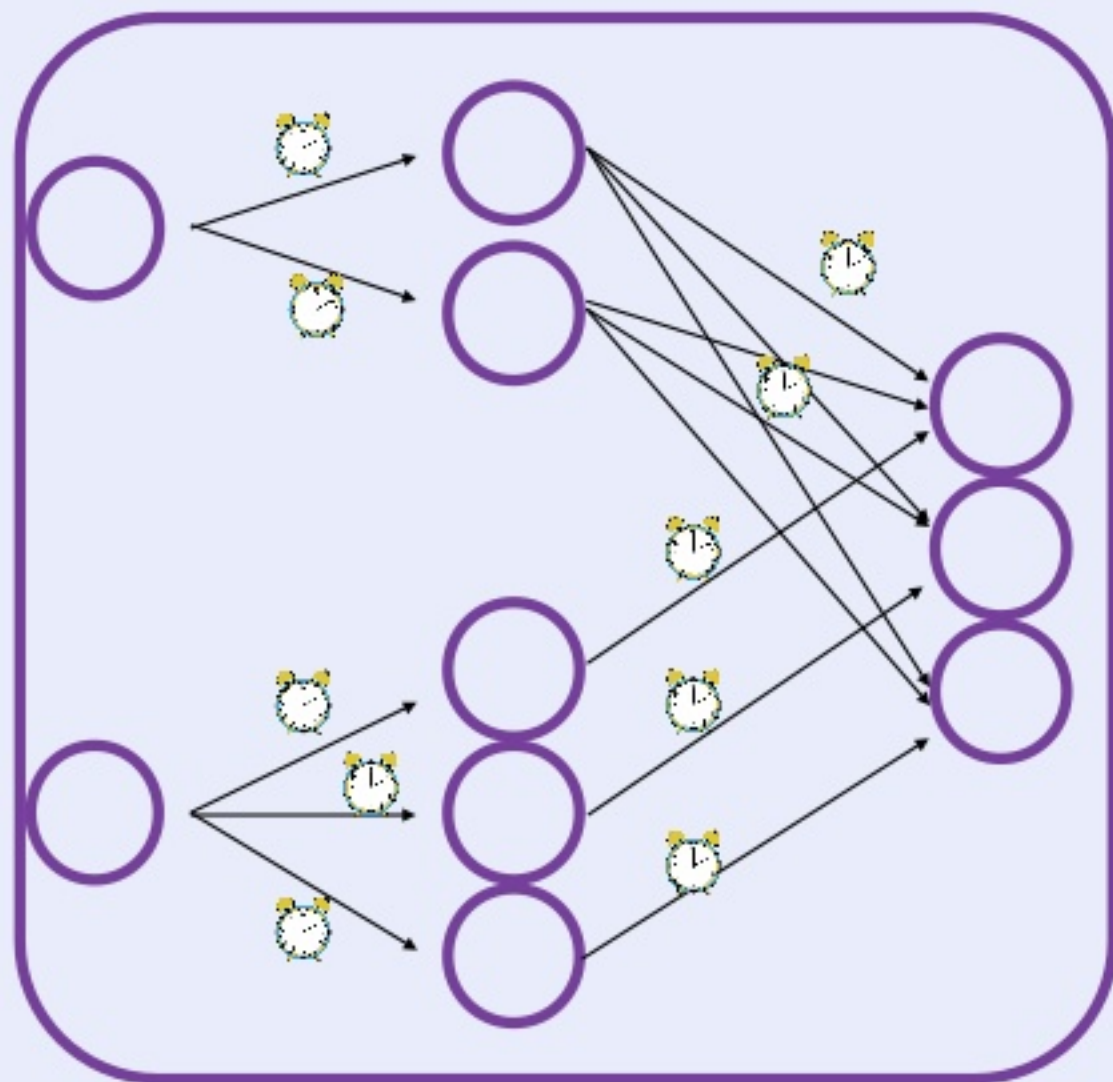


How Event-Time Progress Works



- **Every record carries a timestamp t**
- **$t \geq T$**
- **T can only increment (monotonicity)**

