

Applying Real-time Processing Using Apache Storm

UNDERSTANDING THE COMPONENTS OF STORM



Swetha Kolalapudi

CO-FOUNDER, LOONYCORN

www.loonycorn.com

Overview

Differentiate between real-time processing and batch processing

Understand the components of a Storm topology

Understand how data is represented

Implement a stock price tracker topology

Quant Trading



**Quant traders at a
hedge fund**



**Evaluate stock
performance**



Performance Reports

Fixed frequency

- Quarterly, Monthly, Weekly, Daily etc

Compute metrics over the period

- High, low, average, momentum etc

Build models from historical data

Batch Processing



Monitoring, Alerts, Actions

Monitor the stock price in real-time

Apply trading models to identify actions

Alerts/triggers for specific trades

Real-time Processing

Real-time Processing vs. Batch Processing



Process each event as it happens

Continuous processing

Requires low latency



Process a group of events together

On-demand/At fixed frequency

High latency is acceptable

Real-time Processing vs. Batch Processing



Storm
Heron



Spark
Flink



Hadoop
MapReduce



Apache Storm

Open-sourced originally by Twitter

Real-time processing counterpart for Hadoop

Analytics, transformations, machine learning in real-time

Low latency due to distributed architecture

Components of a Storm Topology

Storm Application



Data is received as a stream

- Stock prices
- Log messages
- Tweets

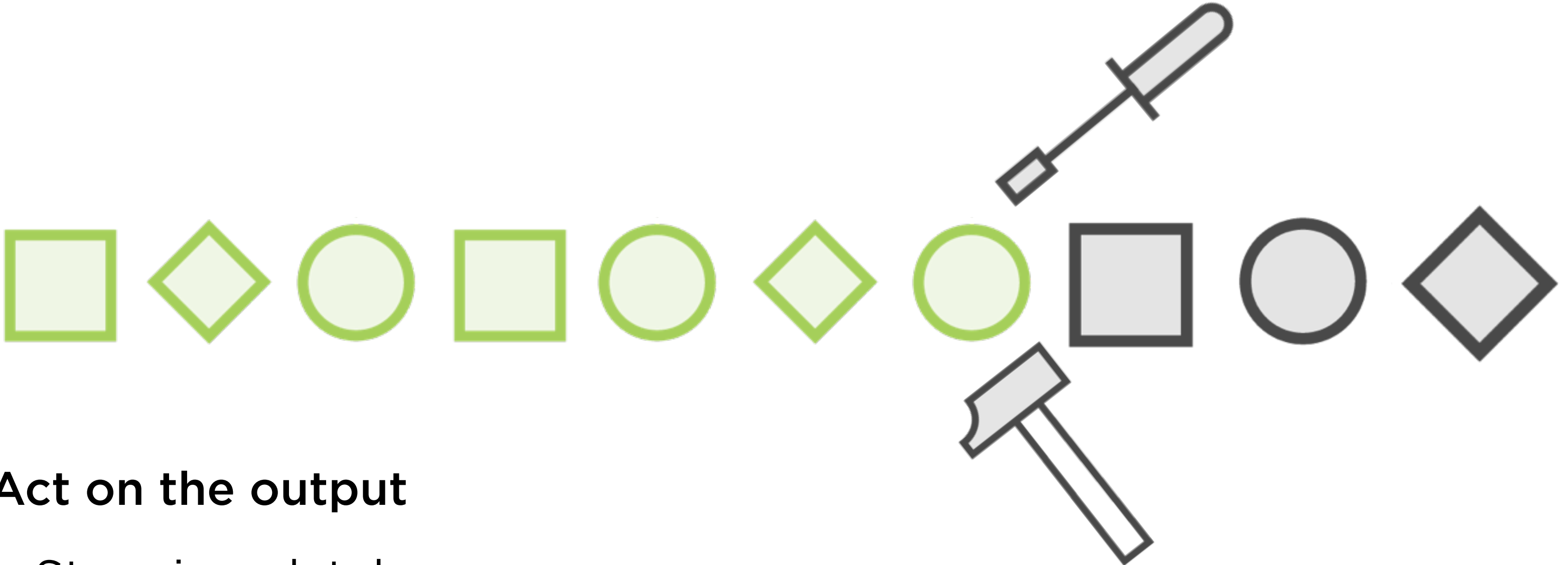
Storm Application



Process the data one entity at a time

- Compare the price to a constant
- Monitor error messages
- Find latest trending hashtags

