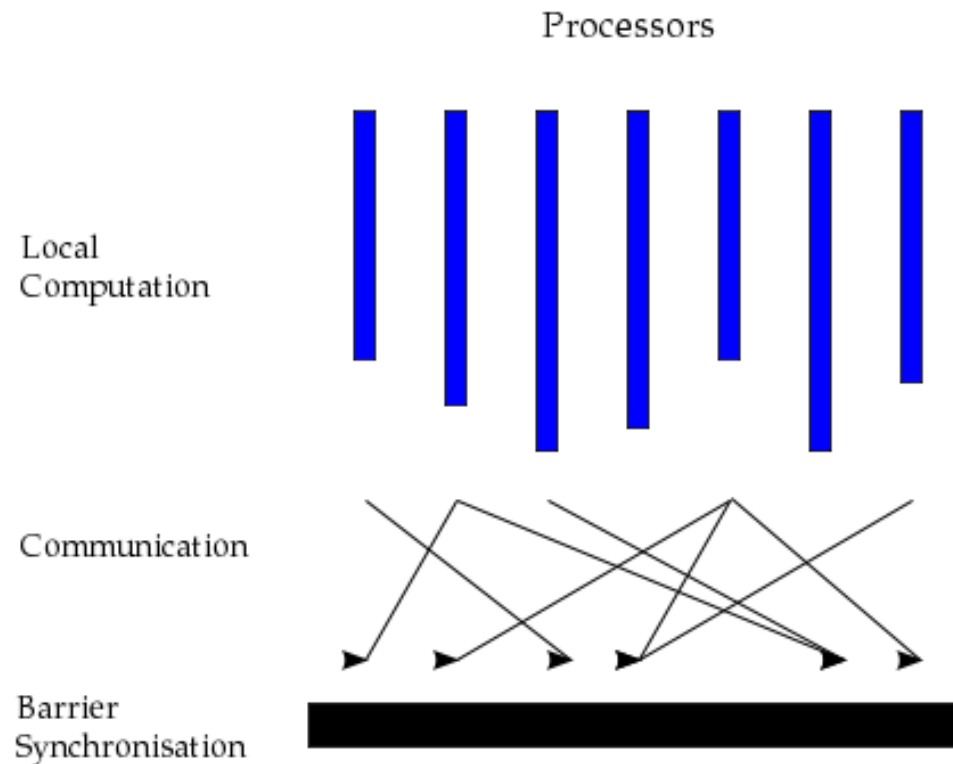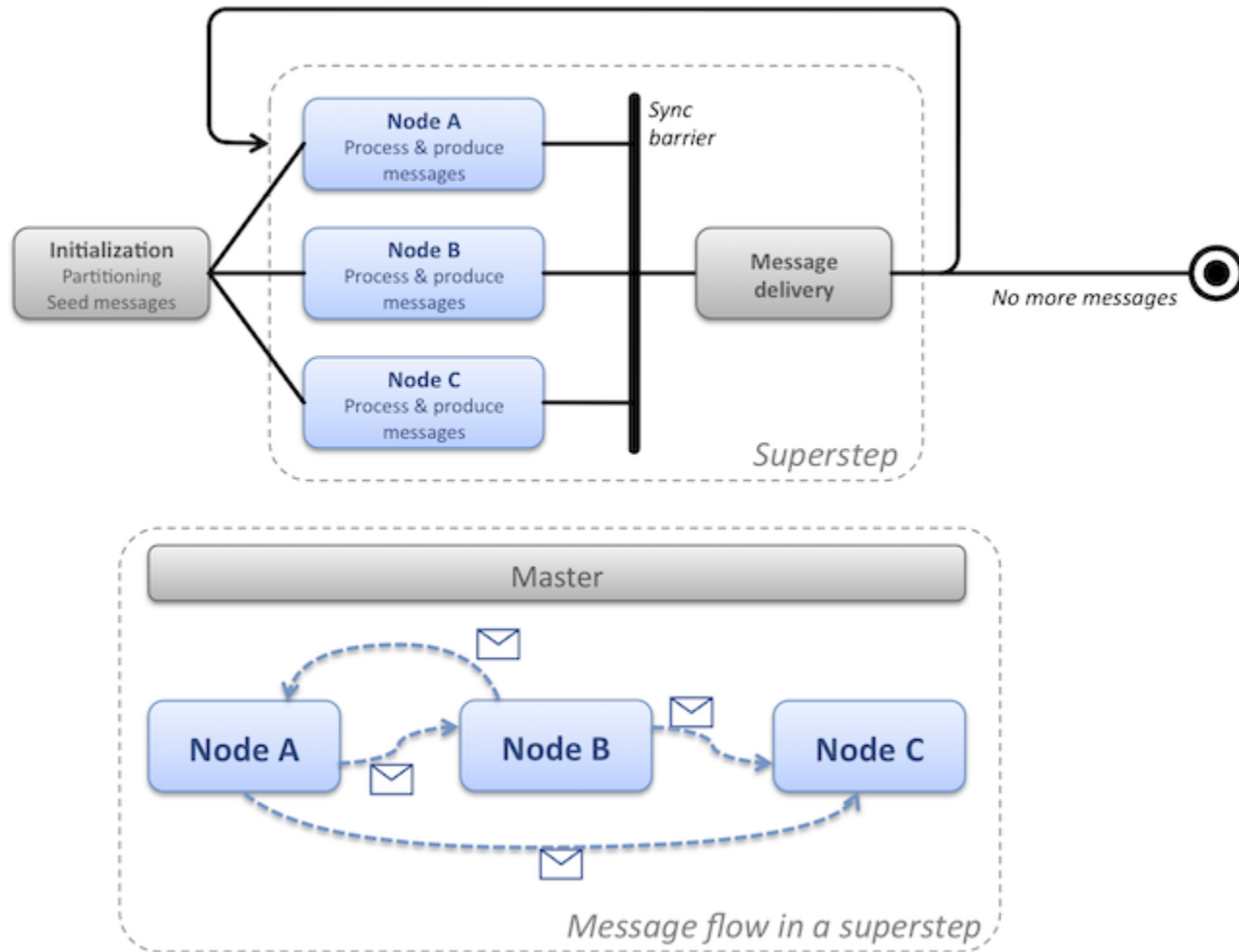# Whatever is left

Alexander Lazovik

a.lazovik@rug.nl
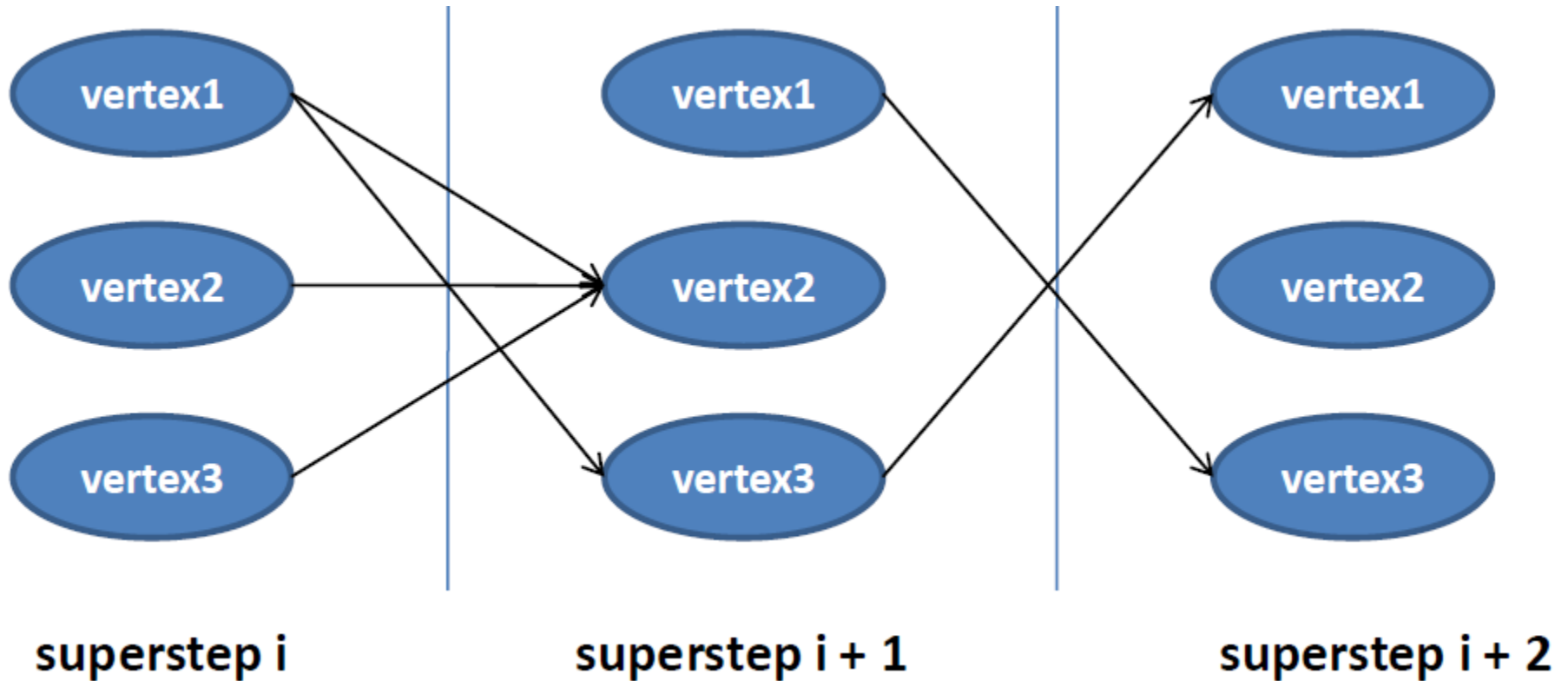
# Bulk Synchronous Parallel (1990)

