

Storm

The Real-Time Layer Your Big Data's Been Missing

Dan Lynn
dan@fullcontact.com
[@danklynn](https://twitter.com/danklynn)



FullContact

Keeps Contact Information Current and Complete

Based in Denver, Colorado



Advisors to FullContact

include some of the leading
pioneers in internet startups.



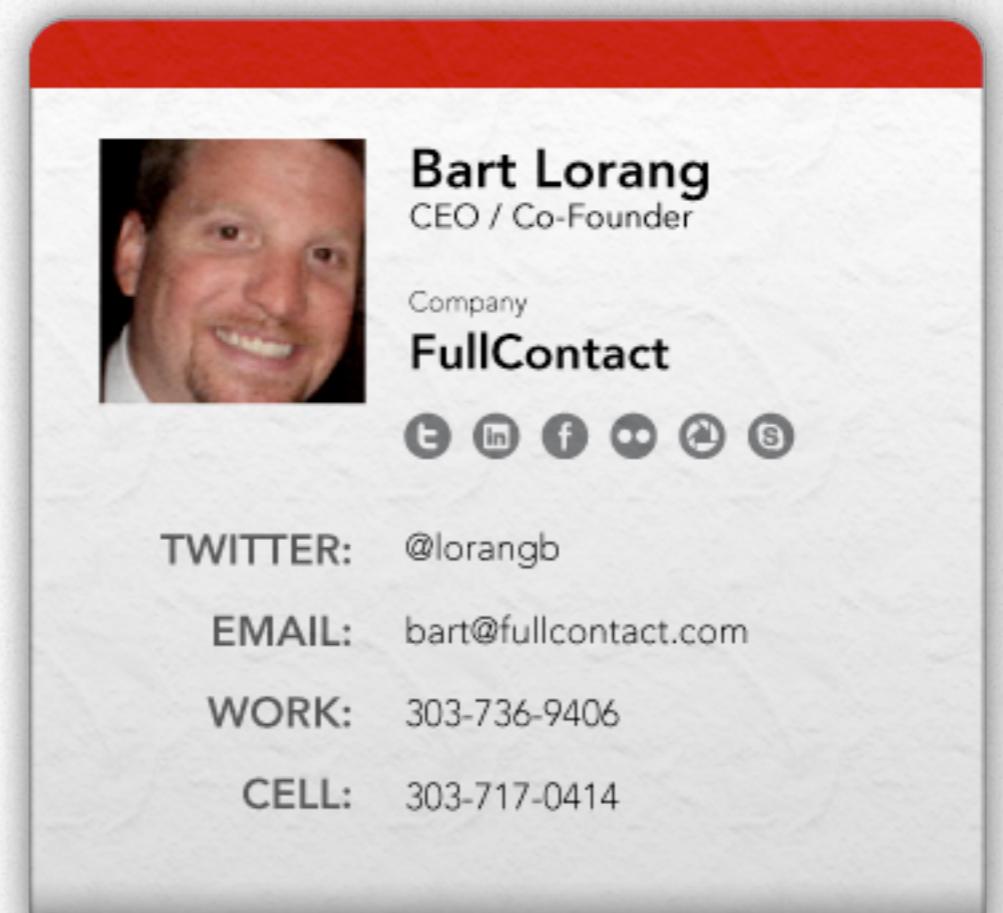
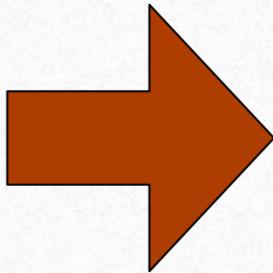
FullContact is

a TechStars 2011 Boulder
Graduate.



CTO & Co-Founder
dan@fullcontact.com
[@danklynn](https://twitter.com/danklynn)

