The background of the slide features a wide-angle photograph of a severe thunderstorm. The sky is filled with heavy, dark grey and black clouds, suggesting a supercell or a very intense convective system. In the distance, a bright horizon line marks the transition from the dark sky to a lighter area where some buildings and trees are visible. The overall mood is somber and dramatic.

Apache Storm

ECE 4150
Spring 2024
Vijay Madisetti



Plan for Class Till End of Semester

- Next Week is Spring Break
- Total of five labs. One more lab will be released.
- March & 27: AWS EMR in Depth Lectures
- Apr 1- Lecture by Industry Visitor (Zoom likely)
- April 3 – Pre-Recorded Zoom lecture (Real-Time Container Scheduling)
- Project Plan: Short paper (use the template of a typical lab to show the various steps of your work on it) on either Spark or Storm, try out some sample code on Spark or Storm and show that you have used it somehow. Lots of tutorials available online, and I will distribute some handouts as well. Due mid April.

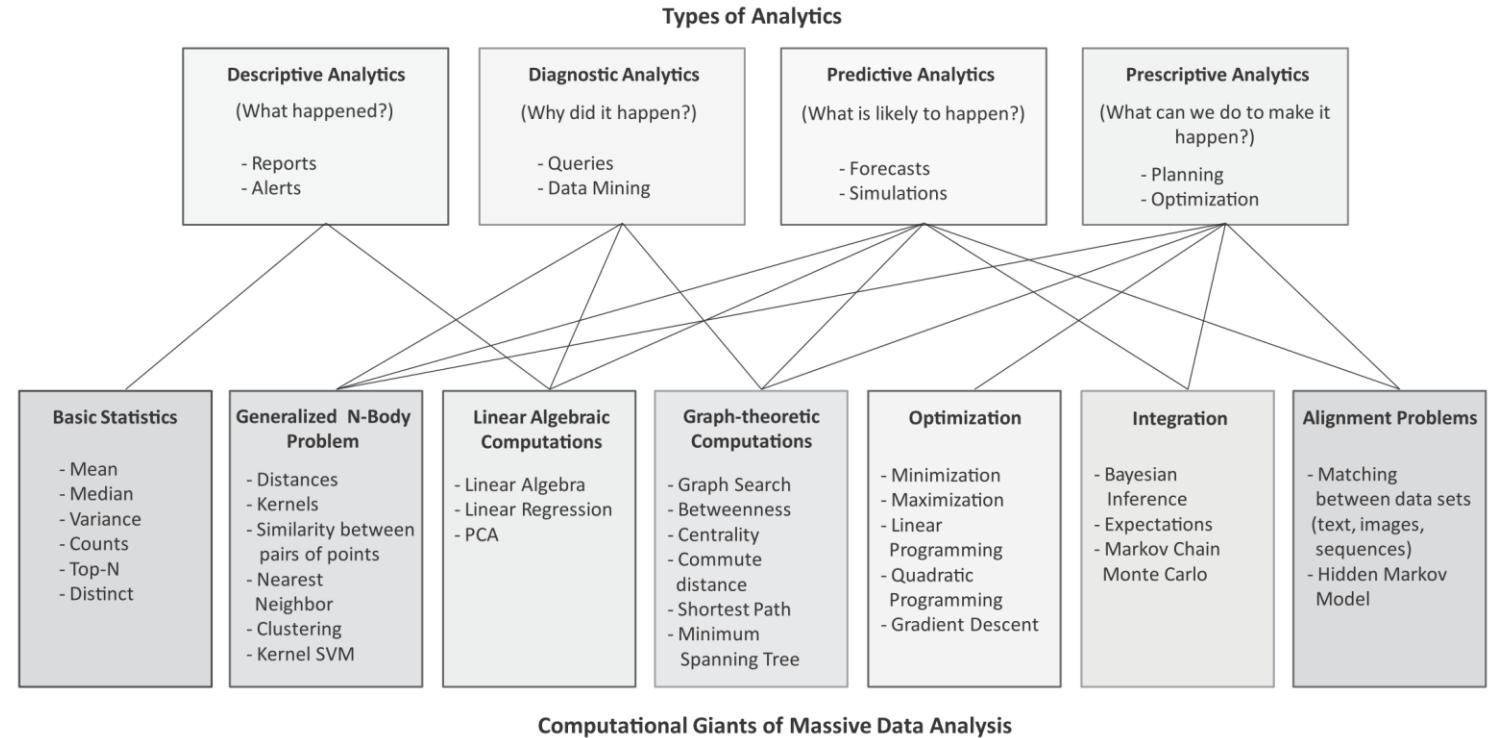
What is Analytics?

- Analytics is a broad term that encompasses the processes, technologies, frameworks and algorithms to extract meaningful insights from data.
- Raw data in itself does not have a meaning until it is contextualized and processed into useful information.
- Analytics is this process of extracting and creating information from raw data by filtering, processing, categorizing, condensing and contextualizing the data.
- This information obtained is then organized and structured to infer knowledge about the system and/or its users, its environment, and its operations and progress towards its objectives, thus making the systems smarter and more efficient.

The Seven Giants

- 
- The National Research Council has done a characterization of computational tasks for massive data analysis (called the seven “giants”).
 - These computational tasks include:
 - (1) Basic Statistics,
 - (2) Generalized N-Body Problems,
 - (3) Linear Algebraic Computations,
 - (4) Graph-Theoretic Computations,
 - (5) Optimization,
 - (6) Integration and
 - (7) Alignment Problems.
 - This characterization of computational tasks aims to provide a taxonomy of tasks that have proved to be useful in data analysis and grouping them roughly according to mathematical structure and computational strategy.

Mapping between types of analytics and computational tasks or ‘giants’





Types of Analytics

- **Descriptive Analytics**

- Descriptive analytics comprises analyzing past data to present it in a summarized form which can be easily interpreted. Descriptive analytics aims to answer - ***What has happened?*** A major portion of analytics done today is descriptive analytics through use of statistics functions such as counts, maximum, minimum, mean, top-N, percentage, for instance. These statistics help in describing patterns in the data and present the data in a summarized form.

- **Diagnostic Analytics**

- Diagnostic analytics comprises analysis of past data to diagnose the reasons as to why certain events happened. Diagnostic analytics aims to answer - ***Why did it happen?*** While descriptive analytics can be useful for summarizing the data by computing various statistics (such as mean, minimum, maximum, variance, or top-N), diagnostic analytics can provide more insights into why certain a fault has occurred based on the patterns in the sensor data for previous faults.