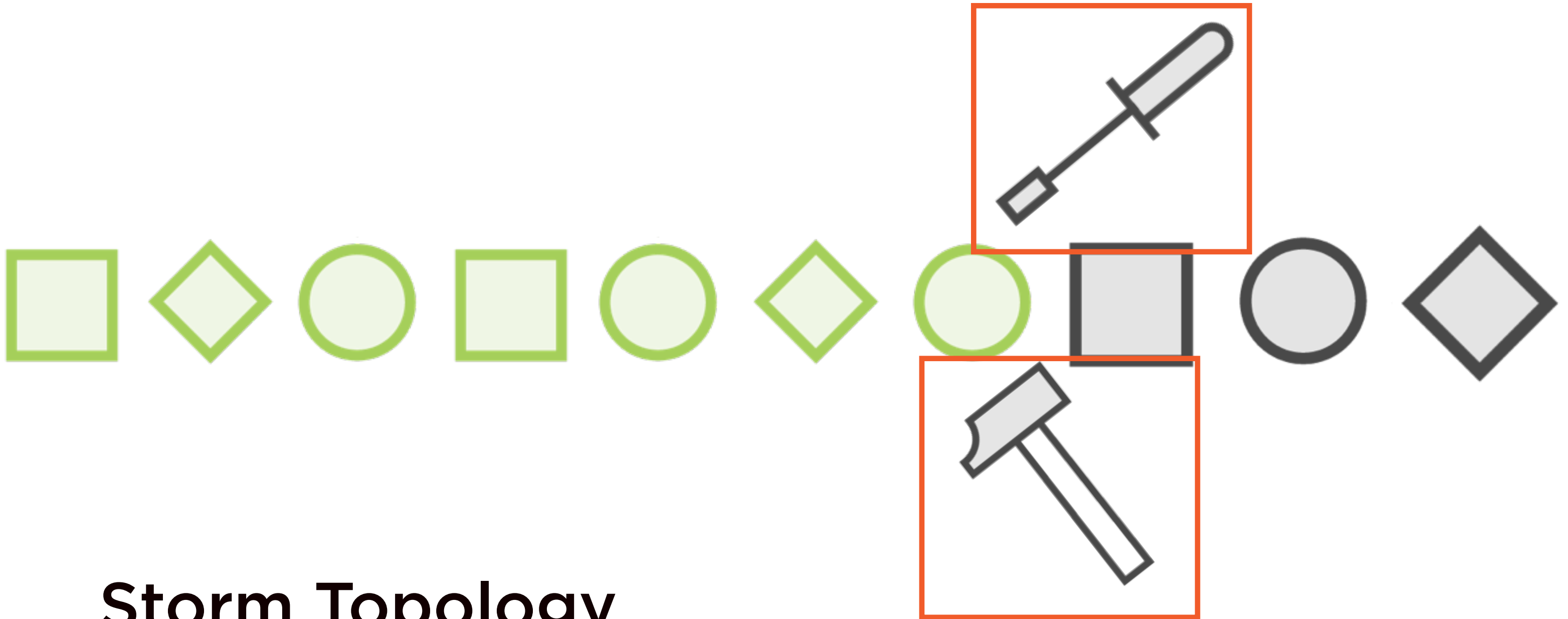
Storm Application



Act on the output

- Store in a database
- Trigger an alert

Storm Application



Storm Topology

Components of a Storm Topology



Spout

Receives data from the source
Emits it to the rest of the
topology



Bolt

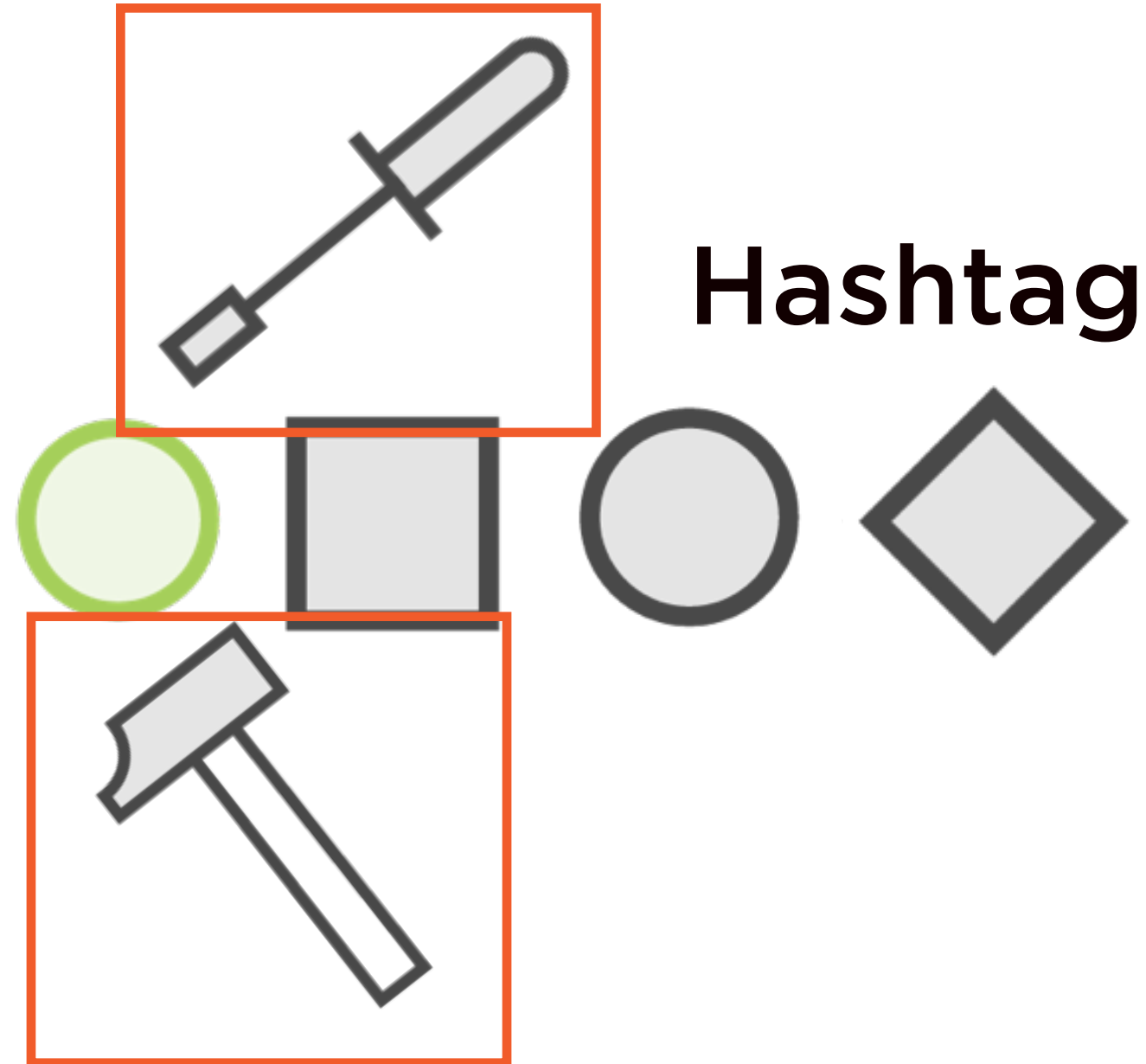
Processes the data

Hashtag Extractor

Tweets

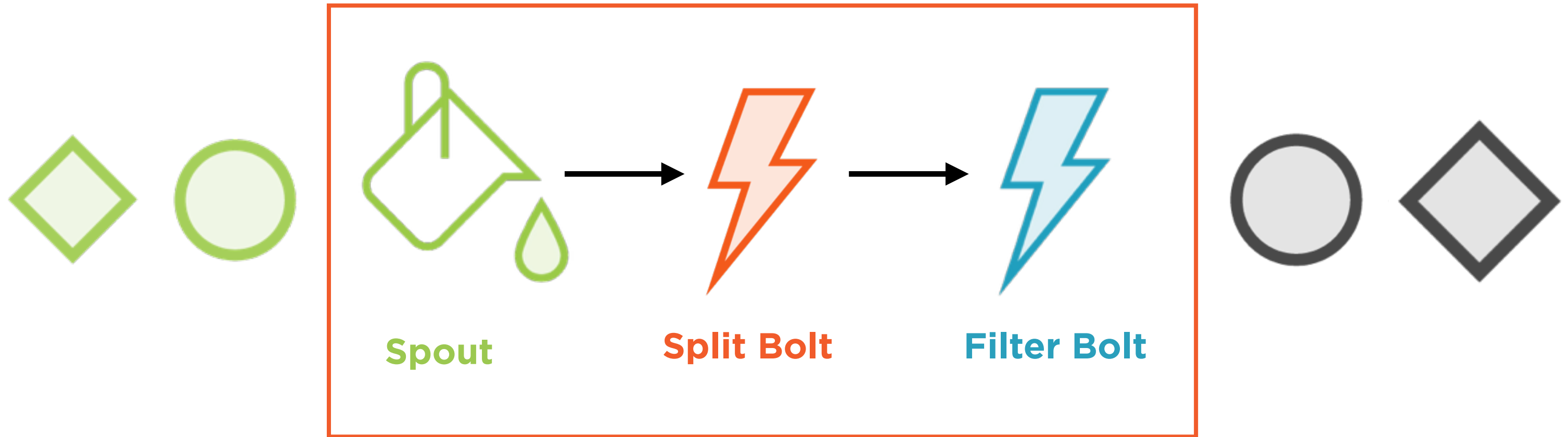