# BSP

# BSP on distributed graphs



superstep i          superstep i + 1          superstep i + 2

# Map/Reduce implementation

▸ *https://highlyscalable.wordpress.com/2012/02/01/mapreduce-patterns/*
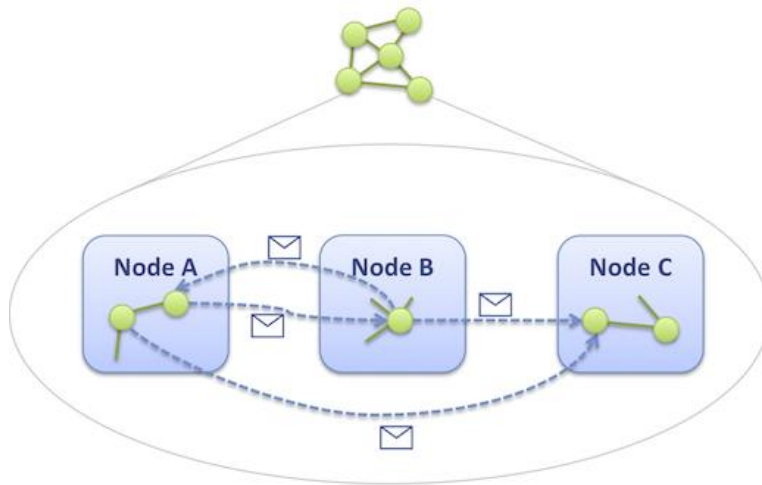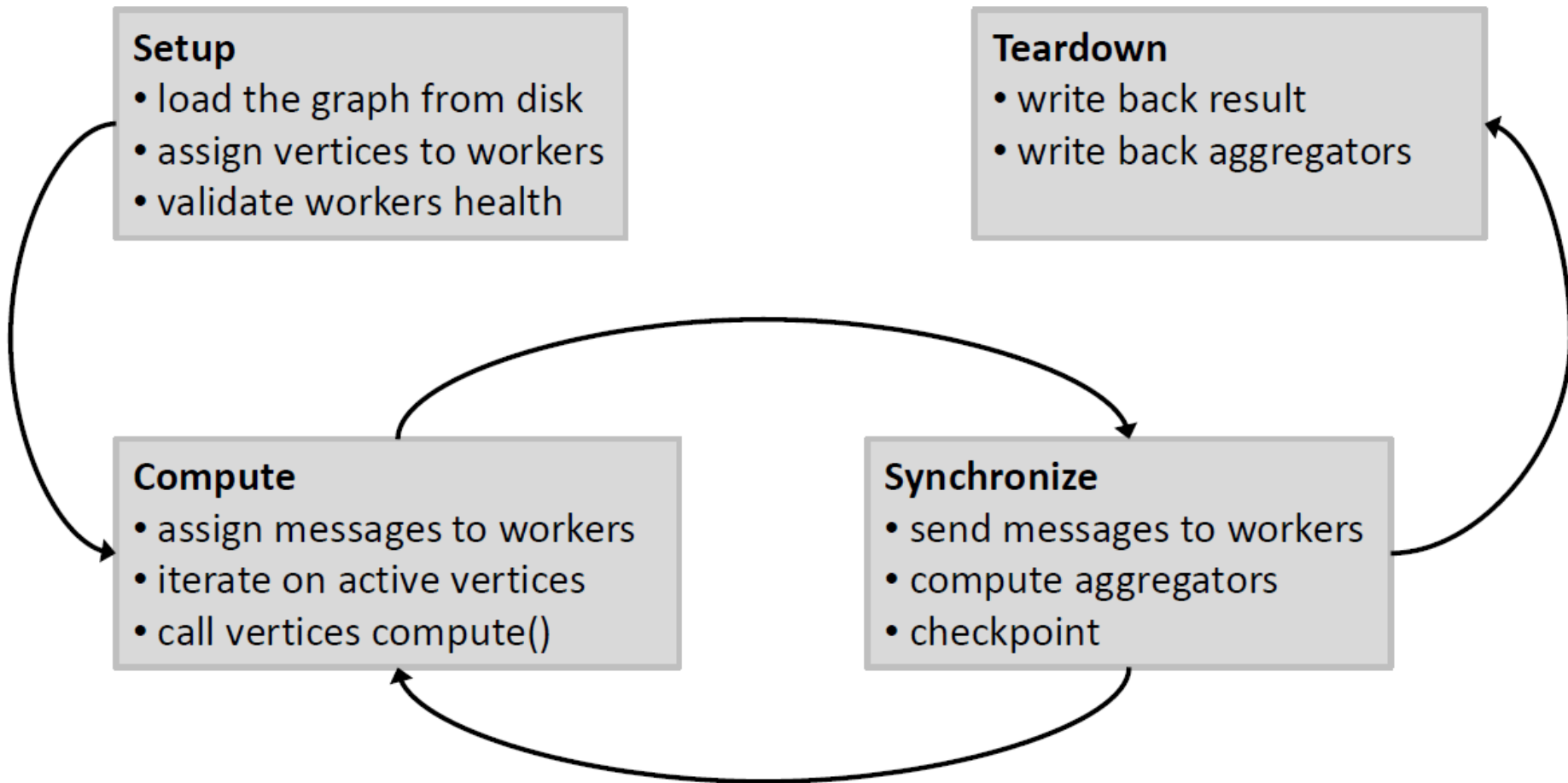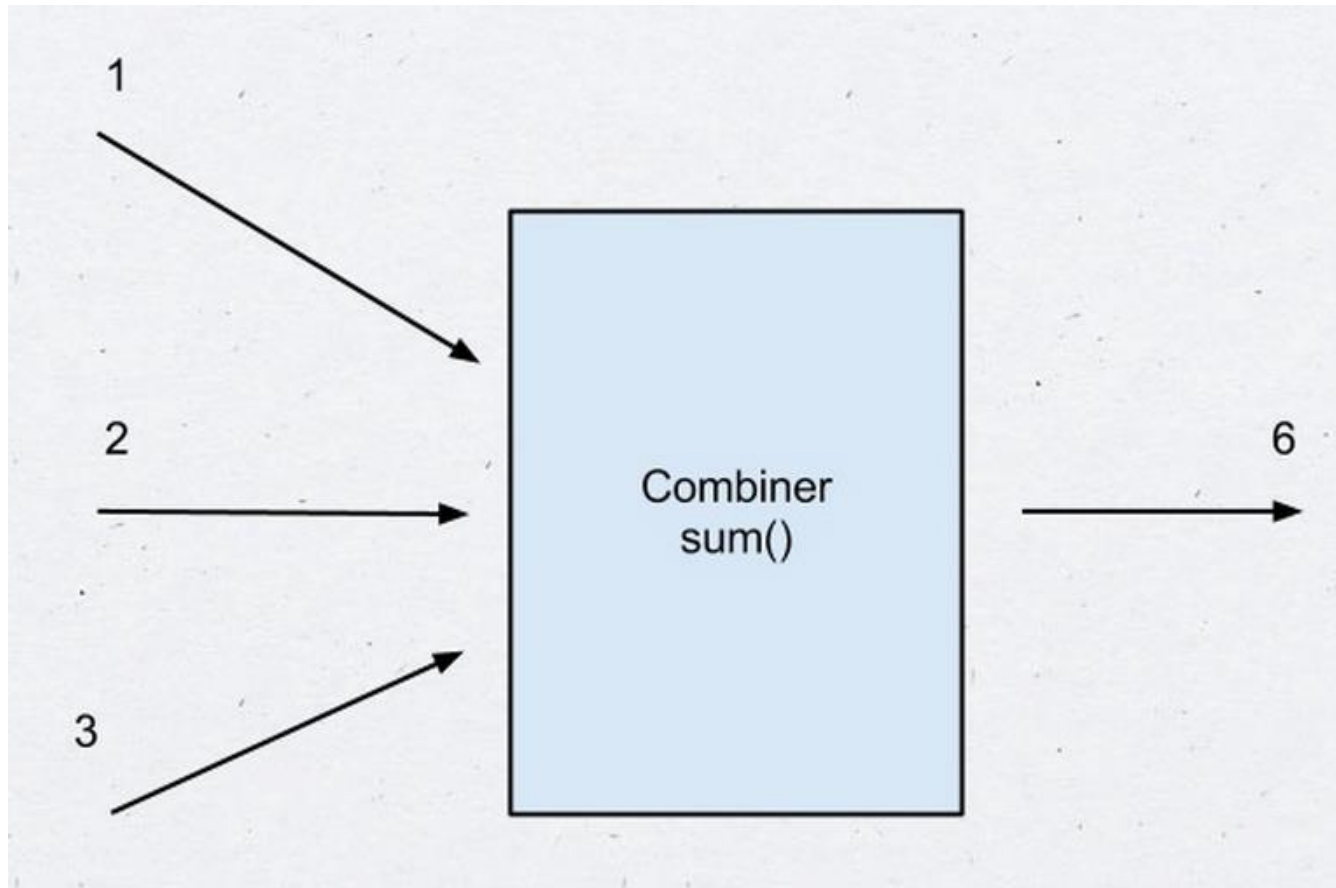
# Pregel: BSP applied to distributed graphs

▸ Apache Giraph: an open source implementation of Pregel

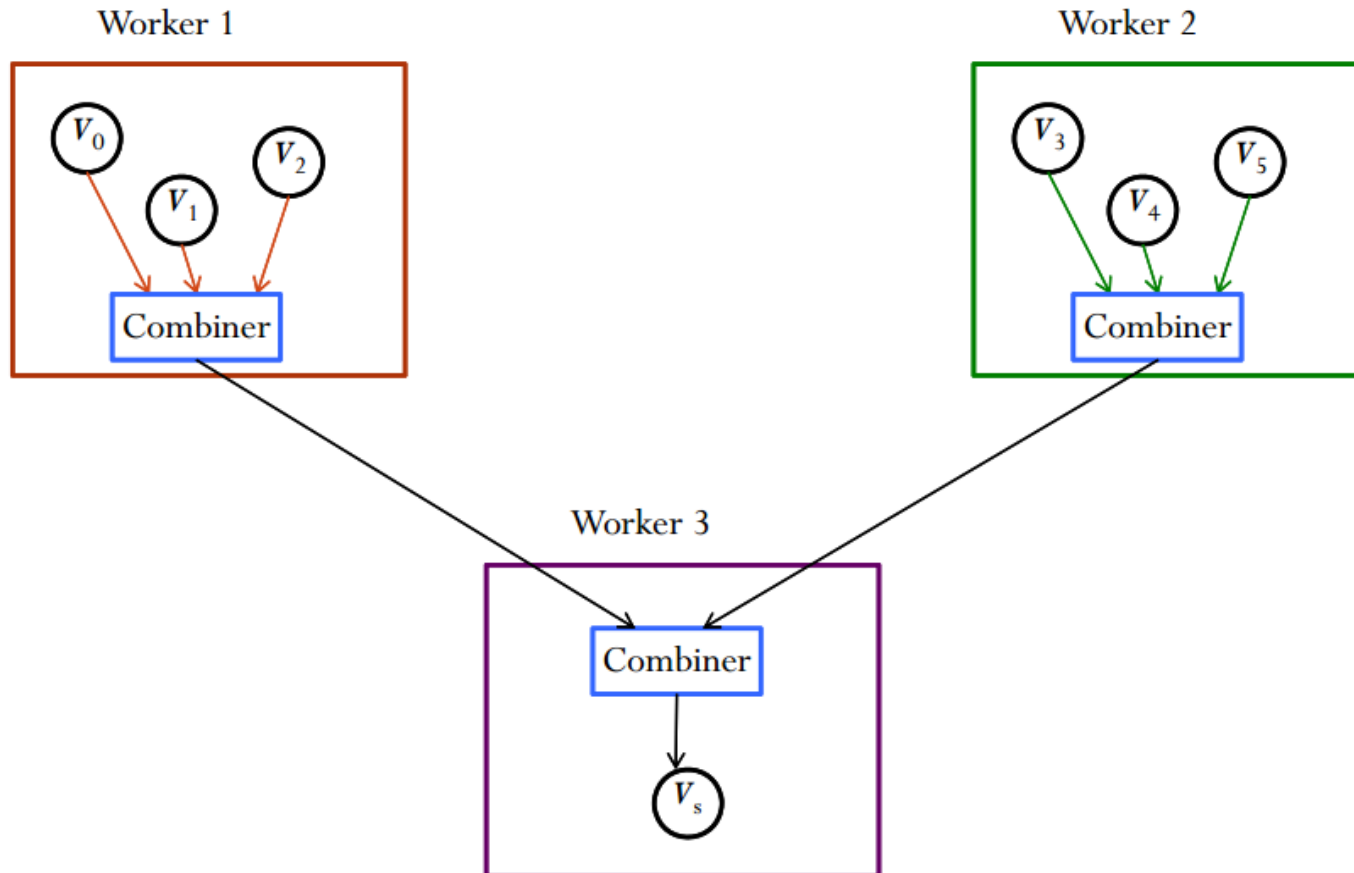- ▸ graph is distributed across several machines
- ▸ computation "node" is a vertex