How Event-Time Progress Works



- **low watermarks** propagate progress metrics along the computational graph
- **progress metric = minimum** watermark across channels
- works correctly as long as monotonicity is maintained

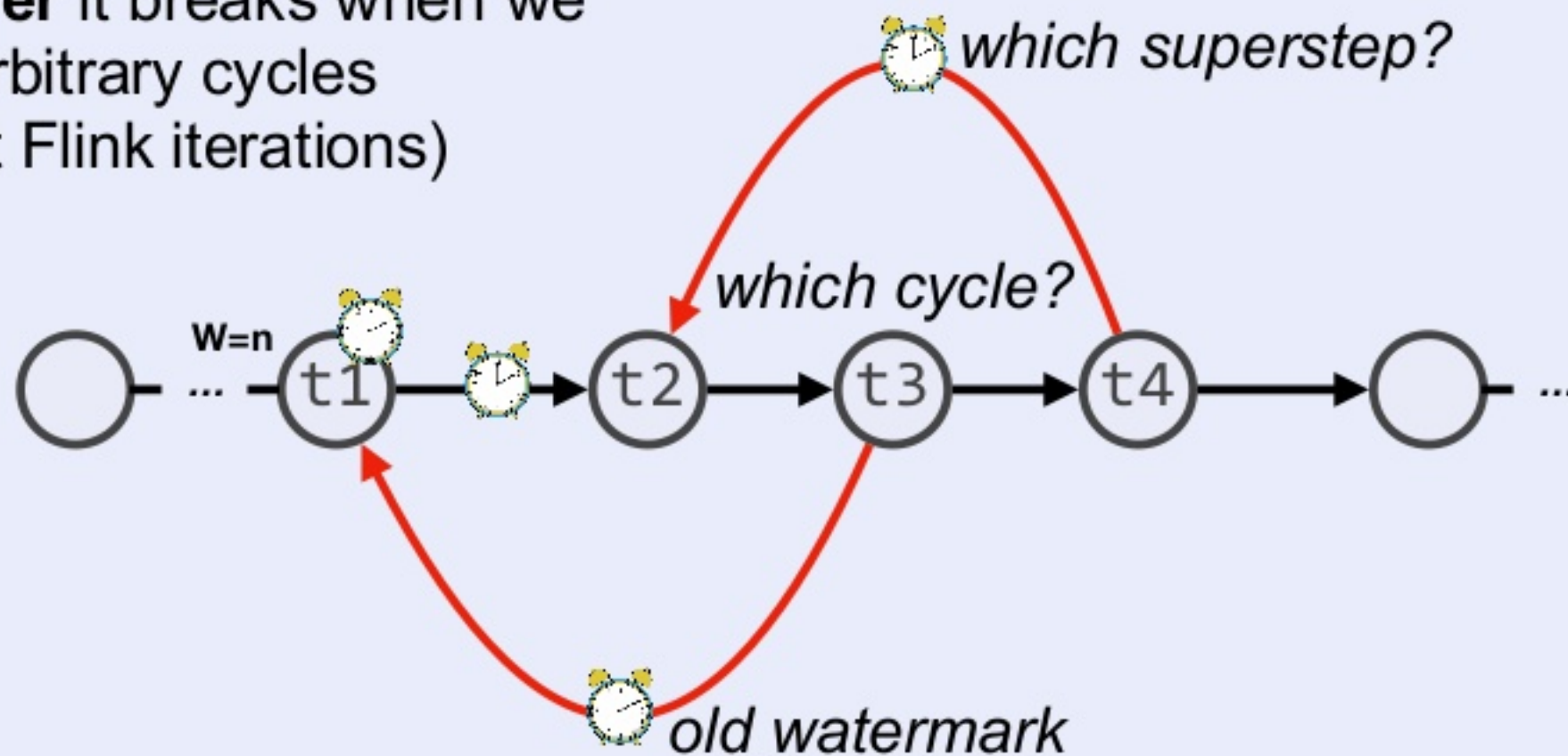
Intuition