Types of Analytics

- **Predictive Analytics**

- Predictive analytics comprises predicting the occurrence of an event or the likely outcome of an event or forecasting the future values using prediction models. Predictive analytics aims to answer - ***What is likely to happen?*** Predictive Analytics is done using predictive models which are trained by existing data. These models learn patterns and trends from the existing data and predict the occurrence of an event or the likely outcome of an event (classification models) or forecast numbers (regression models).

- **Prescriptive Analytics**

- While predictive analytics uses prediction models to predict the likely outcome of an event, prescriptive analytics uses multiple prediction models to predict various outcomes and the best course of action for each outcome. Prescriptive analytics aims to answer - ***What can we do to make it happen?*** Prescriptive Analytics can predict the possible outcomes based on the current choice of actions.

Domain Specific Examples of Big Data

- Web
- Financial
- Healthcare
- Internet of Things
- Environment
- Logistics & Transportation
- Industry
- Retail



Analytics Flow for Big Data

- **Data Collection**
 - Data collection is the first step for any analytics application. Before the data can be analyzed, the data must be collected and ingested into a big data stack. The choice of tools and frameworks for data collection depends on the source of data and the type of data being ingested.
- **Data Preparation**
 - Data can often be dirty and can have various issues that must be resolved before the data can be processed, such as corrupt records, missing values, duplicates, inconsistent abbreviations, inconsistent units, typos, incorrect spellings and incorrect formatting. Data preparation step involves various tasks such as data cleansing, data wrangling or munging, de-duplication, normalization, sampling and filtering.