# Apache Giraph

**Setup**
- load the graph from disk
- assign vertices to workers
- validate workers health

**Teardown**
- write back result
- write back aggregators

**Compute**
- assign messages to workers
- iterate on active vertices
- call vertices compute()

**Synchronize**
- send messages to workers
- compute aggregators
- checkpoint

# Combiners (user-defined)

# Combiners



main goal: reduce network bandwidth / number of supersteps
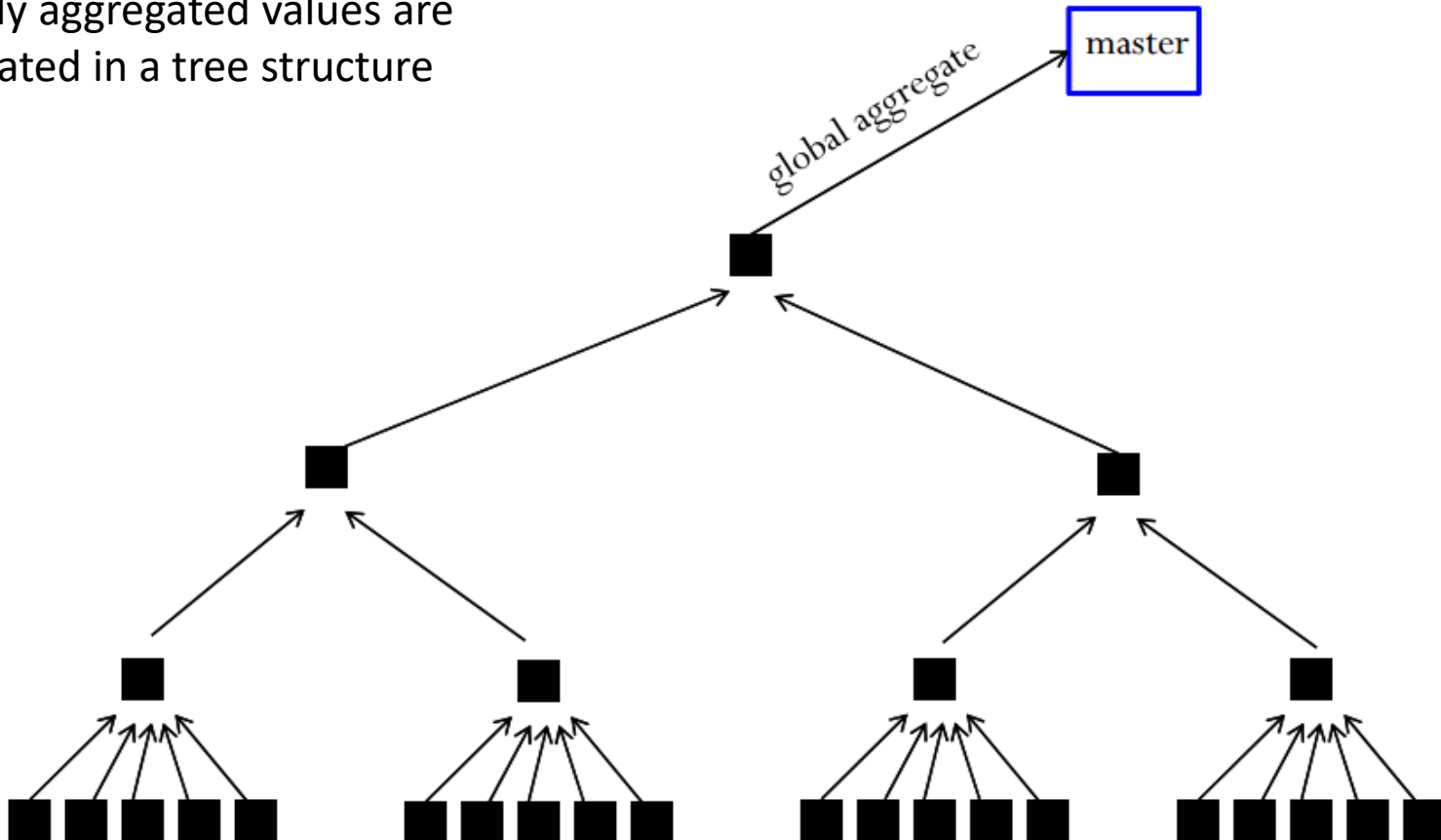
# Aggregators
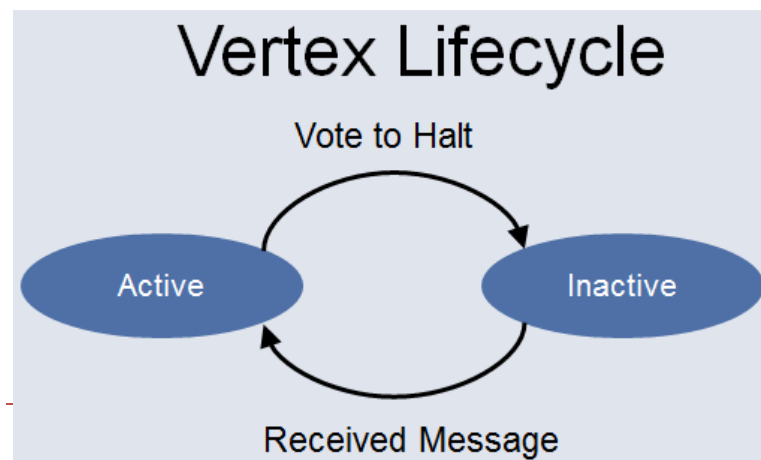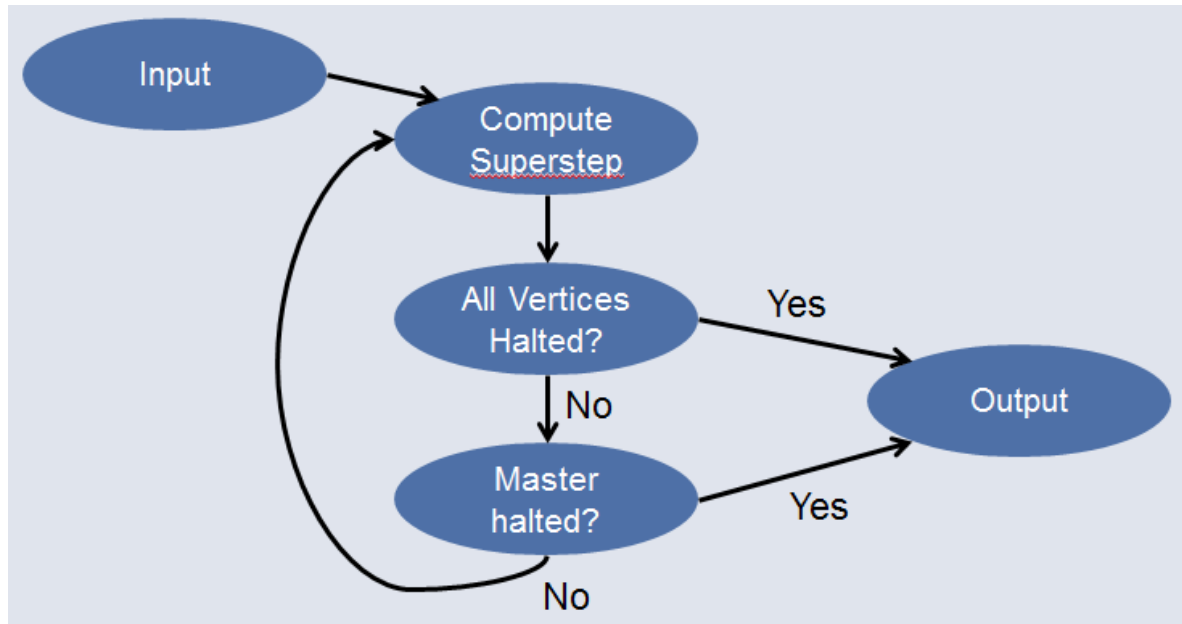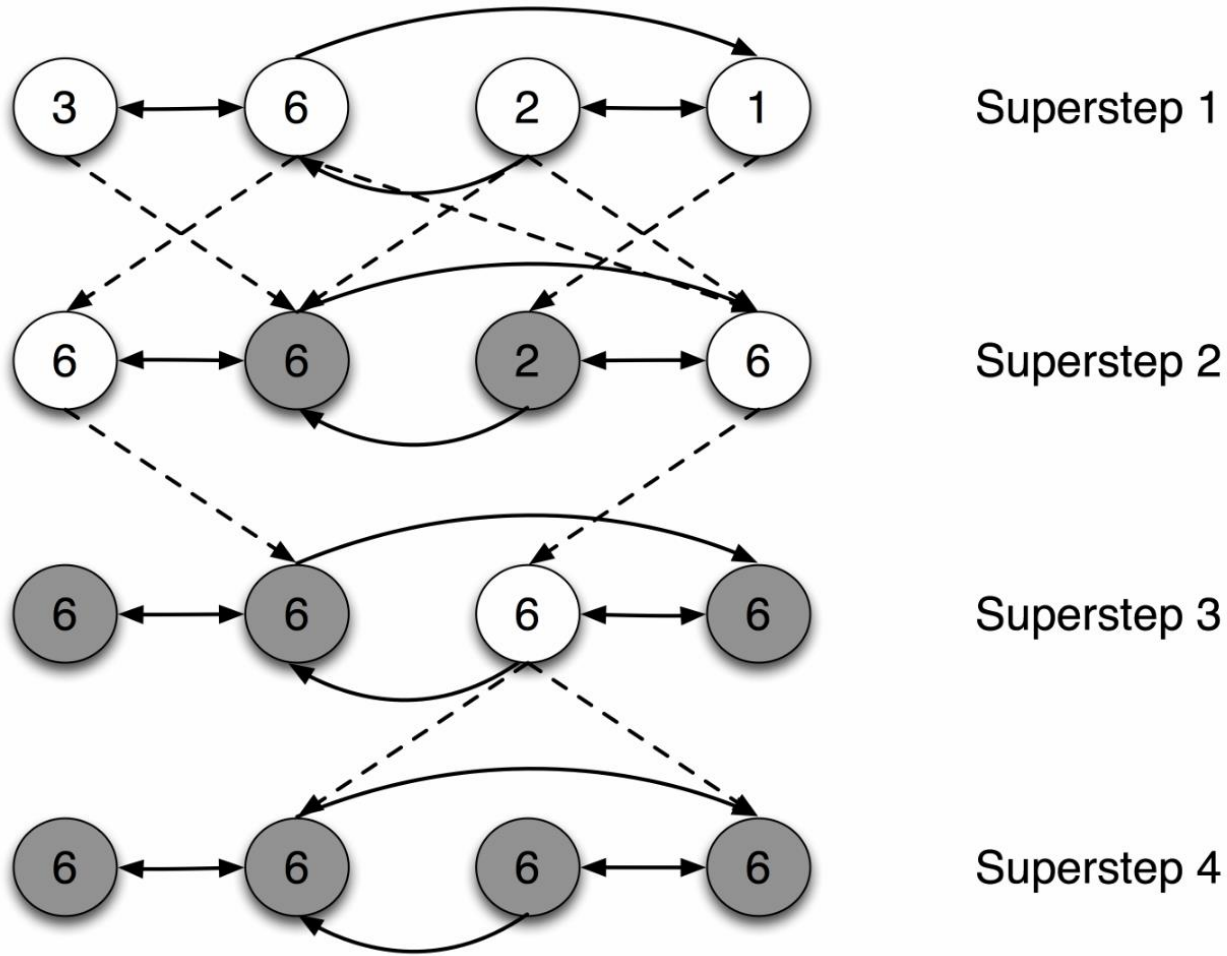
# Aggregators

- each worker aggregates values from its vertices
- partially aggregated values are
  aggregated in a tree structure

# Apache Giraph Lifecycle



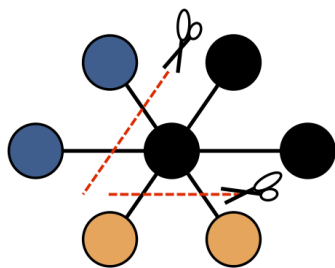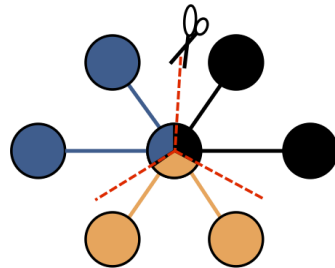Vertex Lifecycle