Hashtags

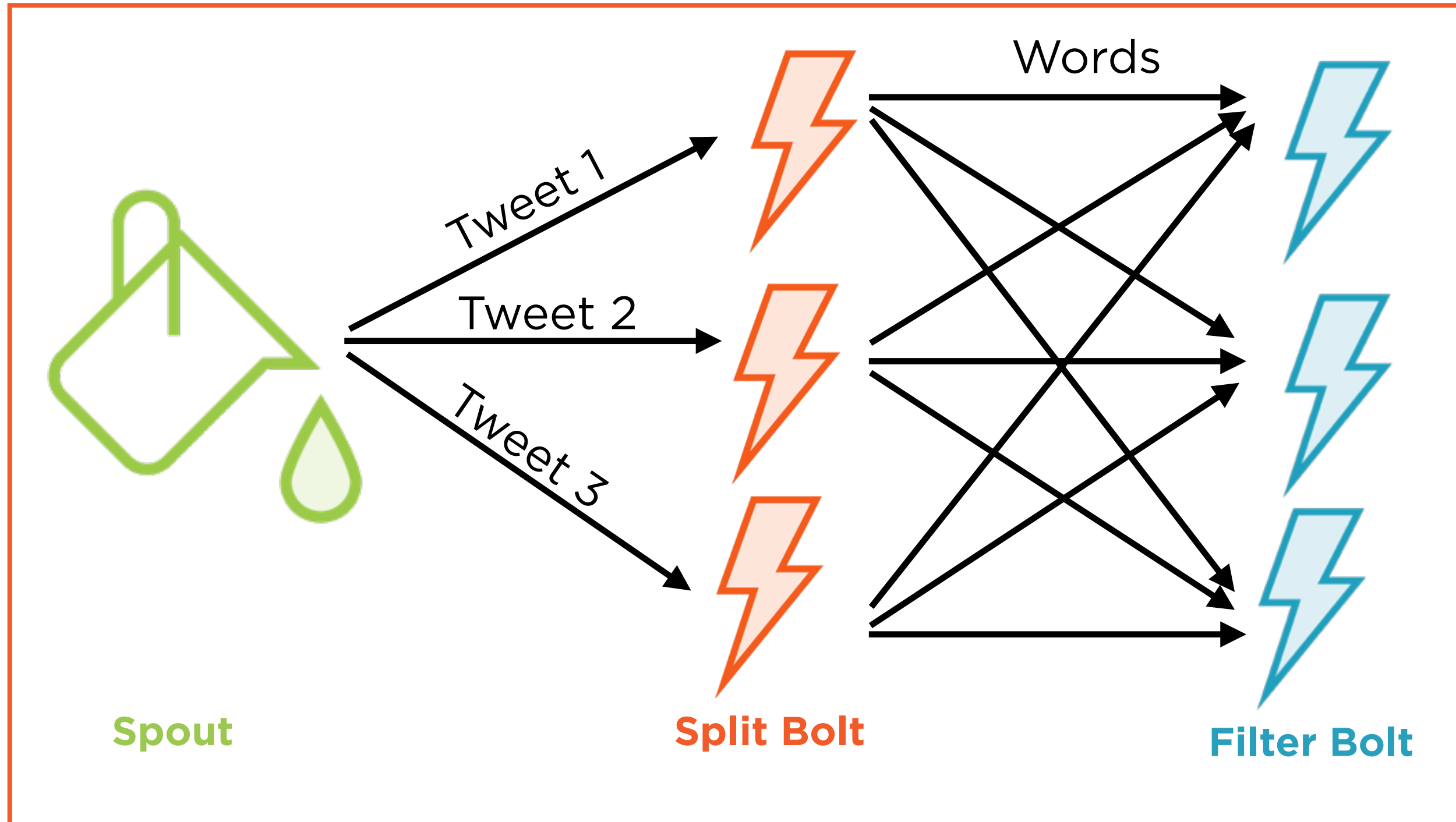


- Split tweets into words
- Filter words which are hashtags

Hashtag Extractor Topology



Hashtag Extractor Topology

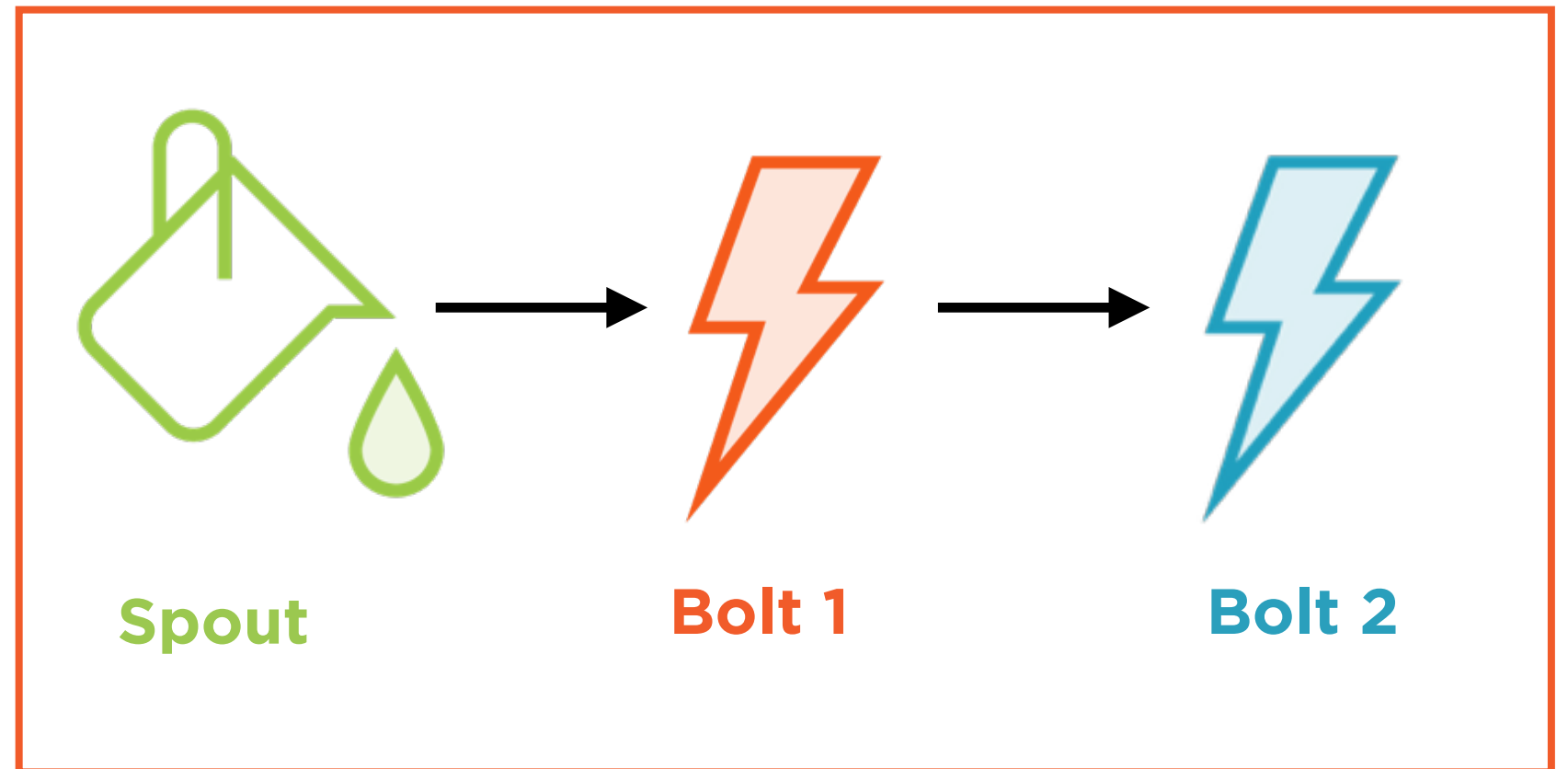


Storm Topology

**Components run by
Java processes**

**Need a central
coordinating process**

**All processes run in a
storm cluster**



Storm Cluster



Local