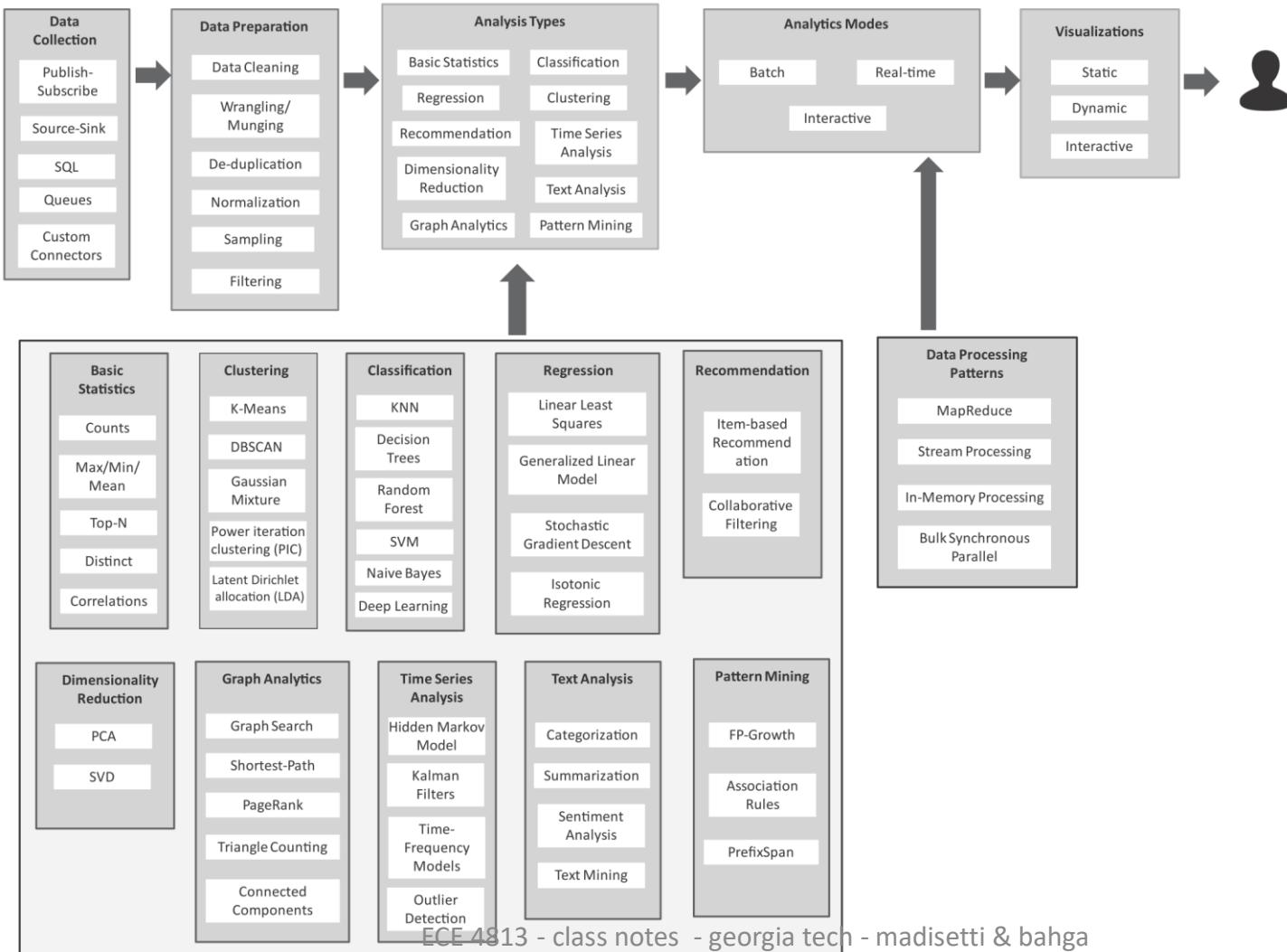


Analytics Flow for Big Data

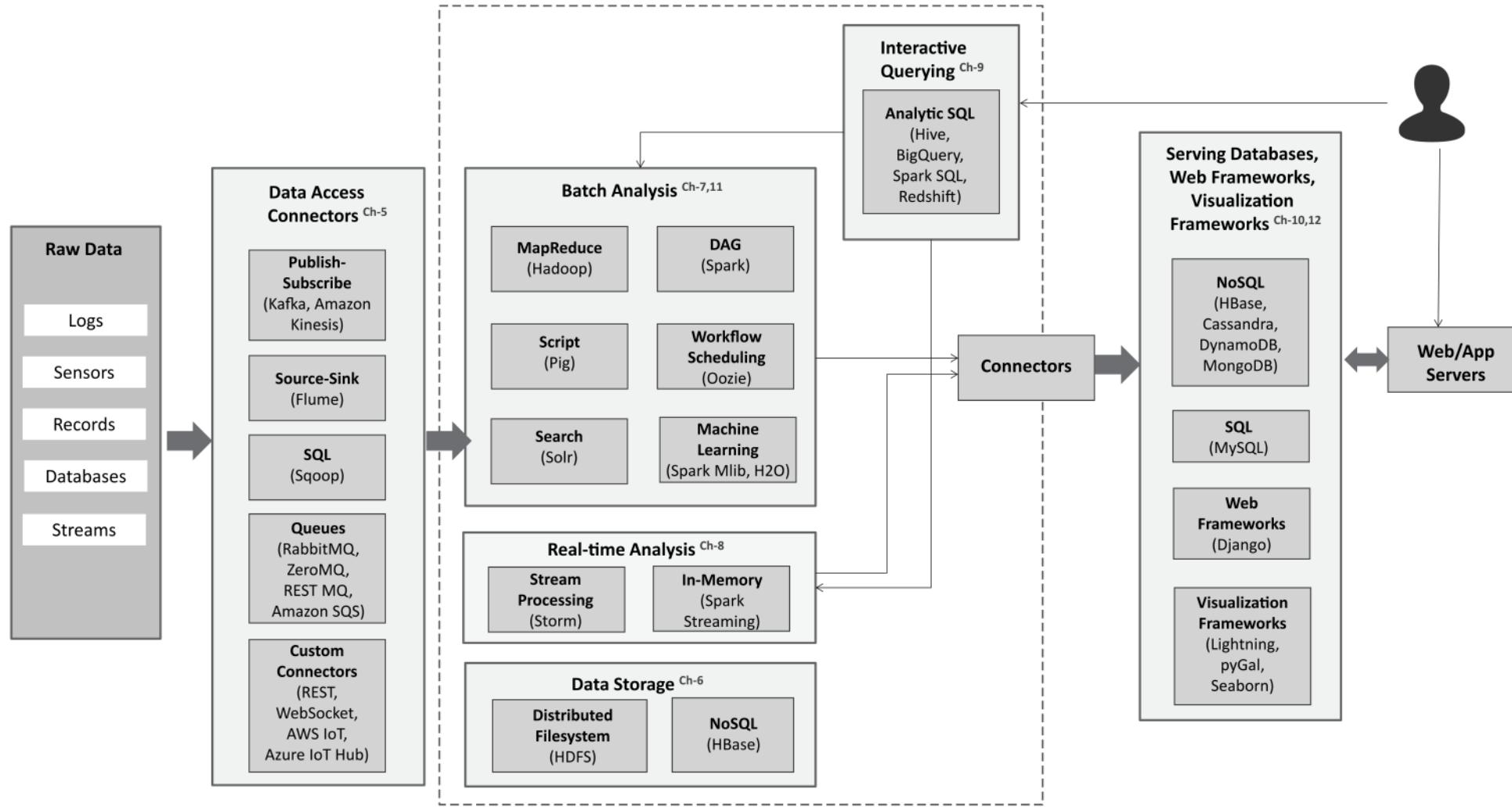
- **Analysis Types**
 - The next step in the analysis flow is to determine the analysis type for the application.
- **Analysis Modes**
 - With the analysis types selected for an application, the next step is to determine the analysis mode, which can be either batch, real-time or interactive. The choice of the mode depends on the requirements of the application.
- **Visualizations**
 - The choice of the visualization tools, serving databases and web frameworks is driven by the requirements of the application. Visualizations can be static, dynamic or interactive.



Analytics Flow for Big Data



Big Data Stack



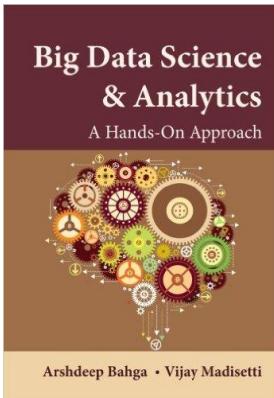
Big Data Stack

- **Raw Data Sources**
 - In any big data analytics application or platform, before the data is processed and analyzed, it must be captured from the raw data sources into the big data systems and frameworks.
- **Data Access Connectors**
 - The Data Access Connectors includes tools and frameworks for collecting and ingesting data from various sources into the big data storage and analytics frameworks. The choice of the data connector is driven by the type of the data source.
- **Data Storage**
 - The data storage block in the big data stack includes distributed filesystems and non-relational (NoSQL) databases, which store the data collected from the raw data sources using the data access connectors.

Big Data Stack

- **Batch Analytics**
 - The batch analytics block in the big data stack includes various frameworks which allow analysis of data in batches such as Hadoop-MapReduce, Pig, Oozie, Pig, Spark, Solr.
- **Real-time Analytics**
 - The real-time analytics block includes the Apache Storm and Spark Streaming frameworks.
- **Interactive Querying**
 - Interactive querying systems allow users to query data by writing statements in SQL-like languages.
- **Serving Databases, Web & Visualization Frameworks**
 - While the various analytics blocks process and analyze the data, the results are stored in serving databases for subsequent tasks of presentation and visualization. These serving databases allow the analyzed data to be queried and presented in the web applications.





Mapping Analytics Flow to Big Data Stack

Data Collection

Analysis Type	Framework (Mode)
Publish-Subscribe	Kafka, Kinesis
Source-Sink	Flume
SQL	Sqoop
Queues	SQS, RabbitMQ, ZeroMQ, RESTMQ
Custom Connectors	REST, WebSocket, MQTT

Data Preparation

Analysis Type	Framework
Data Cleaning	Open Refine
Data Wrangling	Open Refine DataWrangler
De-Duplication	Open Refine, Pig, Hive, Spark SQL
Normalization Sampling, Filtering	MapReduce, Pig, Hive, Spark SQL

Basic Statistics

Analysis Type	Framework (Mode)
Counts, Max, Min, Mean, Top-N, Distinct	Hadoop-MapReduce (Batch), Pig (Batch), Spark (Batch), Spark Streaming (Realtime), Spark SQL (Interactive), Hive (Integrative), Storm (Real-time)
Correlations	Hadoop-MapReduce (Batch), Spark Mlib (Batch)