# Finding max
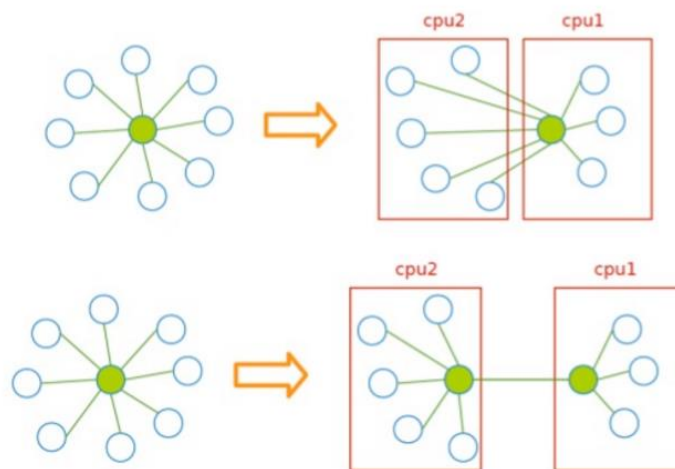
# Spark GraphX: Resilient Distributed Graphs

Edge Cut

Vertex Cut

vertex vs. edge cutting

# Spark GraphX: Resilient Distributed Graphs



```scala
// Vertex collection
class VertexRDD[VD] extends RDD[(VertexId, VD)]

// Edge collection
class EdgeRDD[ED] extends RDD[Edge[ED]]
case class Edge[ED](srcId: VertexId = 0, dstId: VertexId = 0,
                    attr: ED = null.asInstanceOf[ED])

// Edge Triple
class EdgeTriplet[VD, ED] extends Edge[ED]
```
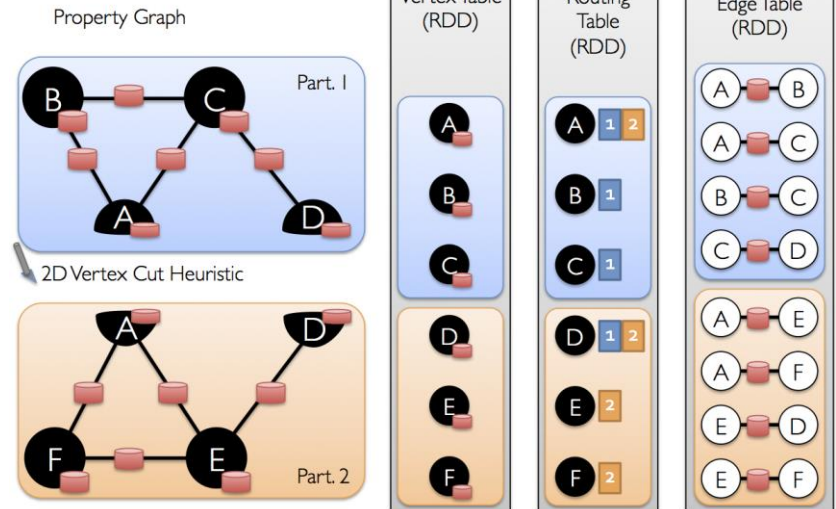
```scala
// VD: the type of the vertex attribute
// ED: the type of the edge attribute
class Graph[VD, ED] {
  val vertices: VertexRDD[VD]
  val edges: EdgeRDD[ED]
}
```

```scala
class Graph[VD, ED] {
  def mapVertices[VD2](map: (VertexId, VD) => VD2): Graph[VD2, ED]

  def mapEdges[ED2](map: Edge[ED] => ED2): Graph[VD, ED2]

  def mapTriplets[ED2](map: EdgeTriplet[VD, ED] => ED2): Graph[VD, ED2]
}
```

# Streaming (unbounded)

▸ You cannot get answers to some questions
  ▸ e.g., what is the average of all elements?
  ▸ you can answer them for a given subset though
    ▸ last 20 elements, last 20 seconds
  ▸ effectively, generating a new stream
    ▸ not necessarily after each 20 elements

▸ Two approaches:
  ▸ Microbatching (combine elements and then use existing tools)
    ▸ e.g., Spark
    ▸ additional buffering may help better distribute load
      ☐ do not have a convincing example (though, you can think of some)
  ▸ Process-each-element
    ▸ e.g., Twitter Storm, Apache Flink
    ▸ response time/availability of the first processed element is better

# Streaming frameworks

### Apache Storm
- True streaming, low latency - lower throughput
- Low level API (Bolts, Spouts) + Trident

### Spark Streaming
- Stream processing on top of batch system, high throughput -  higher latency
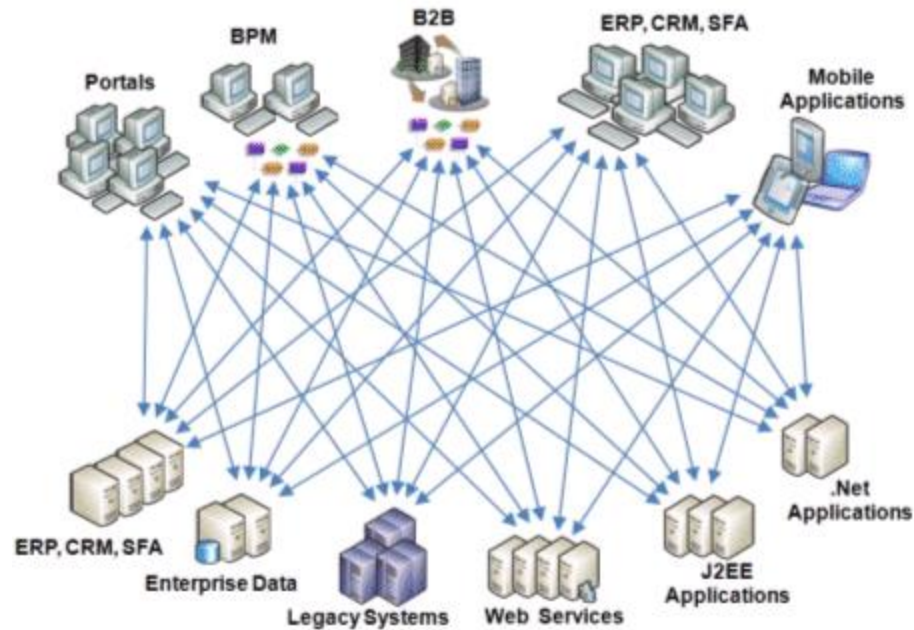- Functional API (DStreams), restricted by batch runtime

### Apache Samza
- True streaming built on top of Apache Kafka, state is first class citizen
- Slightly different stream notion, low level API
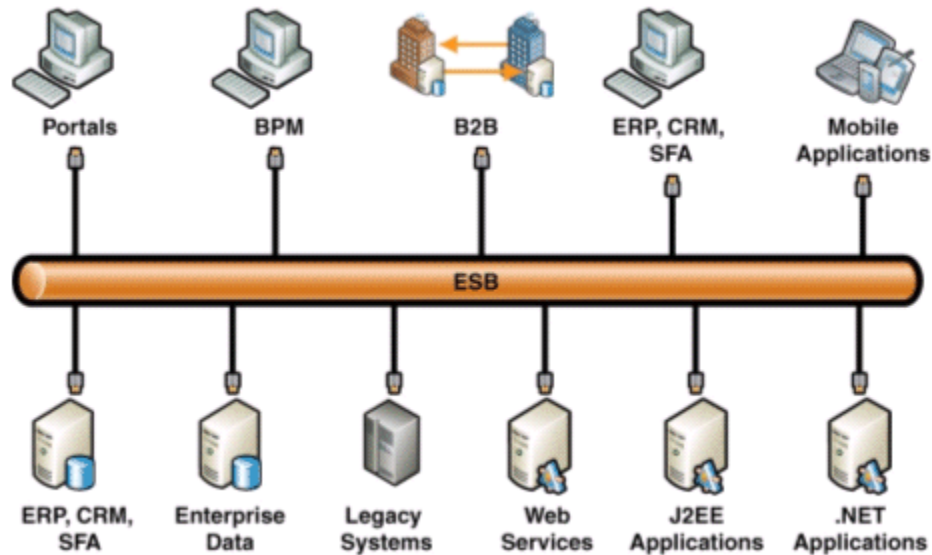
### Apache Flink
- True streaming with adjustable latency-throughput trade-off
- Rich functional API exploiting streaming runtime; e.g. rich windowing semantics

# Enterprise Integration

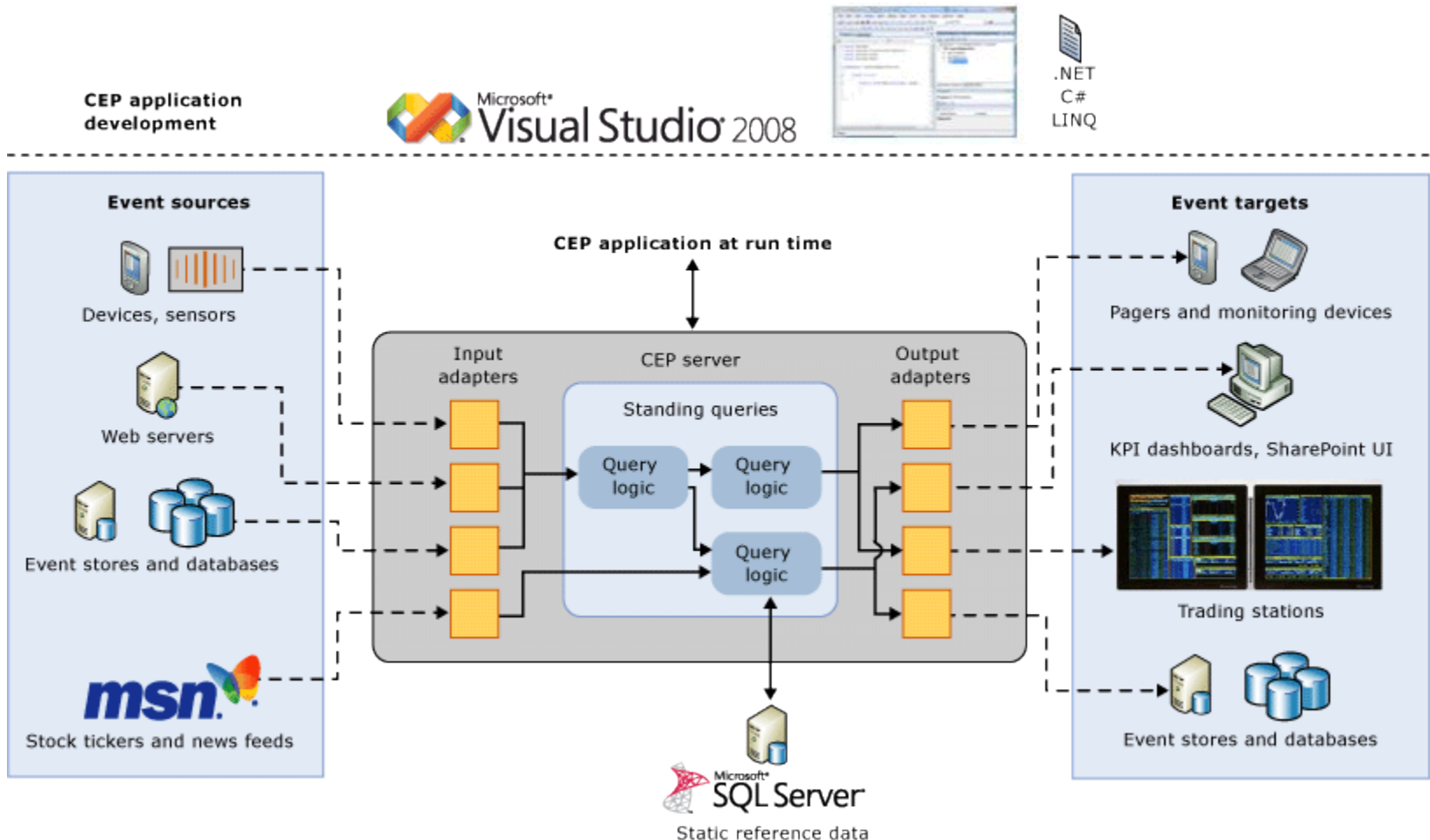

point-to-point communication
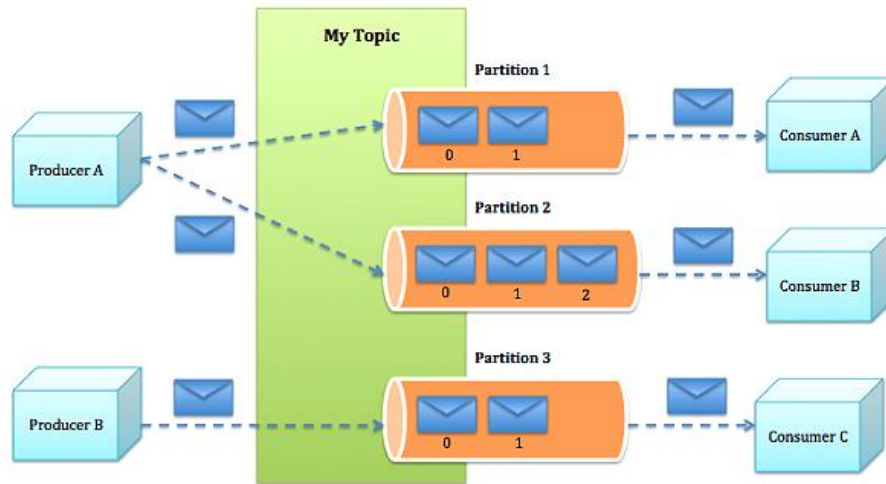
# Enterprise Integration