Iteration Superstep progress bears similarities to **Event-Time** progress

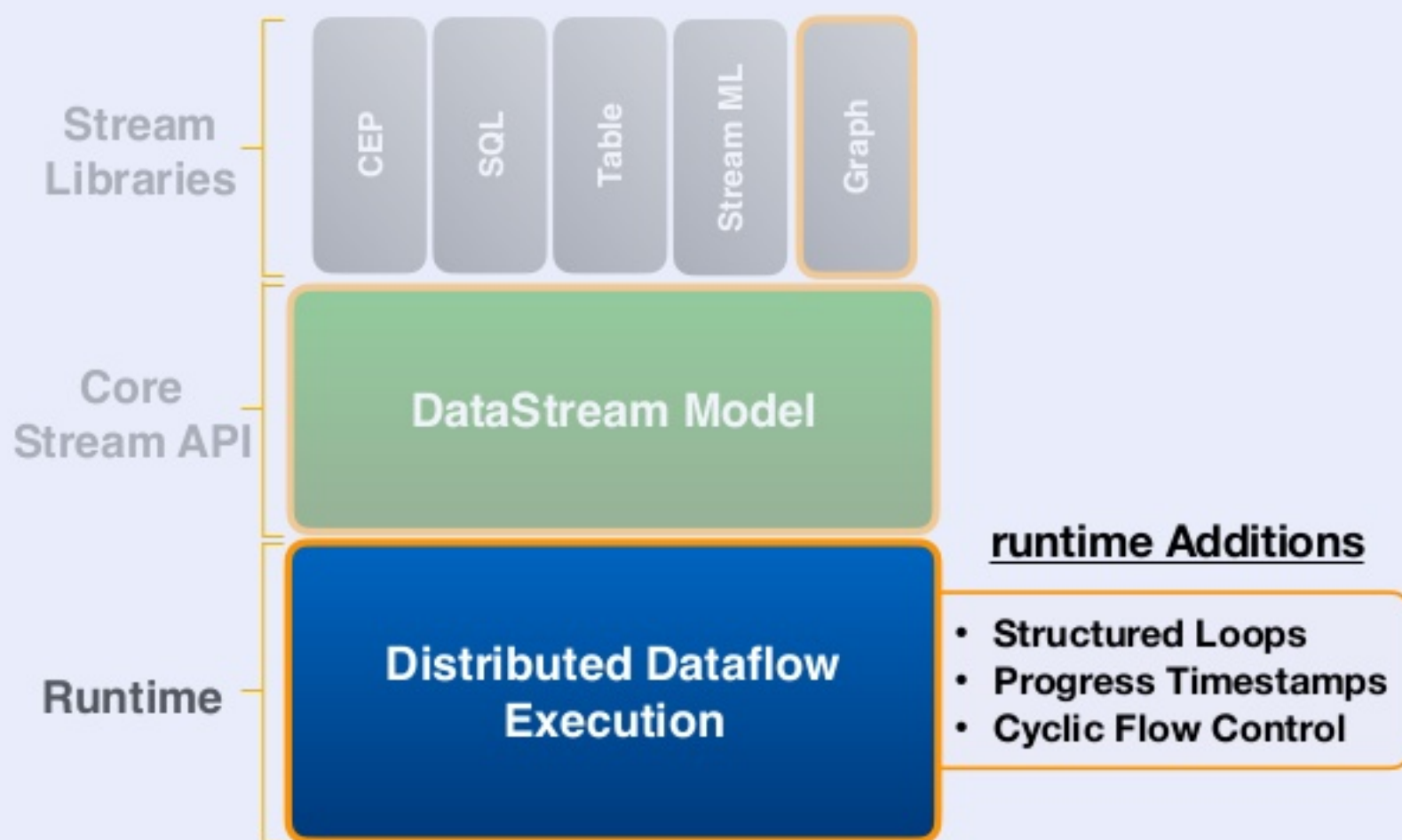
1. Can be used to run computation for out-of-order streams.
2. Need for a decentralised mechanism (no coordinator/scheduler).
3. Only relevant to respective parts of the graph (i.e., sources cannot make use of sink progress)