All processes run locally
Run topologies from IDE

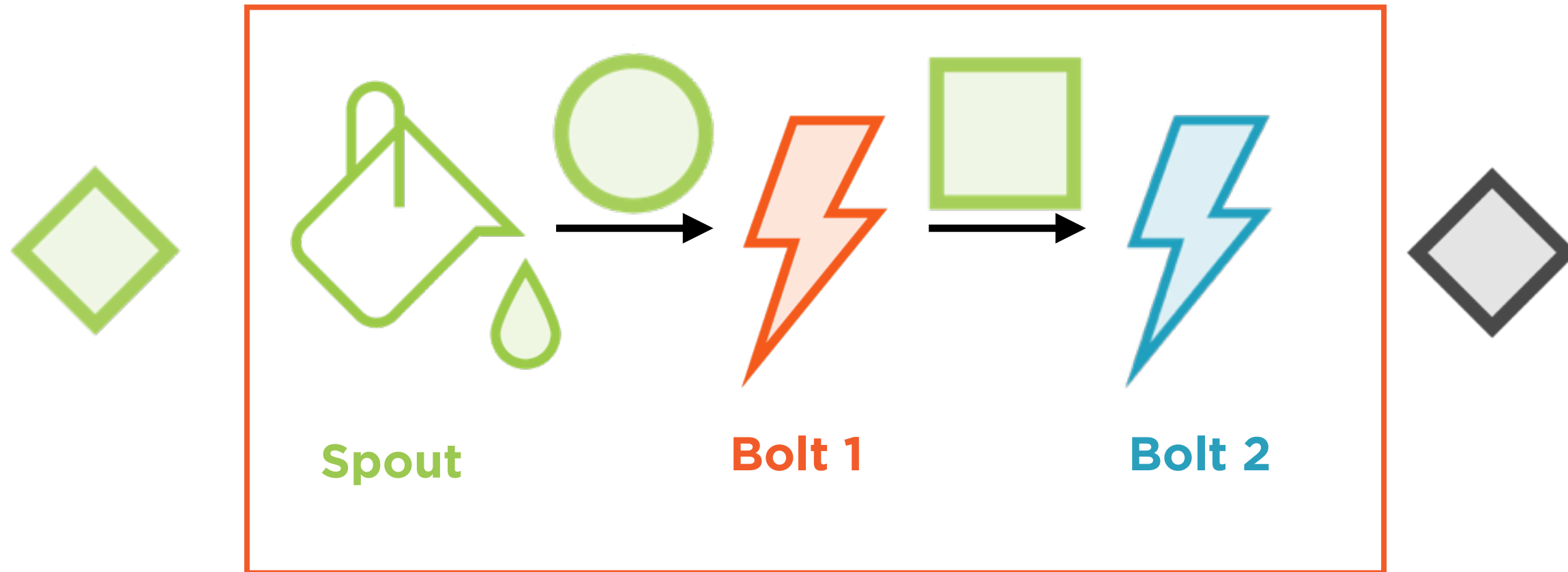


Remote

All processes run on different
machines
Submit a jar to cluster

Representing Data in Storm Topologies

Storm Topology



Components emit **Tuples**

Atomic data points

Named tuples

Storm Tuples

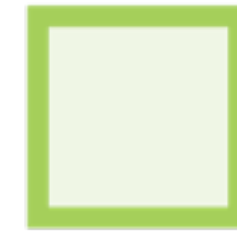


Symbol	Datetime	Price
MSFT	Mar 1, 12 9:00 AM	63

Schema of the tuple

