

An Introduction to Distributed Data Streaming

Elements and Systems

Paris Carbone <parisc@kth.se>
PhD Candidate
KTH Royal Institute of Technology



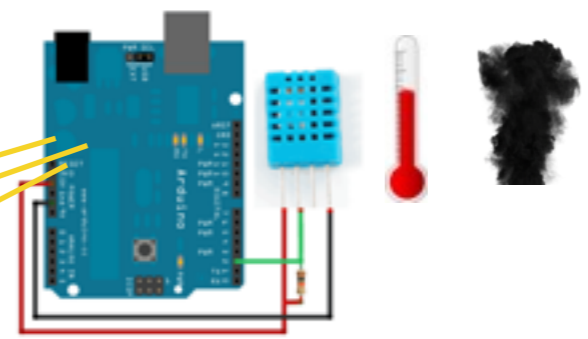
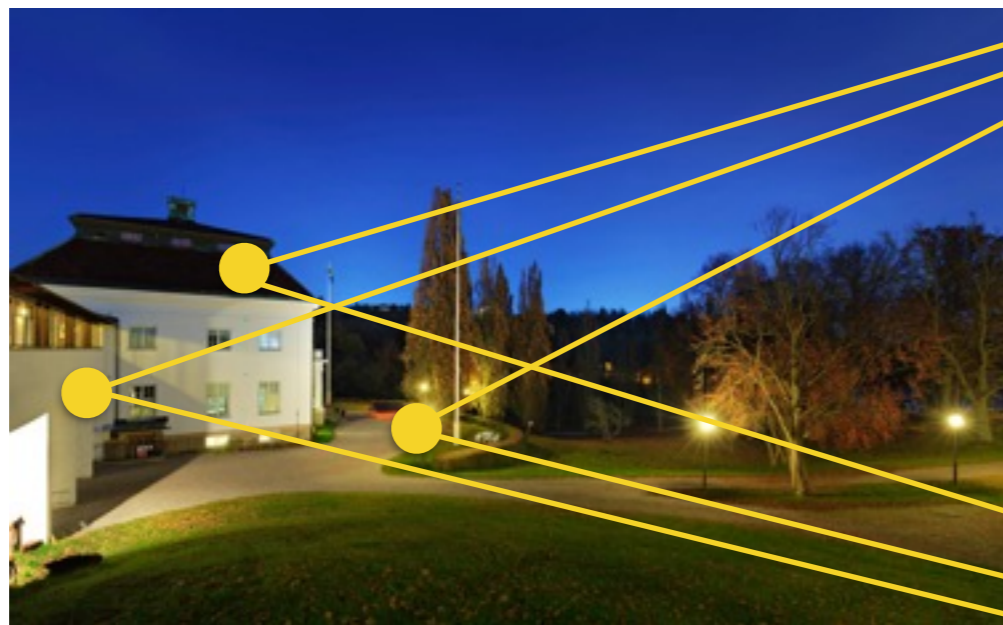




how to avoid this?



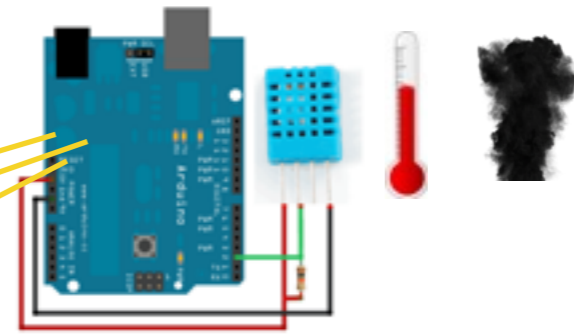
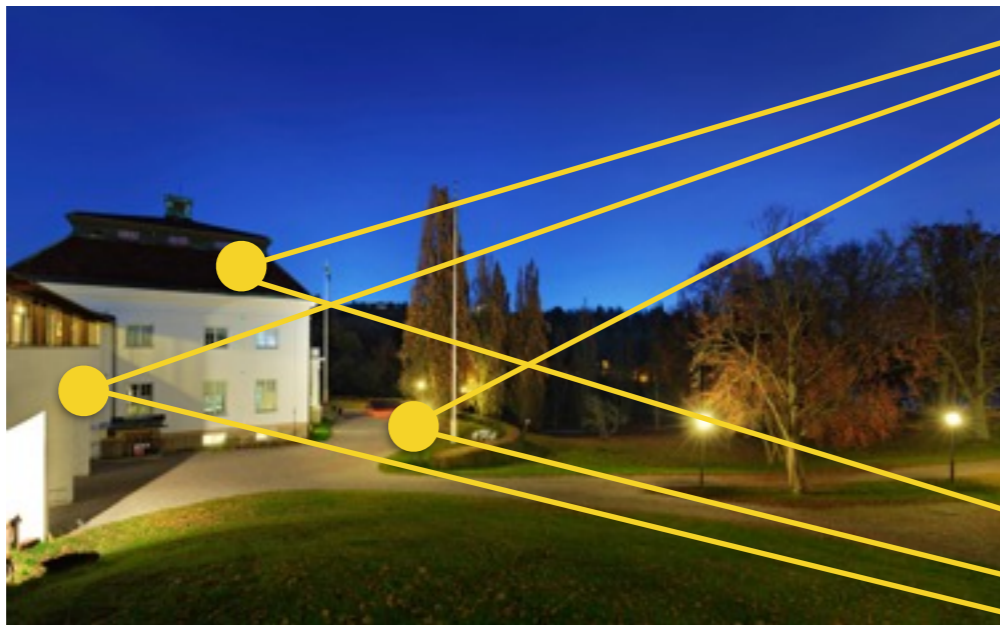
how to avoid this?



Q



how to avoid this?



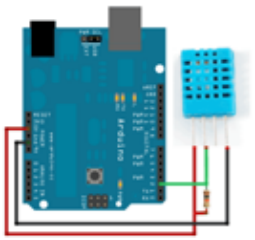
Q



$$Q = \text{thermometer} + \text{smoke plume}$$

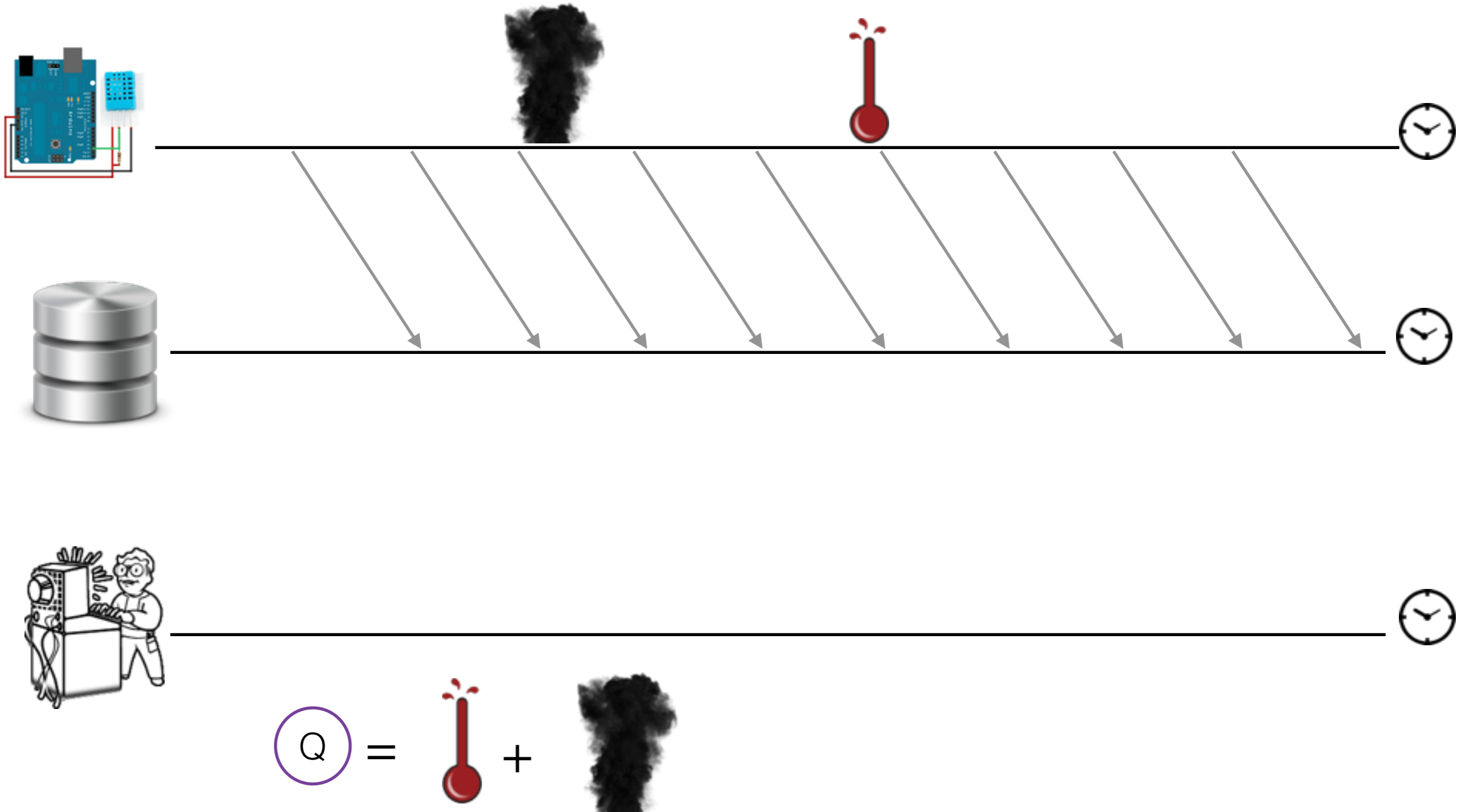
how to avoid this?