Watermarks and Unstructured Loops

- **Low watermarking** is a very powerful mechanism to measure and propagate progress of a **single progress metric**.
- **However** it breaks when we have arbitrary cycles (current Flink iterations)

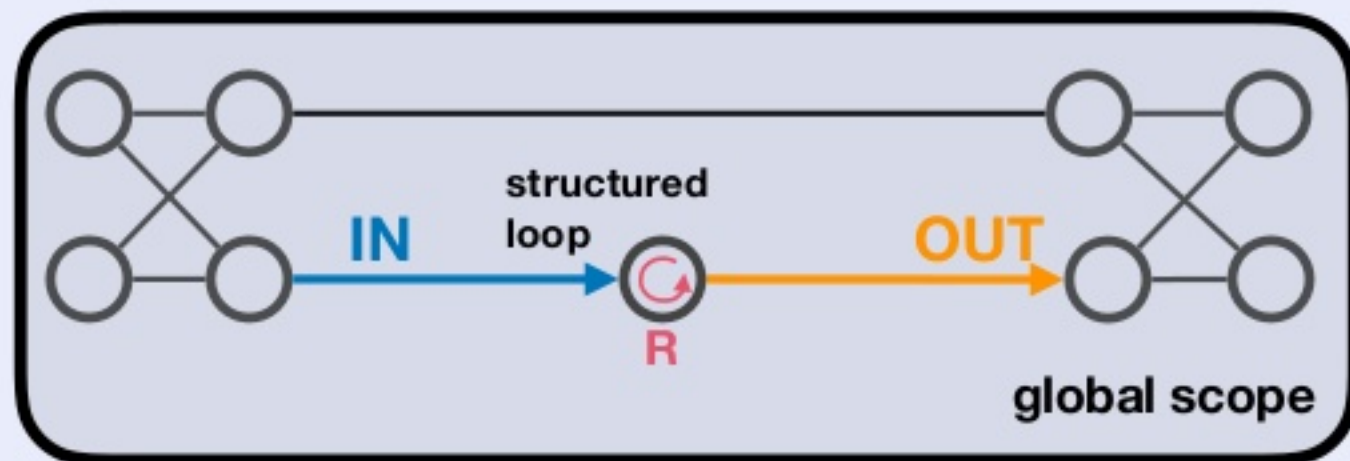


Proposed Extensions



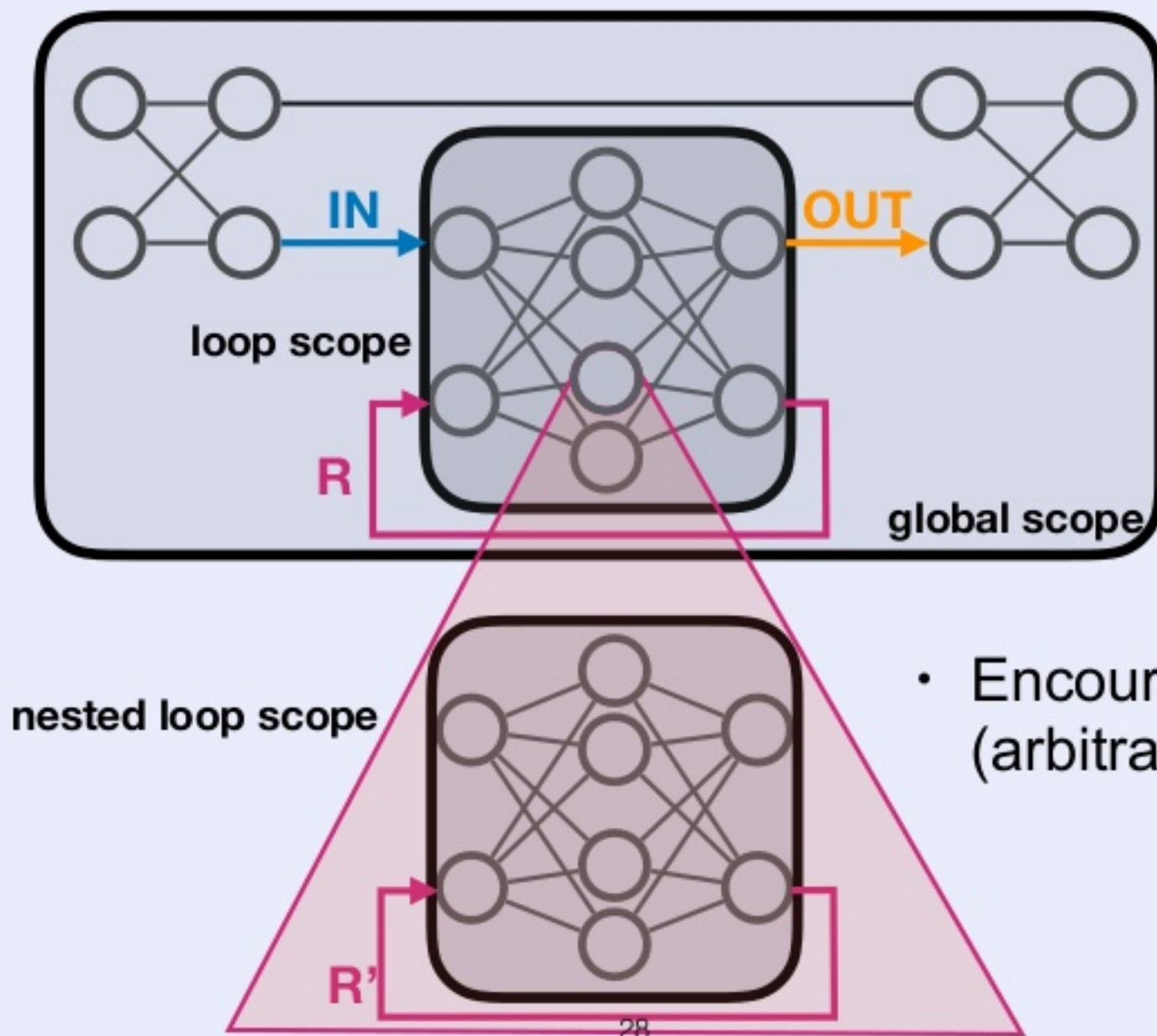
based on Flink 1.6

Structured Loops



- Structure loops are *low level* stream graph primitives.
- They can be seen as '**Iterative Operators**'
- Introduce the notion of “structured programming” for data streaming.
- Each **structured loop** has a **scope**.
- Each **scope** operates on its own **progress metric**.
- **Global scope** operates on **event-time progress** (no changes made).

Structured Loops Composition



- Encourages compositionally (arbitrary nesting)

Using Progress Timestamps

$$P = [\overbrace{p^n}^{p^T}, \underbrace{p^{n-1}, p^{n-2}, \dots, p^0}_{p^{ctx}}]$$