Turn Partial Contacts Into **Full** Contacts



Your data is old

Old data is confusing

First, there was the spreadsheet

Next, there was SQL

Then, there wasn't SQL

Batch Processing



Stale Data

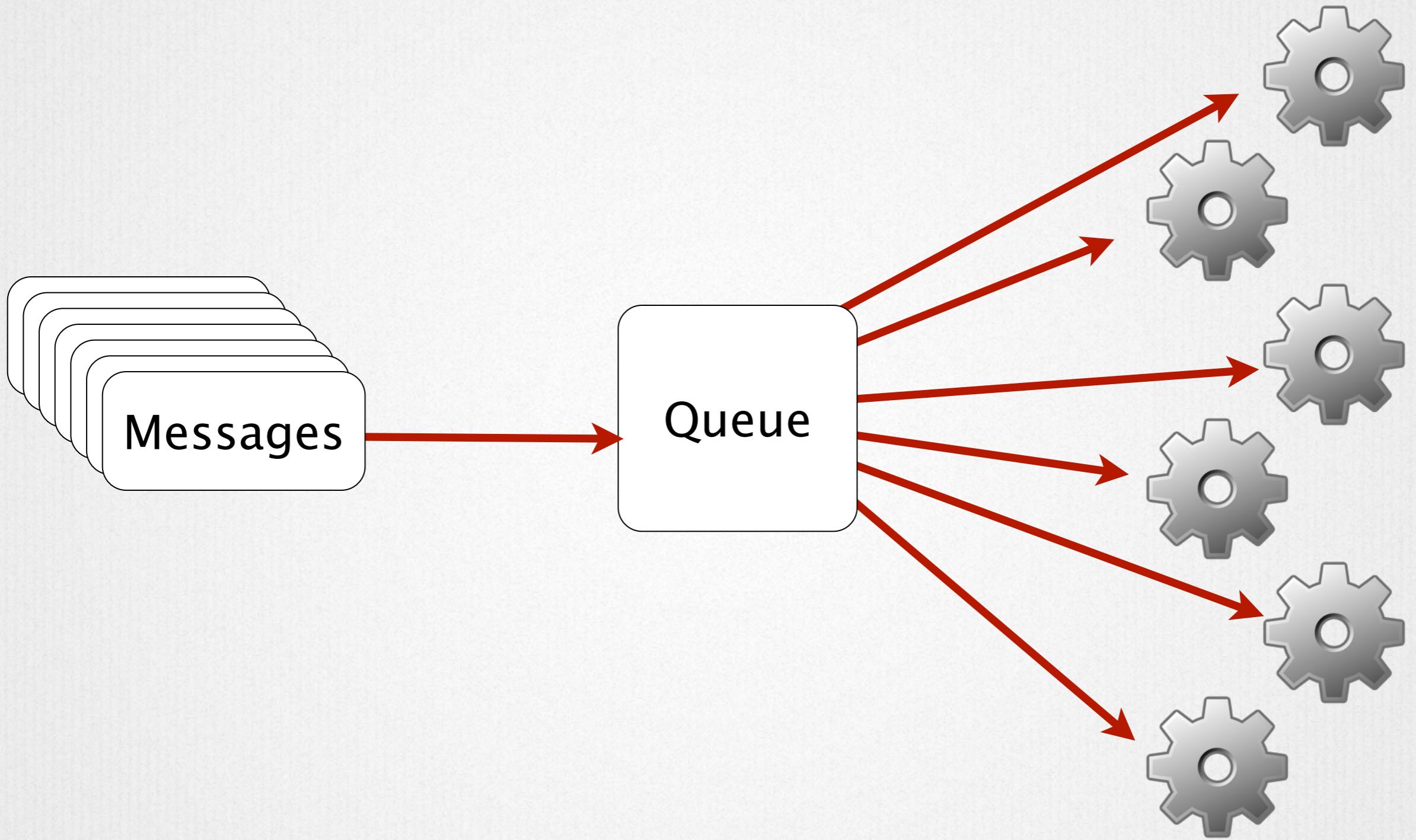


FullContact

Streaming Computation

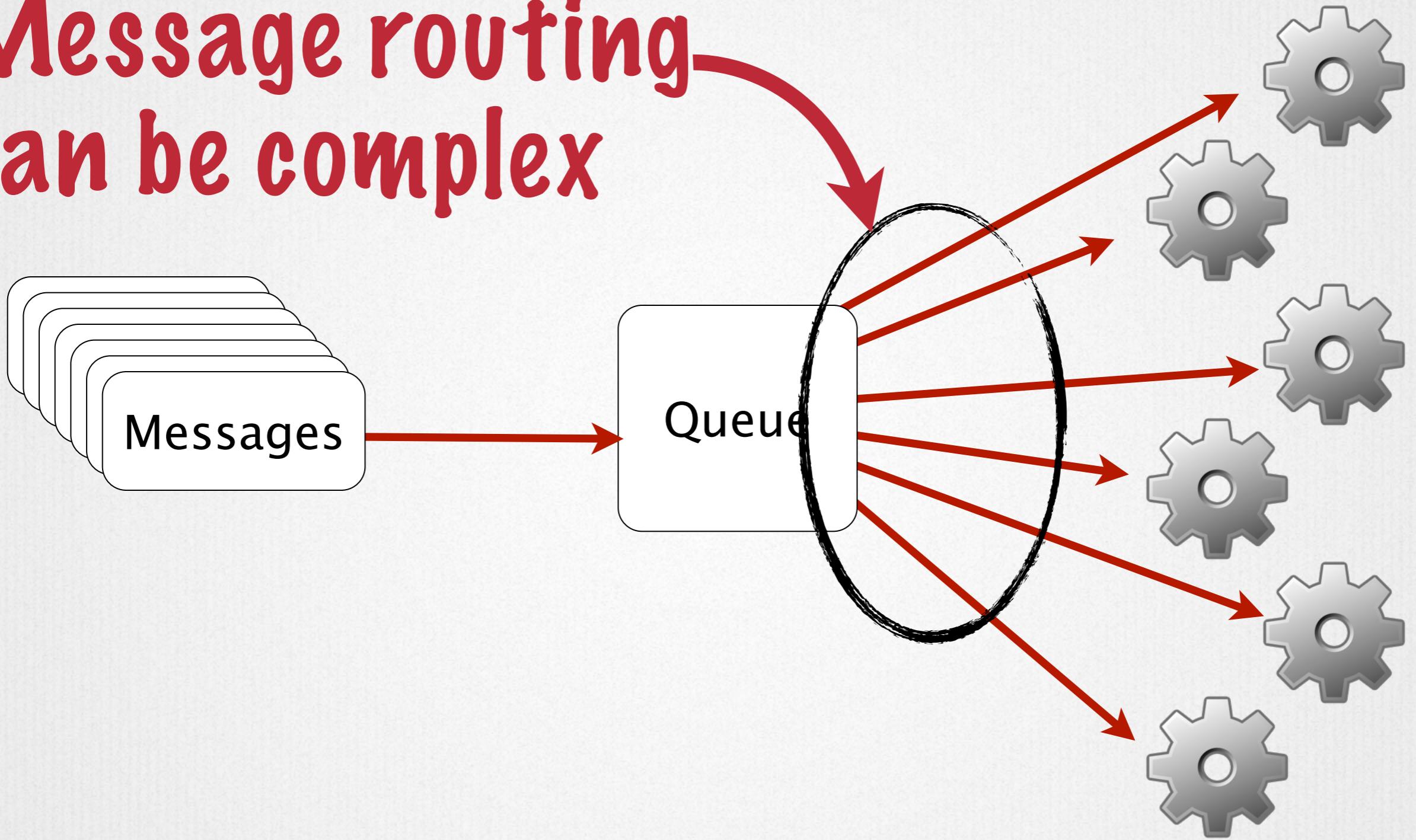
Queues / Workers

Queues / Workers

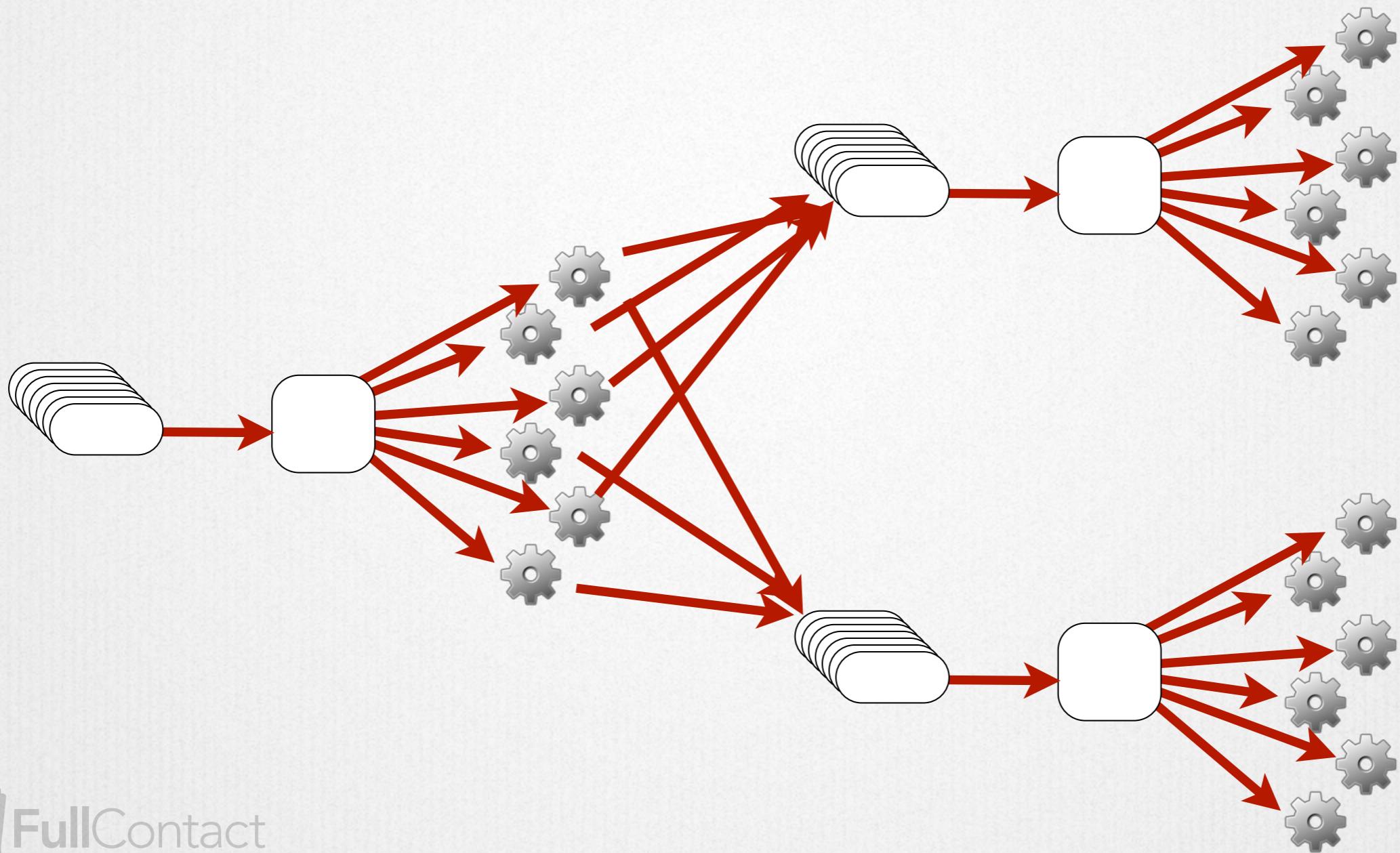