Fields

Storm Tuples

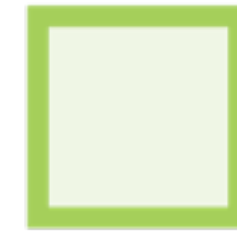


Symbol	Datetime	Price
MSFT	Mar 1, 12 9:00 AM	63

Values

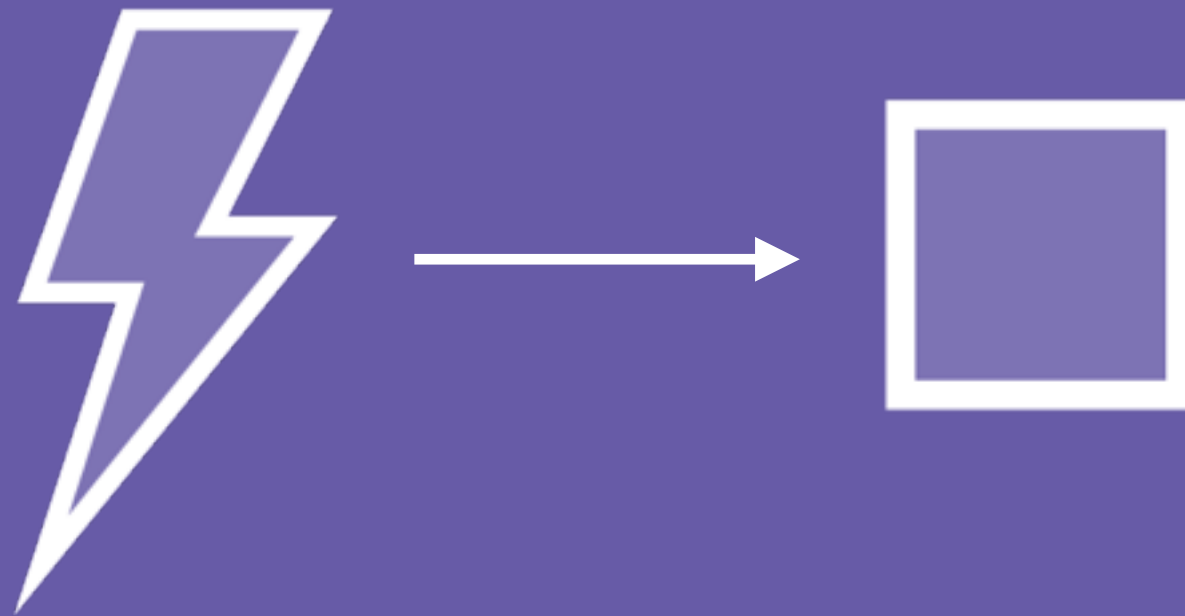
Data in the tuple

Storm Tuples

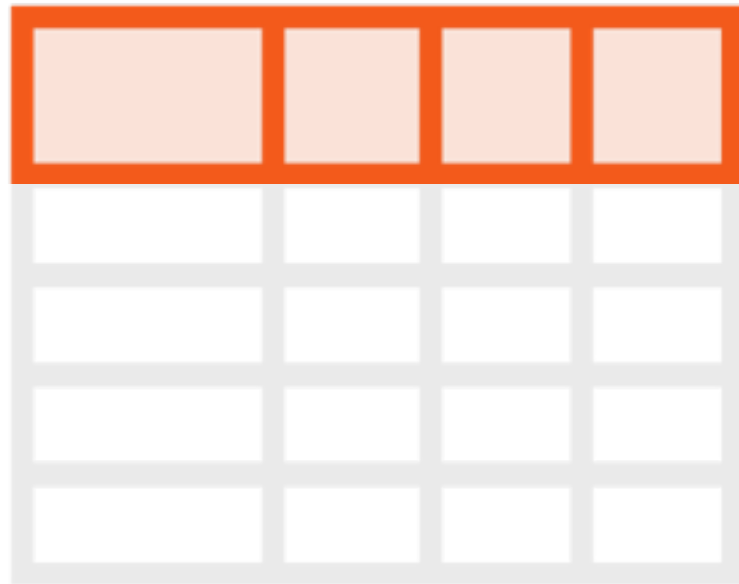


Symbol	Datetime	Price
MSFT	Mar 1, 12 9:00 AM	63

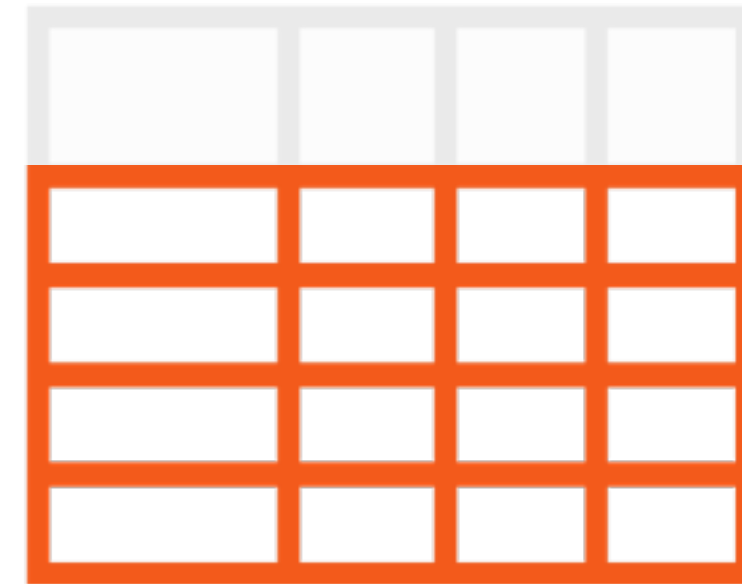
Each storm component emits
Tuples with a **fixed schema**