Motivation

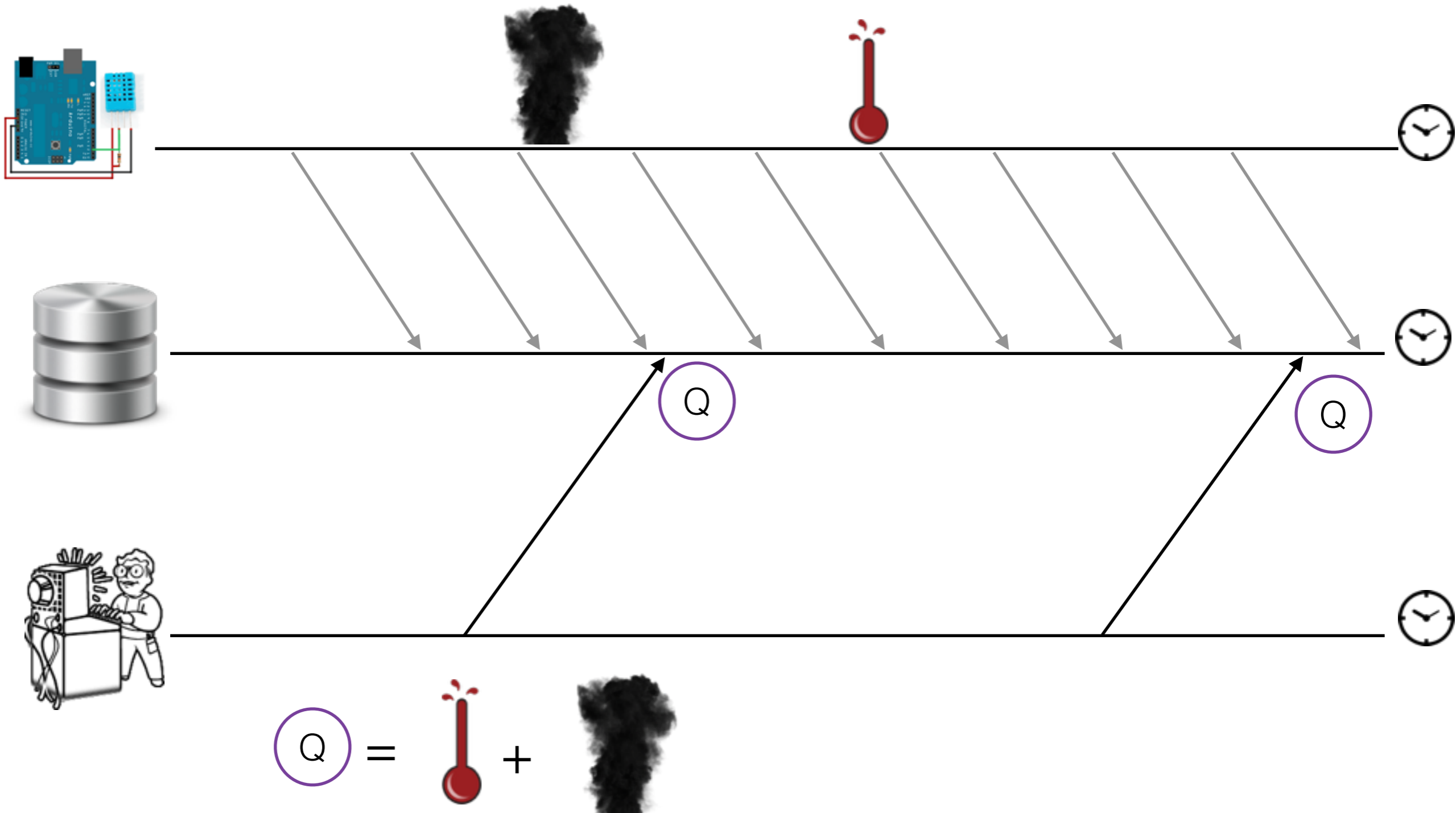


$$Q = \text{thermometer} + \text{smoke}$$

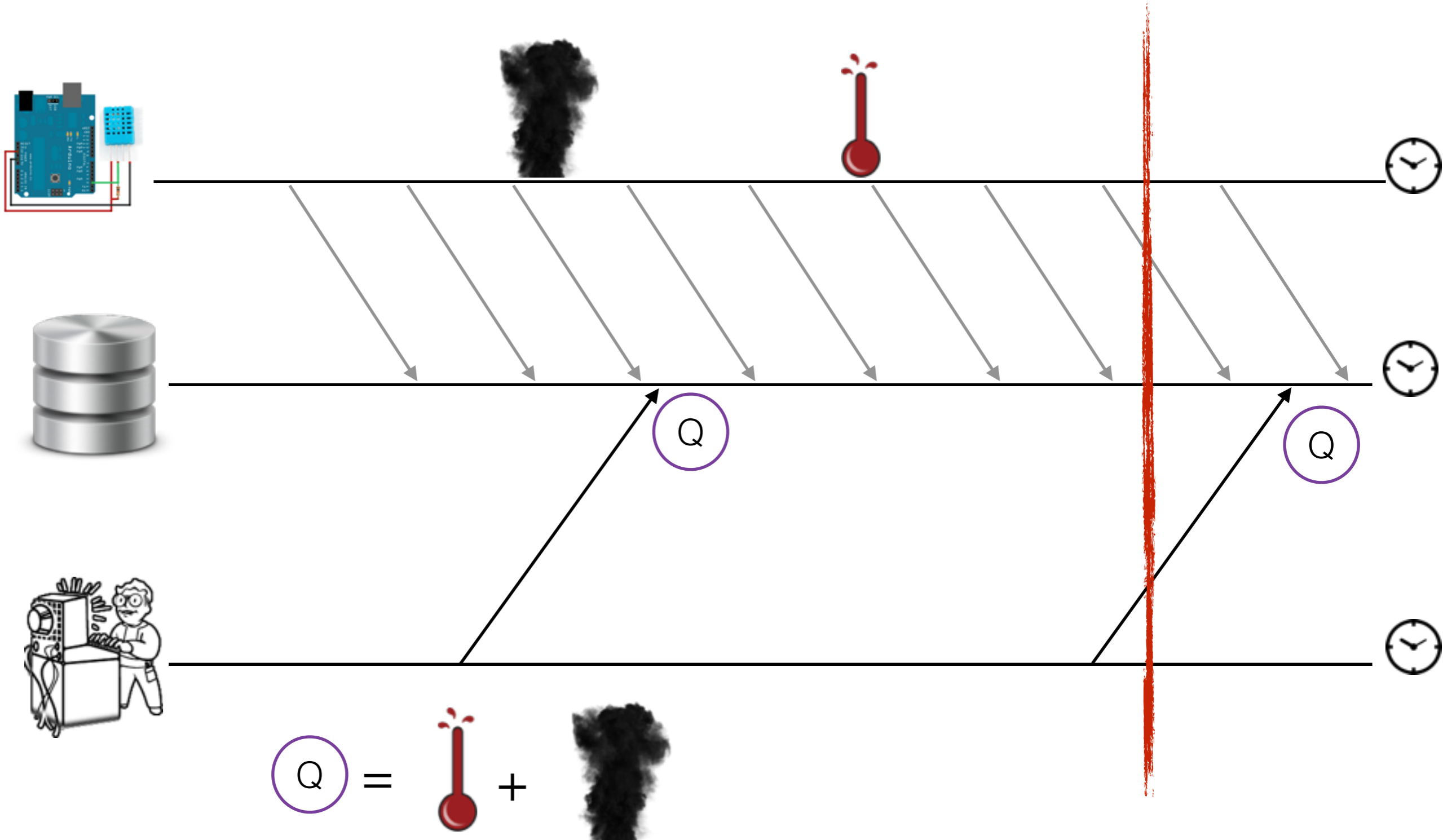
Motivation



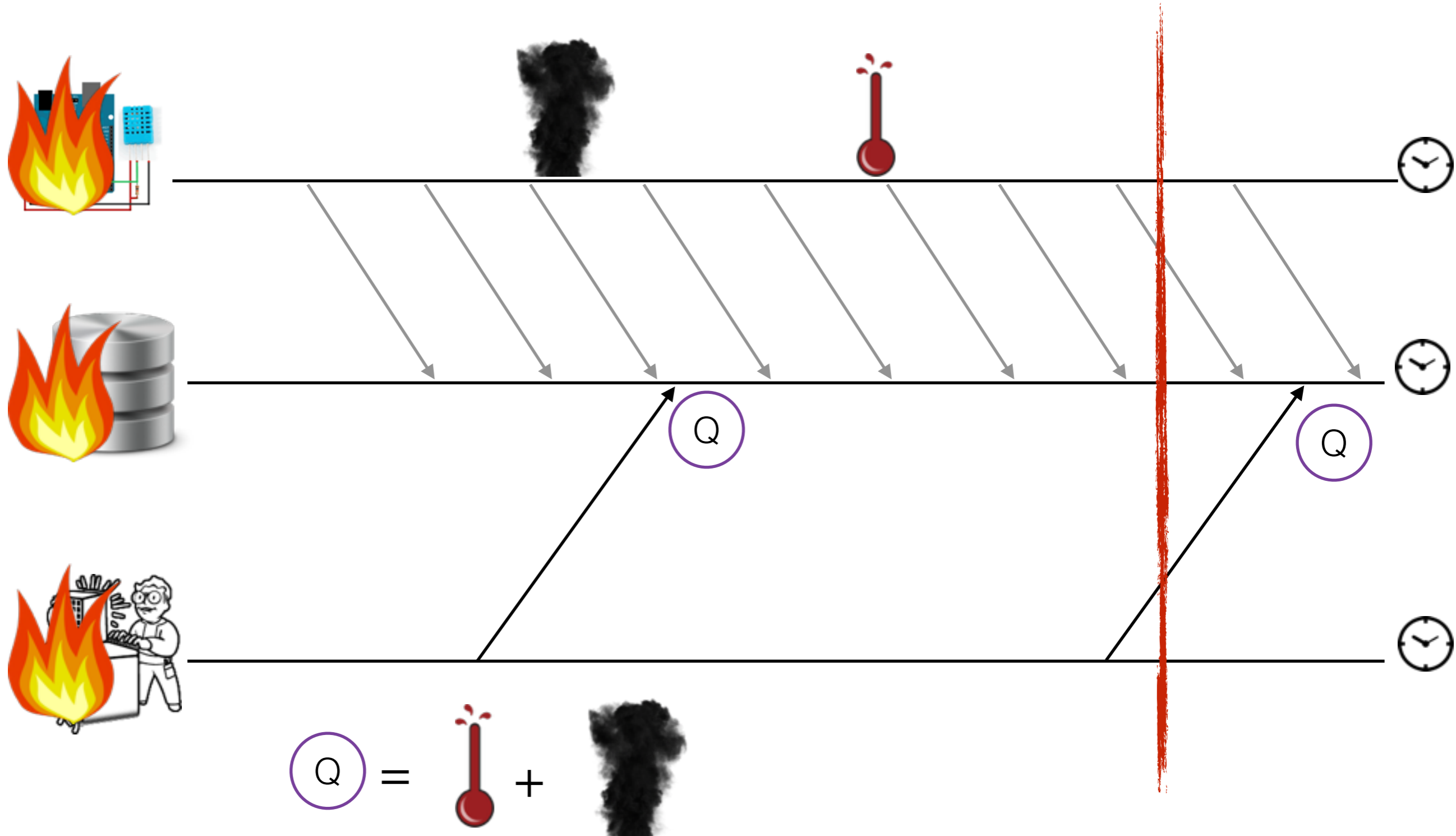
Motivation



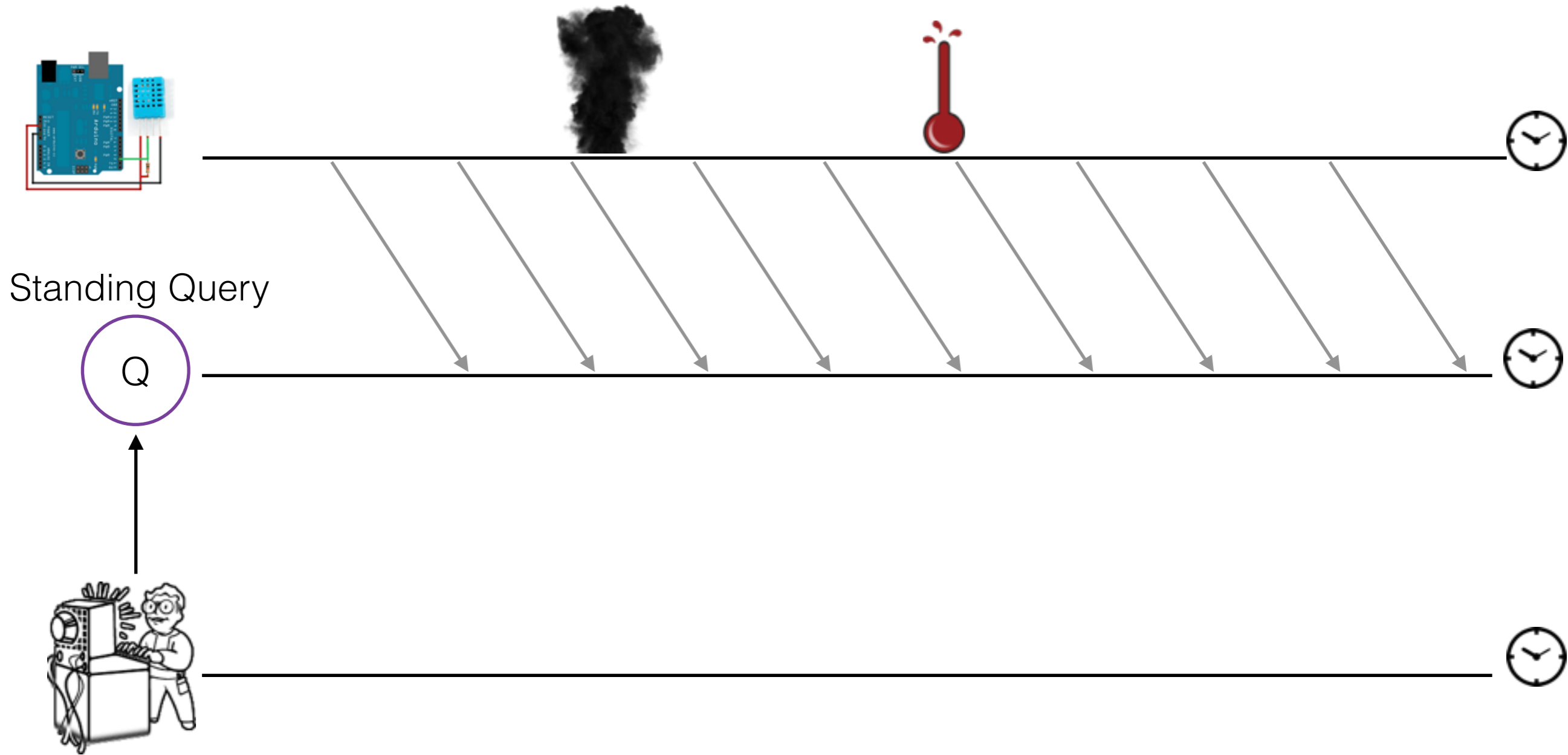
Motivation



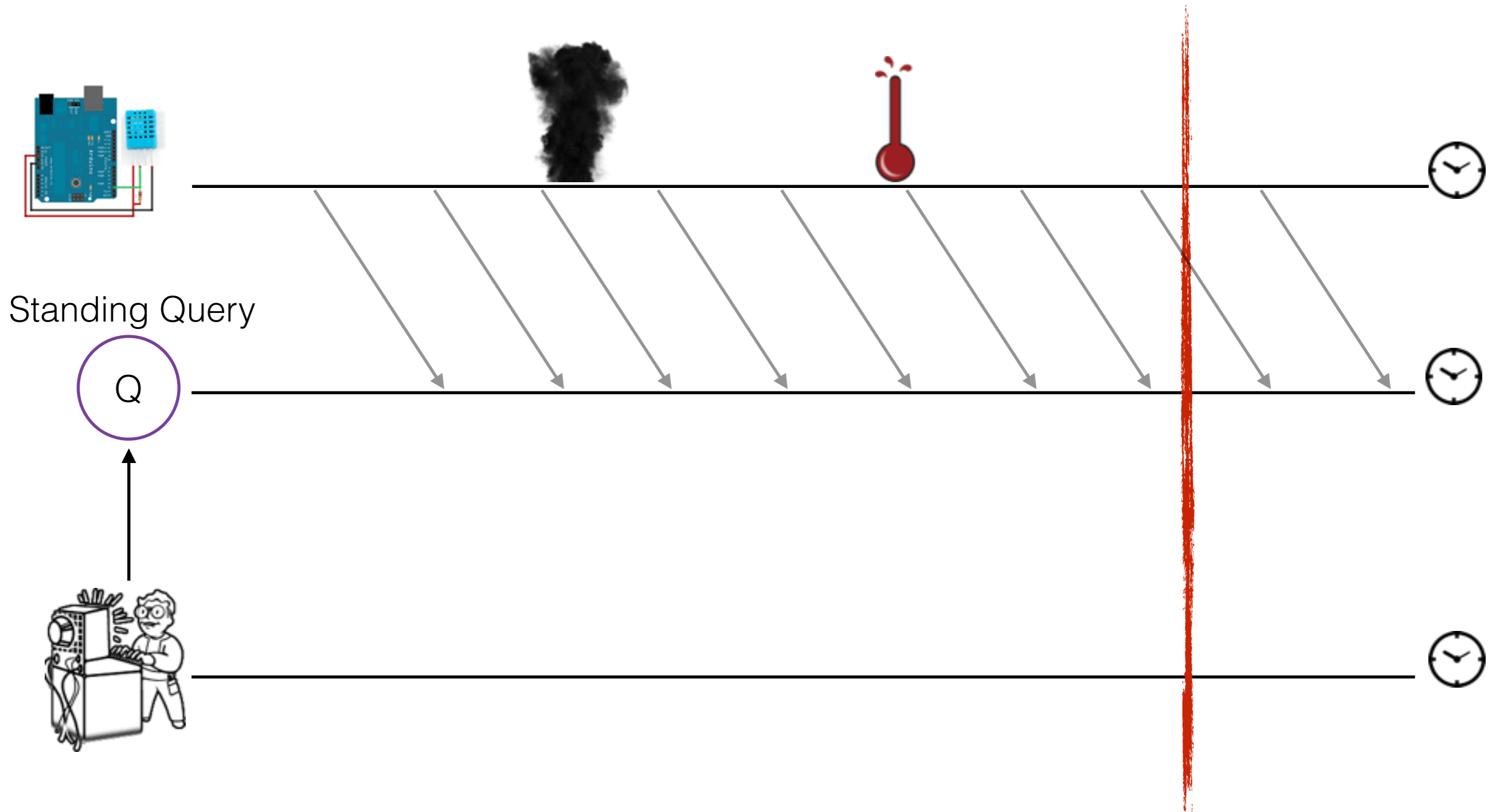
Motivation



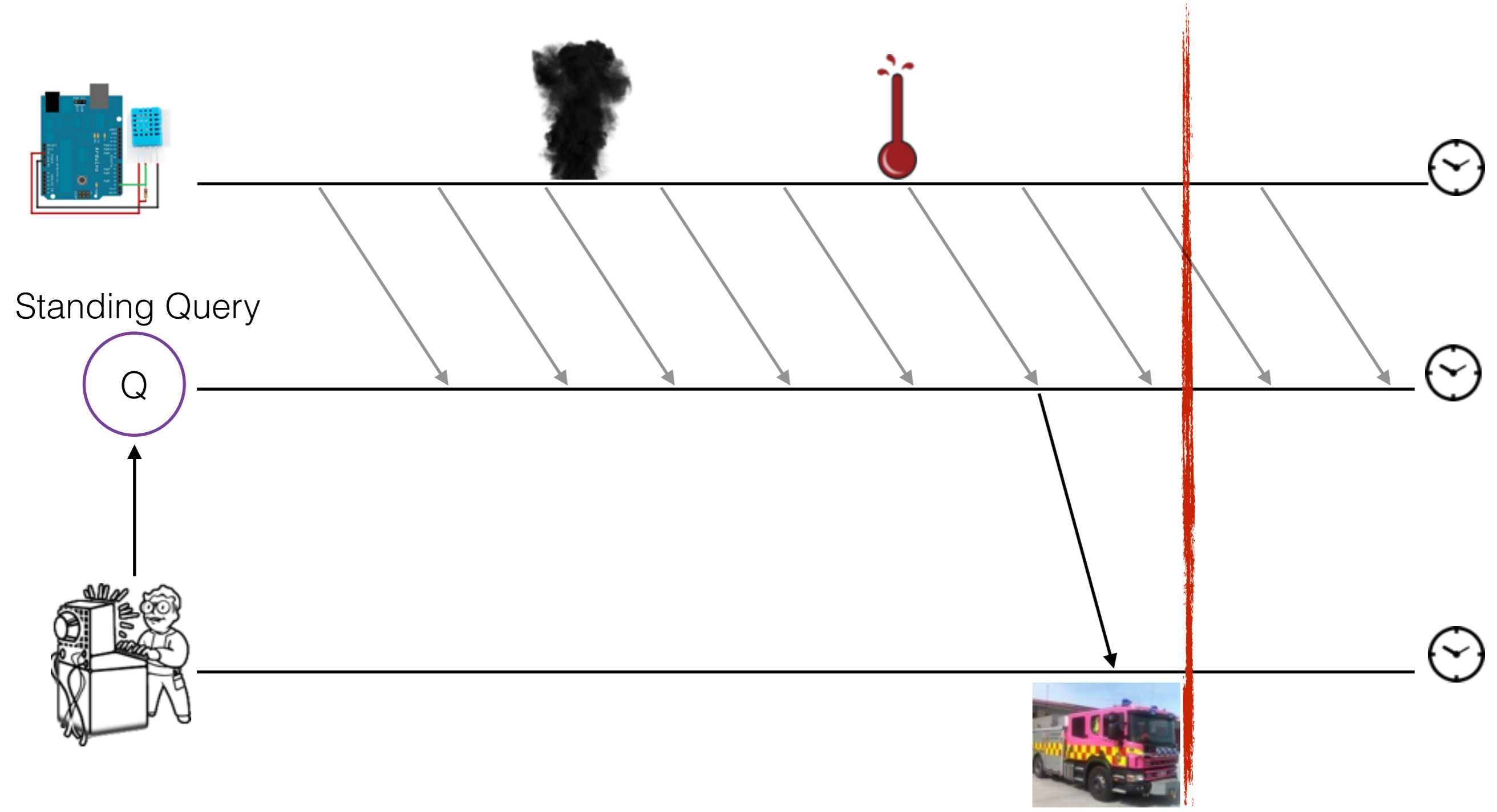
Motivation



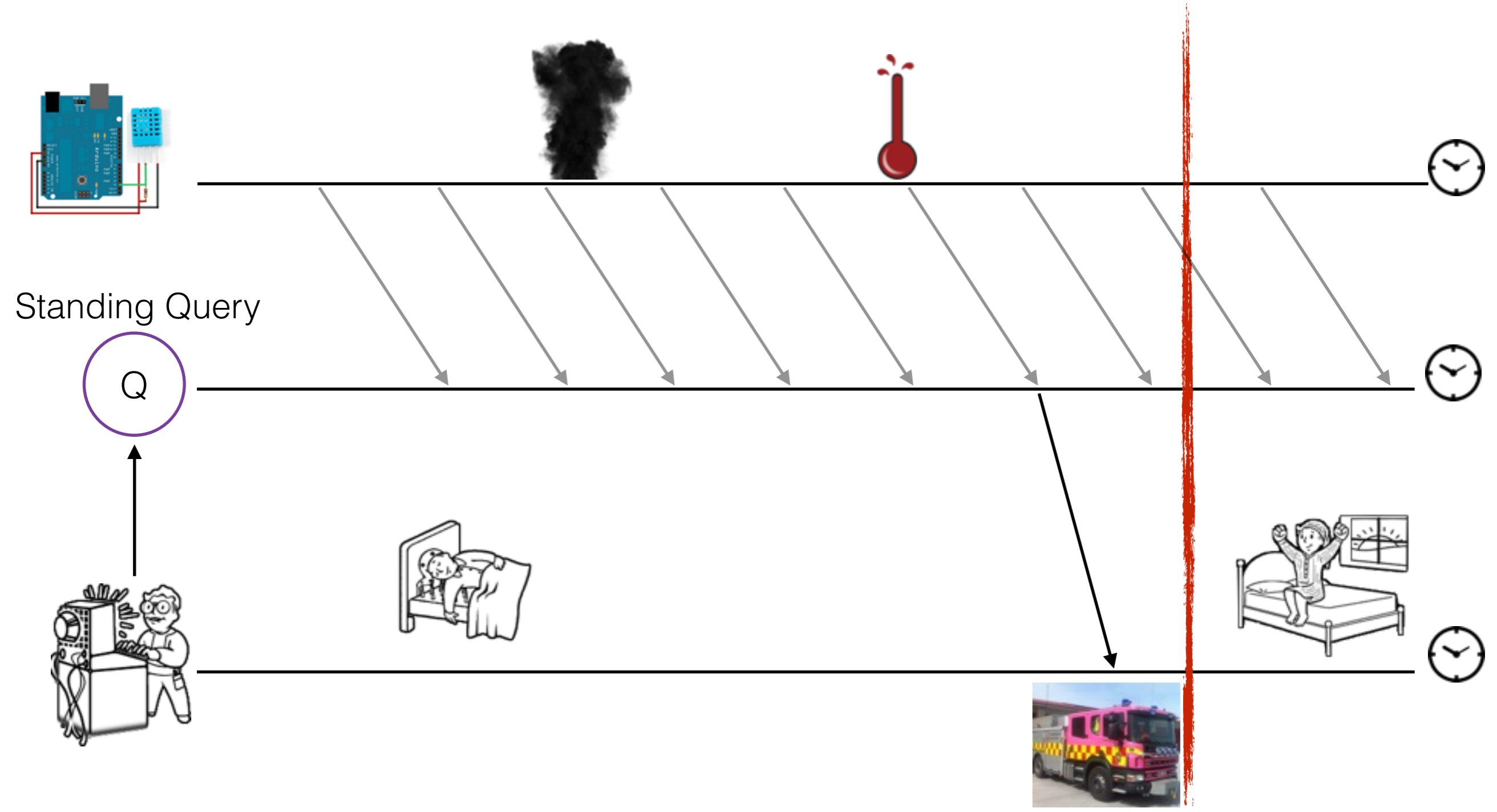
Motivation



Motivation



Motivation

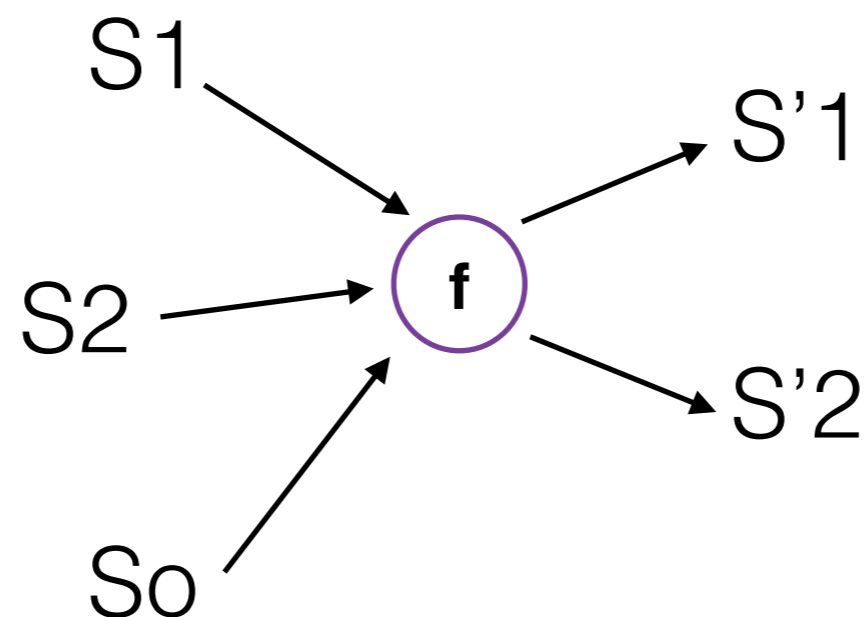


Preliminaries

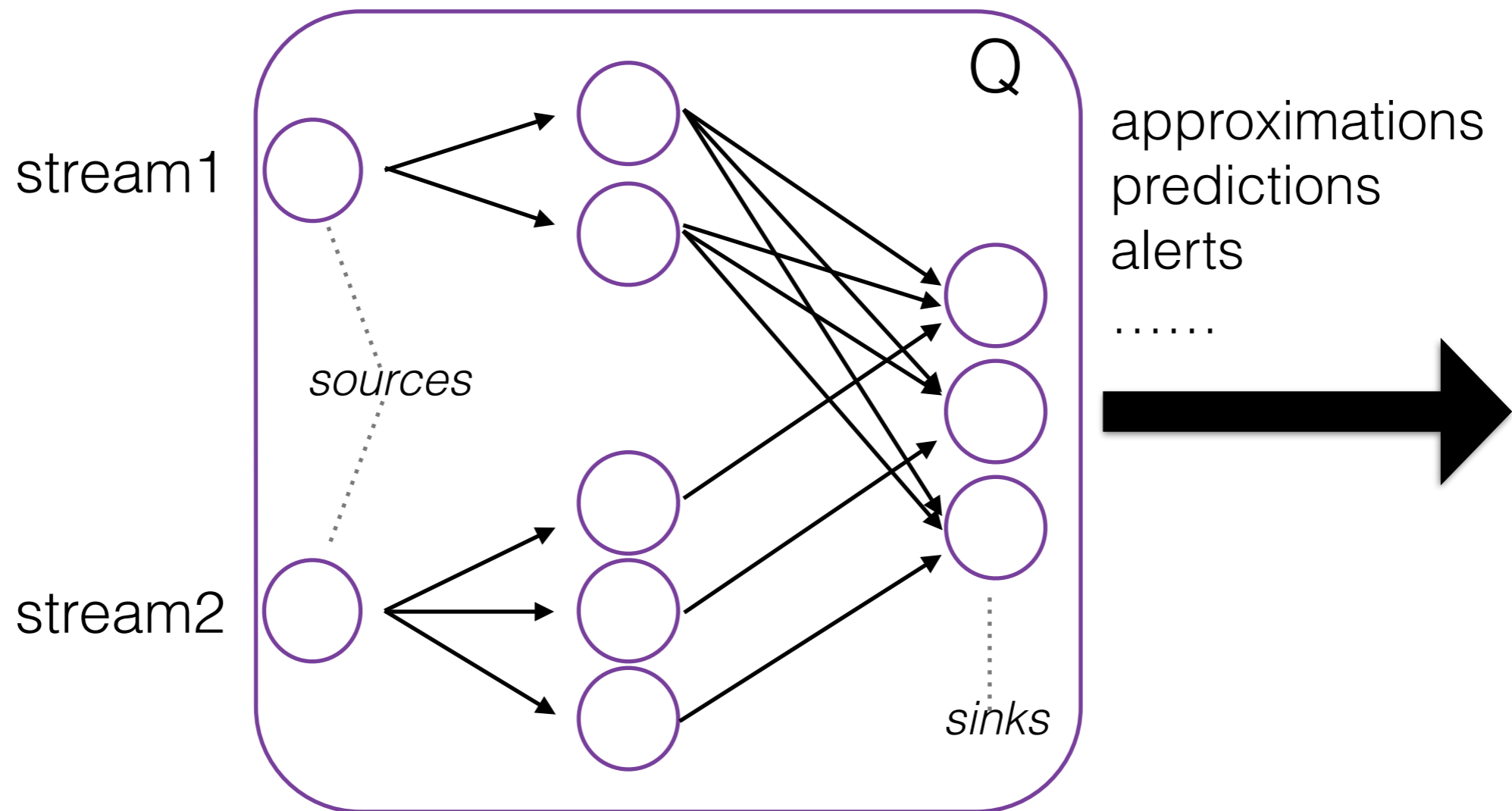
- Data Streaming Paradigm
 - Incoming data is unbound - continuous arrival
 - Standing queries are evaluated continuously
 - Queries operate on the full data stream or on the most recent views of the stream ~ windows

Data Streams Basics

- Events/Tuples : elements of computation - respect a schema
- Data Streams : unbounded sequences of events
- Stream Operators: consume streams and generate new ones.
 - Events are consumed once - no backtracking!



Streaming Pipelines



Core Abstractions

- Windows
- Synopses (summary state)
- Partitioning

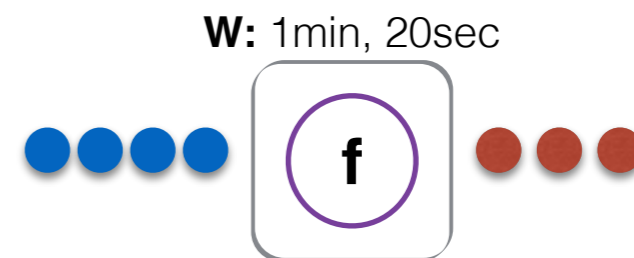
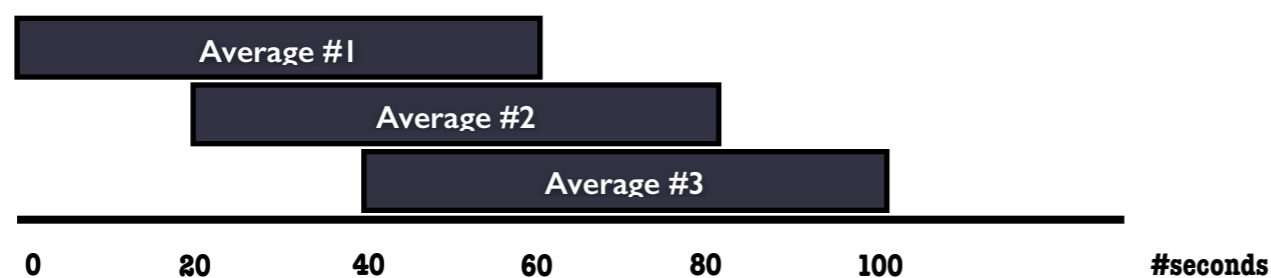
Windows

Discussion

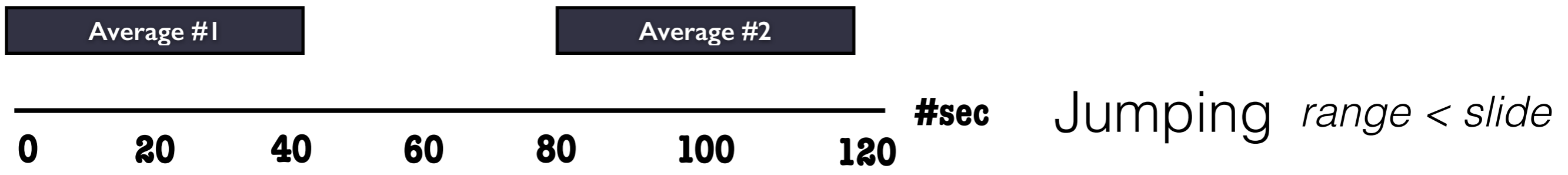
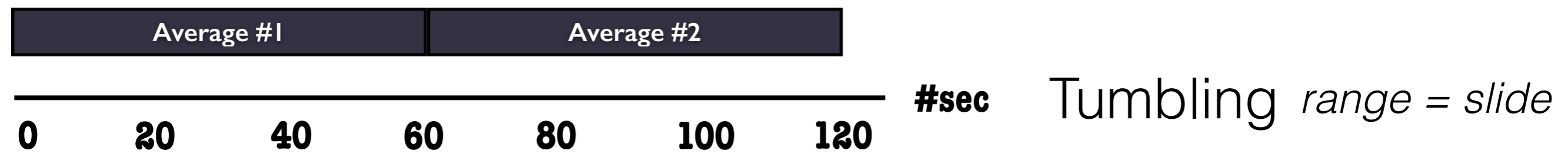
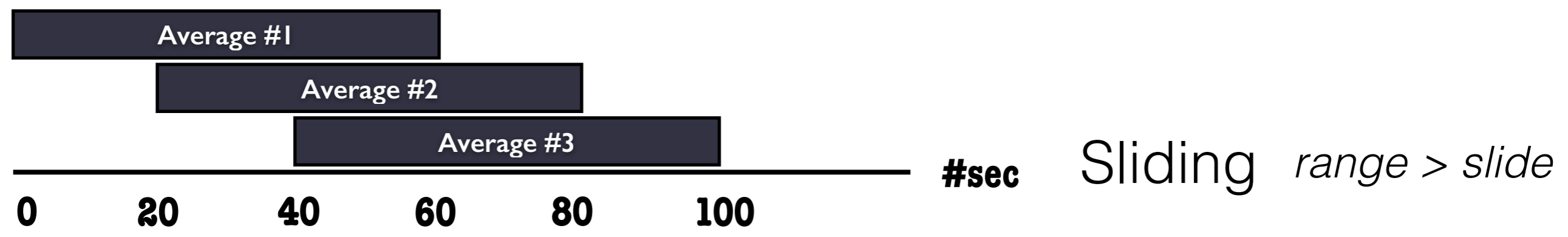
Why do we need windows?

Windows

- We are often interested only in **fresh data**
 - f = “average temperature over the last minute every 20 sec”