Clustering

Analysis Type	Framework (Mode)
K-Means	Hadoop-MapReduce (Batch), Spark Mlib (Batch & Real-time) H2O (Batch)
DBSCAN	Spark (Batch)
Gaussian Mixture	Spark Mlib (Batch)
PIG	Spark Mlib (Batch)
LDA	Spark Mlib (Batch)

Classification

Analysis Type	Framework (Mode)
KNN	Spark Mlib (Batch , Realtime)
Decision Trees	Spark Mlib (Batch, Realtime)
Random Forest	Spark Mlib (Batch , Realtime), H2O (Batch)
SVM	Spark Mlib (Batch , Realtime)
Naïve Bayes	Spark Mlib (Batch, Realtime), H2O (Batch)
Deep Learning	H2O (Batch)

Regression

Analysis Type	Framework (Mode)
Linear Least Squares	Spark Mlib (Batch, Realtime)
Generalized Linear Model	H2O (Batch)
Stochastic Gradient Descent	Spark Mlib (Batch, Realtime)
Isotonic Regression	Spark Mlib (Batch, Realtime)

Mapping Analytics Flow to Big Data Stack

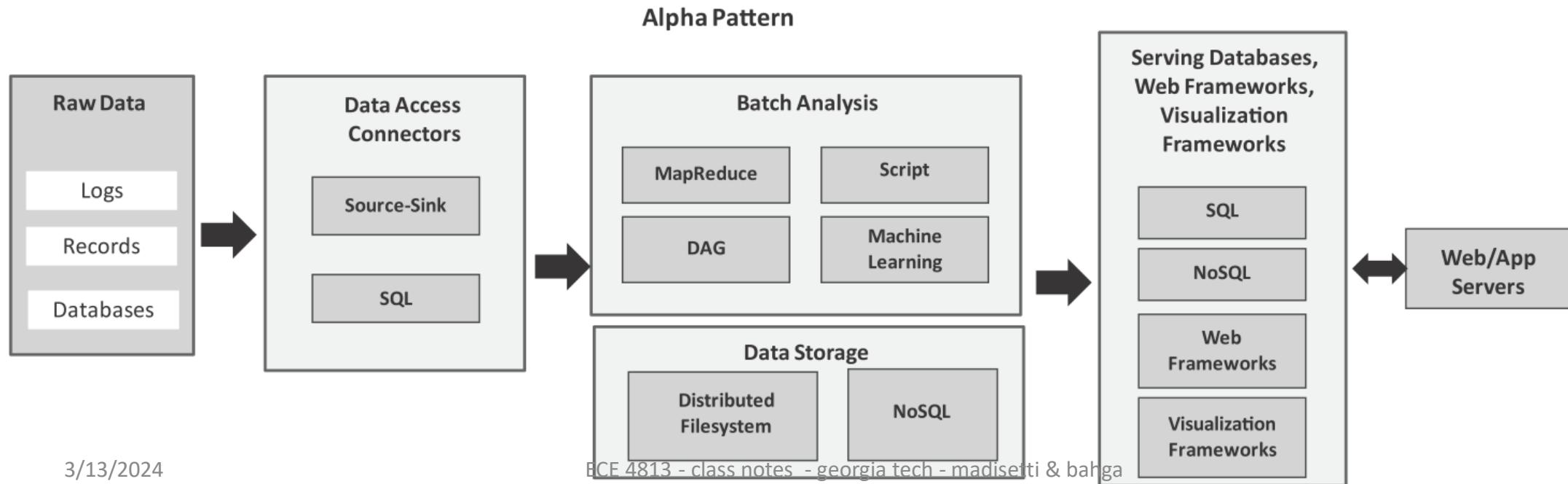
Graph Analytics		Time Series Analysis		Dimensionality Reduction		Recommendation	
Analysis Type	Framework (Mode)	Analysis Type	Framework (Mode)	Analysis Type	Framework (Mode)	Analysis Type	Framework (Mode)
Graph Search	Spark GraphX (Batch)	Kalman Filter	Spark (Realtime)	SVD	Spark Mlib (Batch)	Item-bases Recommendation	Spark Mlib (Batch)
Shortest-Path	Spark GraphX (Batch)	Time Frequency Models	Spark (Realtime)	PCA	Spark Mlib (Batch), H2O (Batch)	Collaborative Filtering	Spark Mlib (Batch)
PageRank	Spark GraphX (Batch)	Outlier Detection	Storm (Realtime), Spark (Batch, Realtime)				
Triangle Counting	Spark GraphX (Batch)						
Connected Components	Spark GraphX (Batch)						

Text Analysis		Pattern Mining		Visualization	
Analysis Type	Framework (Mode)	Analysis Type	Framework (Mode)	Analysis Type	Framework (Mode)
Categorization	Hadoop-MapReduce (Batch), Storm (Realtime), Spark (Batch, Realtime)	FP-Growth	Spark Mlib (Batch)	Web Frameworks	Django, Flask
Summarization	Spark (Batch)	Association Rules	Spark Mlib (Batch)	SQL Databases	MySQL
Sentiment Analysis	Storm (Realtime), Spark (Batch, Realtime)	PrefixSpan	Spark Mlib (Batch)	NoSQL Databases	Hbase, DynamoDB, Cassandra, MongoDB
Text Mining	Storm (Realtime), Spark (Batch, Realtime)			Visualization Frameworks	Lightning, pyGal, Seaborn