Queues / Workers

Message routing
can be complex



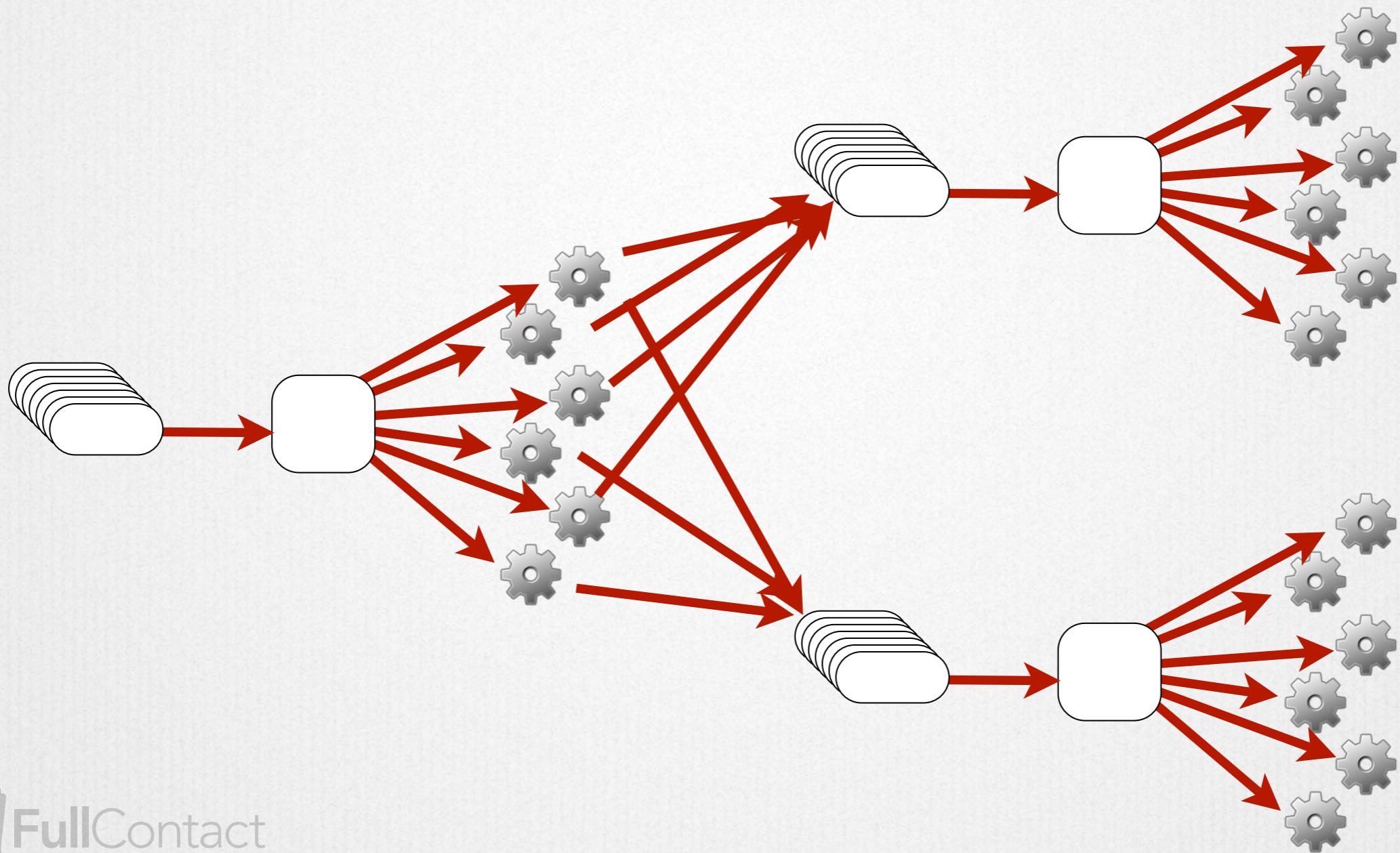
Queues / Workers



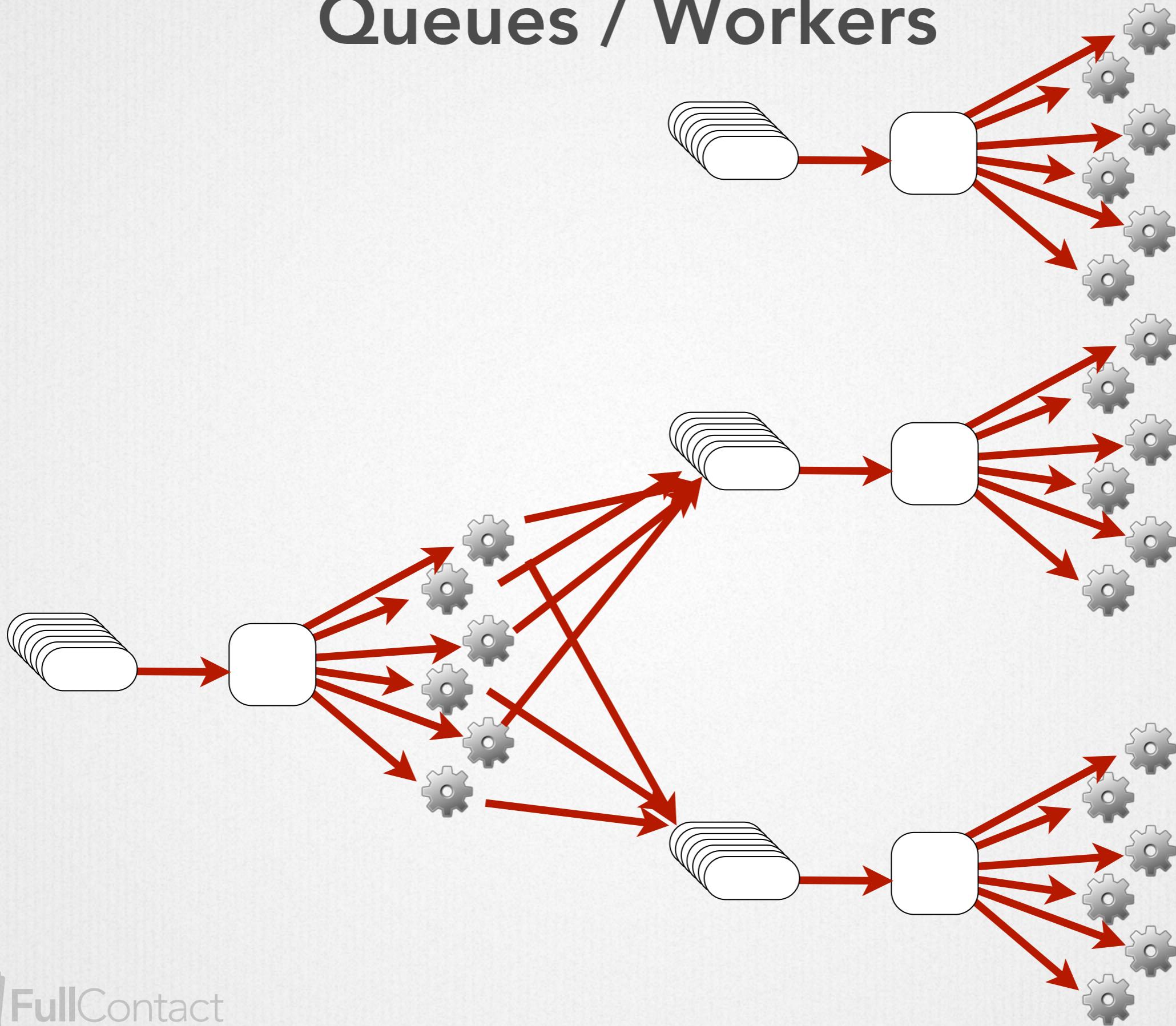
Brittle

Hard to scale

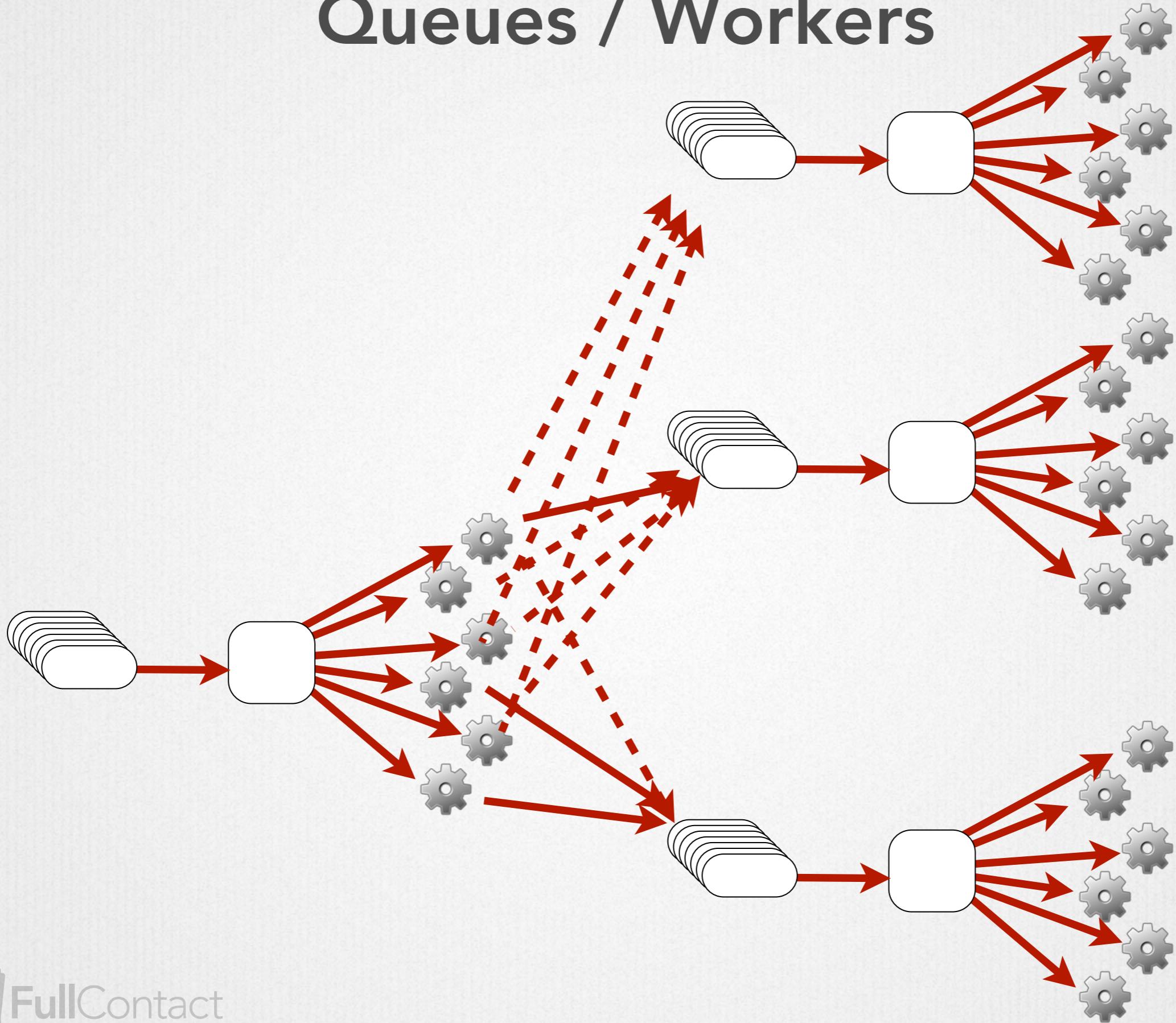
Queues / Workers



Queues / Workers

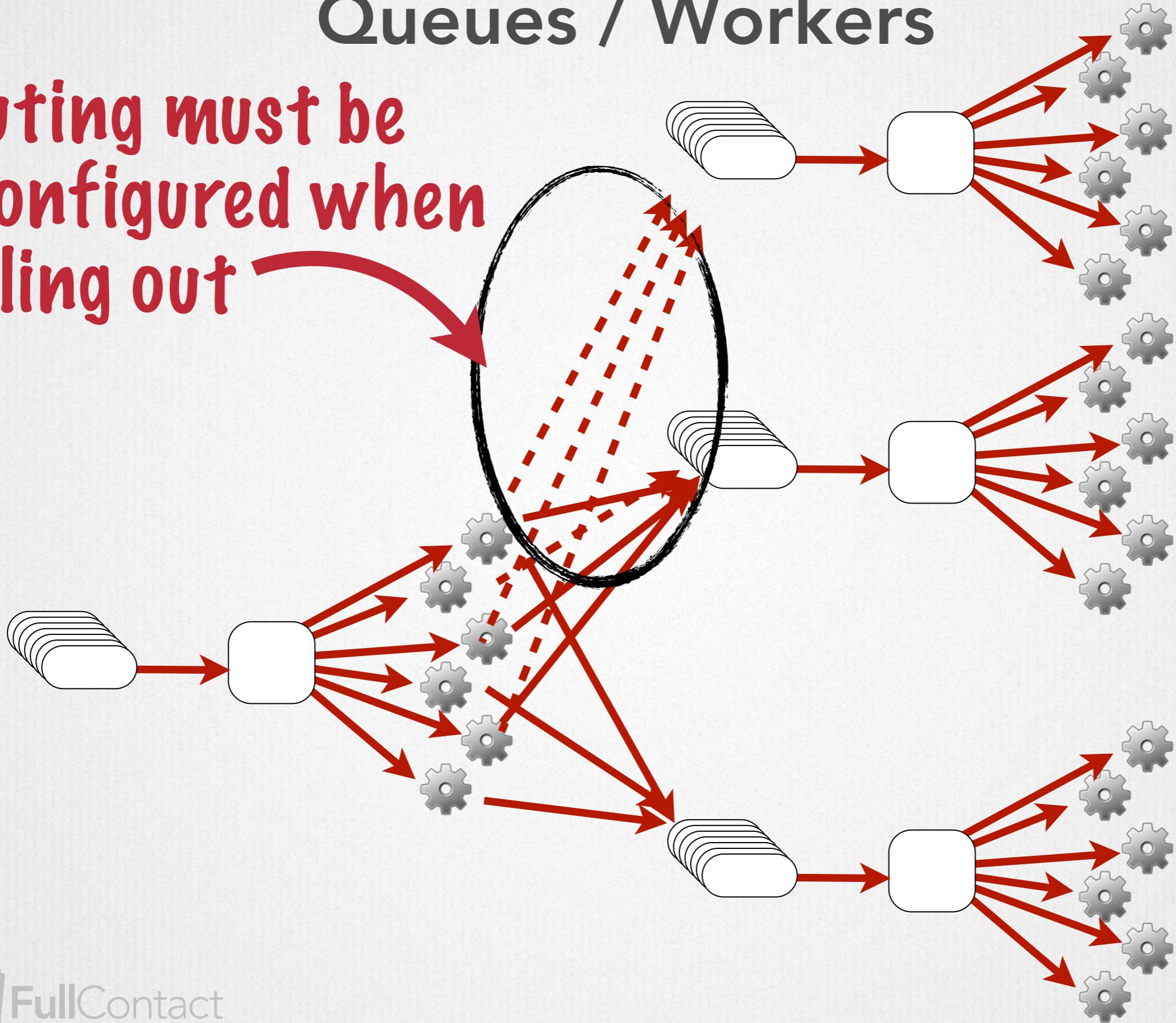


Queues / Workers



Queues / Workers

Routing must be
reconfigured when
scaling out



Storm

Storm