- **Range:** Most data stream processing systems allow window operations on the most recent history (eg. 1 minute, 1000 tuples)
- **Slide:** The frequency/granularity f is evaluated on a given range



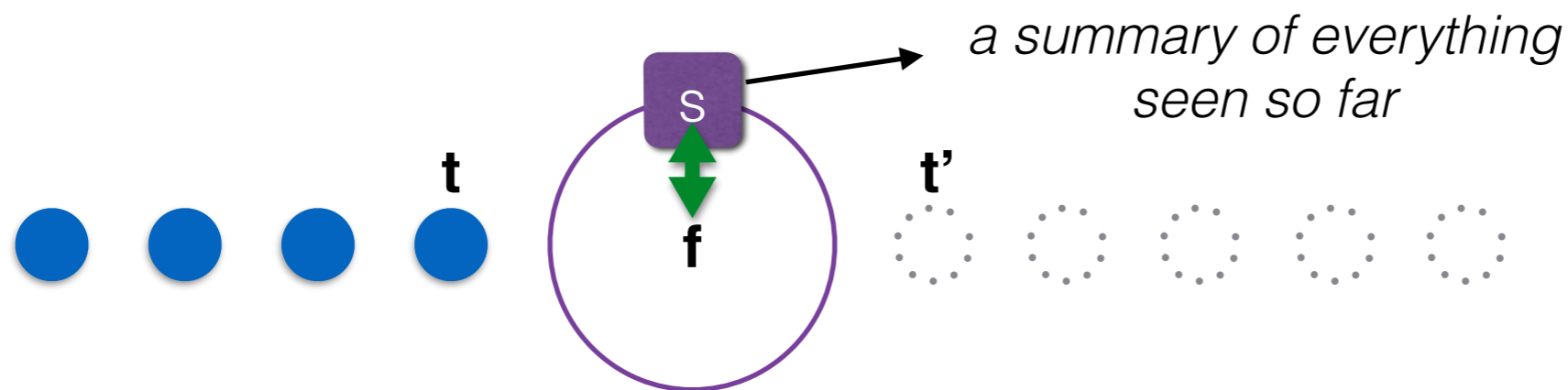
Window Types



Synopses

We cannot infinitely store all events seen

- **Synopsis:** A summary of an infinite stream
 - It is in principle any streaming operator state
- Examples: samples, histograms, sketches, state machines...



1. process t , s
2. update s
3. produce t'

What about window synopses?

Synopses-Aggregations

- **Discussion** - Rolling Aggregations
- Propose a synopsis, $s=?$ when
 - $f= \max$
 - $f= \text{ArithmeticMean}$
 - $f= \text{stDev}$

Synopses-Approximations

- **Discussion** - *Approximate Results*
- Propose a synopsis, $s=?$ when
 - $f=$ uniform random sample of k records over the whole stream
 - $f=$ filter distinct records over windows of 1000 records with a 5% error

Synopses-ML and Graphs

- Examples of cool synopses to check out
 - **Sparsifiers/Spanners** - approximating graph properties such as shortest paths
 - **Change detectors** - detecting concept drift
 - **Incremental decision trees** - continuous stream training and classification

Partitioning

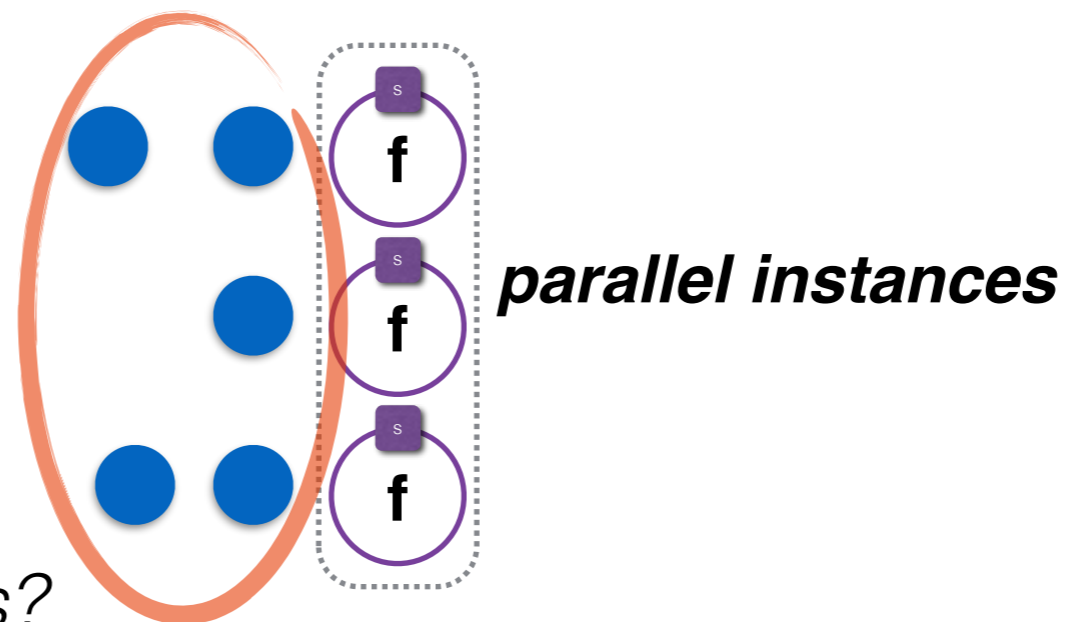
- One stream operator is not enough



- **Data** might be too large to process

- e.g. very high input rate, too many stream sources

- **State** could possibly not fit in memory

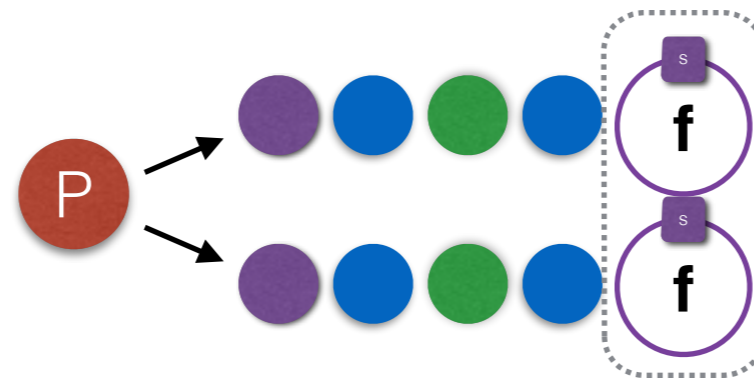


How do we partition the input streams?