enterprise service bus

*consider integration across organizations*
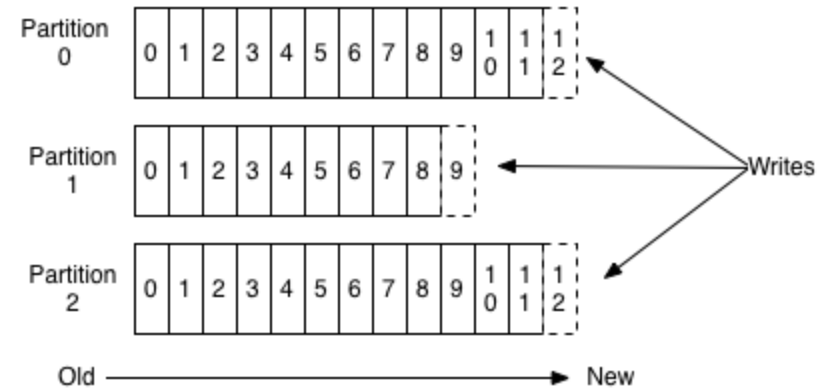
# Complex Event Processing (CEP)

# Kafka
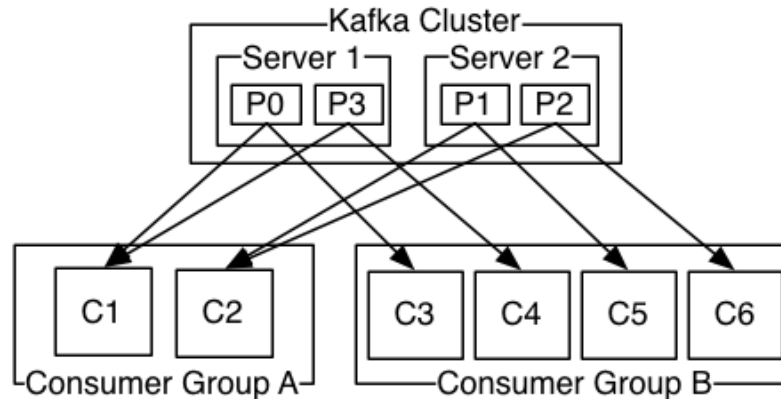

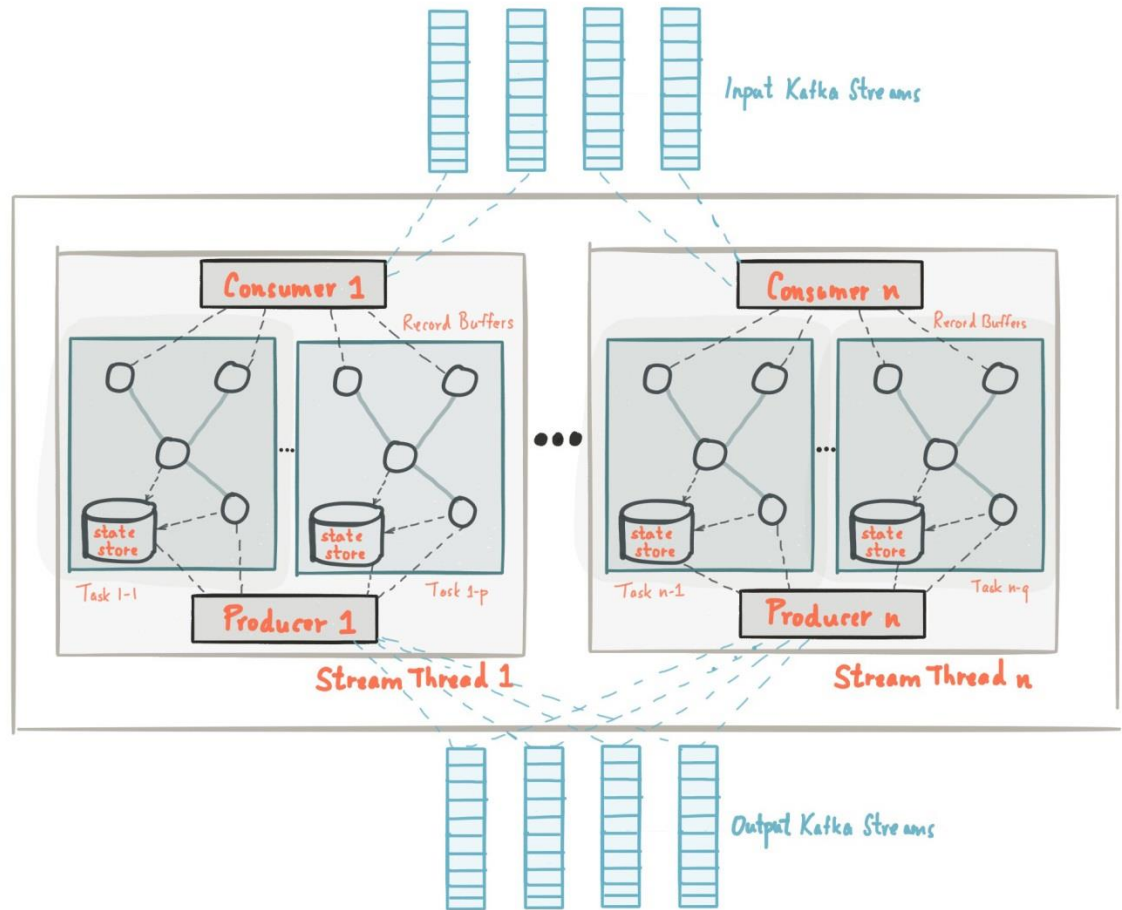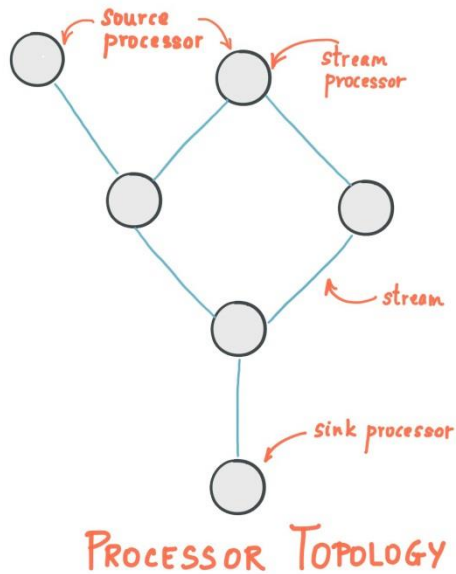
- partitions are immutable sequences
    - and stored (on disk) for a configurable retention period
- producer is responsible to chose topic/partition
- within a partition, order is preserved
- partitions are replicated over several servers (one acting as a leader)
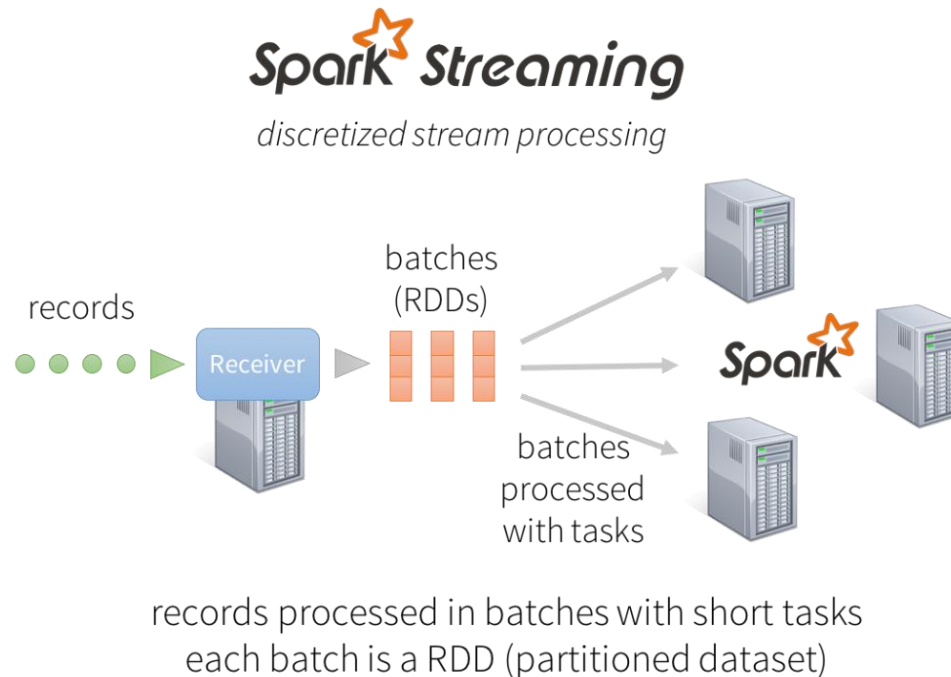
# Kafka



- each record is delivered to one consumer from each consumer group
    - one consumer group => load balancing
    - consumer group per consumer => broadcasting
- each consumer is assigned a partition exclusively
    - one partition => max one consumer per consumer group

# Kafka Streaming API



map, filter, reduceByKey, …

# Spark Microbatching approach



Spark Streaming
*discretized stream processing*

records → Receiver → batches (RDDs) → Spark

batches processed with tasks

records processed in batches with short tasks
each batch is a RDD (partitioned dataset)

▸ DStream is a stream of RDDs

  ▸ each element (RDD) representing one batch

# Windows and Sliding

Sensor → , 9, 6, 8, 4, 7, 3, 8, 4, 2, 1, 3, 2, ⟹ out

Sensor → , 9, 6, 8, 4, 7, 3, 8, 4, 2, 1, 3, 2, ⟹

rolling sum → , 57, 48, 42, 34, 30, 23, 20, 12, 8, 6, 5, 2, → out

# Windows and Sliding

# Tumbling vs hopping vs sliding windows

▸ Tumbling

  ▸ non-overlapping windows

▸ Hopping

  ▸ overlapping windows with a fixed jump-ahead period

▸ Sliding

  ▸ overlapping "continuous" window
  ▸ all possible windows of a given size

# Reactive Streams

▶ (Akka) Reactive Streams

  ▶ good for streaming-like APIs

```
public interface Subscriber<T> {
    public void onSubscribe(Subscription s);
    public void onNext(T t);
    public void onError(Throwable t);
    public void onComplete();
}
```