Analytics Patterns

- **Alpha Pattern**

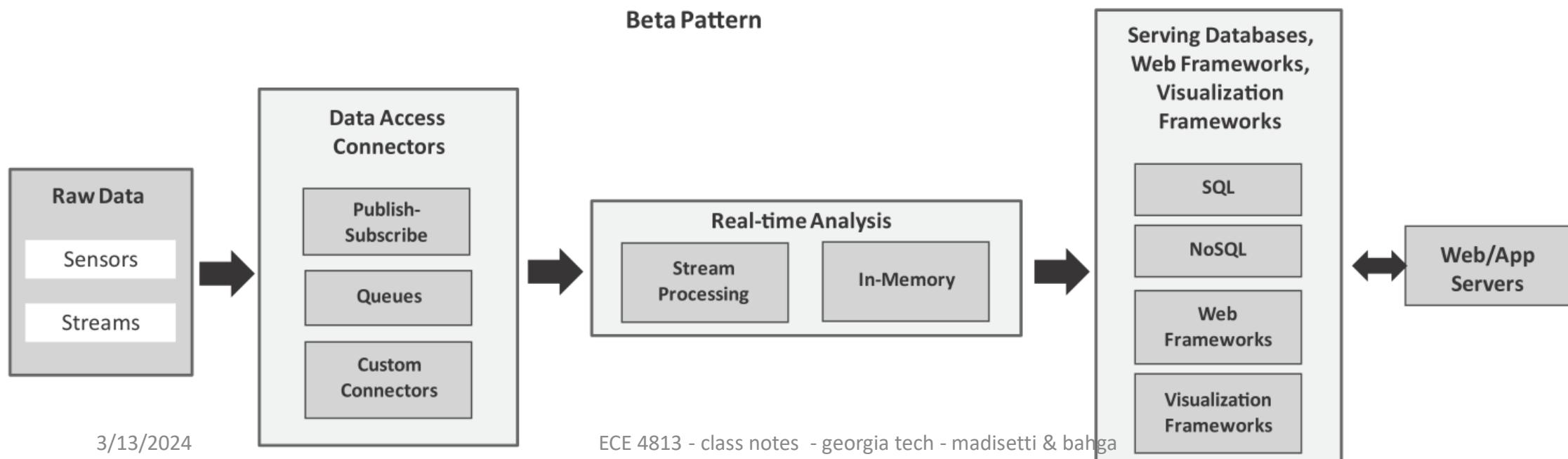
- This pattern can be used for ingesting large volumes of data into a distributed filesystem (such as HDFS) or a NoSQL database (such as HBase) using source-sink connectors (such as Flume) and SQL connectors (such as Sqoop). After the data is moved to the stack, the data can be analyzed in batch mode with batch analysis frameworks including MapReduce (using Hadoop), scripting frameworks (such as Pig), distributed acyclic graph frameworks (such as Spark), machine learning frameworks (such as Spark MLlib). The analysis results are stored either in relational or non-relational databases.





Analytics Patterns

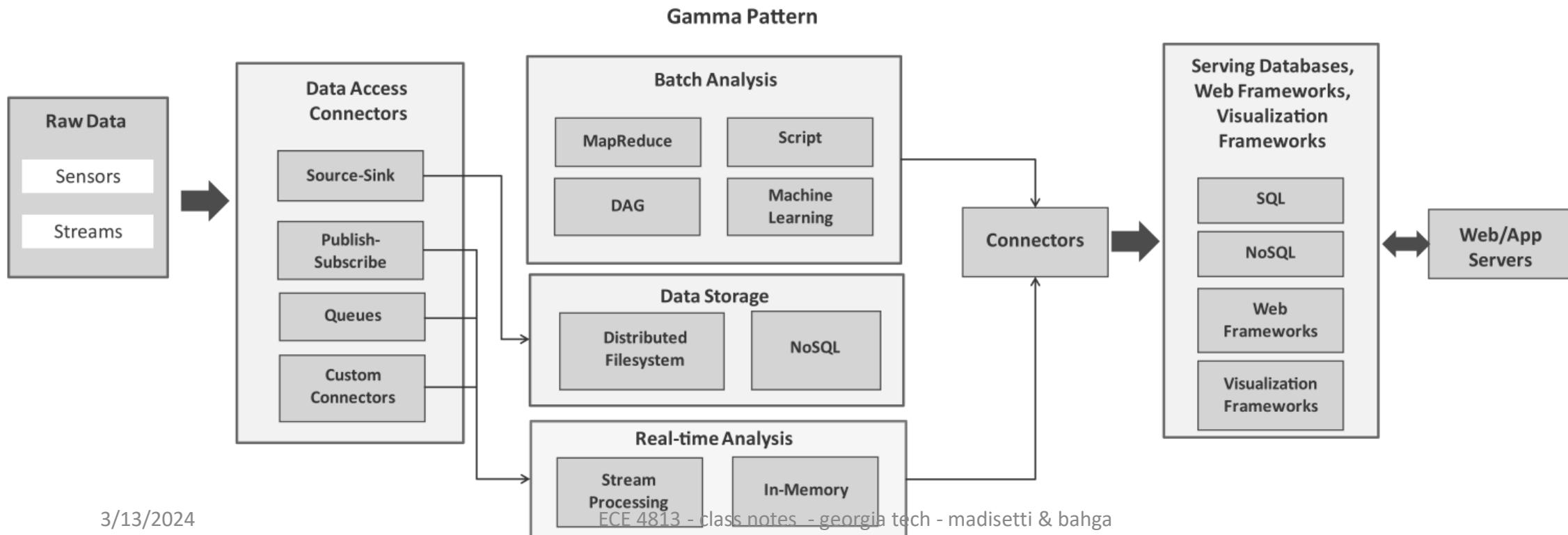
- **Beta Pattern:**
 - This pattern can be used for ingesting streaming data using publish-subscribe messaging frameworks, queues and custom connectors. For real-time analysis, we can use stream processing frameworks (such as Storm) or in-memory processing frameworks (such as Spark).