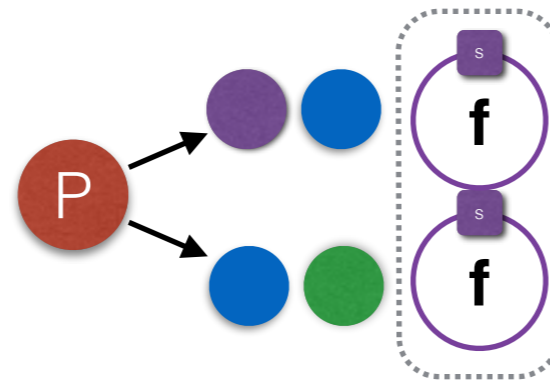
Partitioning

- Partitioning defines how we allocate events to each parallel instance. Typical partitioners are:

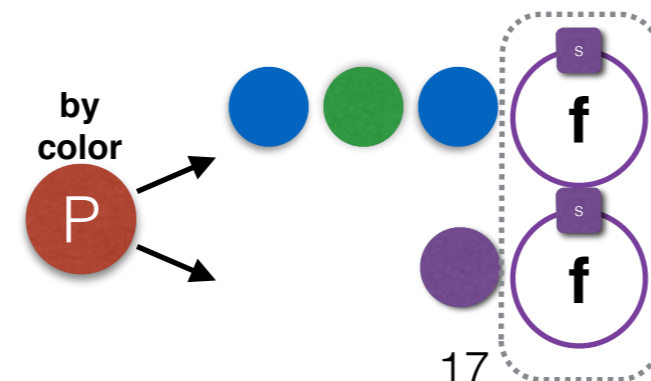
- Broadcast



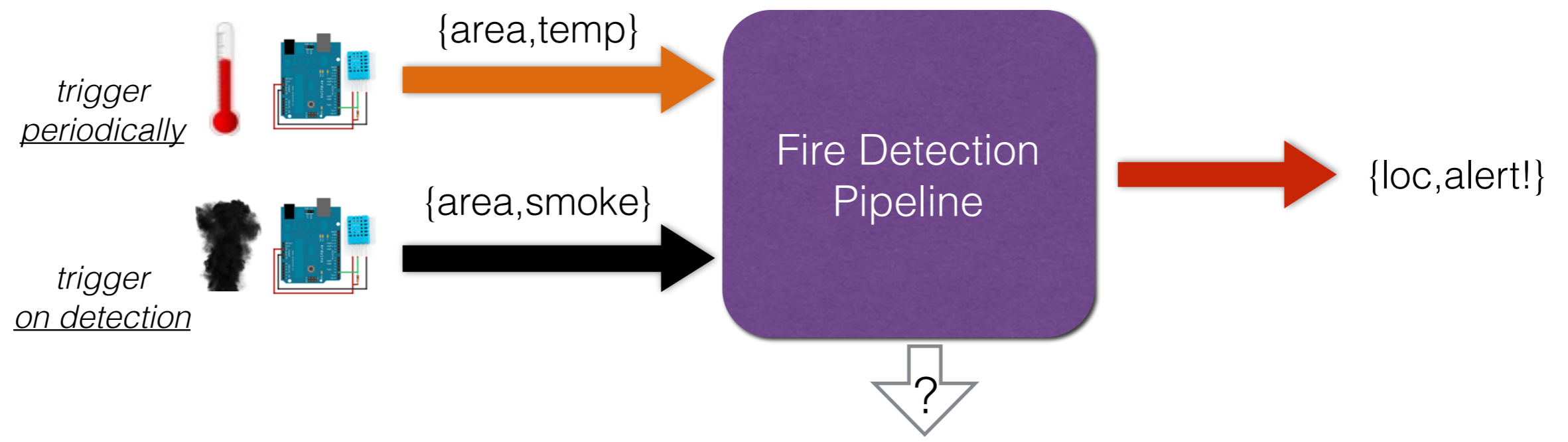
- Shuffle



- Key-based



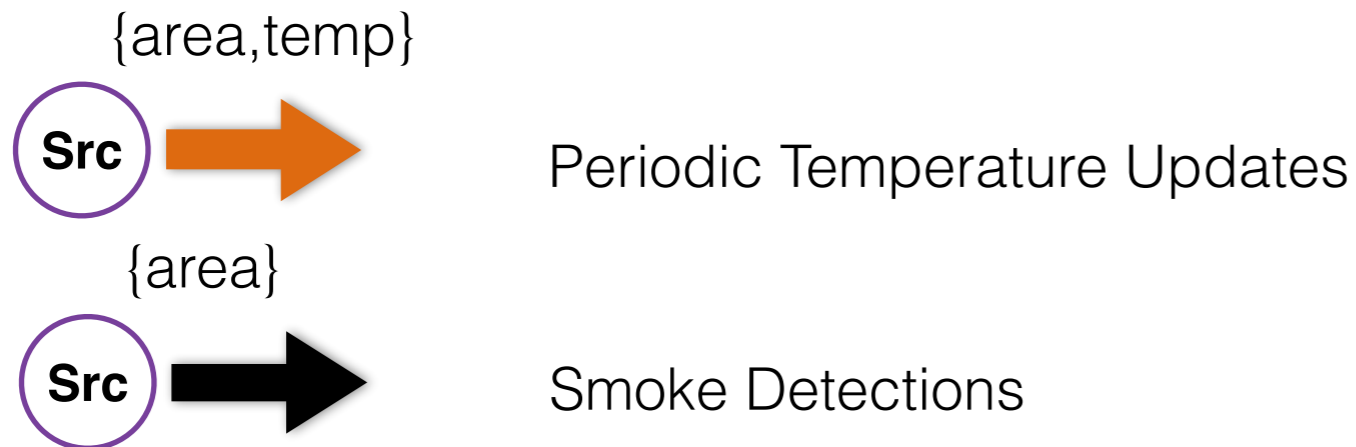
Putting Everything Together



- operators
 - synopses
 - windows
 - partitioning

Operators

Sensor Data Sources



Rolling Arithmetic Mean of Temperatures

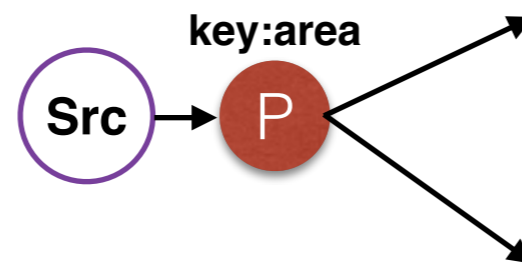


State Machine-based Fire Alarm



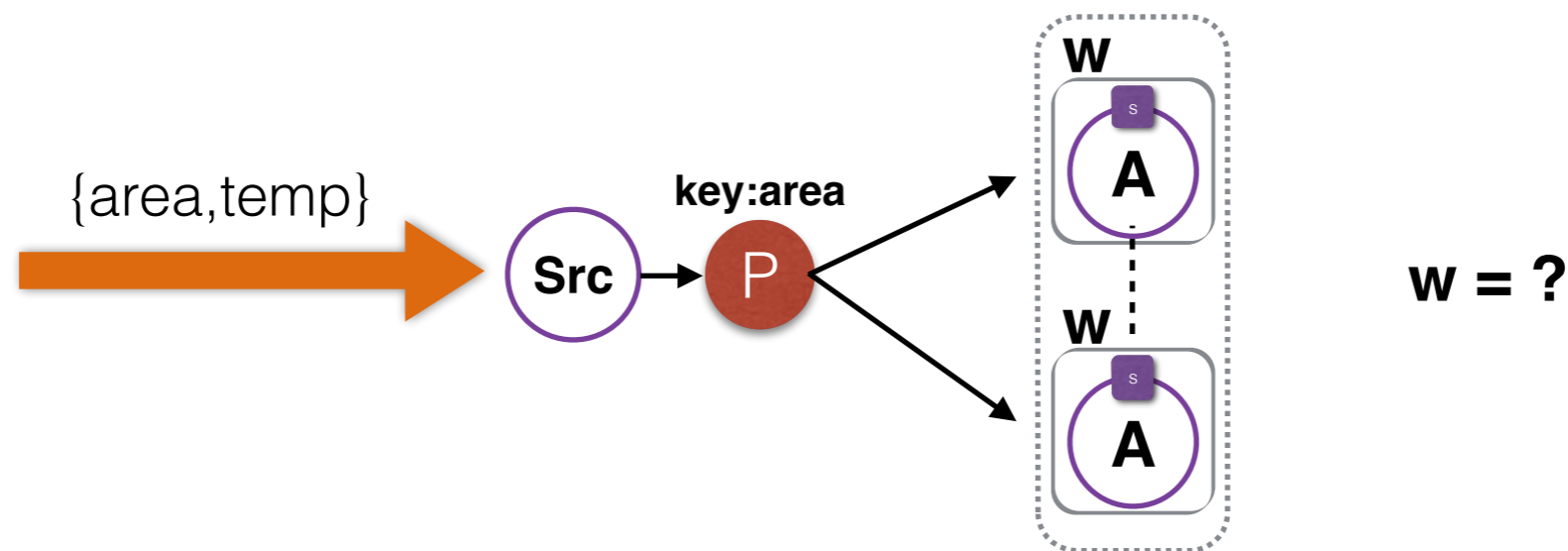
Partitioning

- We are only interested in correlating smoke and high temperature within the same area
- Events carry area information so we can partition our computation by **area**



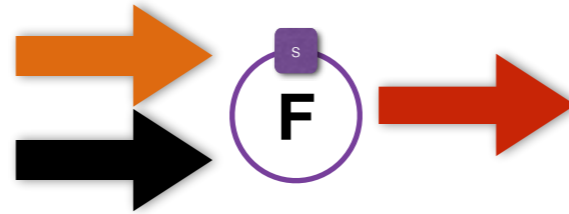
Windowing

- Individual sensor data could be potentially faulty
- We need to gather data from all temperature sensors of an area and produce an average
- We want fresh average temperatures



The Fire Alarm

The Fire Alarm

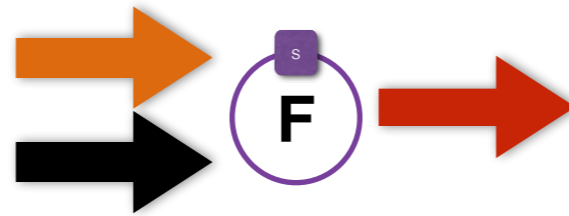


The Fire Alarm

T : avgTemp > 40

T : avgTemp < 40

S : Smoke

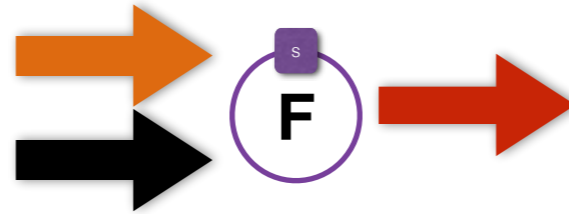


The Fire Alarm

T : avgTemp > 40

T : avgTemp < 40

S : Smoke



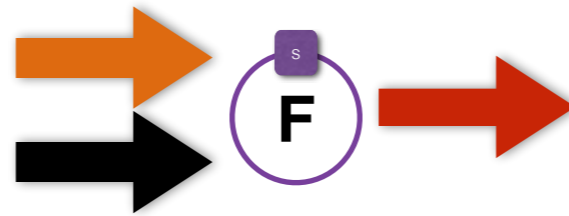
...TTTSTTSTTT...

The Fire Alarm

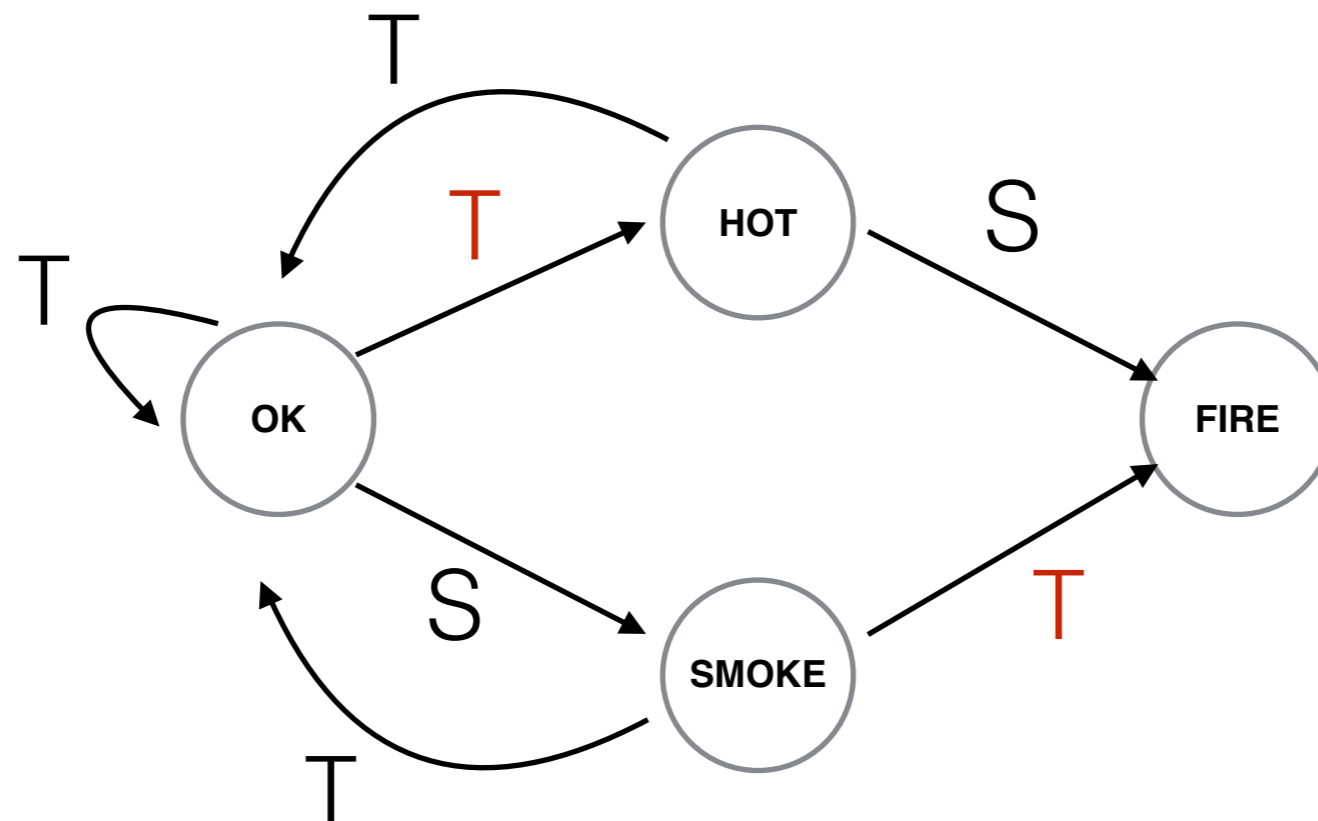
T : avgTemp > 40

T : avgTemp < 40

S : Smoke



...TTTSTTSTTTT...