Distributed and fault-tolerant real-time computation



Storm

Distributed and fault-tolerant real-time computation

Storm

Distributed and fault-tolerant real-time computation

Storm

Distributed and fault-tolerant real-time computation

Key Concepts

Tuples

Ordered list of elements

Tuples

Ordered list of elements

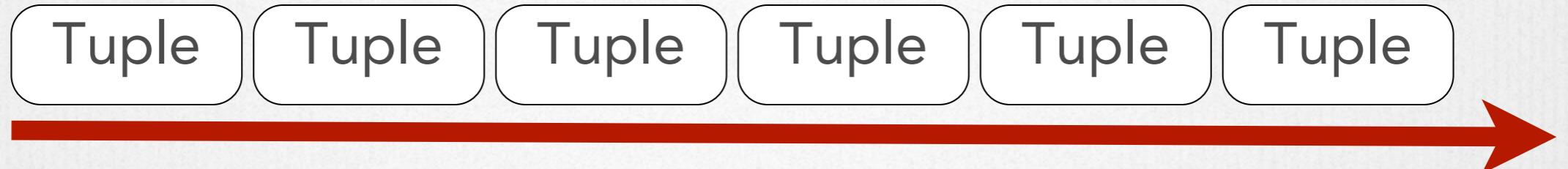
```
("search-01384", "e:dan@fullcontact.com")
```