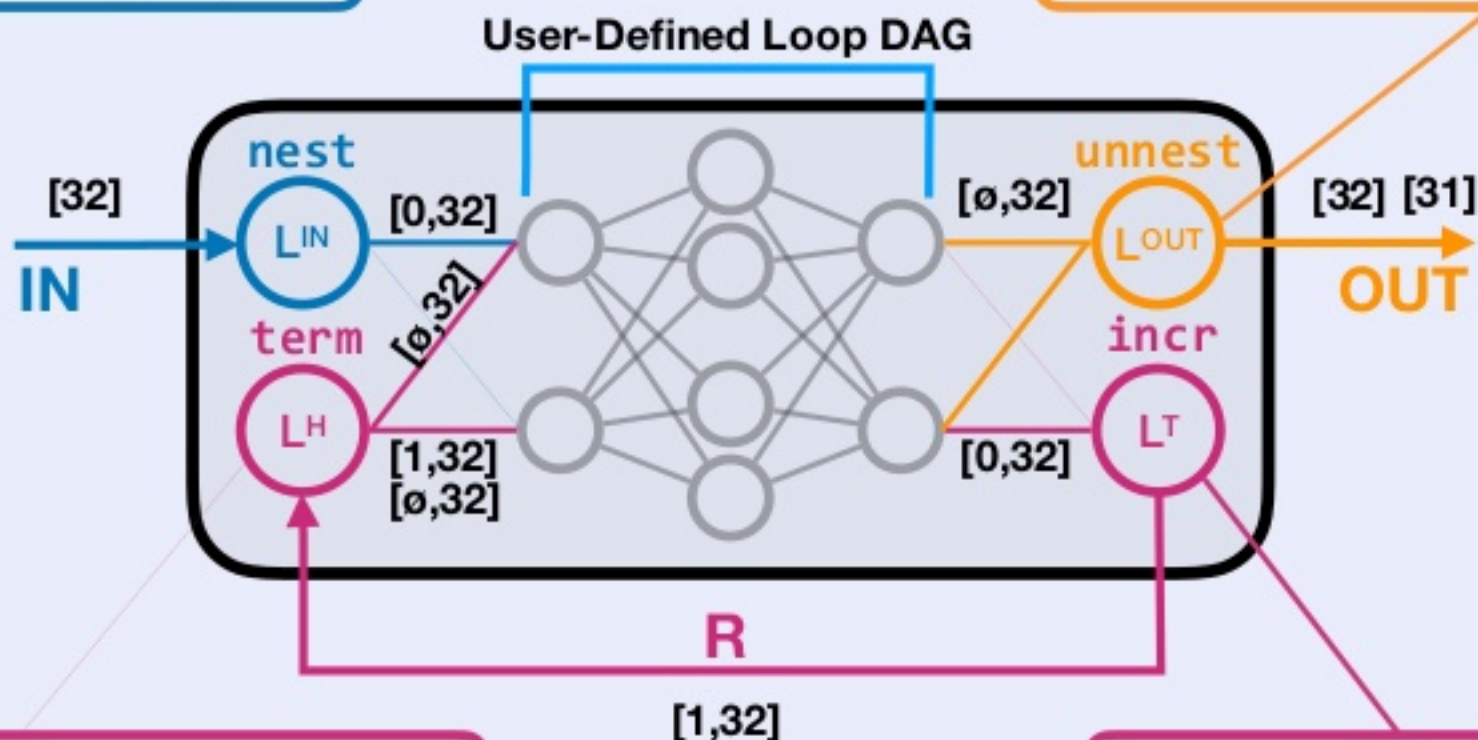
Operation	Implem.	Example
$progress: \mathbb{N}^n \rightarrow \mathbb{N}$	$head(P)$	p_n
$unnest: \mathbb{N}^n \rightarrow \mathbb{N}^{n-1}$	$tail(P)$	$[p_{n-1}, \dots, p_1]$
$nest: \mathbb{N}^n \rightarrow \mathbb{N}^{n+1}$	$0 :: P$	$[0, p_n, \dots, p_1]$
$incr: \mathbb{N}^n \rightarrow \mathbb{N}^n$	$head(P)+1$ $:: tail(P)$	$[p_n+1, \dots, p_1]$

- Nests **progress timestamps**
- (e.g., last superstep p^T per window P^{ctx})
- How do we guarantee monotonic progress metrics?