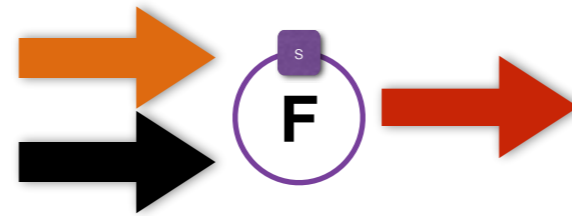


The Fire Alarm

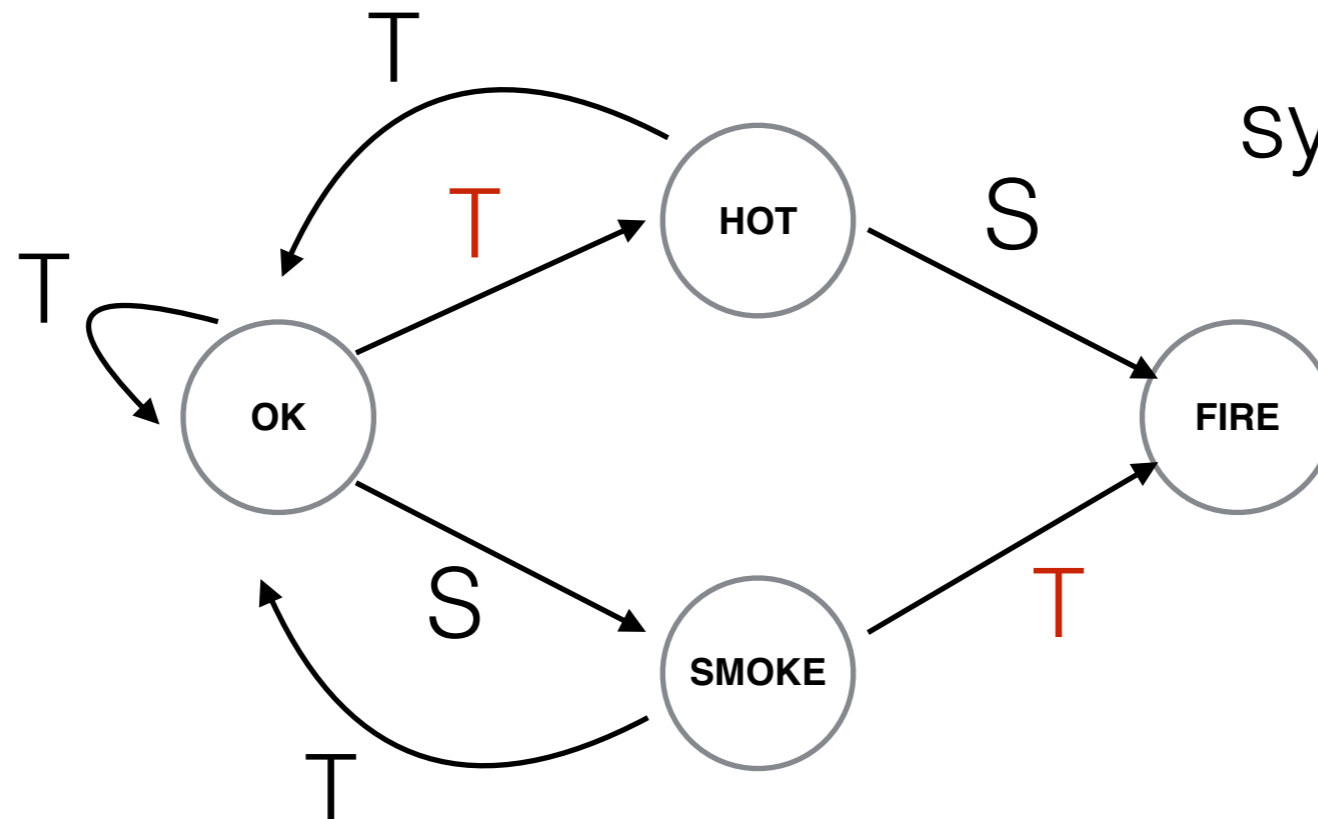
T : avgTemp > 40

T : avgTemp < 40

S : Smoke

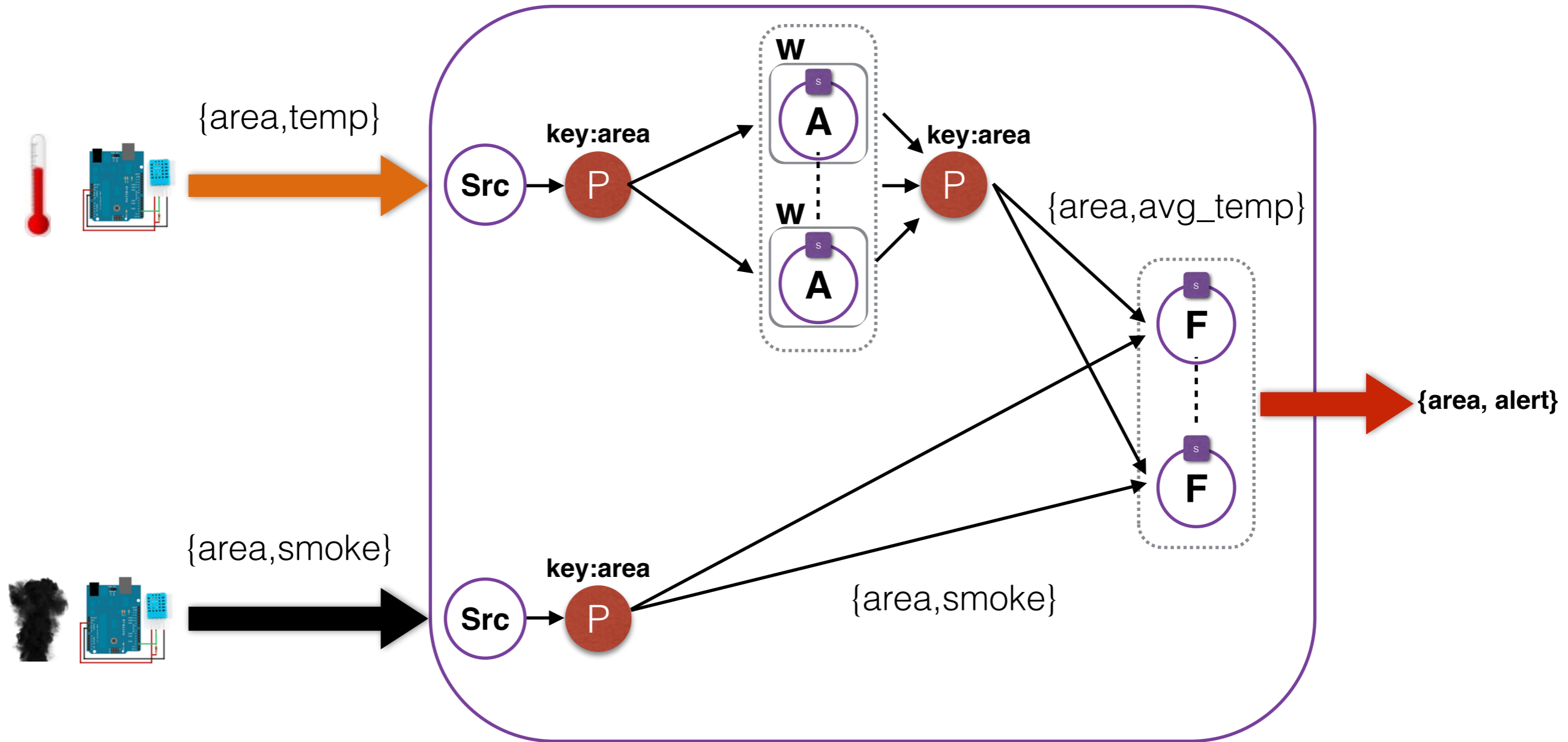


...TTTSTTSTTT...



synopsis= 1 state

Putting Everything Together



Systems: The Big Picture

Proprietary

Google
DataFlow



IBM
Infosphere



Microsoft
Azure



Open Source



Flink



Samza



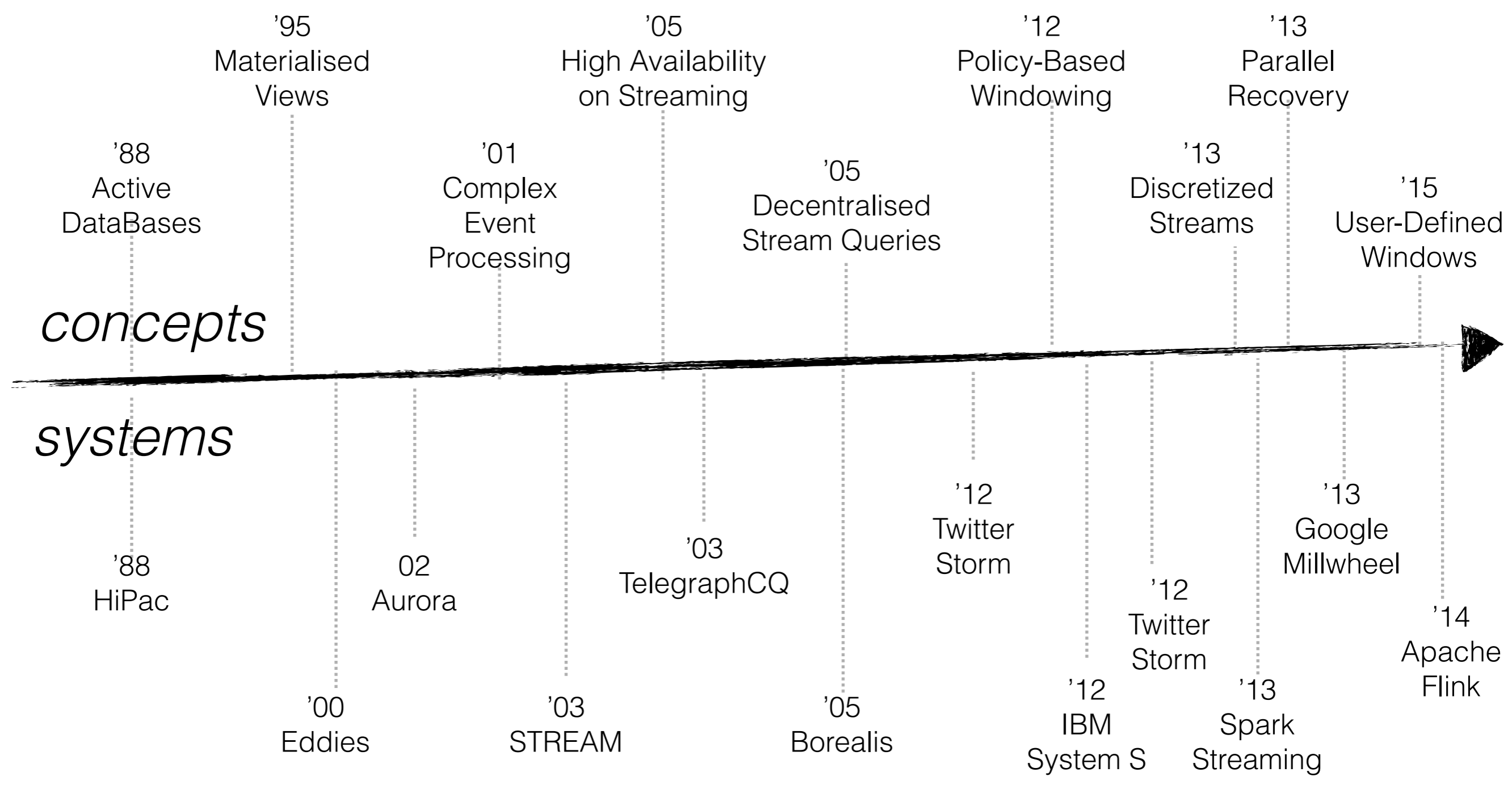
Spark



Storm



Evolution



Programming Models

Compositional

Declarative



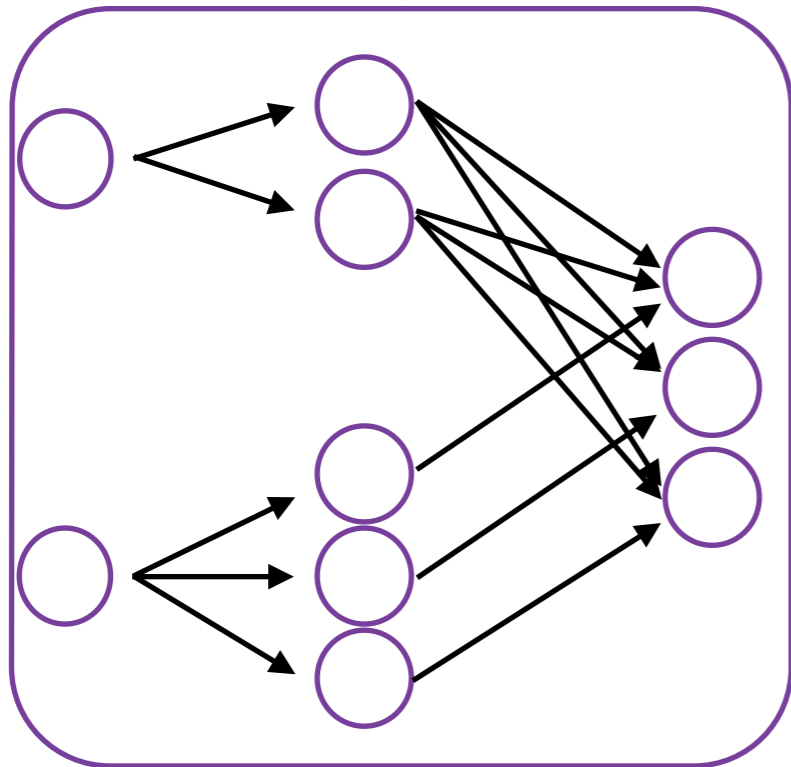
- Offer basic building blocks for composing custom operators and topologies
- Advanced behaviour such as windowing is often missing
- Custom Optimisation

- Expose a high-level API
- Operators are higher order functions on abstract data stream types
- Advanced behaviour such as windowing is supported
- Self-Optimisation

Programming Model Types



DStream, DataStream,
PCollection...



samza

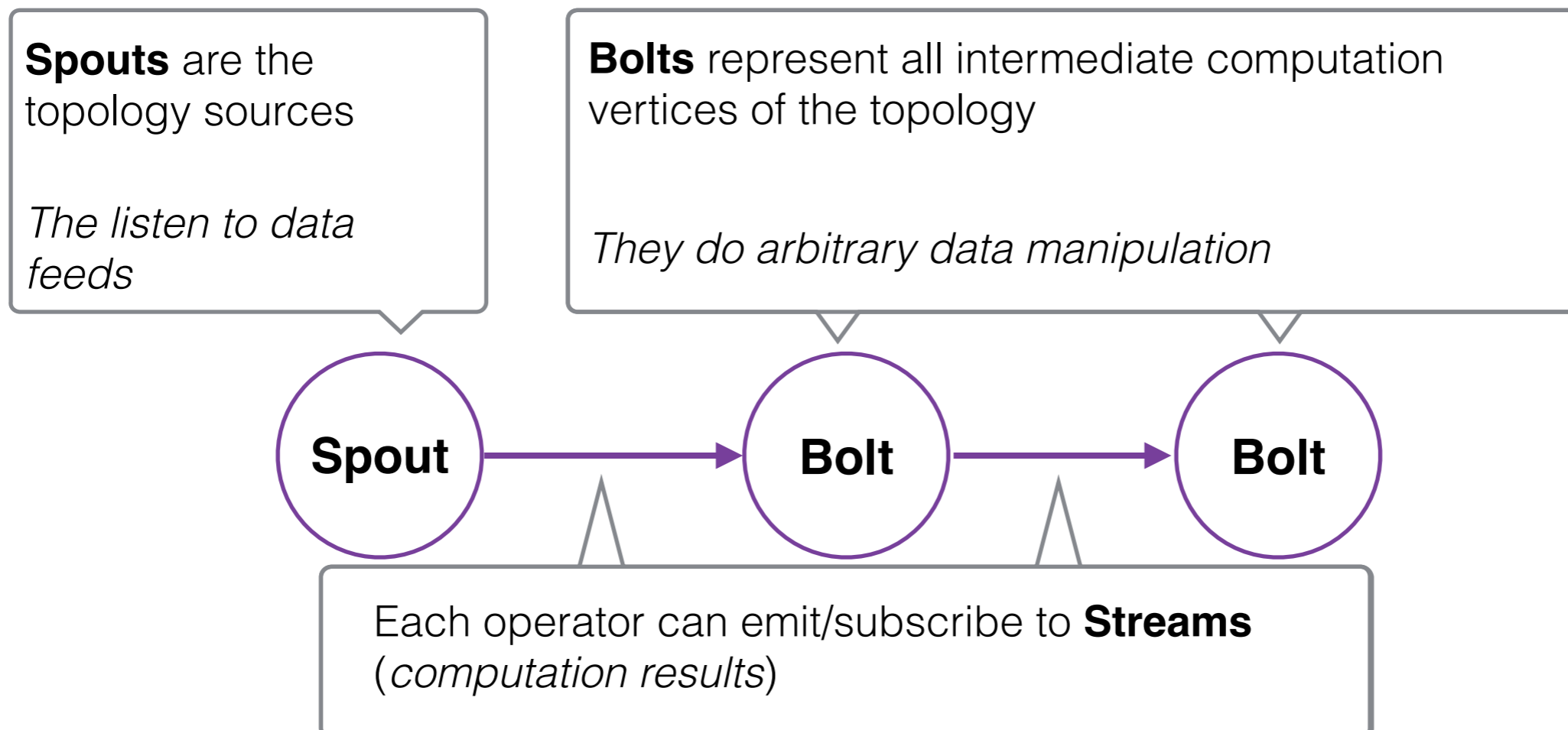


- Transformations abstract operator details
- Suitable for engineers and data analysts

- Direct access to the execution graph / topology
- Suitable for engineers

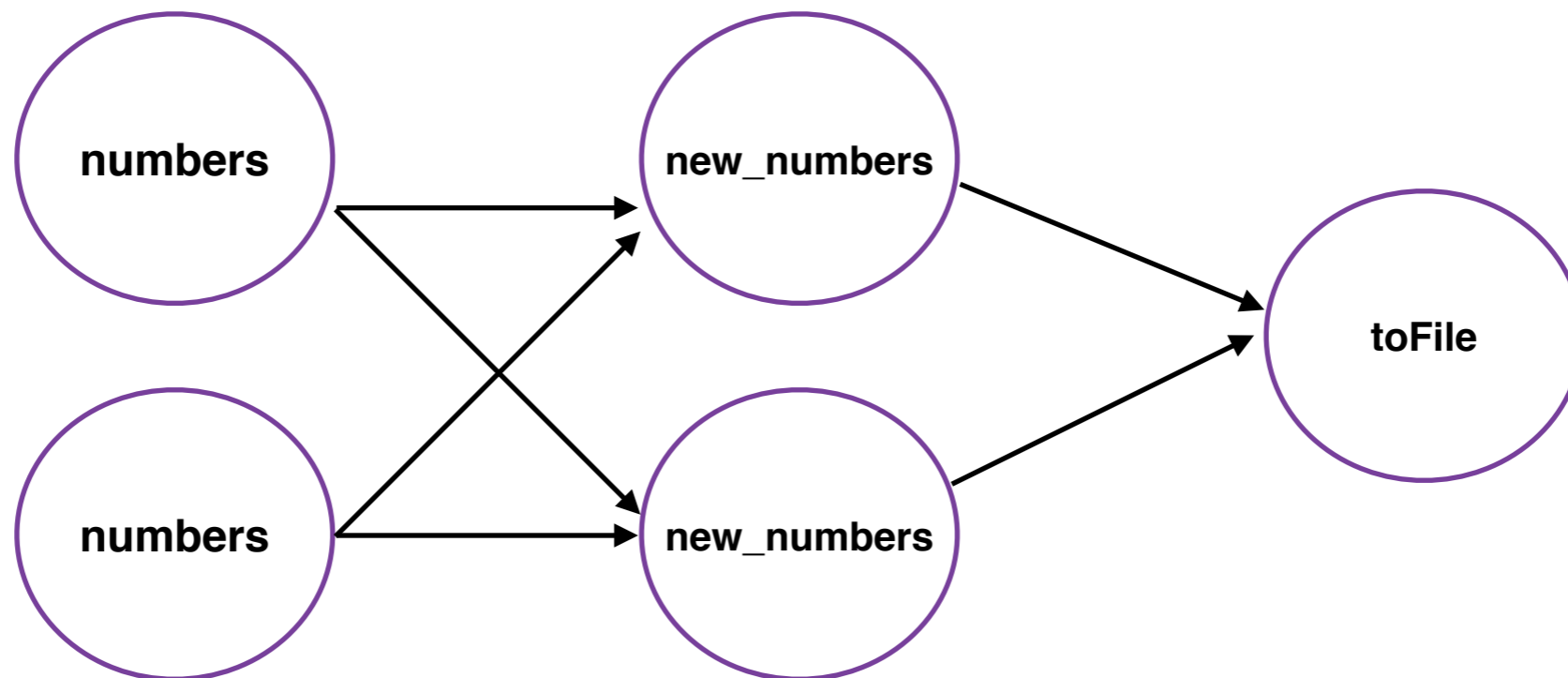
Standing Queries with Apache Storm

- Step 1: Implement input (**Spouts**) and intermediate operators (**Bolts**)
- Step 2: Construct a **Topology** by combining operators