Streams

Unbounded sequence of tuples

Streams

Unbounded sequence of tuples



Spouts

Source of streams



FullContact

Spouts

Source of streams



FullContact

Spouts

Source of streams



Tuple Tuple Tuple Tuple Tuple Tuple



FullContact

Spouts can talk with

Spouts can talk with

- Queues

Spouts can talk with

- Queues
- Web logs

Spouts can talk with

- Queues
- Web logs
- API calls

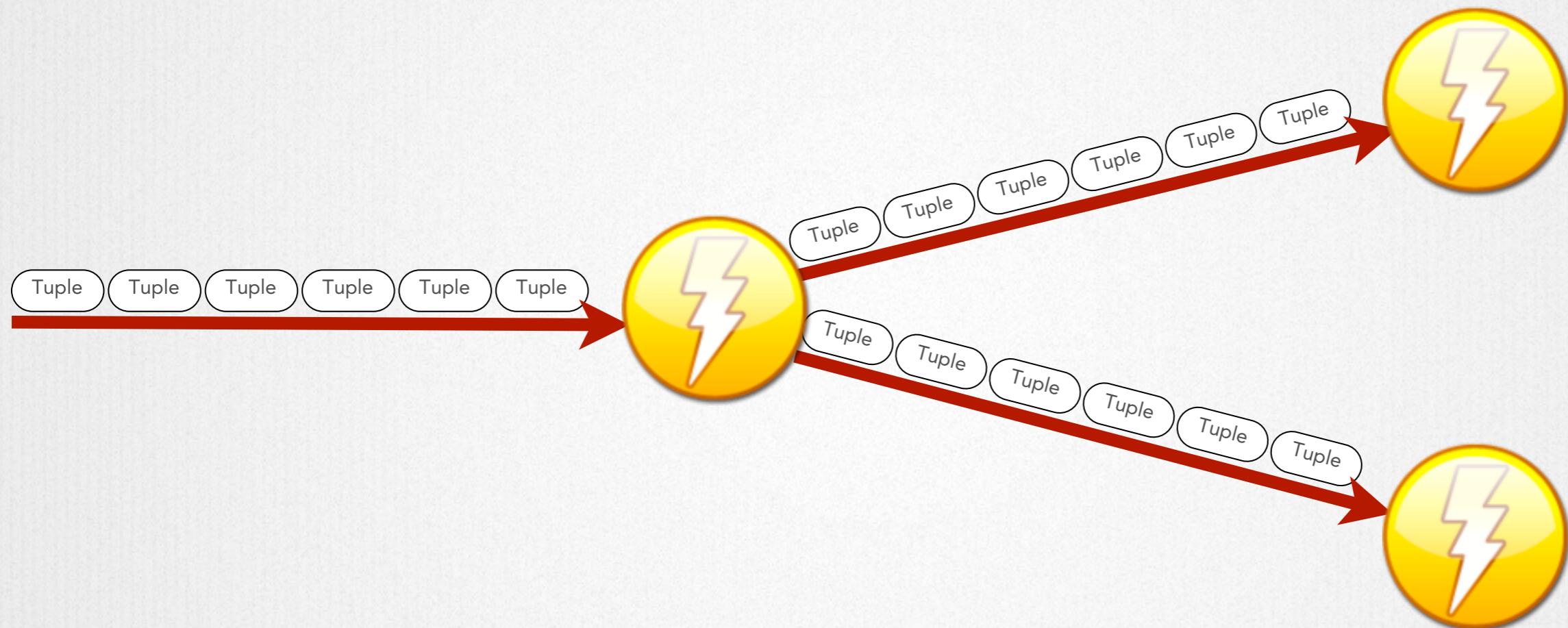
Spouts can talk with

- Queues
- Web logs
- API calls
- Event data

Bolts

Process tuples and create new streams

Bolts



Bolts

Bolts

- Apply functions / transforms

Bolts

- Apply functions / transforms
- Filter

Bolts

- Apply functions / transforms
- Filter
- Aggregation

Bolts

- Apply functions / transforms
- Filter
- Aggregation
- Streaming joins

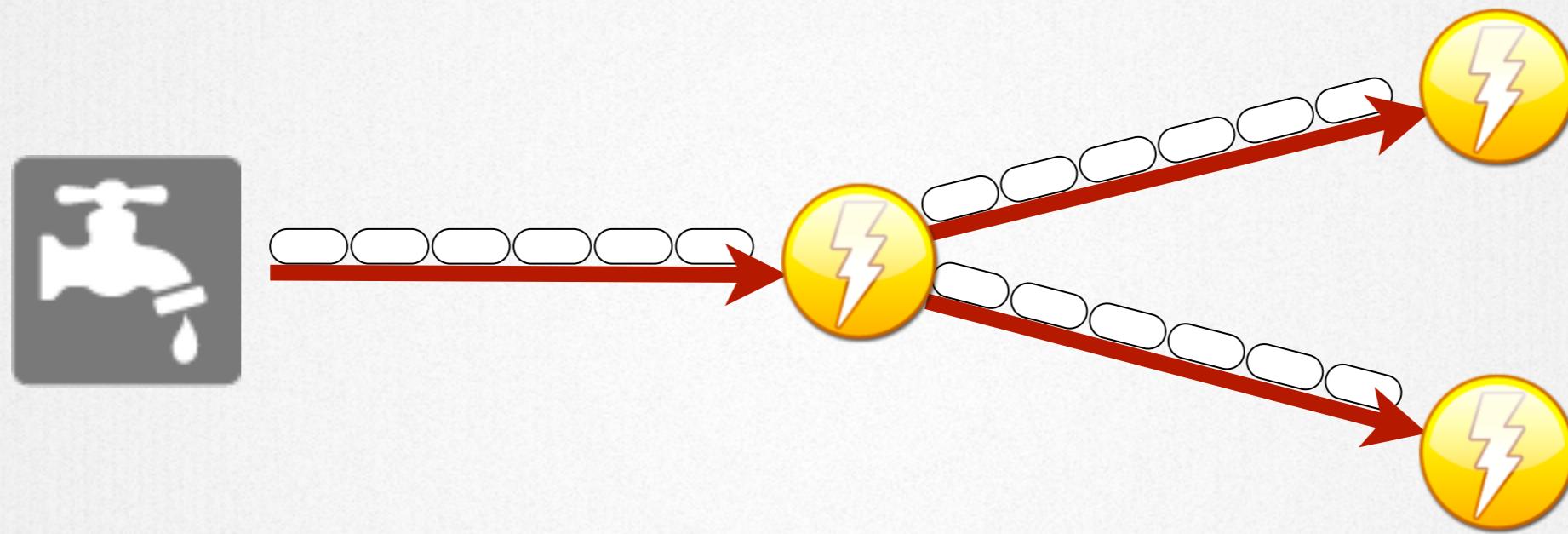
Bolts

- Apply functions / transforms
- Filter
- Aggregation
- Streaming joins
- Access DBs, APIs, etc...