▶ Goal:

let the consumer explicitly tell the producer

"I am ready to accept new data"
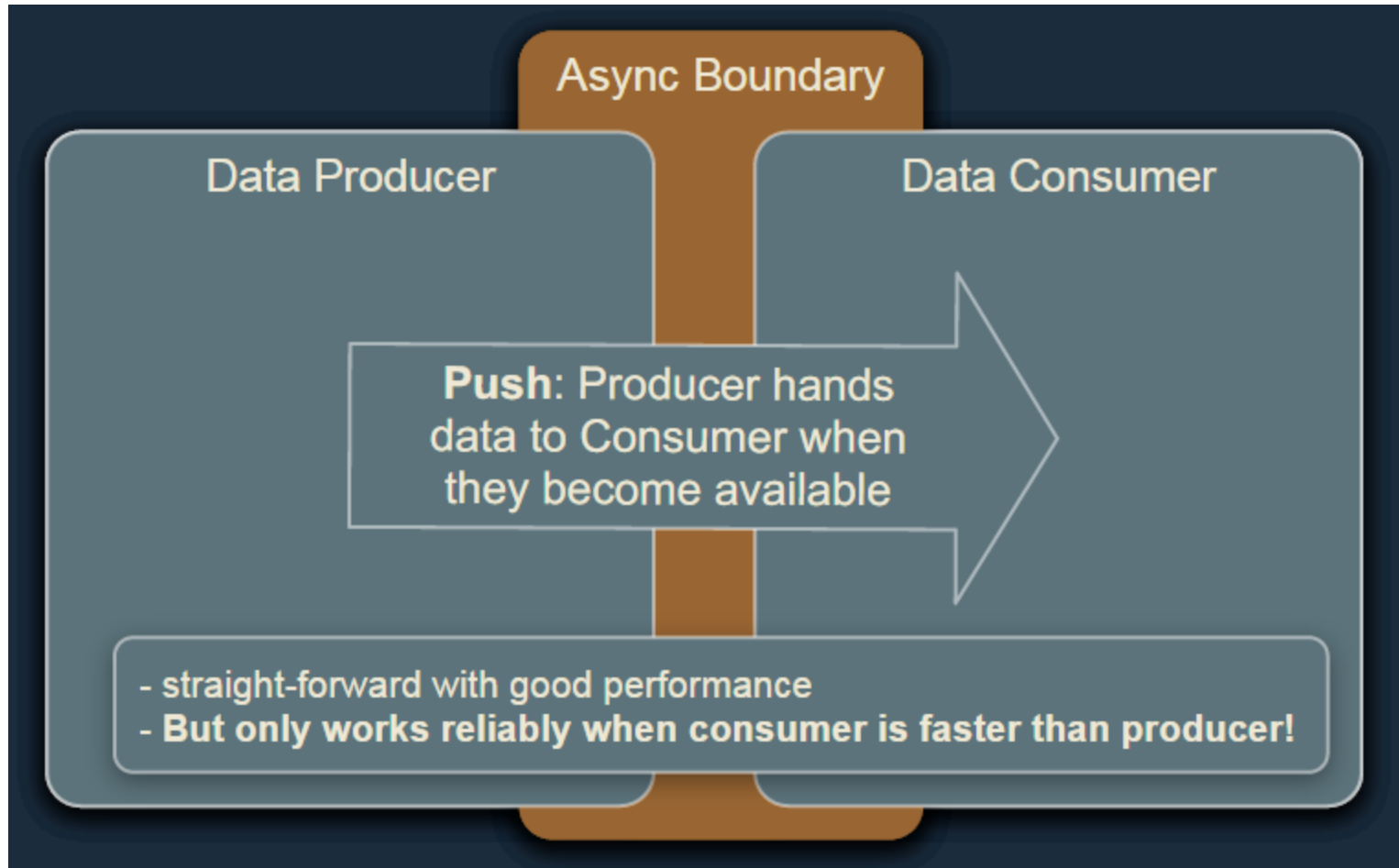
# Slides from Mathias Doenitz (spray.io, akka-http)
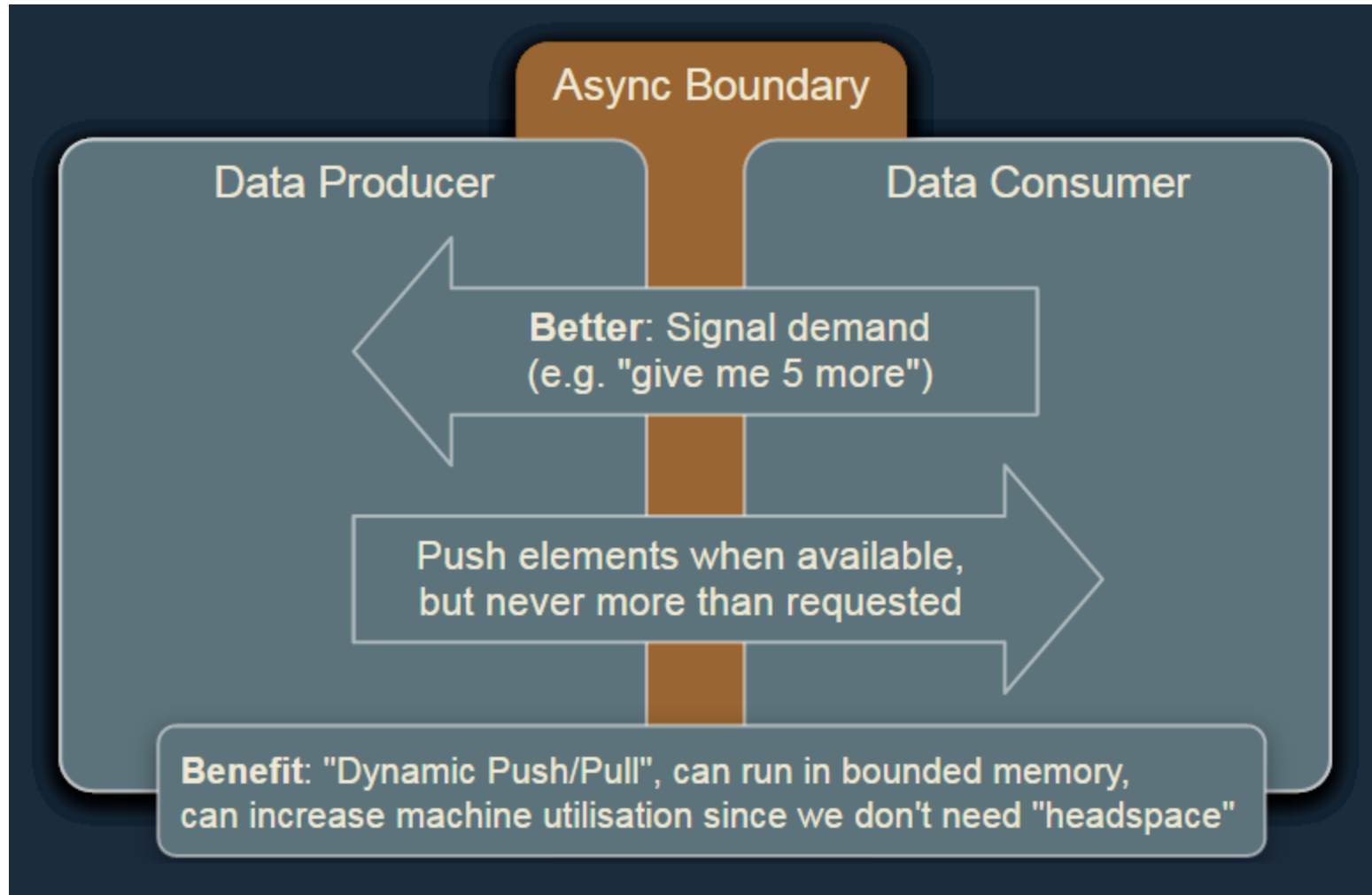
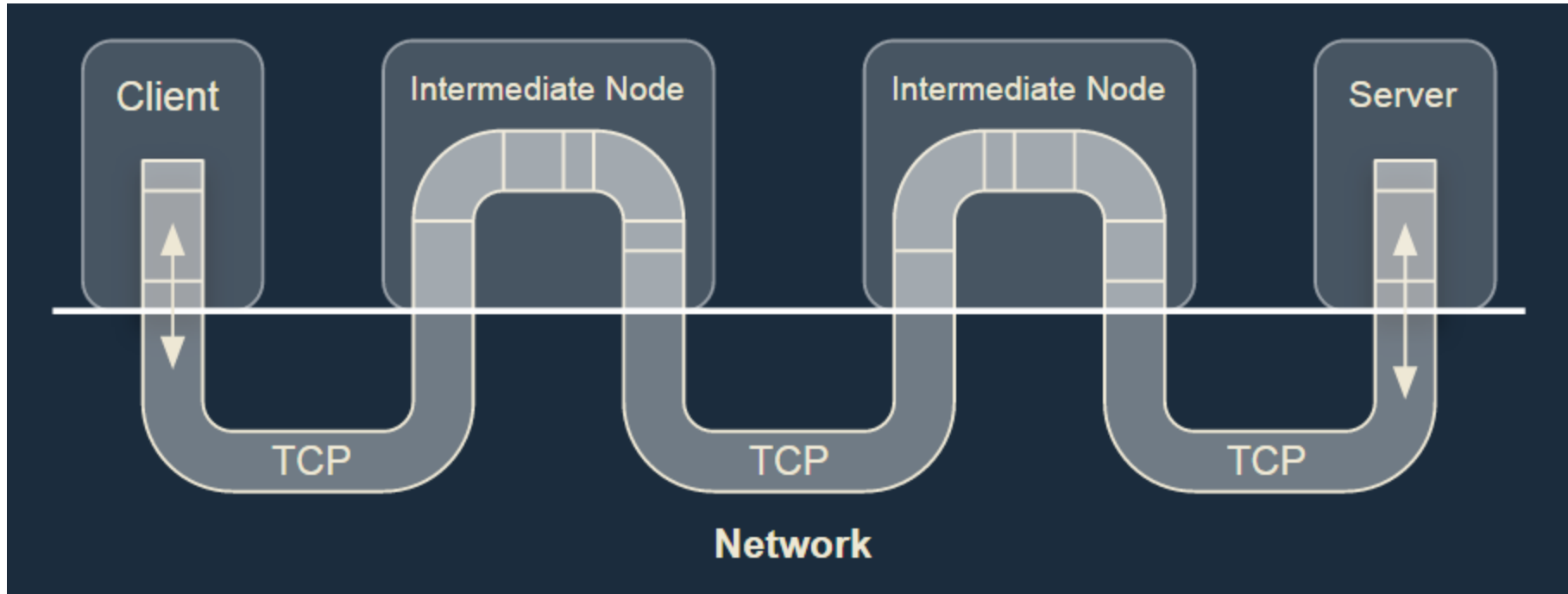# Slides from Mathias Doenitz (spray.io, akka-http)