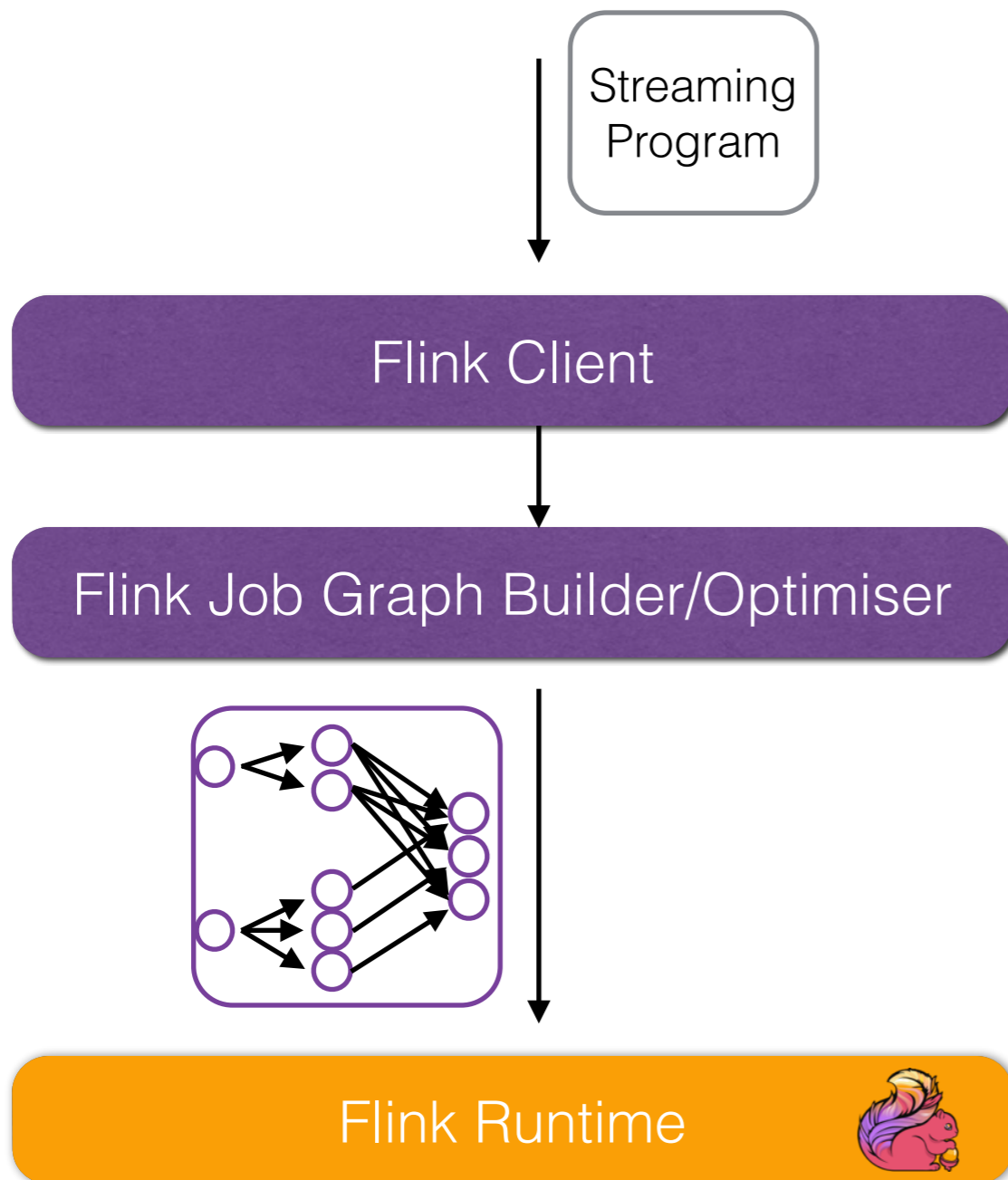


Example: Topology Definition

```
TopologyBuilder builder = new TopologyBuilder()
builder.setSpout("numbers" new NumberGenerator(), 2);
builder.setBolt("new_numbers", new DoubleAndTripleBolt(), 2)
    .shuffleGrouping("numbers");
builder.setBolt("toFile", new DumpToFileBolt(), 1);
    .allGrouping("new_numbers");
```



Standing Queries with Apache Flink



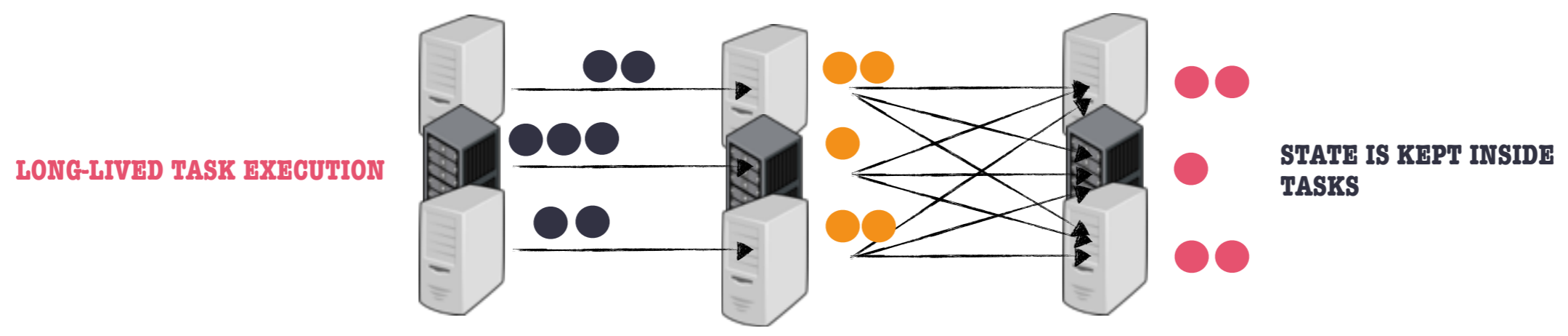
```
val lines: DataStream[String] = env.fromSocketStream(...)
lines.flatMap {line => line.split(" ")
        .map(word => Word(word,1))}
    .window(Time.of(5, SECONDS)) .every(Time.of(1, SECONDS))
    .groupBy("word") .sum("frequency")
    .print()
```

- Operator fusion
- Window Pre-aggregates

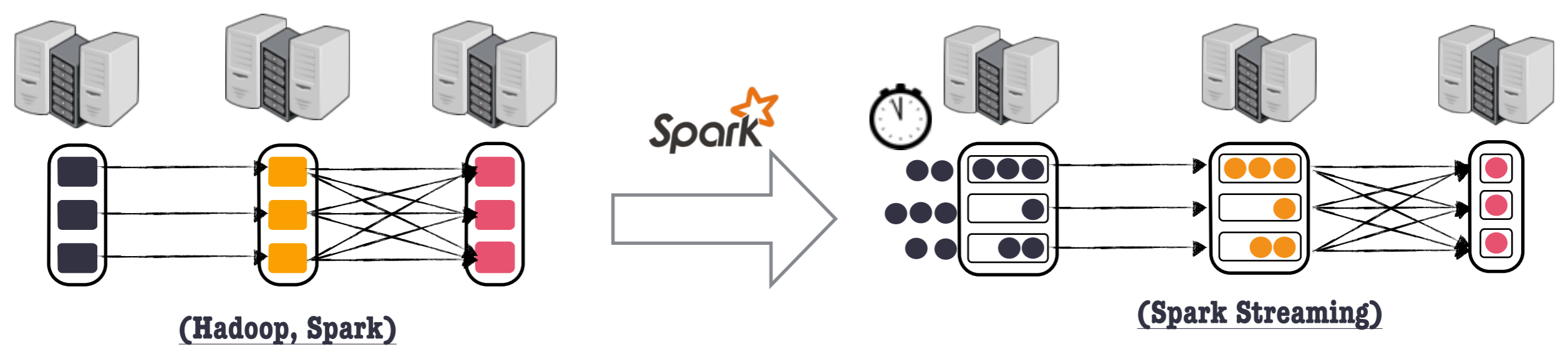
- Deploy Long Running Tasks
- Monitor Execution

Distributed Stream Execution Paradigms

1) Real Streaming (Distributed Data Flow)   



2) Batched Execution



Windows in Action



```
val ssc = new StreamingContext(conf, Seconds(1))
```

```
dstream.window(SECONDS(10), SECONDS(2))
```



```
stream.window(Count.of(100))  
.every(Time.of(10, SECONDS))
```

range

slide

- DStreams are already partitioned in time windows
- Only time windows supported
- Windows decomposed into policies
- Policies can be user-defined too

Windows on Storm?

```
@Override
public void execute(Tuple tuple) {
    if (TupleHelpers.isTickTuple(tuple)) {
        LOG.info("Received tick tuple, triggering emit of current window counts");
        emitCurrentWindowCounts();
    }
    else {
        countObjAndAck(tuple);
    }
}

private void emitCurrentWindowCounts() {
    Map<Object, Long> counts = counter.getCountsThenAdvanceWindow();
    ...
    emit(counts, actualWindowLengthInSeconds);
}

private void emit(Map<Object, Long> counts) {
    for (Entry<Object, Long> entry : counts.entrySet()) {
        Object obj = entry.getKey();
        Long count = entry.getValue();
        collector.emit(new Values(obj, count));
    }
}

private void countObjAndAck(Tuple tuple) {
    Object obj = tuple.getValue(0);
    counter.incrementCount(obj);
    collector.ack(tuple);
}
```

[src-http://www.michael-noll.com/blog/2013/01/18/implementing-real-time-trending-topics-in-storm/](http://www.michael-noll.com/blog/2013/01/18/implementing-real-time-trending-topics-in-storm/)

Partitioning in Action



forward()
shuffle()
broadcast()
keyBy()

partitionCustom()

full control



shuffleGrouping()
allGrouping()
fieldsGrouping()

customGrouping()



repartition(num)
reduceByKey()
updateStateByKey()

no fine-grained control

Synopses in Action

implementing a rolling max per key



```
public static class MaxPerKey extends BaseRichBolt {
    private int Map<Integer,Integer> maxVals = ...
    private OutputCollectorBase _collector;
    @Override
    public void execute(Tuple tuple) {
        int key = tuple.getInteger(1);
        int _max = maxVals.containsKey(key)
            ? maxVals.get(tuple.getInteger(1)) : 0;
        if(_max < tuple.getInteger(2)){
            _max = tuple.getInteger(2);
            maxVals.put(tuple.getInteger(1), _max)
            _collector.emit(tuple, _max);
        }
    }
}
```

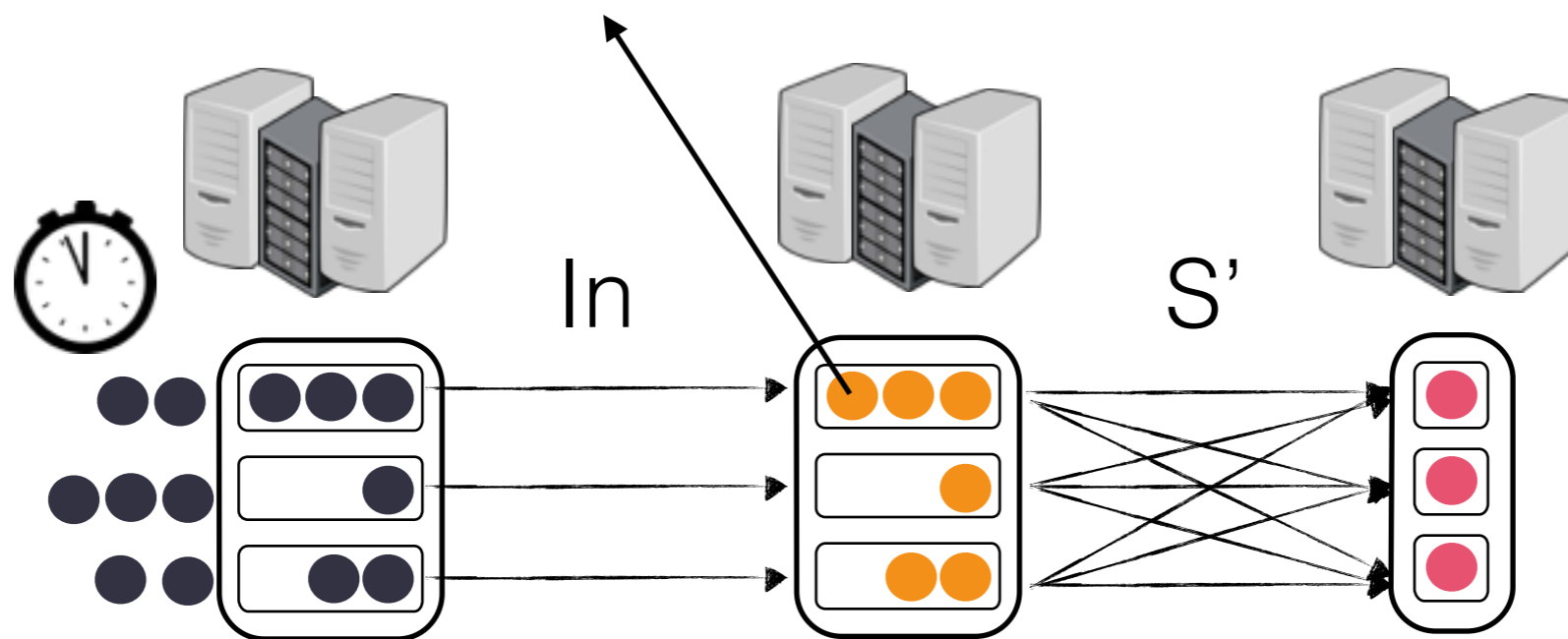


```
stream.keyBy("key").flatMapWithState(
    (in, state: Option[Int]) => state match {
        | | | case Some(val) => if (val < in.val) (List(in), in.val)
        | | | else (List(),Some(state))
        | | | case None => (List(in.val), Some(in.val))
    })
```

State in Spark?

- Streams are partitioned into small batches
- There is practically no state kept in workers (stateless)
- How do we keep state??

`dstream.updateStateByKey(...)` *put new states in output RDD*



(Spark Streaming)