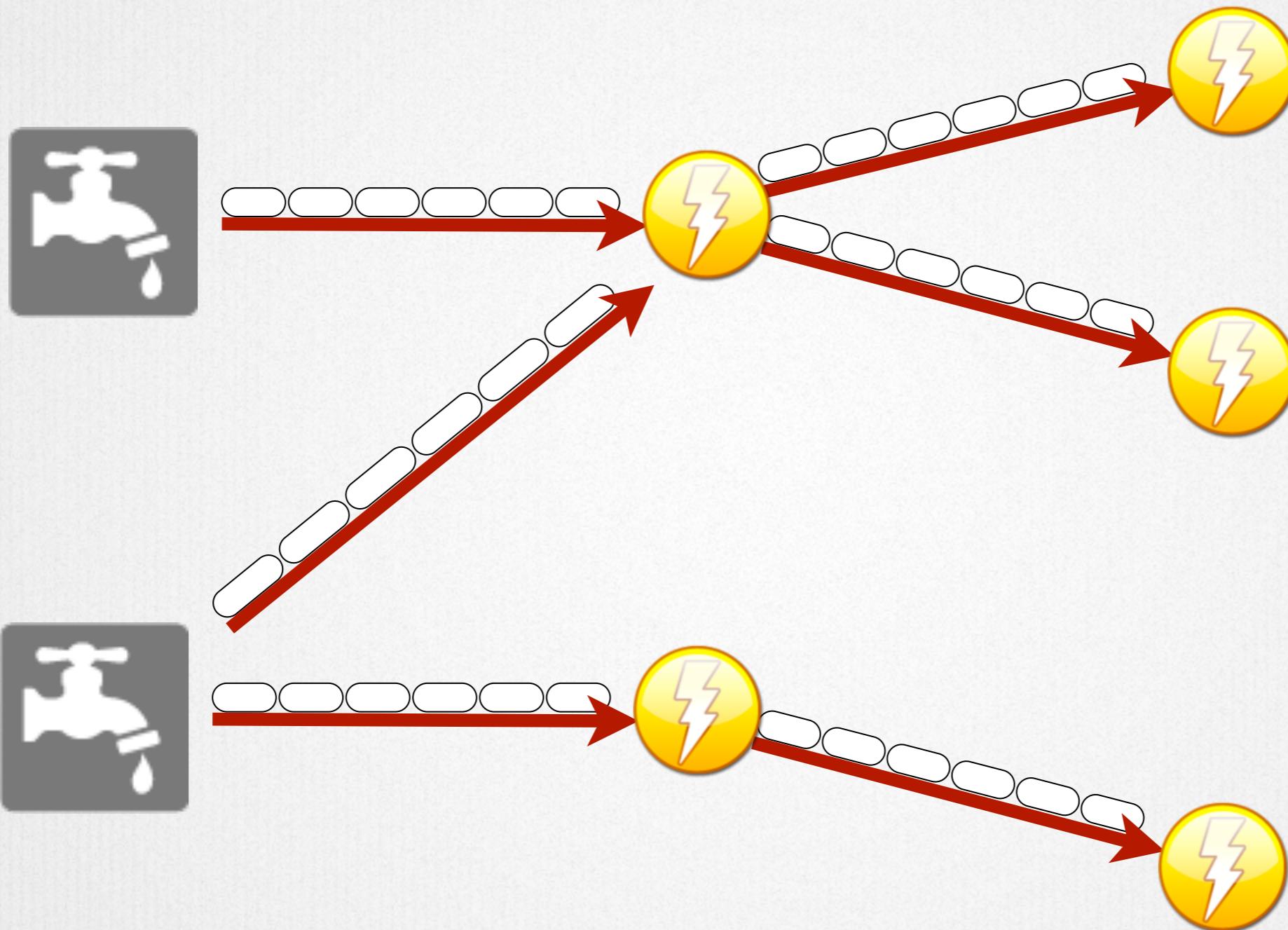
Topologies

A directed graph of Spouts and Bolts

This is a Topology



This is also a topology



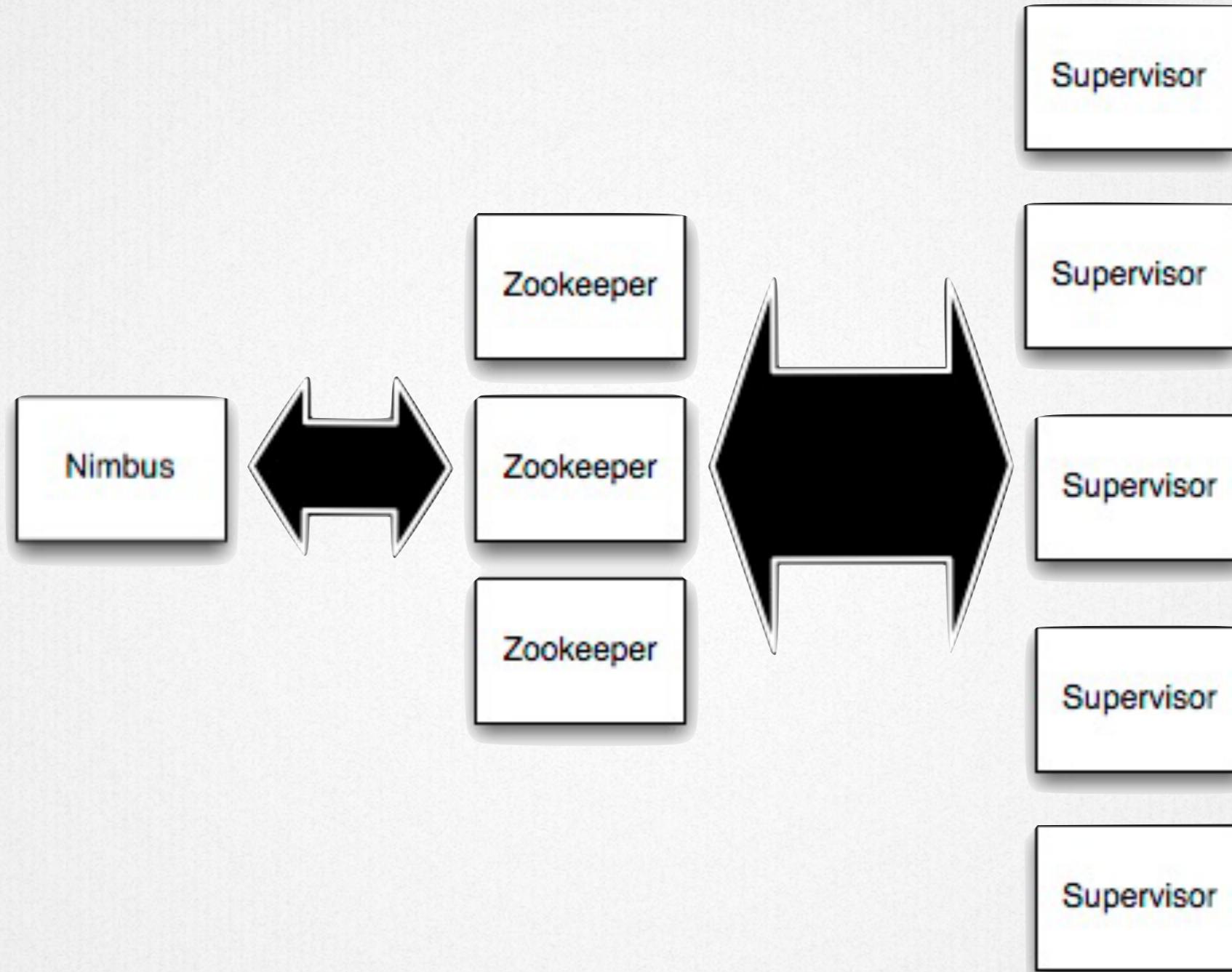
Tasks

Processes which execute Streams or Bolts

Running a Topology

```
$ storm jar my-code.jar com.example.MyTopology arg1 arg2
```

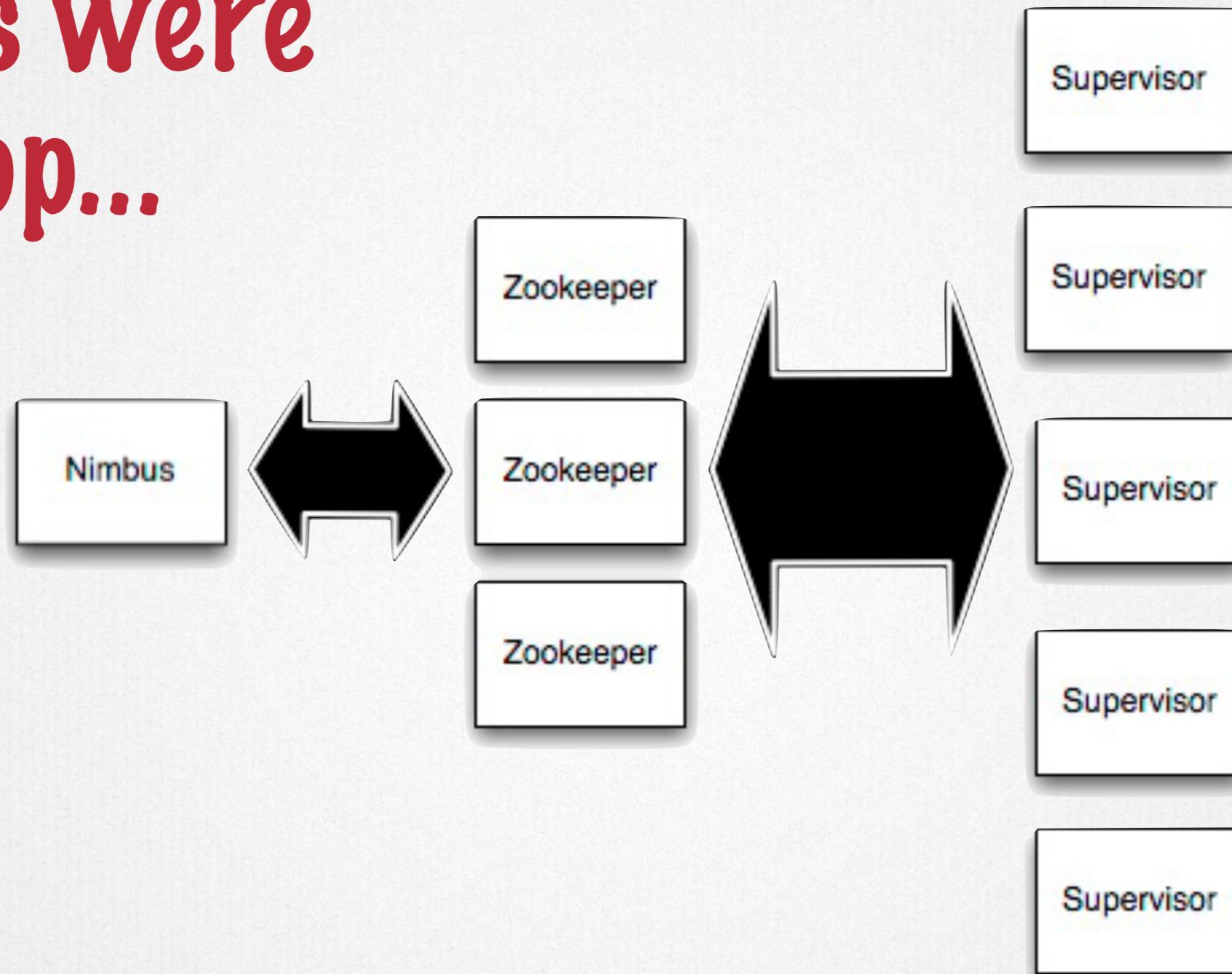
Storm Cluster



Nathan Marz

Storm Cluster

If this were
Hadoop...