Setting up a Storm Component

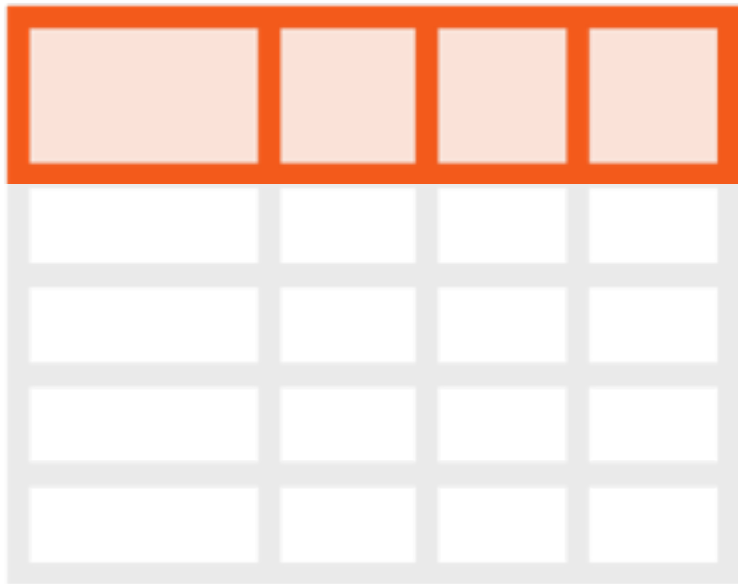


Declare Fields

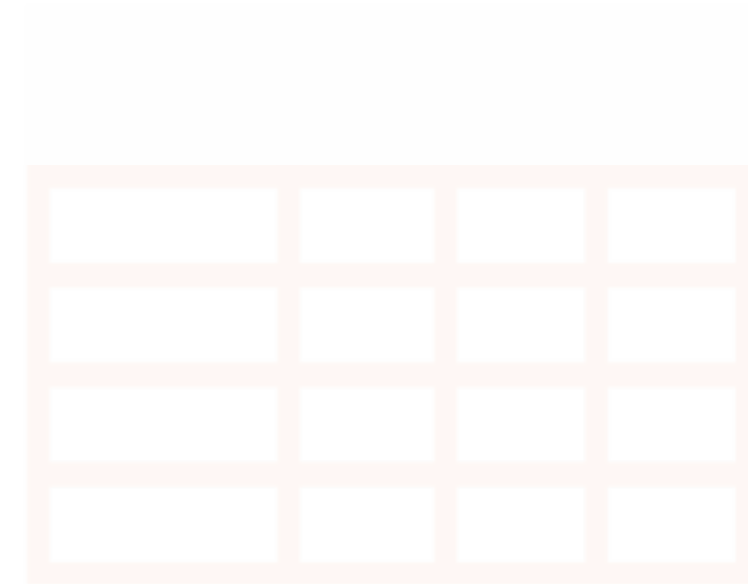


Emit Values

Setting up a Storm Component



Declare Fields



Emit Values

```
public void declareOutputFields(OutputFieldsDeclarer declarer){  
    declarer.declare(new Fields("symbol", "date", "price"));}
```

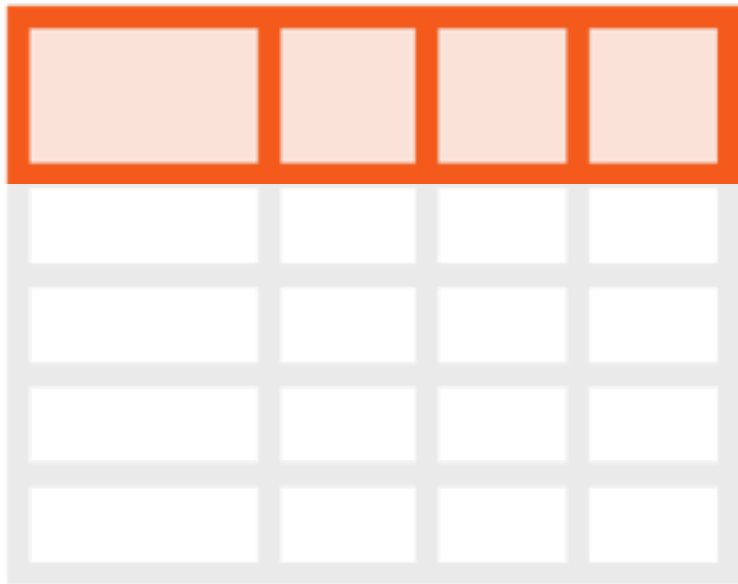
Declare Fields

Implement this method in every storm component

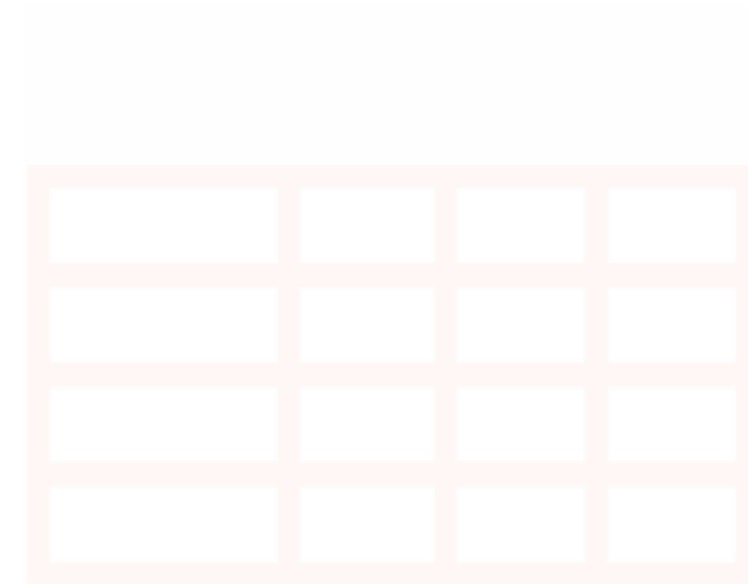
```
public void declareOutputFields(OutputFieldsDeclarer declarer){  
    declarer.declare(new Fields("symbol", "date", "price"));}
```