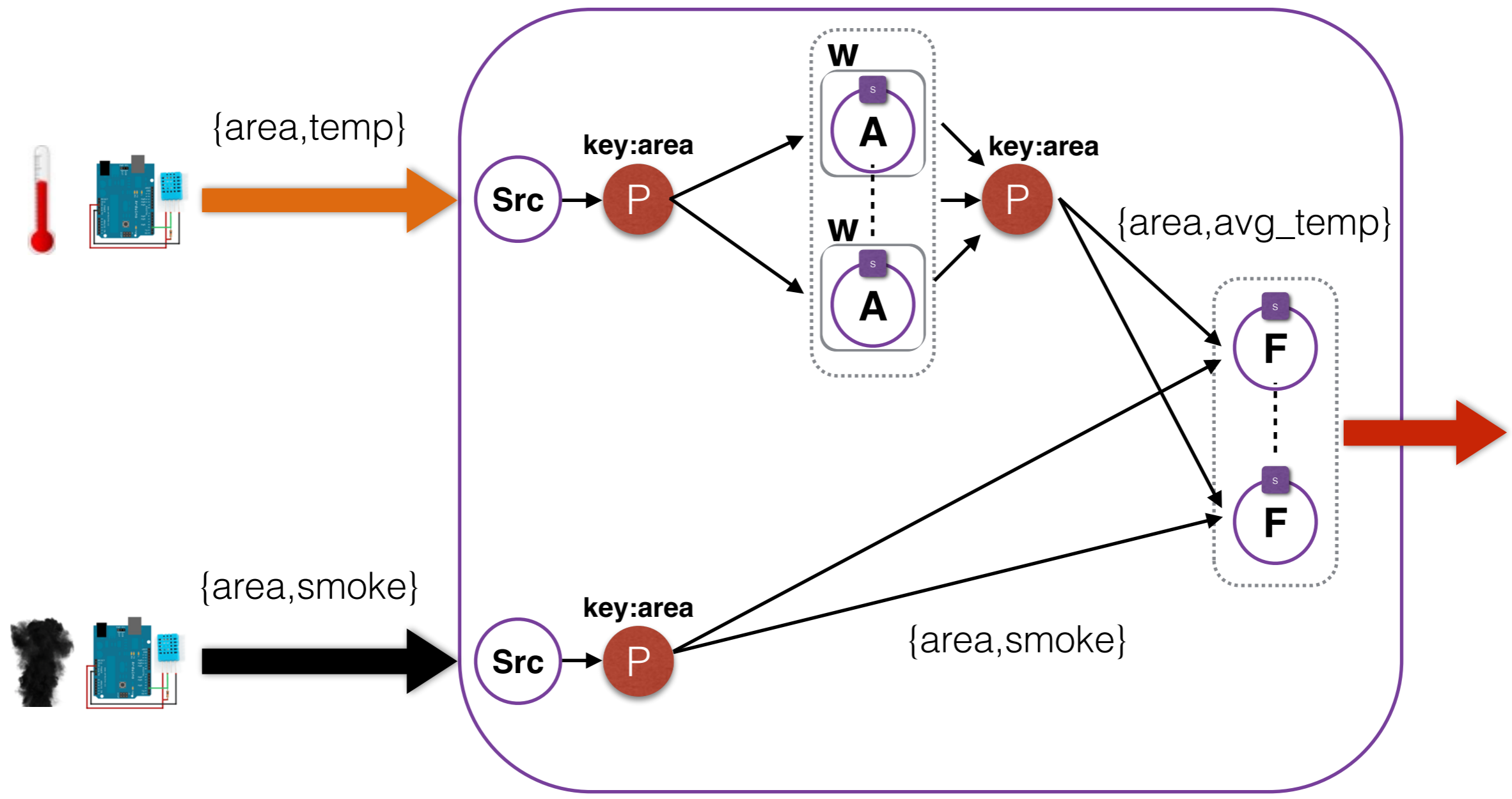
Implementing the alarm in Flink

```
val temperatures = env.socketTextStream(...).keyBy("area")
val smokes = env.socketTextStream(...).keyBy("area")

val avgTemp = temperatures
| .window(Time.of(60, SECONDS)
| .every(Time.of(5, SECONDS)
| .mapWindow(_avgTemp).flatten().keyBy("area")

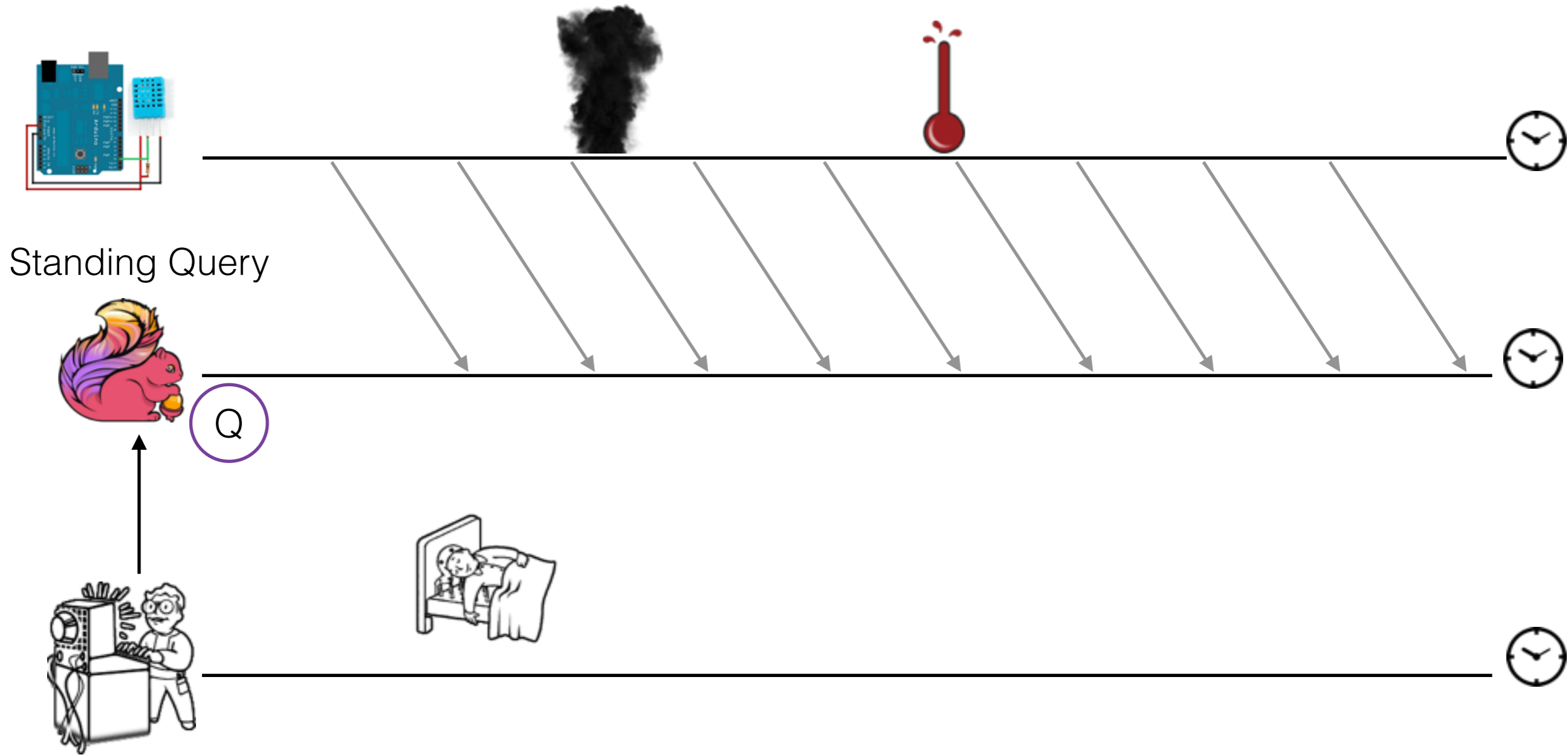
avgTemp.connect(smokes).flatMap(sm_alarm).print()
```

So everything works

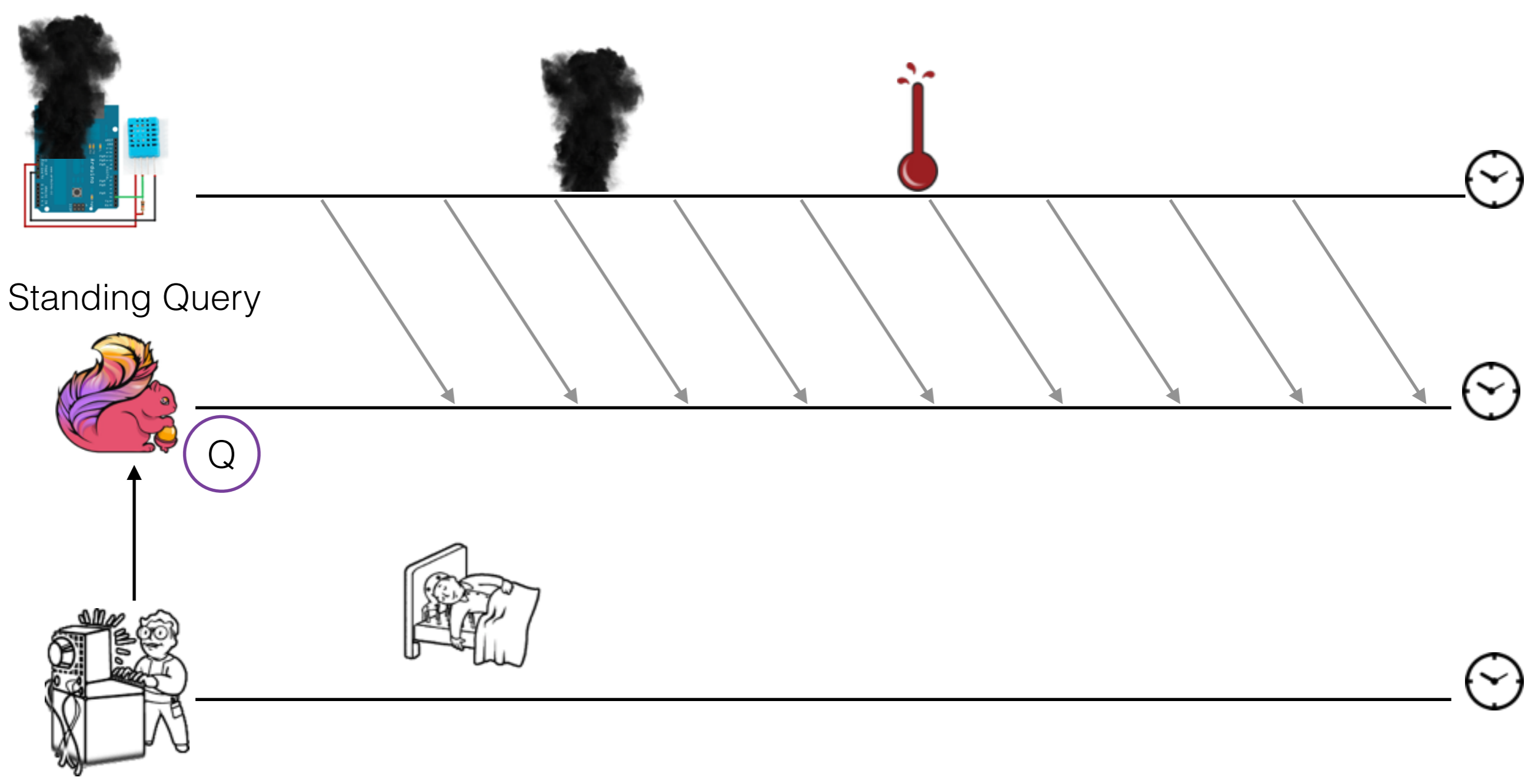


or...