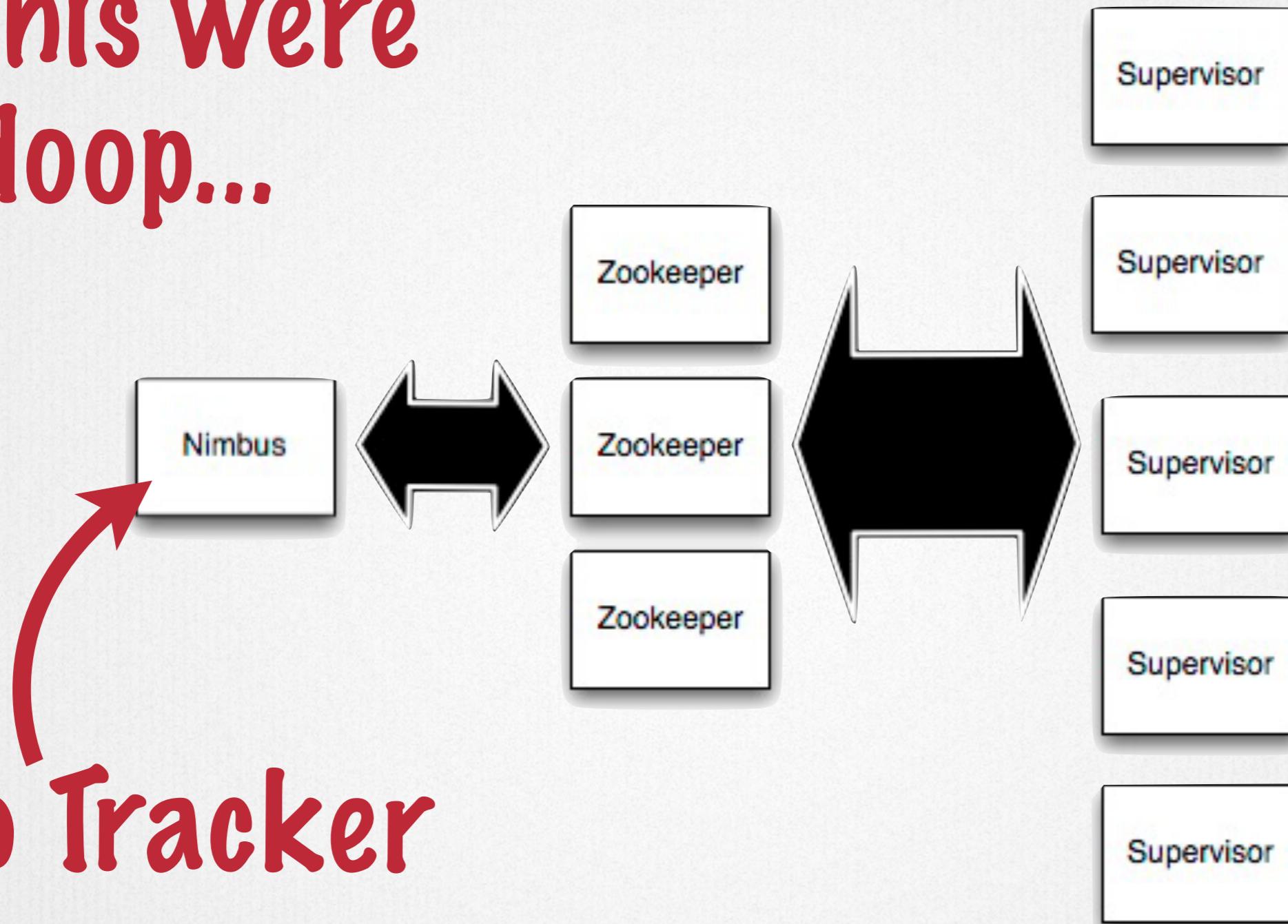


Nathan Marz

Storm Cluster

If this were
Hadoop...

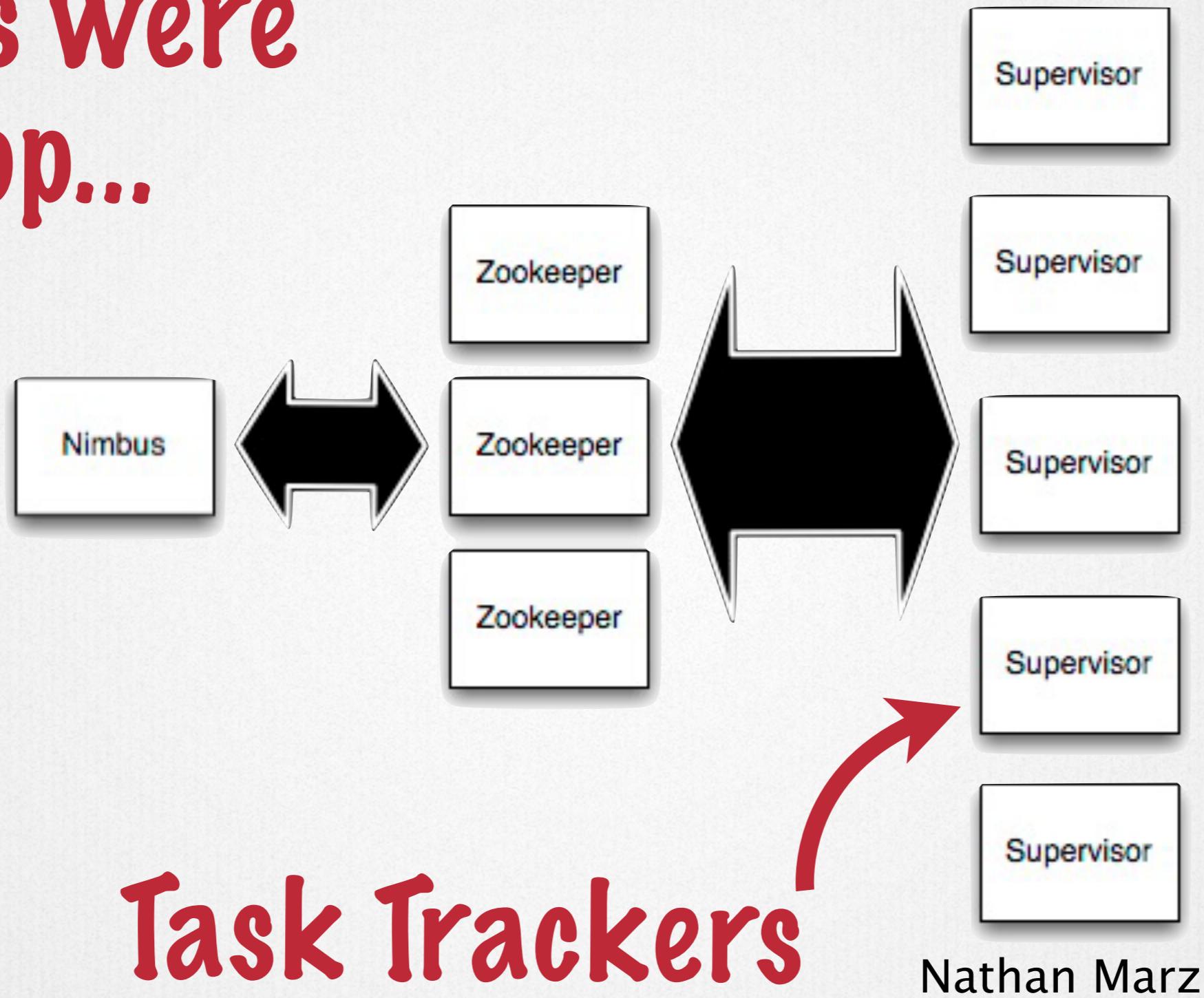
Job Tracker



Nathan Marz

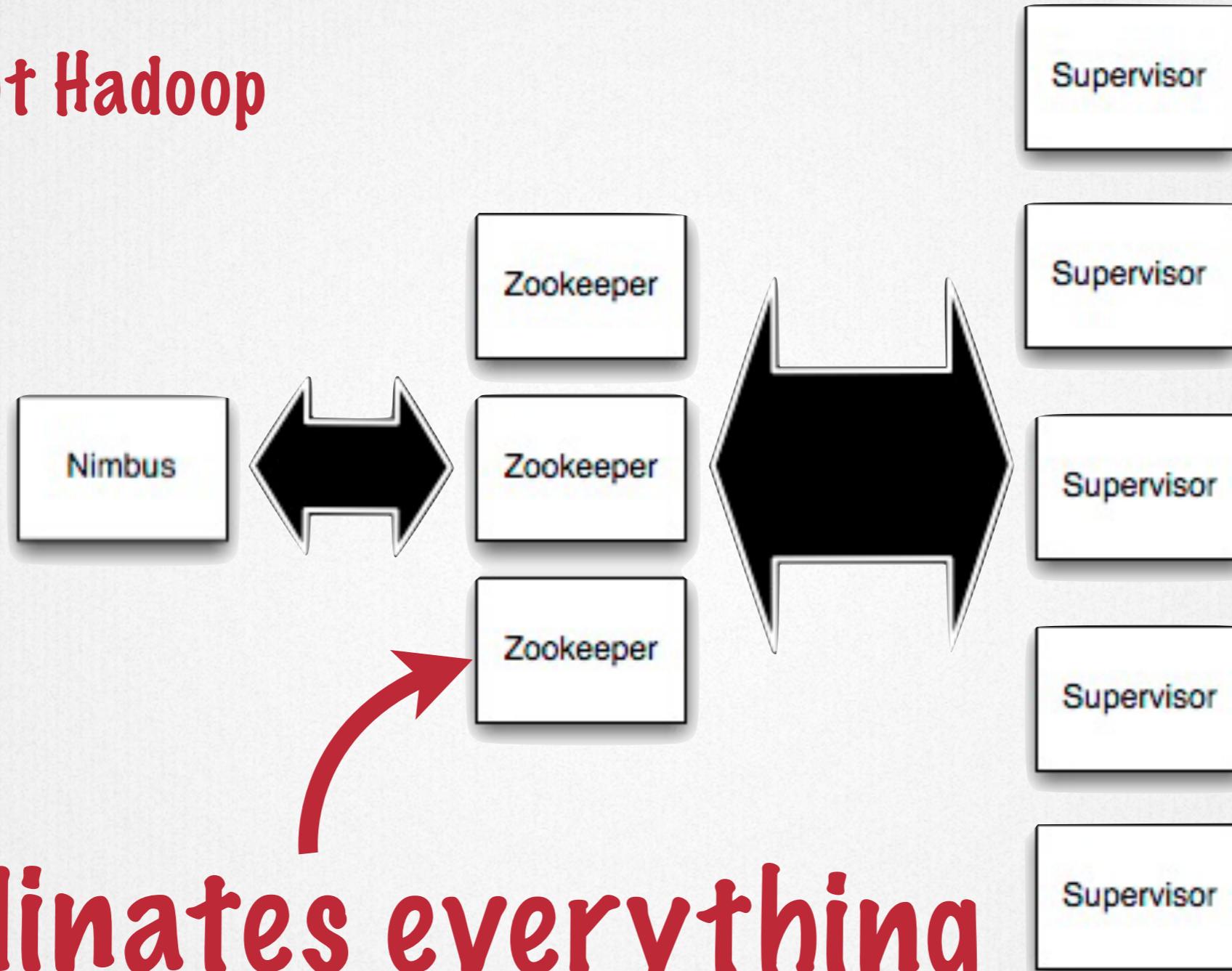
Storm Cluster

If this were
Hadoop...