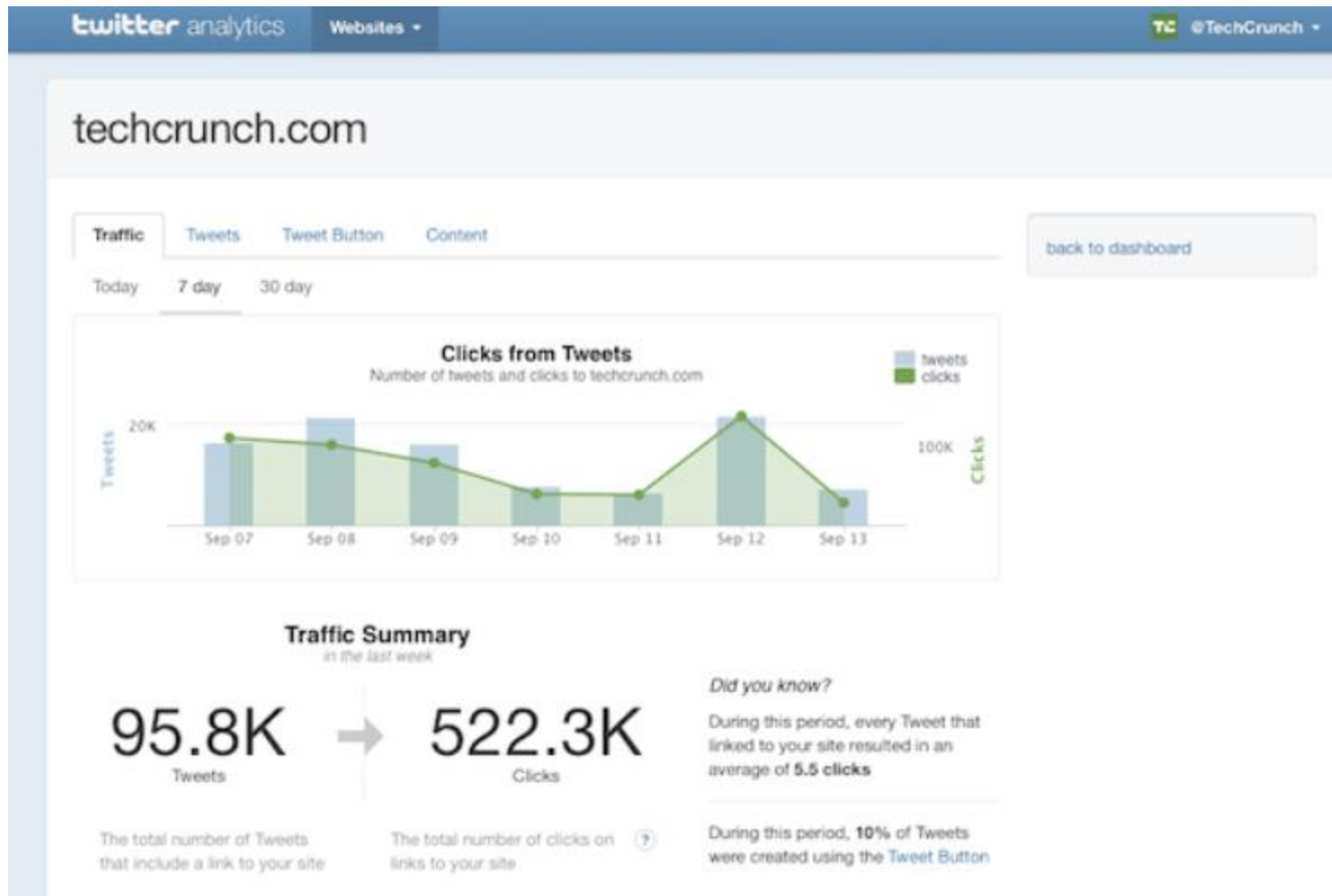# Slides from Mathias Doenitz
# (spray.io, akka-http)

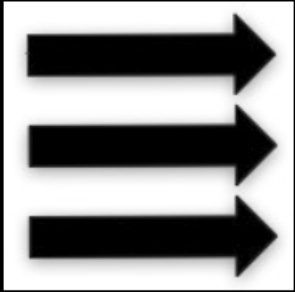# Slides from Mathias Doenitz (spray.io, akka-http)



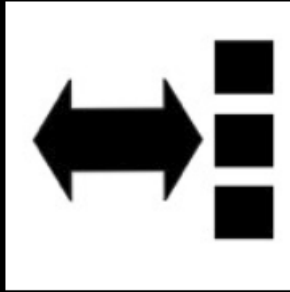full resource utilization if used across the whole pipeline!
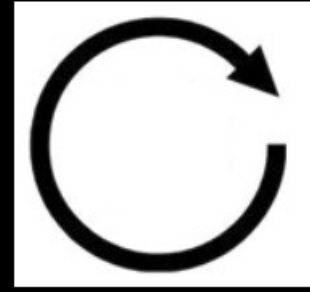
# Storm @ Twitter
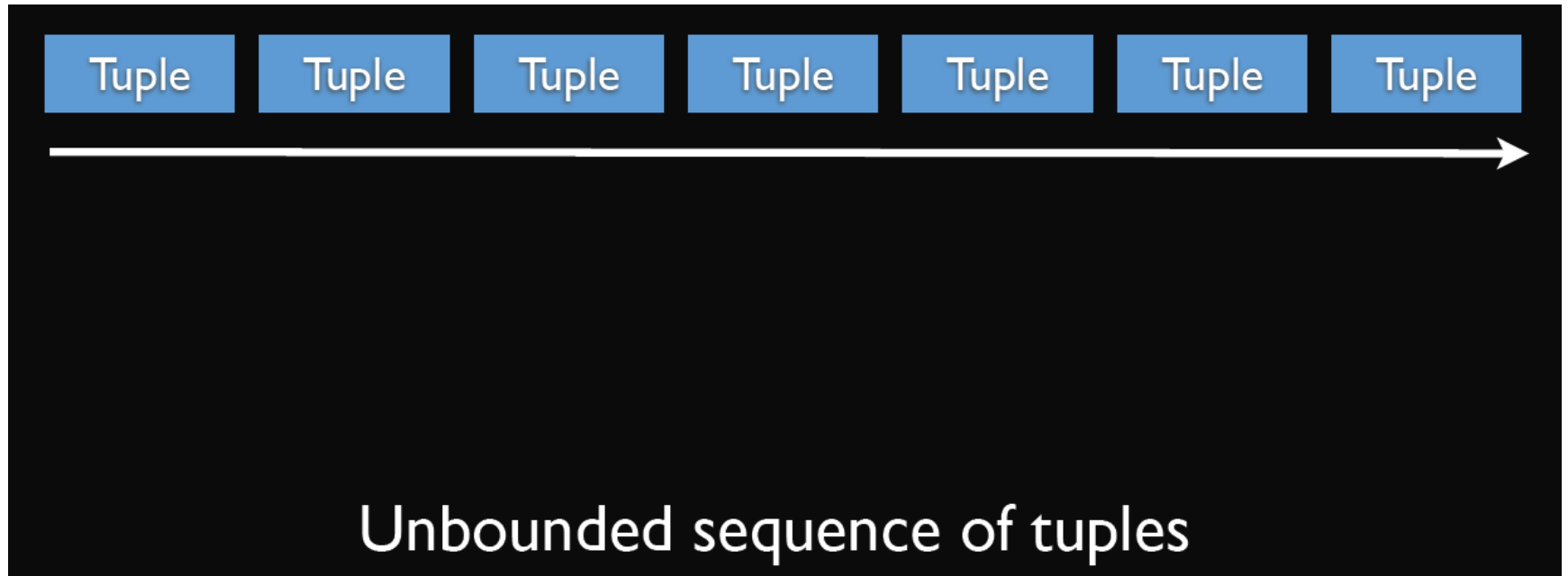
# Use cases