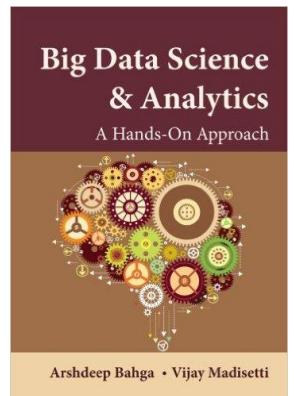


Analytics Patterns

- **Gamma Pattern:**

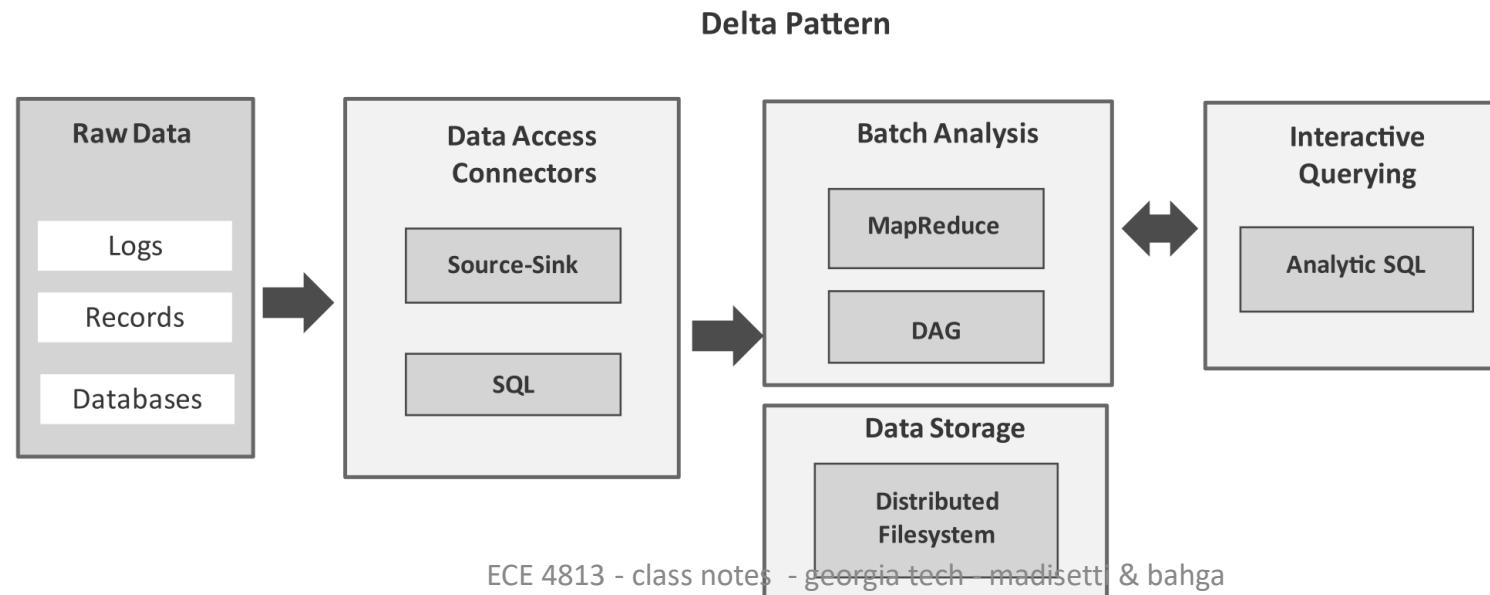
- This pattern is meant for ingesting streaming data into a big data stack and analyzing the data both in real-time and in batch modes.





Analytics Patterns

- **Delta Pattern:**
 - This pattern uses source-sink connectors (such as Flume) or SQL connectors (such as Sqoop) to ingest bulk data into the big data stack. After the data is moved to a distributed filesystem, you can use interactive querying frameworks (such as Hive or Spark SQL) for querying data with SQL-like queries in an interactive mode.



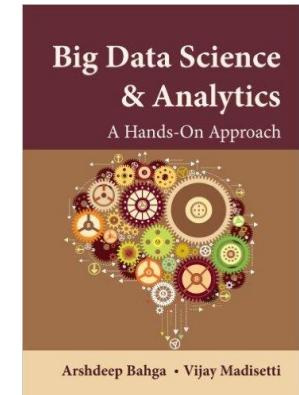
Mapping of blocks in analytics patterns to AWS and Azure services

Data Access Connector	AWS Service	Azure Service
Publish-Subscribe	AWS Kinesis	Azure Event Hubs
Source-Sink	Flume on EMR	Flume on HDInsight
SQL	Sqoop on EMR	Sqoop on HDInsight
Queues	AWS SQS	Azure Queue Service
Custom Connectors	AWS IoT	Azure IoT Hub

Real-time Analysis	AWS Service	Azure Service
Stream Processing	Storm cluster on AWS EC2	Storm on Azure HDInsight, Azure Stream Analytics
In-Memory Processing	Spark on AWS EMR	Spark on Azure HDInsight

Data Storage	AWS Service	Azure Service
Distributed Filesystem	HDFS on AWS EMR	HDFS on Azure HDInsight
NoSQL	AWS DynamoDB	Azure DocumentDB

Serving Databases, Web & Visualization Frameworks	AWS Service	Azure Service
SQL	AWS RDS	Azure SQL DB
NoSQL	AWS DynamoDB	Azure DocumentDB
Web Framework	Django on AWS EC2 instance	Django on Azure Virtual Machines instance
Visualization Framework	Lightning, Pygal, Seaborn, on AWS Virtual Machines instance	Lightning, Pygal, Seaborn, on Azure Virtual Machines instance, Azure Power BI



Complexity Levels for Analytics Patterns

	Analytics Patterns			
Levels	Alpha	Beta	Gamma	Delta
Functionality	Batch Processing	Real time processing	Batch & Real-time processing	Interactive querying
Performance	Process large volumes of data with within timescales of few minutes	Process streaming data on the timescales of few milliseconds to seconds	Combines batch and real-time processing for streaming data, with timescales of milliseconds to seconds for real-time, and few minutes for batch	Interactively query large volumes of data with within timescales of few milliseconds to seconds
Scalability	Scalable to thousands of nodes that can store and process several Petabytes of data	Scalable to thousands of nodes that can process very high throughput streaming data	Scalable to thousands of nodes that can process very high throughput streaming data	Scalable to thousands of nodes that can store and process several Petabytes of data
Fault Tolerance	<ul style="list-style-type: none"> Distributed, fault tolerant and highly available architectures Data replication, automatic failovers 			
Security	<ul style="list-style-type: none"> Apache Ranger: Authorization, authentication, auditing, data encryption, security administration Apache Knox Gateway: REST API gateway, LDAP and Active Directory Authentication, Federation/SSO, Authorization, Auditing 			



Nathan Marz
@nathanmarz

I'm working on a new kind of stream processing system. If this sounds interesting to you, ping me. I want to learn your use cases.