Storm Cluster

But it's not Hadoop



Coordinates everything

Nathan Marz

Example: Streaming Word Count

Streaming Word Count

```
TopologyBuilder builder = new TopologyBuilder();

builder.setSpout("sentences", new RandomSentenceSpout(), 5);
builder.setBolt("split", new SplitSentence(), 8)
    .shuffleGrouping("sentences");
builder.setBolt("count", new WordCount(), 12)
    .fieldsGrouping("split", new Fields("word"));
```

Streaming Word Count

```
TopologyBuilder builder = new TopologyBuilder();

builder.setSpout("sentences", new RandomSentenceSpout(), 5);
builder.setBolt("split", new SplitSentence(), 8)
    .shuffleGrouping("sentences");
builder.setBolt("count", new WordCount(), 12)
    .fieldsGrouping("split", new Fields("word"));
```

Streaming Word Count

```
public static class SplitSentence extends ShellBolt implements IRichBolt {  
  
    public SplitSentence() {  
        super("python", "splitsentence.py");  
    }  
  
    @Override  
    public void declareOutputFields(OutputFieldsDeclarer declarer) {  
        declarer.declare(new Fields("word"));  
    }  
  
    @Override  
    public Map<String, Object> getComponentConfiguration() {  
        return null;  
    }  
}
```

SplitSentence.java

Streaming Word Count

```
public static class SplitSentence extends ShellBolt implements IRichBolt {  
  
    public SplitSentence() {  
        super("python", "splitsentence.py");  
    }  
  
    @Override  
    public void declareOutputFields(OutputFieldsDeclarer declarer) {  
        declarer.declare(new Fields("word"));  
    }  
  
    @Override  
    public void execute(Tuple tuple, BasicOutputCollector collector) {  
        String sentence = tuple.getString(0);  
        String[] words = sentence.split(" ");  
        for (String word : words) {  
            collector.emit(word);  
        }  
    }  
}  
  
import storm  
public class SplitSentenceBolt(storm.BasicBolt):  
    def process(self, tup):  
        words = tup.values[0].split(" ")  
        for word in words:  
            storm.emit([word])
```

SplitSentence.java

splitsentence.py

Streaming Word Count

```
public static class SplitSentence extends ShellBolt implements IRichBolt {  
  
    public SplitSentence() {  
        super("python", "splitsentence.py");  
    }  
  
    @Override  
    public void declareOutputFields(OutputFieldsDeclarer declarer) {  
        declarer.declare(new Fields("word"));  
    }  
  
    @Override  
    public Map<String, Object> getComponentConfiguration() {  
        return null;  
    }  
}
```

SplitSentence.java

Streaming Word Count

```
TopologyBuilder builder = new TopologyBuilder();

builder.setSpout("sentences", new RandomSentenceSpout(), 5);
builder.setBolt("split", new SplitSentence(), 8)
    .shuffleGrouping("sentences");
builder.setBolt("count", new WordCount() 12)
    .fieldsGrouping("split", new Fields("word"));
```

java

Streaming Word Count

```
public static class WordCount extends BaseBasicBolt {  
    Map<String, Integer> counts = new HashMap<String, Integer>();  
  
    @Override  
    public void execute(Tuple tuple, BasicOutputCollector collector) {  
        String word = tuple.getString(0);  
        Integer count = counts.get(word);  
        if(count==null) count = 0;  
        count++;  
        counts.put(word, count);  
        collector.emit(new Values(word, count));  
    }  
  
    @Override  
    public void declareOutputFields(OutputFieldsDeclarer declarer) {  
        declarer.declare(new Fields("word", "count"));  
    }  
}
```

WordCount.java

Streaming Word Count

```
TopologyBuilder builder = new TopologyBuilder();

builder.setSpout("sentences", new RandomSentenceSpout(), 5);
builder.setBolt("split", new SplitSentence(), 8)
    .shuffleGrouping("sentences");
builder.setBolt("count", new WordCount(), 12)
    .fieldsGrouping("split", new Fields("word"));
```



java

Groupings control how tuples are routed



FullContact

Shuffle grouping

Tuples are randomly distributed across all of the tasks running the bolt

Fields grouping

Groups tuples by specific named fields and routes them to the same task

Analogous to Hadoop's
partitioning behavior

Fields grouping

Groups tuples by specific named fields and routes
them to the same task

Distributed RPC

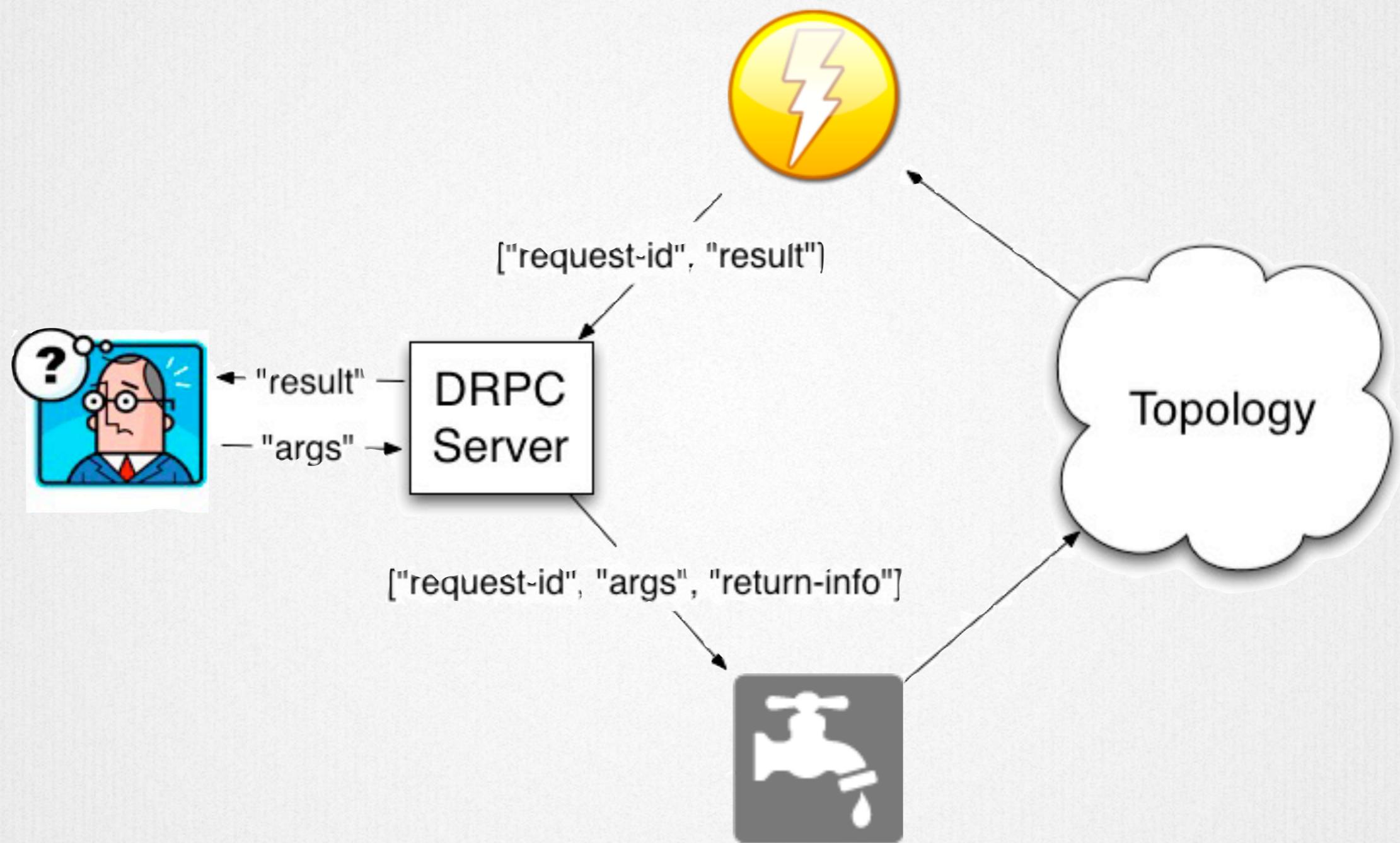
Before Distributed RPC,
time-sensitive queries relied on a
pre-computed index



FullContact

What if you didn't need an index?

Distributed RPC



Try it out!

Huge thanks to Nathan Marz - @nathanmarz

<http://github.com/nathanmarz/storm>

<https://github.com/nathanmarz/storm-starter>

@stormprocessor

Questions?

dan@fullcontact.com



FullContact