Structured Loop Embeddings

- Initializes Scope's Progress Metric for both records and watermarks (nest)

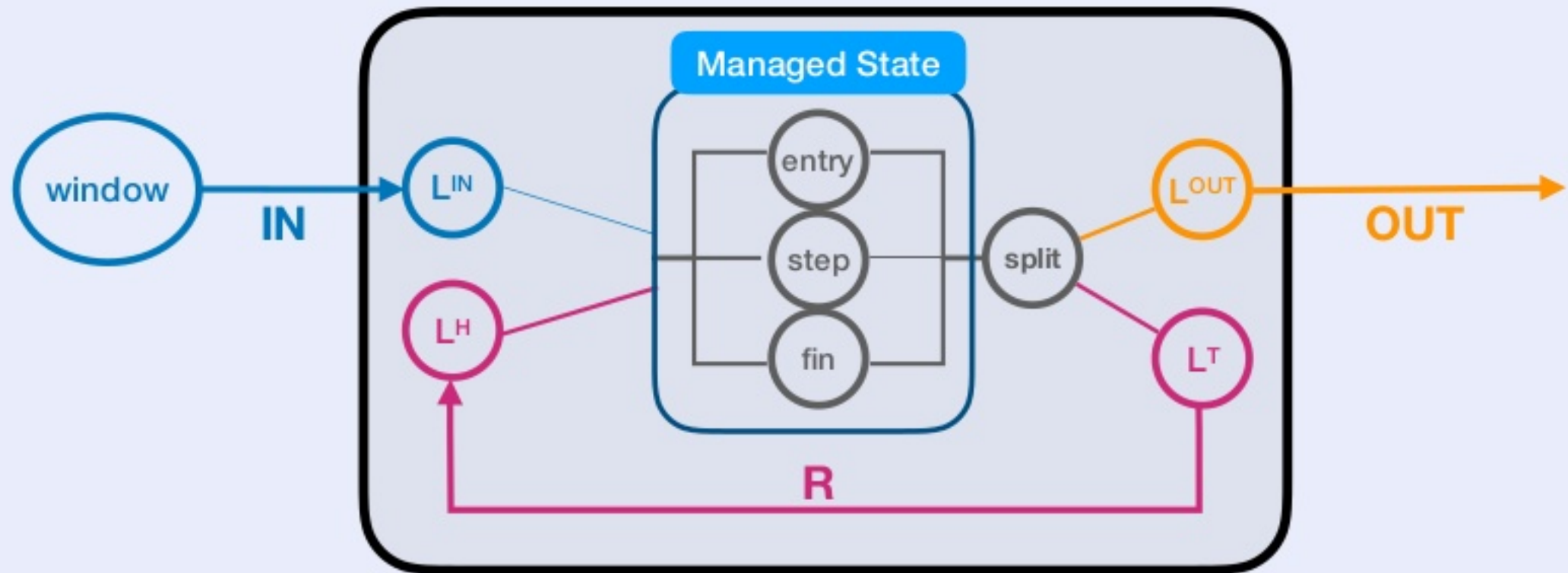
- Triggers final output of an iterative operator
- Removes inner progress metric (unnest)



- Forwards Records and Watermarks
- Evaluates and Disseminates Termination (\emptyset)

- Increments Progress Metric (incr)