Declare Fields

Setting up a Storm Component



Declare Fields



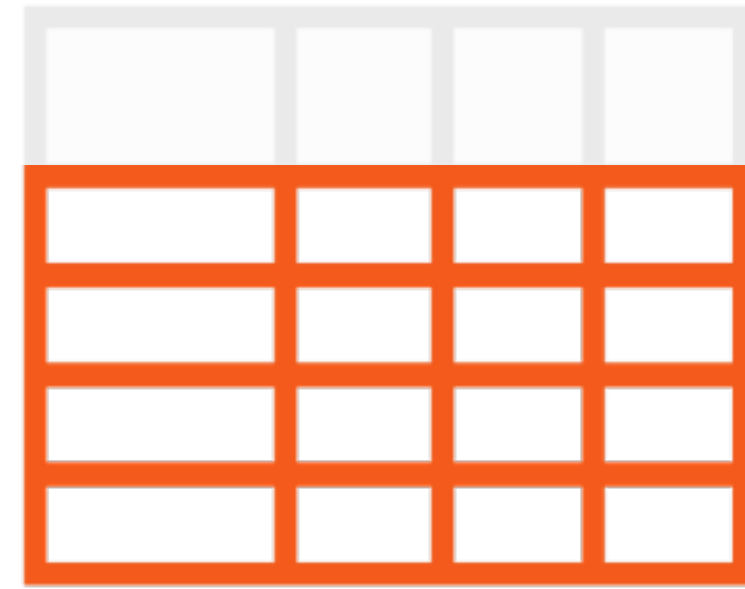
Emit Values

Setting up a Storm Component



A diagram illustrating the declaration of fields in a Storm component. It consists of a 5x4 grid of cells. The top row is highlighted with a thick orange border, representing the field names. The remaining four rows have thin gray borders, representing the data structure.

Declare Fields



A diagram illustrating the emission of values from a Storm component. It consists of a 5x4 grid of cells. The top row has a thin gray border, representing the field names. The remaining four rows have thick orange borders, representing the data values being emitted.

Emit Values

```
collector.emit(new Values("MSFT", "Mar 12 9:00 AM", 62.5))
```

Emit Values

All components use a collector object to emit values


```
collector.emit(new Values("MSFT", "Mar 12 9:00 AM", 62.5))
```