2:17 PM - Dec 14, 2010

2

13

Creator of Apache Storm

RELATED POSTS

[New – AWS PrivateLink for AWS Services: Kinesis, Service Catalog, EC2 Systems Manager, Amazon EC2 APIs, and ELB APIs in your VPC](#)

AWS CloudFormation Feature Updates: Support for Amazon Athena and Coverage Updates for Amazon S3, Amazon RDS, Amazon Kinesis and Amazon CloudWatch

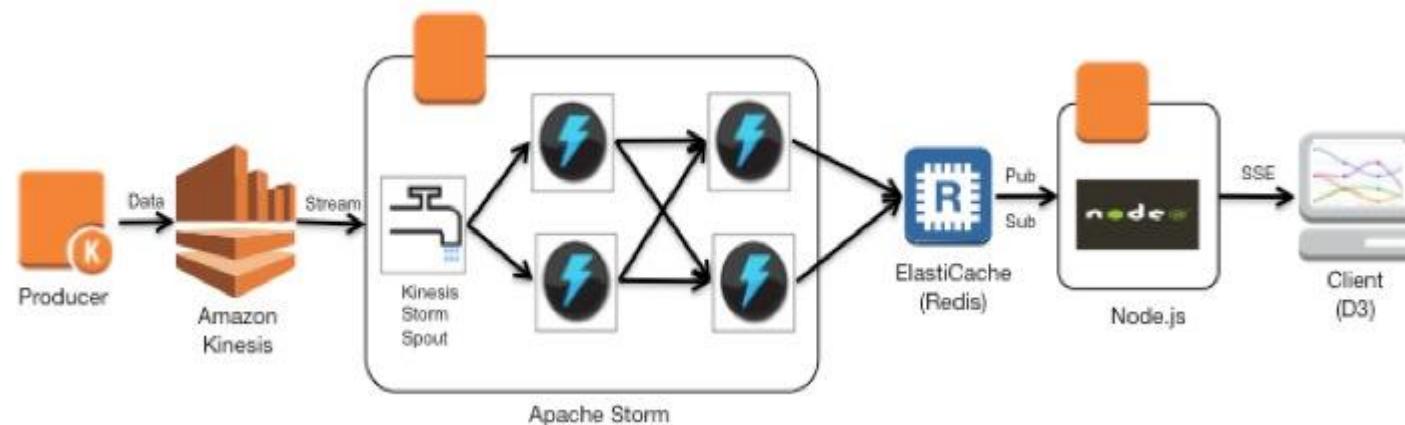
Preprocessing Data in Amazon Kinesis Analytics with AWS Lambda

AWS CloudFormation Supports Amazon Kinesis Analytics Applications

New: Server-Side Encryption for Amazon Kinesis Streams

Perform Stream and Am

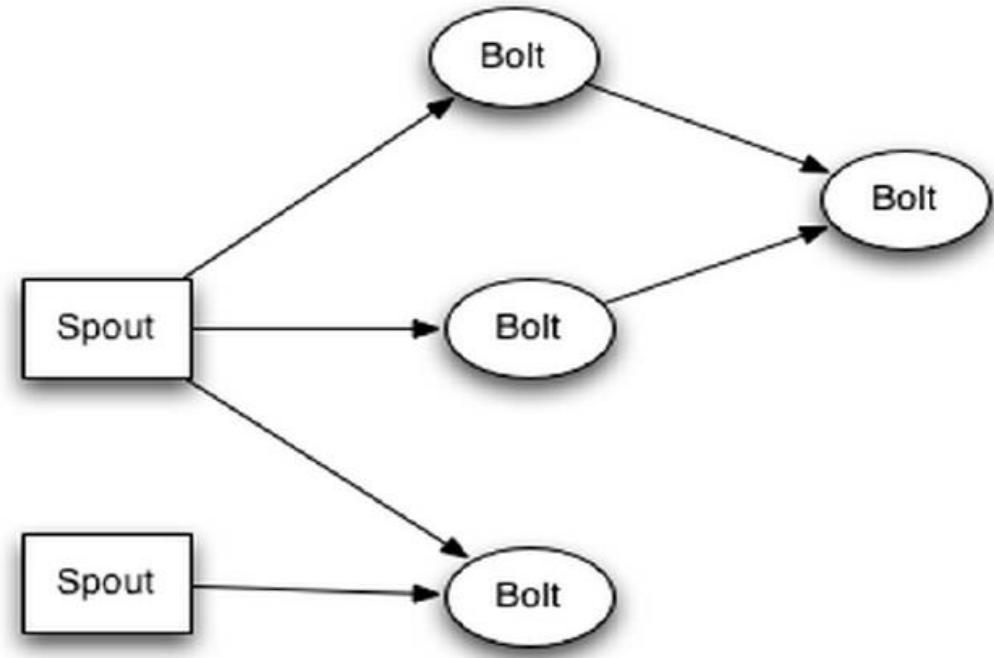
<https://aws.amazon.com/blogs/big-data/implement-a-real-time-sliding-window-application-using-amazon-kinesis-and-apache-storm/>



Storm Processing Flow

- **Topology**

- Graph of computation – can have cycles
- Network of Spouts and Bolts
- Spouts and bolts execute as many tasks across the cluster



<http://storm.apache.org/releases/current/Tutorial.html>

Use Case - Trending topics on Twitter

9:21-9:29 PM

"Take your broken heart, make it into art" -- **#CarrieFisher** via **#MerylStreep**

find someone who looks at you the way everyone looks at **meryl streep**

Can't stop watching this lady having **NONE** of **Meryl Streep**. 😂😂😂
#GoldenGlobes 🎉

Meryl Streep's speech at the **#GoldenGlobes** 🎉 is definitely one of the biggest highlights of the night.

Celebrities couldn't contain their love for **Meryl Streep** after **#GoldenGlobes** 🎉 speech huff.to/2j8ajZ

Goodbye **#MerylStreep** you just lost 60 million fans. Get use to movie flops.



9:30 PM

Meryl Streep
Moonlight
#HiddenFences
Casey Affleck
Donald Glover
#ThingsGhostsProbablyDo
Emma Stone
Brad Pitt
#ILaughAboutItBecause
#ThingsNerdsSay