Example: Window Iteration

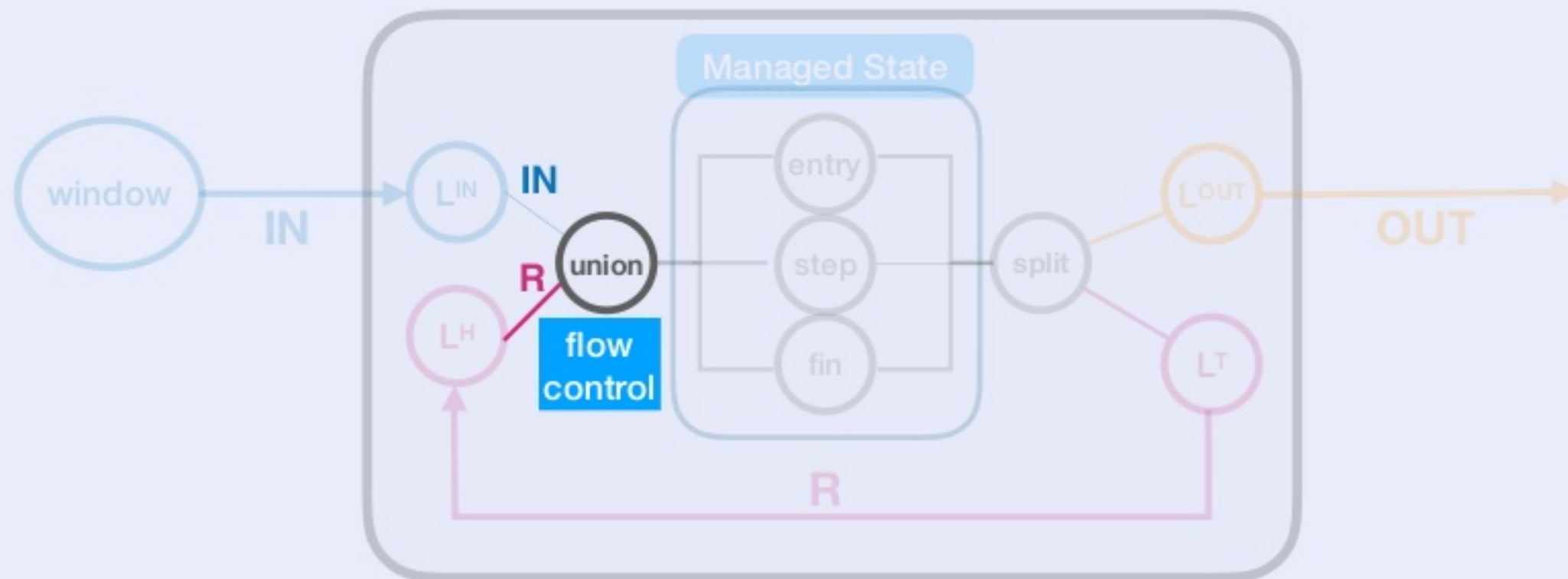


Further Challenges

what keeps us busy

- **Cyclic Flow Control**
 - avoid deadlocks
 - encourage iteration completion
- **State Management Integration**
 - one more level of indirection (mapvalue per key)

Cyclic Flow Control



- (Consumed Buffers < Threshold): Round-Robin between **IN** and **R**
- (Consumed Buffers \geq Threshold): Exclusive Ingestion of **R**

FLIP-15

- Introduces **Scoping** and **Nesting** (for Structured Loops)
- Addresses **Flow Control** Alternatives
- Discusses Computation-agnostic Loop Termination

<https://cwiki.apache.org/confluence/pages/viewpage.action?pageId=66853132>

Research Team and Sponsors

- **Marius Melzer (TUDresden)**
- **Asterios Katsifodimos (TUDelft)**
- **Vasiliki Kalavri (ETH)**
- **Seif Haridi (KTH)**
- **Pramod Bhatotia (Un.Edinburgh)**



STREAMLINE. CONSTRUCTIVE
SCIENCE



SWEDISH FOUNDATION for
STRATEGIC RESEARCH

References

- **Scalable and Reliable Data Stream Processing - Paris Carbone (2018)**
- Inflight Progress Tracking for Stream Processing Systems with Cyclic Dataflows - Marius Melzer (2017)
- Naiad: a timely dataflow system - Murray, McSherry et al. (2013)
- The dataflow model: a practical approach to balancing correctness, latency, and cost in massive-scale, unbounded, out-of-order data processing - Tyler Akidau et al. (2015)
- Out-of-order processing: a new architecture for high-performance stream systems - Li, Maier et al. (2008)
- Flexible time management in data stream systems - Srivastava, Widom (2004)

Stream Loops on Flink

Reinventing the wheel for the streaming era



Paris Carbone

Systems Researcher@KTH/SICS <parisc@kth.se>
Committee@Flink <senorcarbone@apache.org>