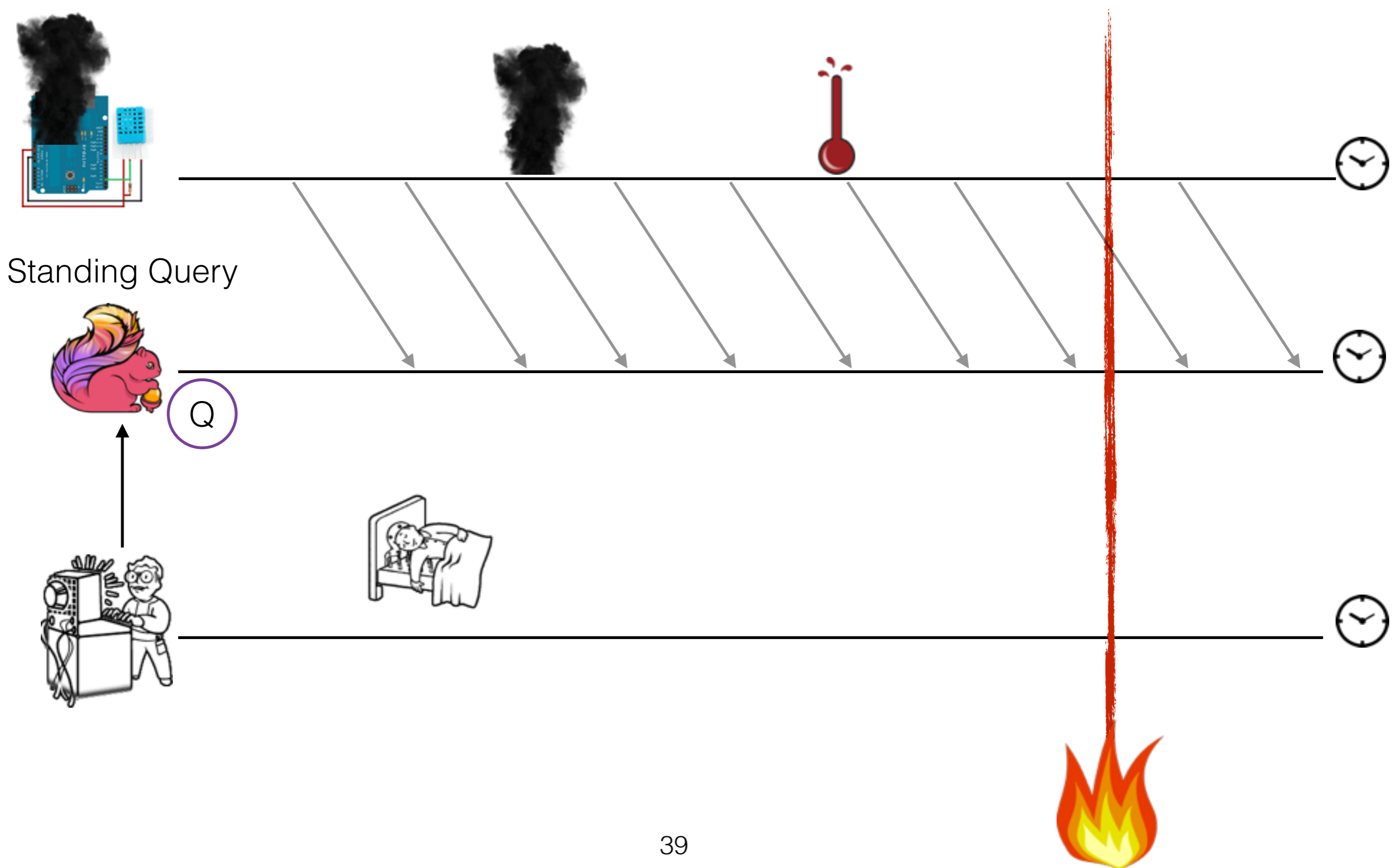
Unreliable Sources



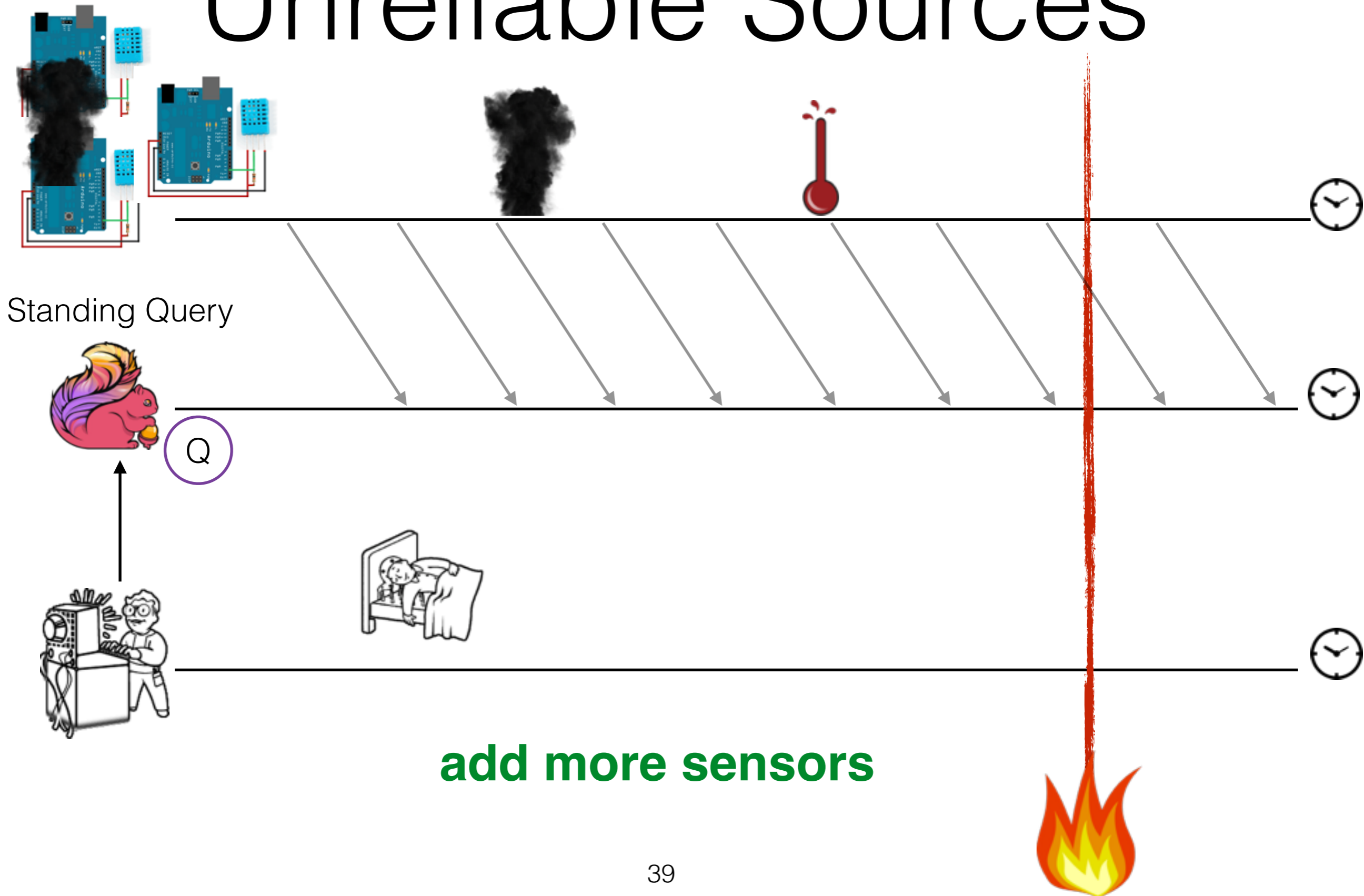
Unreliable Sources



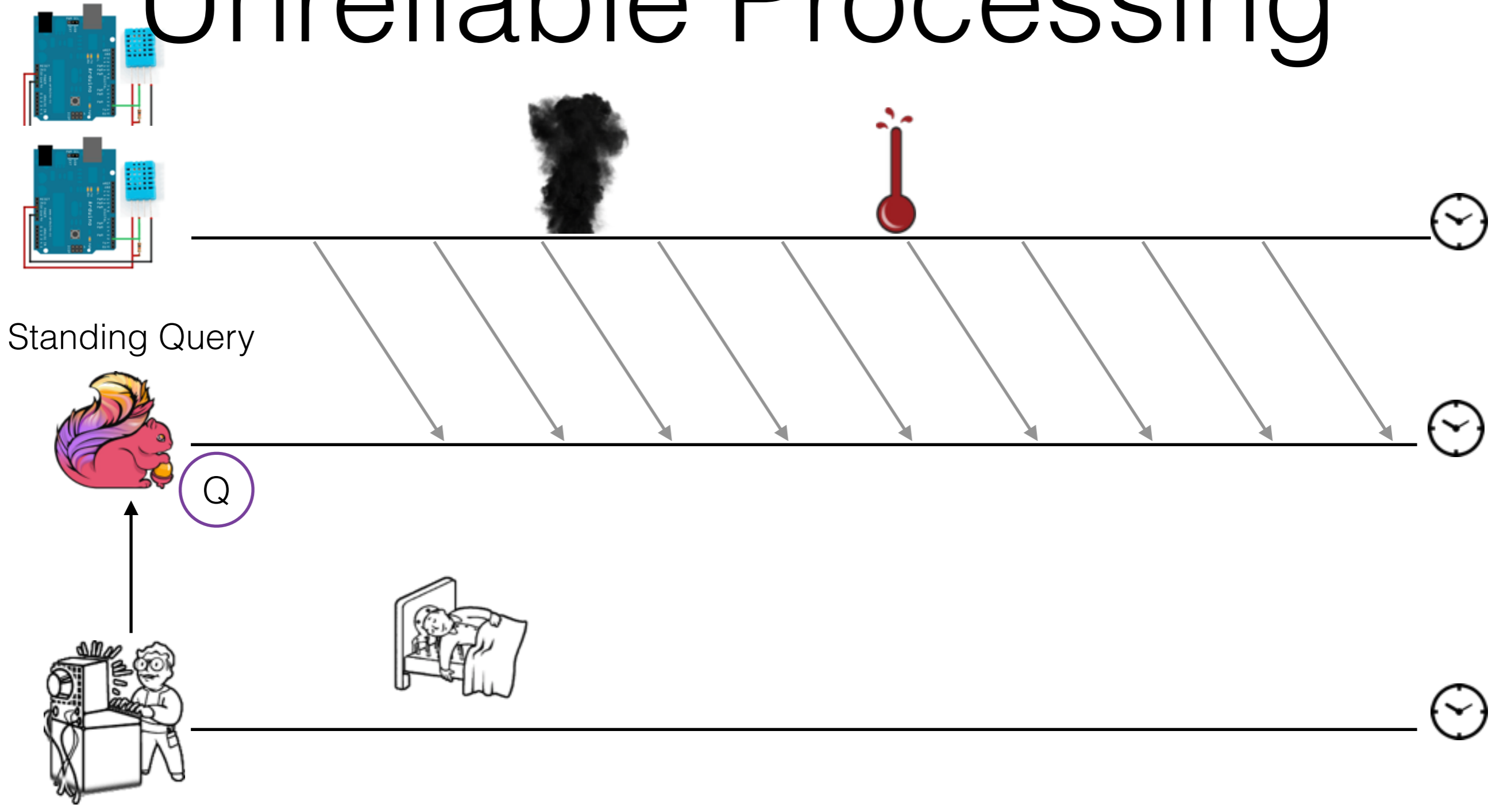
Unreliable Sources



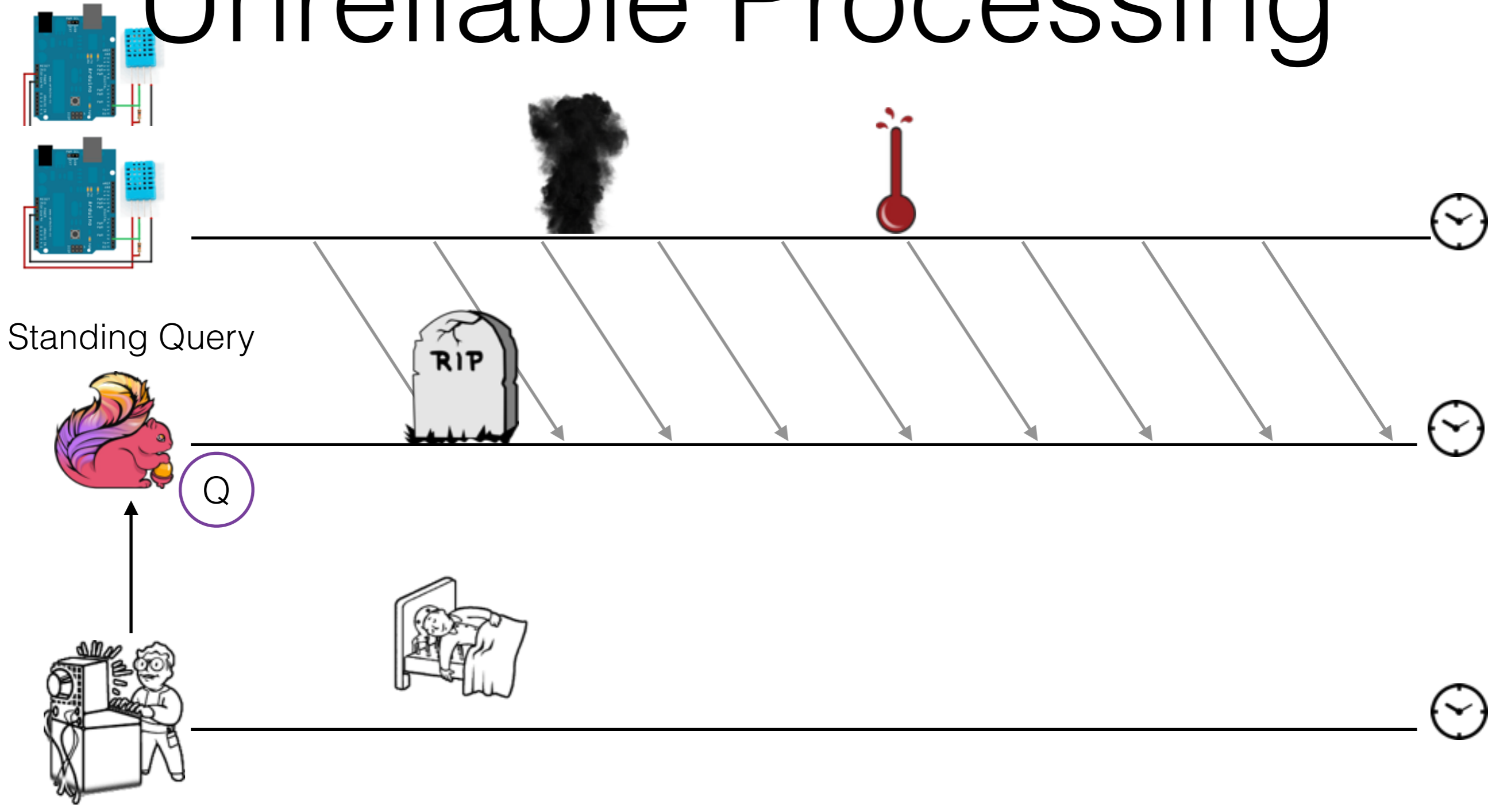
Unreliable Sources



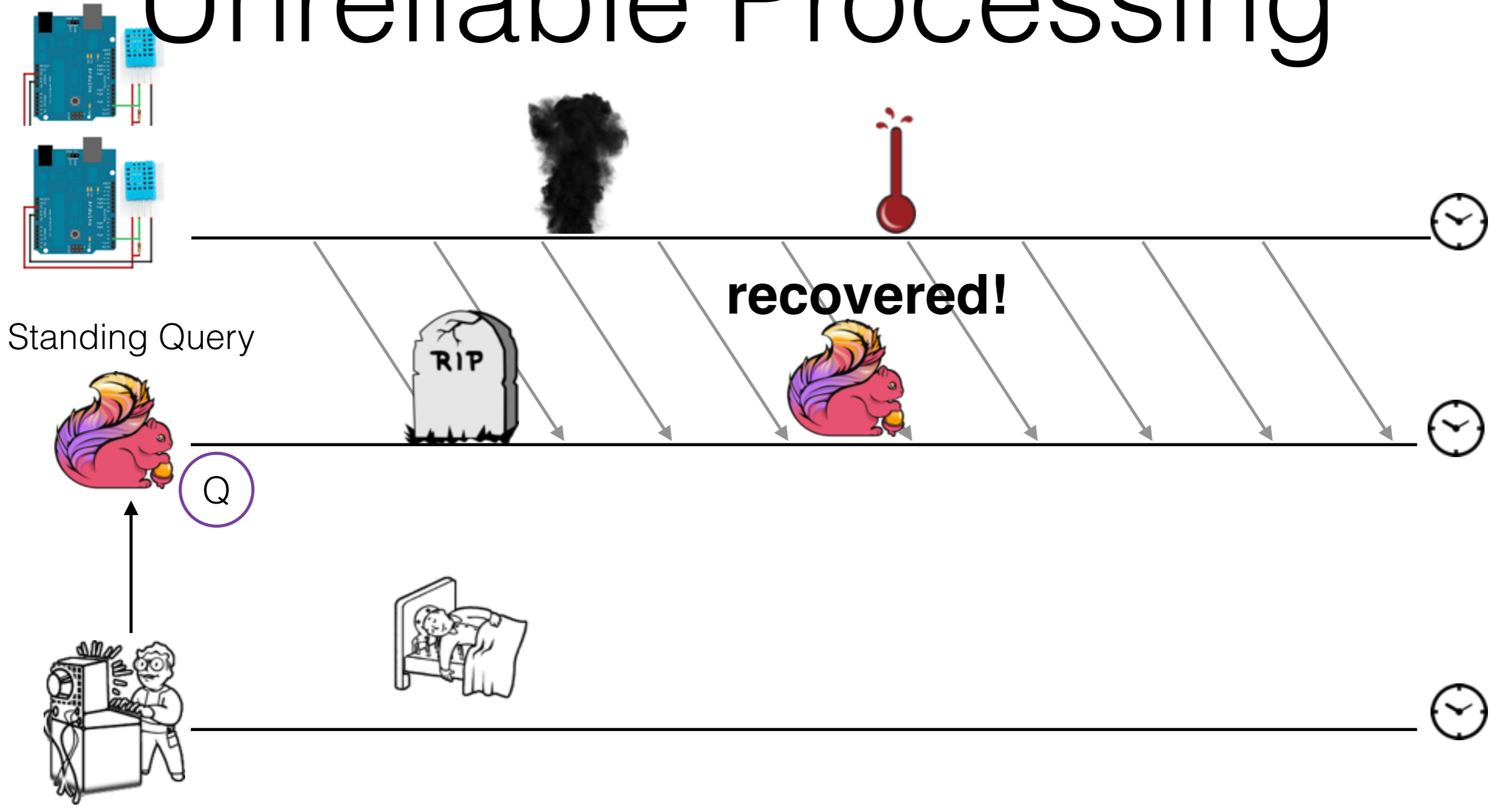
Unreliable Processing



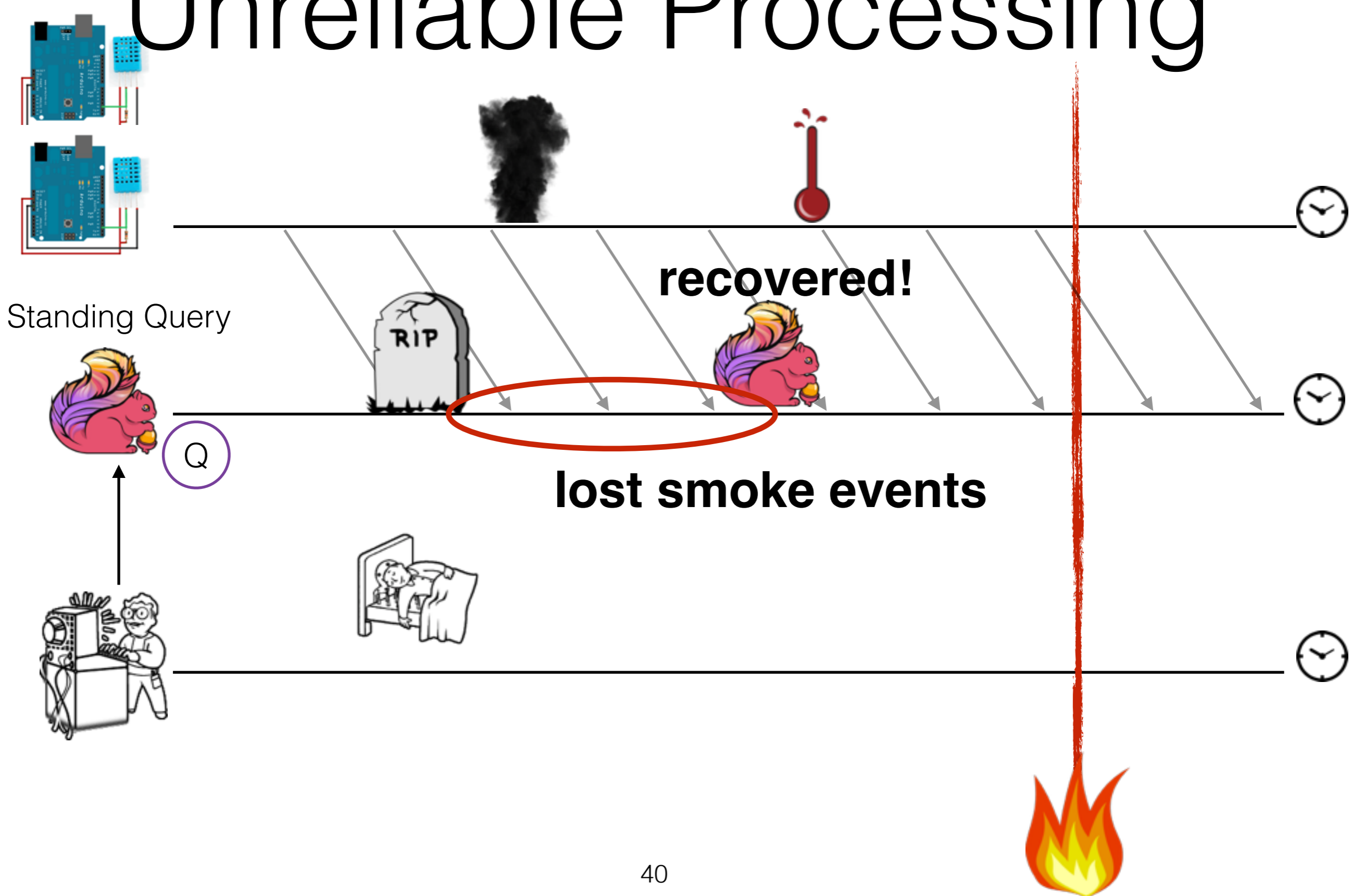
Unreliable Processing



Unreliable Processing



Unreliable Processing



Resilient Brokers



Main Features

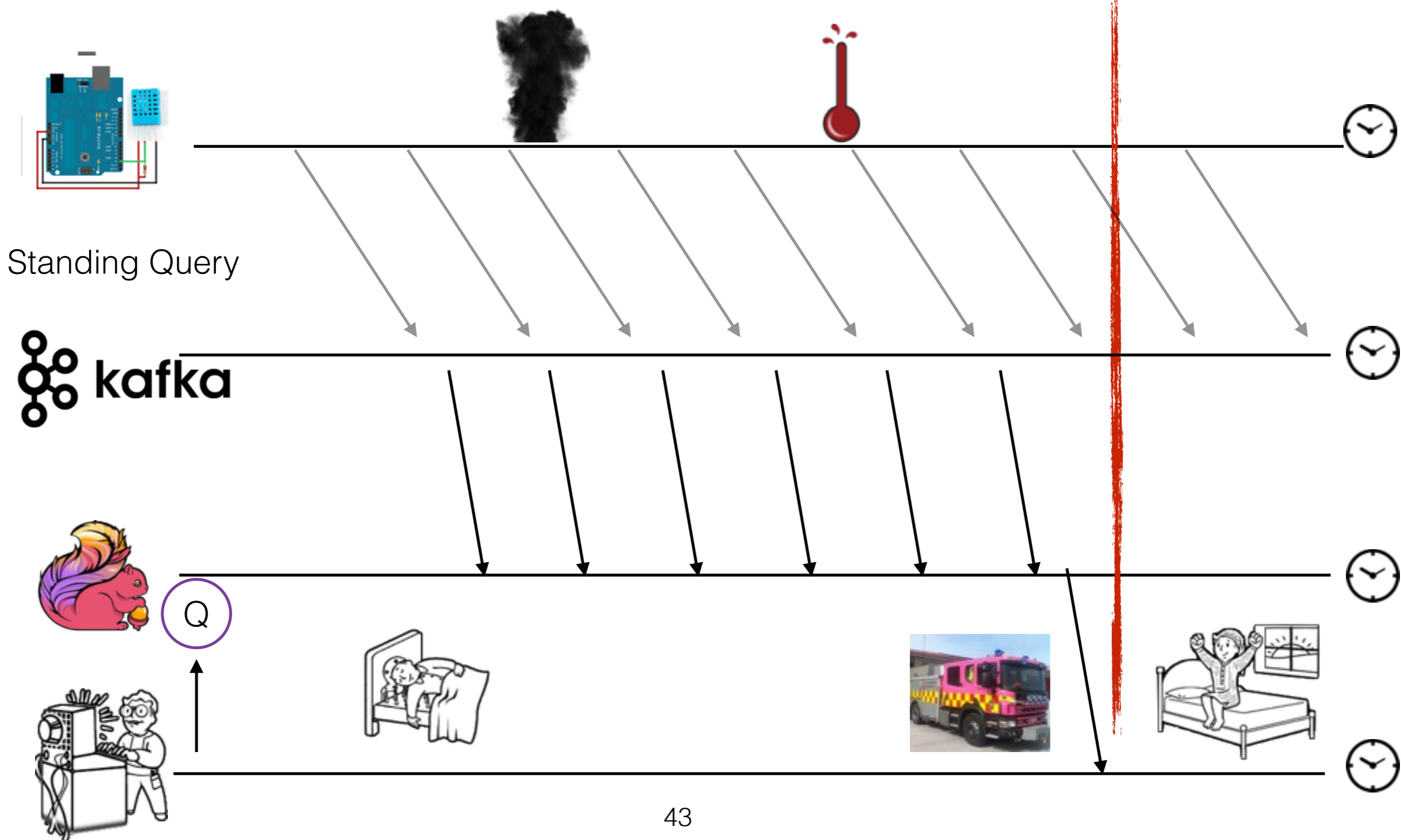
- Topic-based partitioned queues
- Strongly consistent offset mapping to records

Processing Guarantees

- Kafka solves the source consistency problem
- How about the rest of the states of the computation ? (e.g. alert operator state)
- Each system offers different guarantees

	Guarantees	Technique
Storm	at least once	event dependency tracking
Spark	exactly once	source upstream backup
Flink	exactly once	periodic snapshots

Mission Accomplished



Research Topics at KTH/SICS

- Exactly-Once-Output Guarantees
- State management and auto-scaling
- Streaming ML pipelines
- Streaming Graphs