Emit Values

Provide Objects in the same order as declared Fields

Setting up a Storm Component

Declare Fields

Emit Values

Setting up a Storm Component

Declare Fields

Emit Values

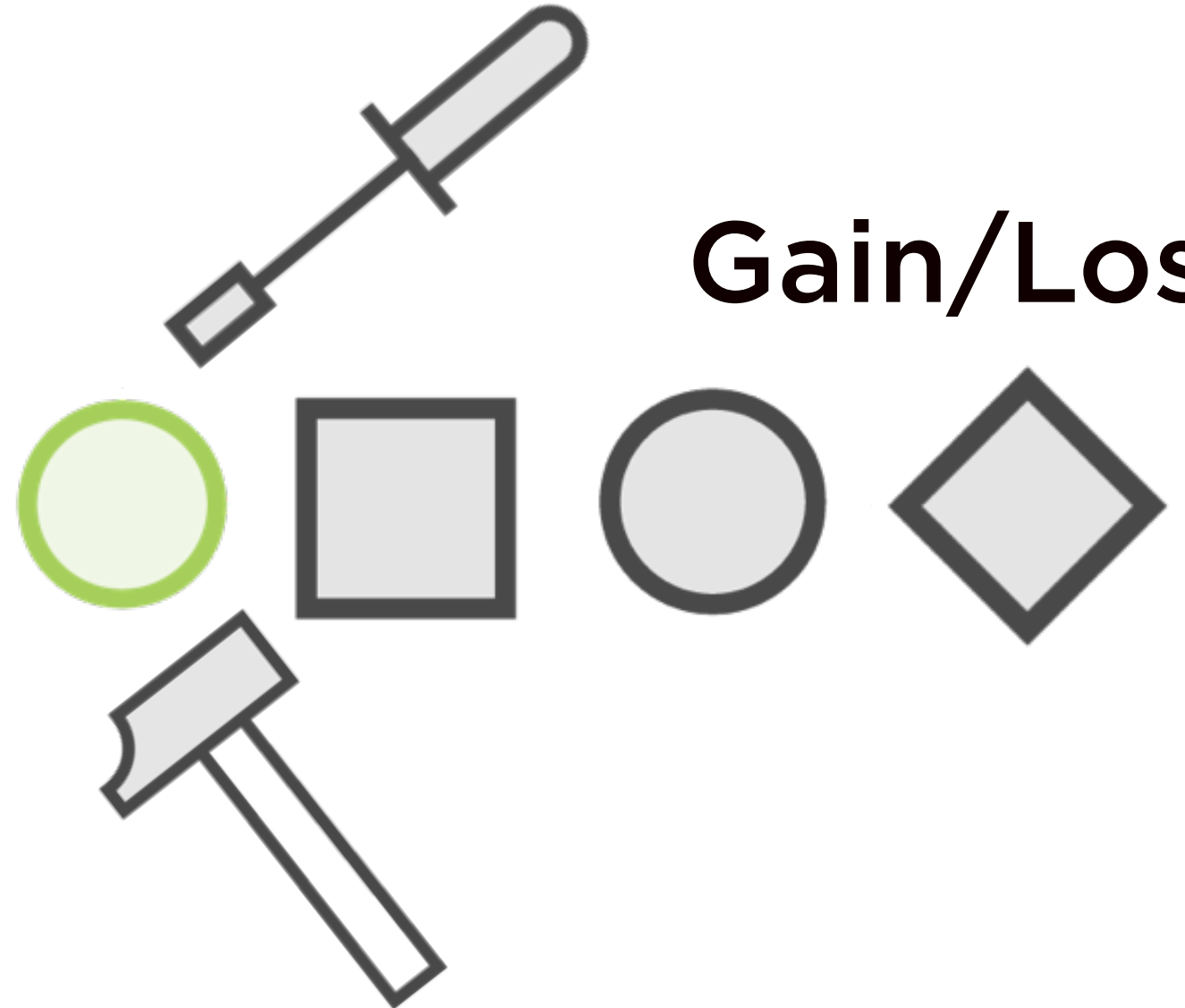
Building a Stock Price Tracker

Stock Price Tracker

Stock Prices



Gain/Loss



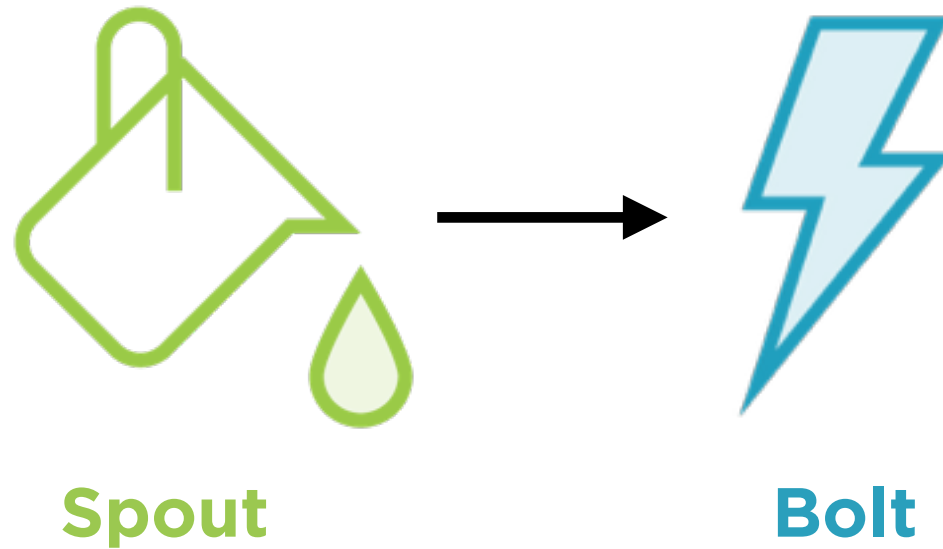
- Get stock prices from Yahoo Finance
- If $\text{price} > \text{previous close}$, **Gain** else **Loss**

Stock Price Tracker

Topology



**Yahoo
Finance**



Gain/Loss
(Written to file)

Stock Price Tracker

Set up a project

Use Maven to add Storm as a dependency

Set up the Bolt

Extends the
BaseBasicBolt abstract
class

Set up the Spout

Extends the
BaseRichSpout abstract
class

Run the topology

Build the topology
Submit to a local
cluster

Stock Price Tracker

Set up a project

Use Maven to add Storm
as a dependency



Demo

Set up a project

- Add the Storm dependency

Stock Price Tracker

Set up a project

Use Maven to add Storm as a dependency

Set up the Spout

Extends the
BaseRichSpout abstract
class

Demo

Set up the Spout

- Extend the BaseRichSpout class
- Connect to Yahoo Finance API
- Get stock price, previous close and send to Bolt

```
public void open(Map conf, TopologyContext context,  
                SpoutOutputCollector collector) {}  
  
public void nextTuple() {}  
  
public void declareOutputFields(OutputFieldsDeclarer declarer) {}
```

Implementing the Spout

```
public void open(Map conf, TopologyContext context,  
                SpoutOutputCollector collector) {}
```

```
public void nextTuple() {}
```

```
public void declareOutputFields(OutputFieldsDeclarer declarer) {}
```

Initializing the Spout

```
public void open(Map conf, TopologyContext context,  
                 SpoutOutputCollector collector) {}
```

```
public void nextTuple() {}
```

```
public void declareOutputFields(OutputFieldsDeclarer declarer) {}
```

Initializing the Spout

Topology configuration

```
public void open(Map conf, TopologyContext context,  
                  SpoutOutputCollector collector) {}
```

```
public void nextTuple() {}
```

```
public void declareOutputFields(OutputFieldsDeclarer declarer) {}
```

Initializing the Spout
Component name, id etc

```
public void open(Map conf, TopologyContext context,  
                 SpoutOutputCollector collector) {}
```

```
public void nextTuple() {}
```

```
public void declareOutputFields(OutputFieldsDeclarer declarer) {}
```

Initializing the Spout

Used to communicate between components

- Emit data, success/failure


```
public void open(Map conf, TopologyContext context,  
                SpoutOutputCollector collector) {}
```

```
public void nextTuple() {}
```

```
public void declareOutputFields(OutputFieldsDeclarer declarer) {}
```

Initializing the Spout

```
public void open(Map conf, TopologyContext context,  
                  SpoutOutputCollector collector) {}
```

```
public void nextTuple() {}
```

```
public void declareOutputFields(OutputFieldsDeclarer declarer) {}
```

Declare Fields

```
public void open(Map conf, TopologyContext context,  
                  SpoutOutputCollector collector) {}
```

```
public void nextTuple() {}
```

```
public void declareOutputFields(OutputFieldsDeclarer declarer) {}
```

Emit Tuples

Called continuously as long as the spout runs

Stock Price Tracker

Set up a project

Use Maven to add Storm as a dependency

Set up the Bolt

Extends the
BaseBasicBolt abstract
class

Set up the Spout

Extends the
BaseRichSpout abstract
class

Demo

Set up the Bolt

- Extend the BaseBasicBolt class
- Compute Gain/Loss signal
- Write data to file

```
public void prepare(Map stormConf, TopologyContext context) {}  
public void execute(Tuple input, BasicOutputCollector collector) {}  
public void declareOutputFields(OutputFieldsDeclarer declarer) {}  
public void cleanup() {}
```

Implementing the Bolt

```
public void prepare(Map stormConf, TopologyContext context) {}  
public void execute(Tuple input, BasicOutputCollector collector) {}  
public void declareOutputFields(OutputFieldsDeclarer declarer) {}  
public void cleanup() {}
```

Initializing the Bolt

```
public void prepare(Map stormConf, TopologyContext context) {}  
public void execute(Tuple input, BasicOutputCollector collector) {}  
public void declareOutputFields(OutputFieldsDeclarer declarer) {}  
public void cleanup() {}
```

Declare Fields


```
public void prepare(Map stormConf, TopologyContext context) {}  
public void execute(Tuple input, BasicOutputCollector collector) {}  
public void declareOutputFields(OutputFieldsDeclarer declarer) {}  
public void cleanup() {}
```

Process Input Tuple and Emit Values

```
public void prepare(Map stormConf, TopologyContext context) {}  
public void execute(Tuple input, BasicOutputCollector collector) {}  
public void declareOutputFields(OutputFieldsDeclarer declarer) {}  
public void cleanup() {}
```

Cleanup on Bolt Shutdown

Stock Price Tracker

Set up a project

Use Maven to add Storm as a dependency

Set up the Bolt

Extends the
BaseBasicBolt abstract
class

Set up the Spout

Extends the
BaseRichSpout abstract
class

Run the topology

Build the topology
Submit to a local
cluster

Running the Topology

Build the topology

Add the spout and bolt
and connect them

Submit to a local cluster

Configure the topology

Specify the output
file directory



Demo

Running the topology

Summary

Differentiate between real-time processing and batch processing

Understand the components of a Storm topology

Understand how data is represented

Implement a stock price tracker topology