Stream processing

Distributed RPC

Continuous computation

# Concepts

- Streams
- Spouts
- Bolts
- Topologies

# Streams



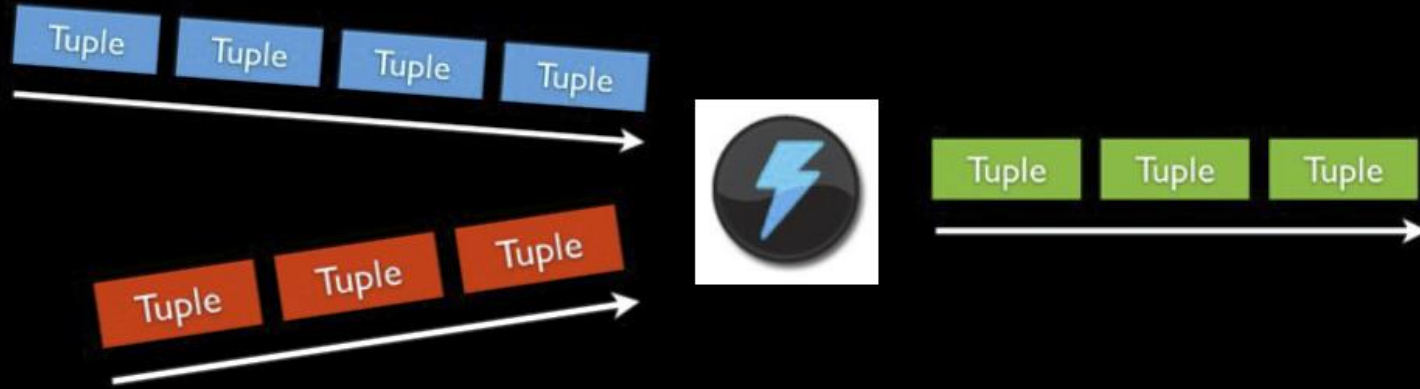Unbounded sequence of tuples

# Spouts
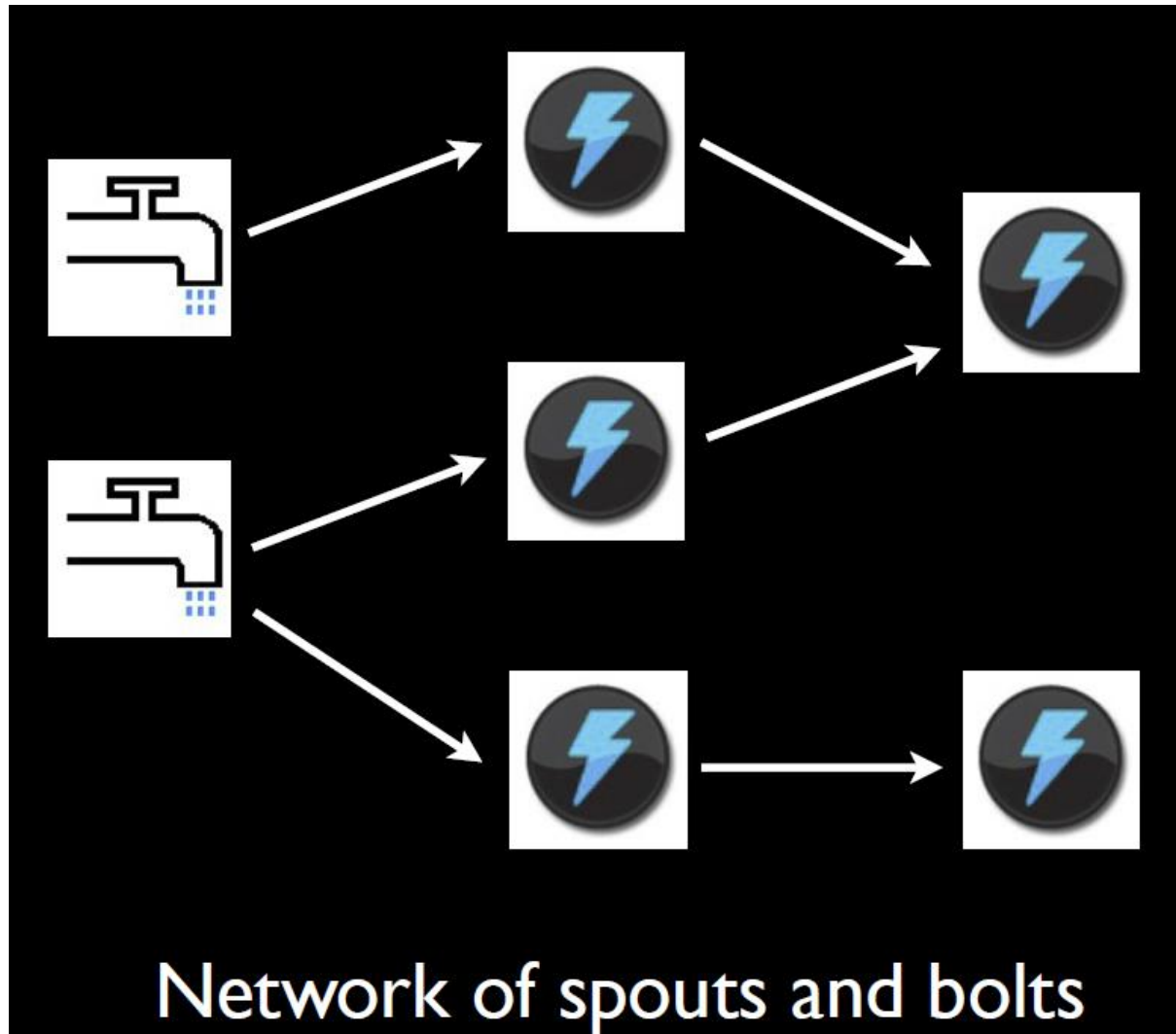


Source of streams

# Bolts



Processes input streams and produces new streams

# Bolts



- Functions
- Filters
- Aggregation
- Joins
- Talk to databases

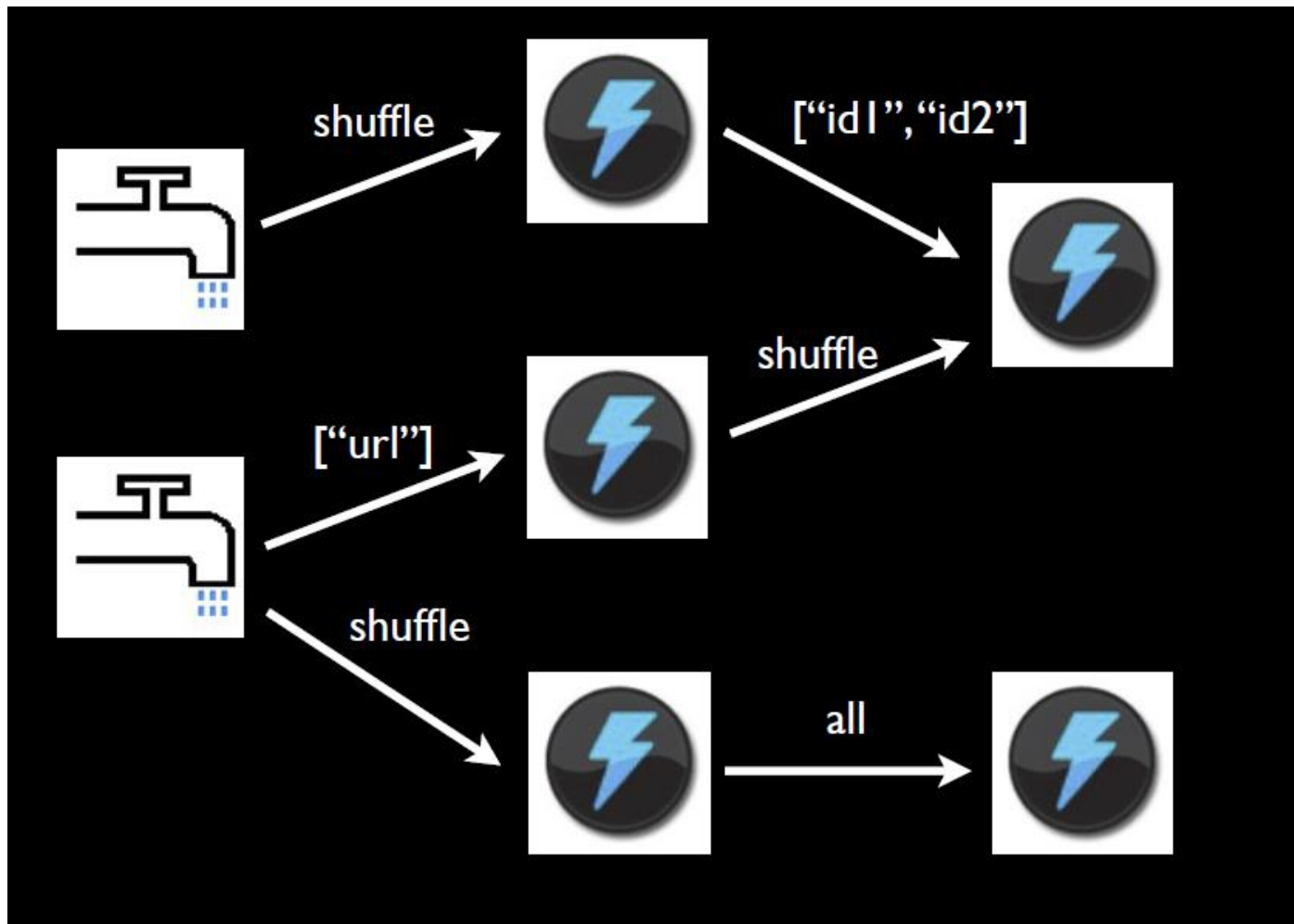# Topology



Network of spouts and bolts

# Tasks



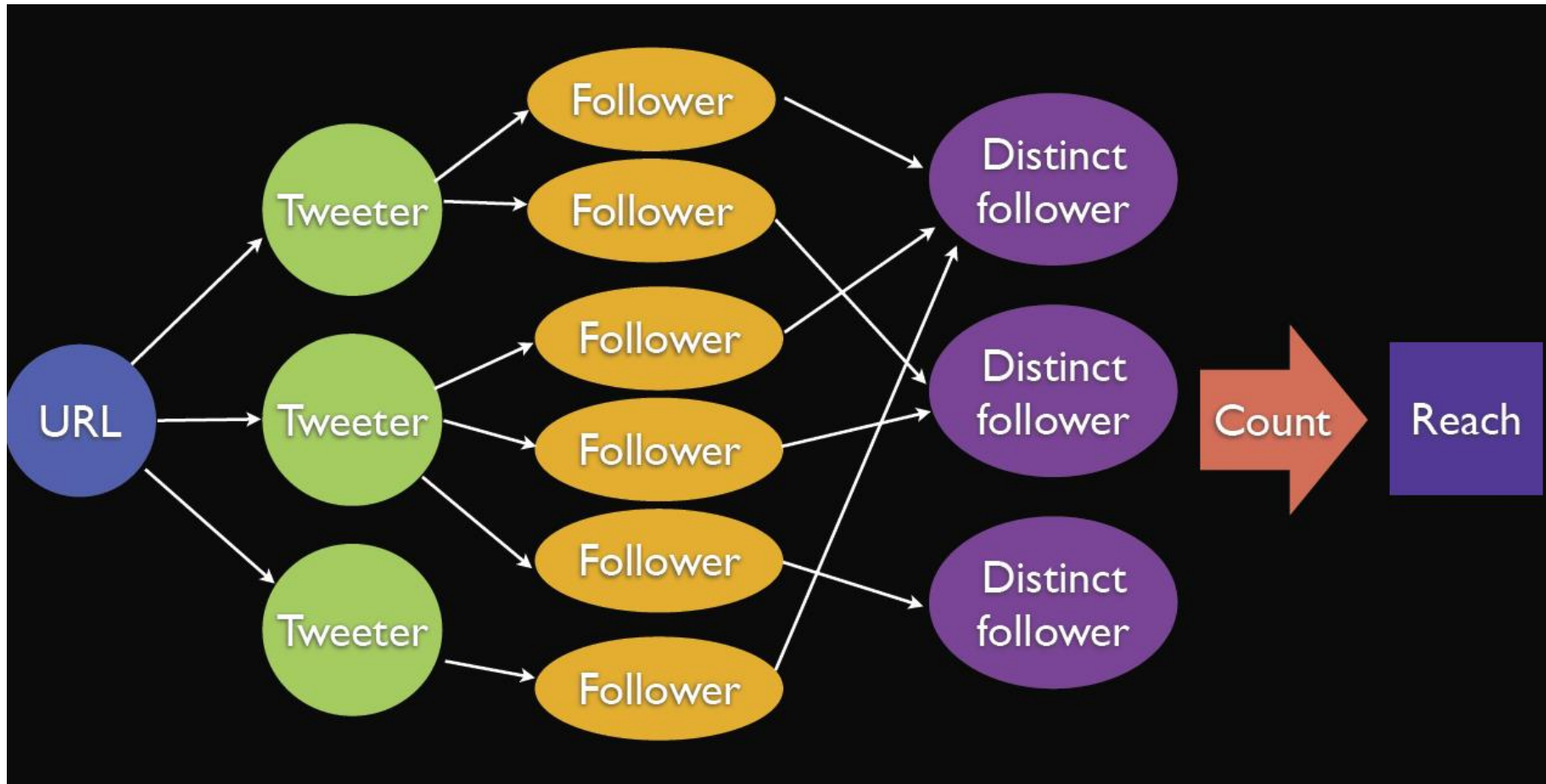Spouts and bolts execute as many tasks across the cluster

# Stream grouping

- **Shuffle grouping:** pick a random task

- **Fields grouping:** consistent hashing on a subset of tuple fields

- **All grouping:** send to all tasks

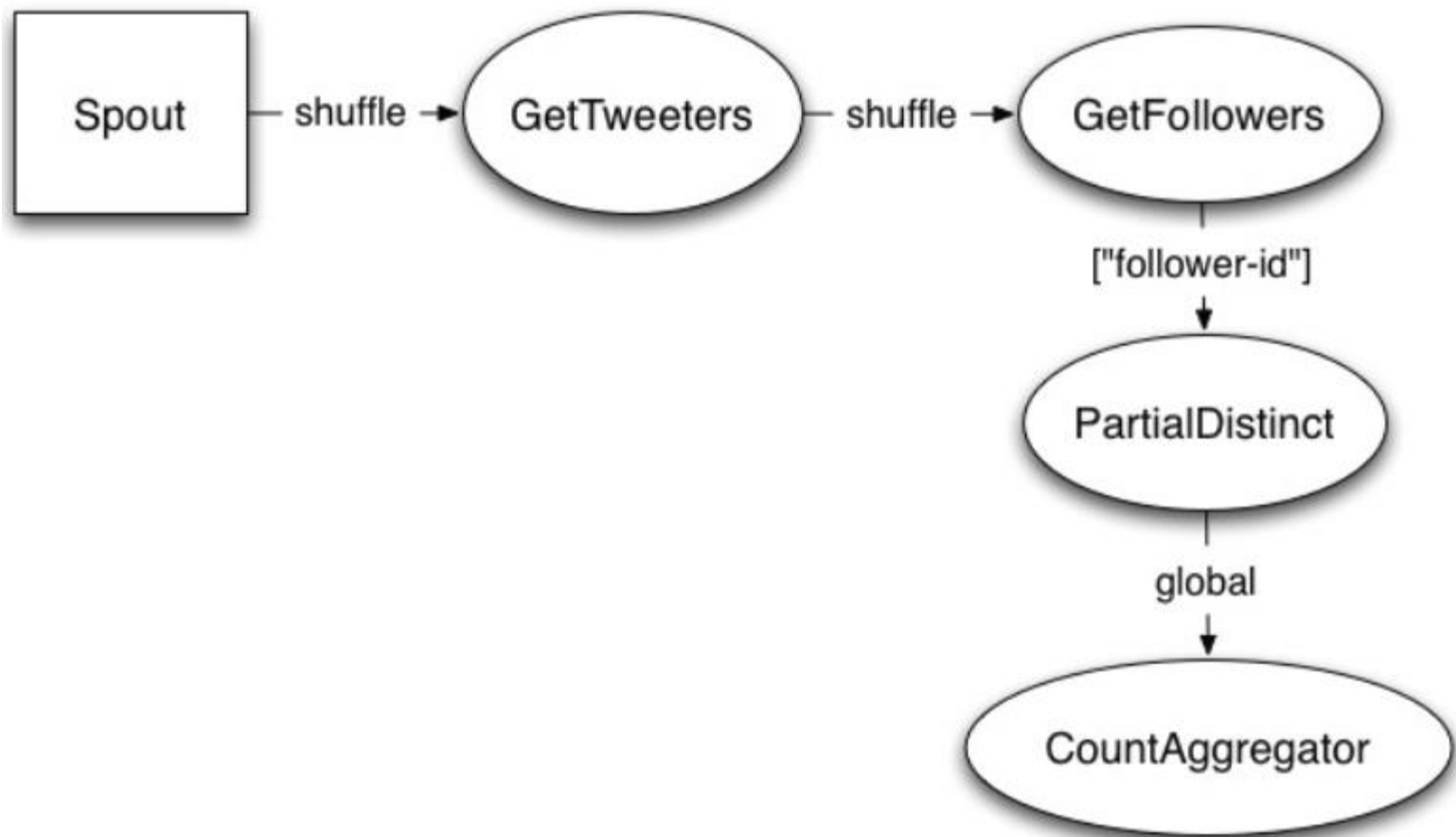- **Global grouping:** pick task with lowest id

# Filtering and grouping

# Sample Application

# Filtering and grouping

# Cluster Coordination



Used for cluster coordination