Top keywords

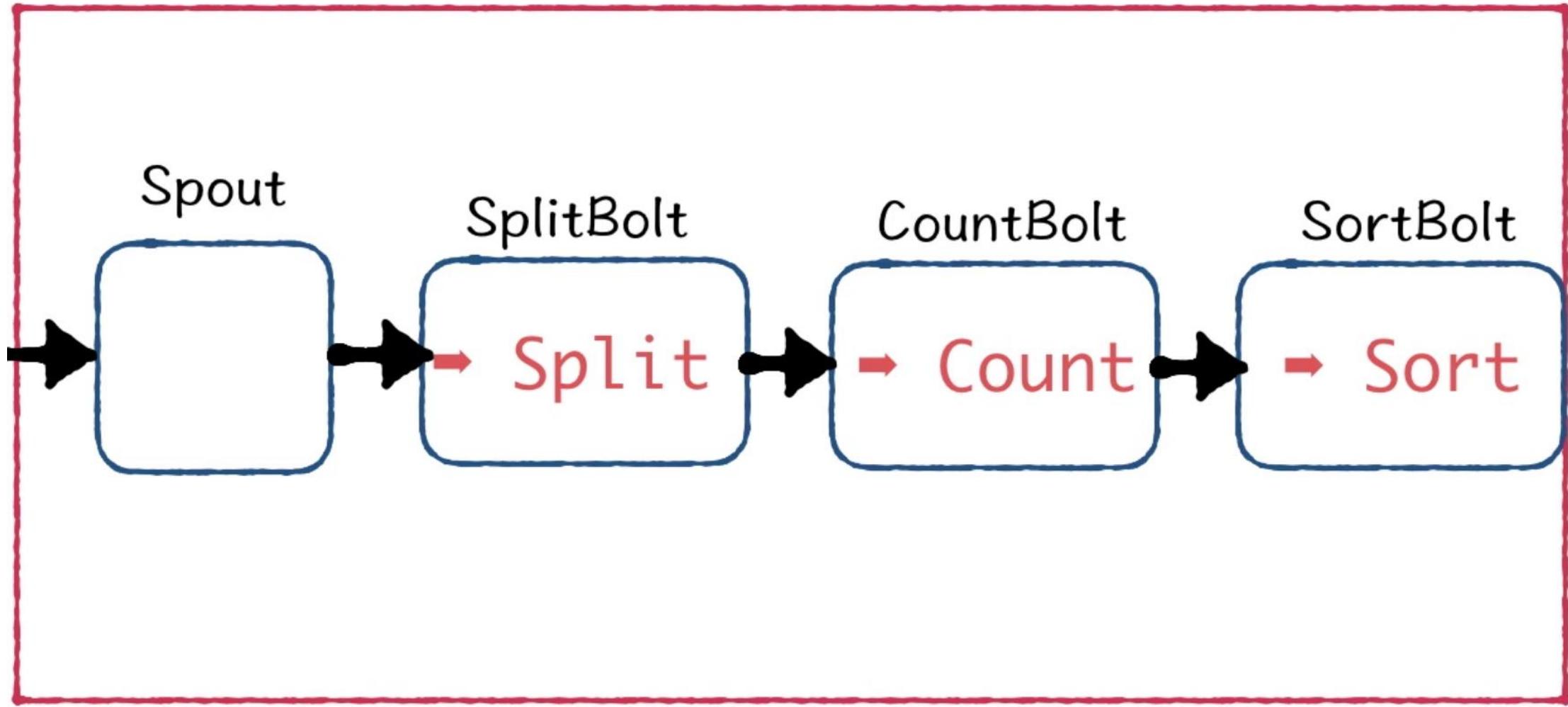
Algorithm Flow

App

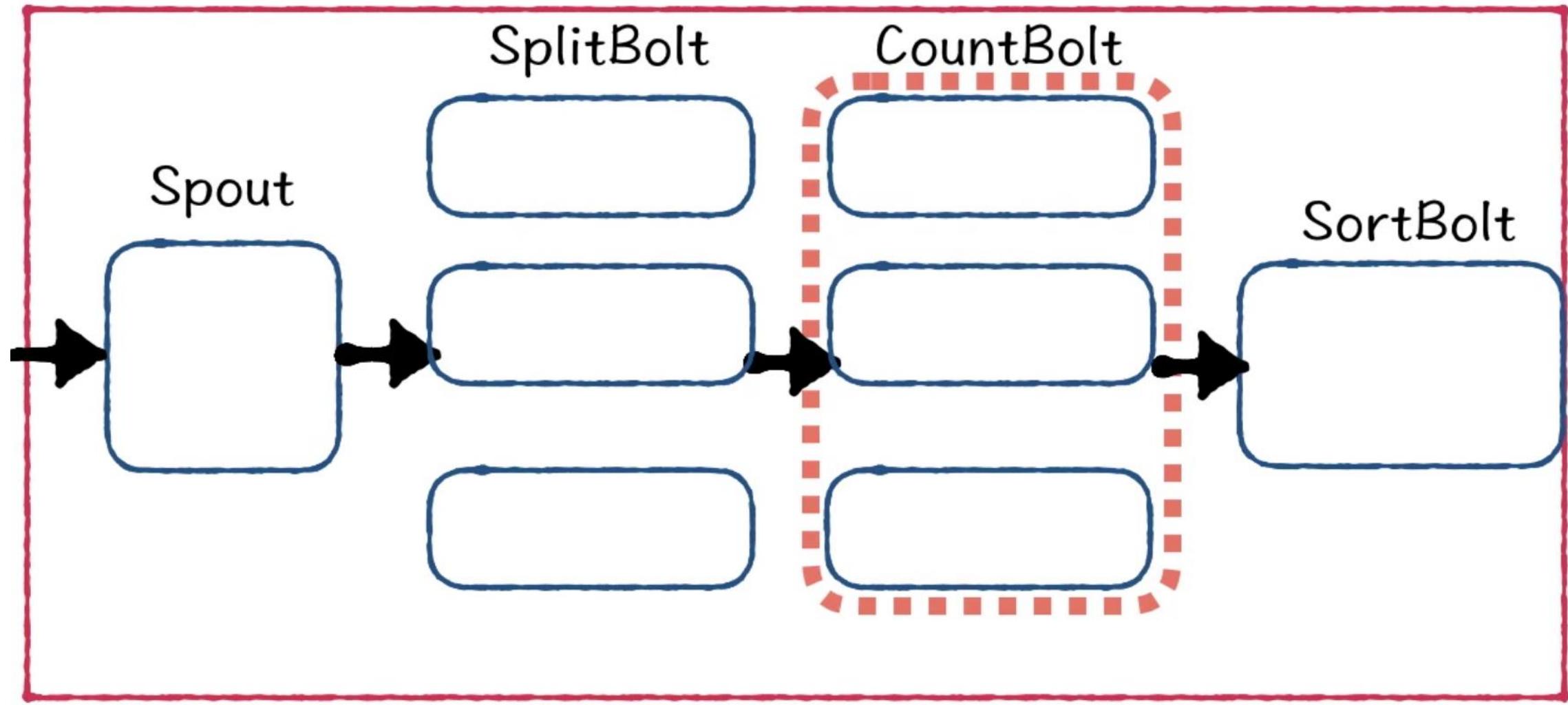
MerylStreep 15000
GoldenGlobes 10000
CarrieFisher 2000

- Split tweets into words
- Find counts for each word in last 10 mins
- Sort words by count

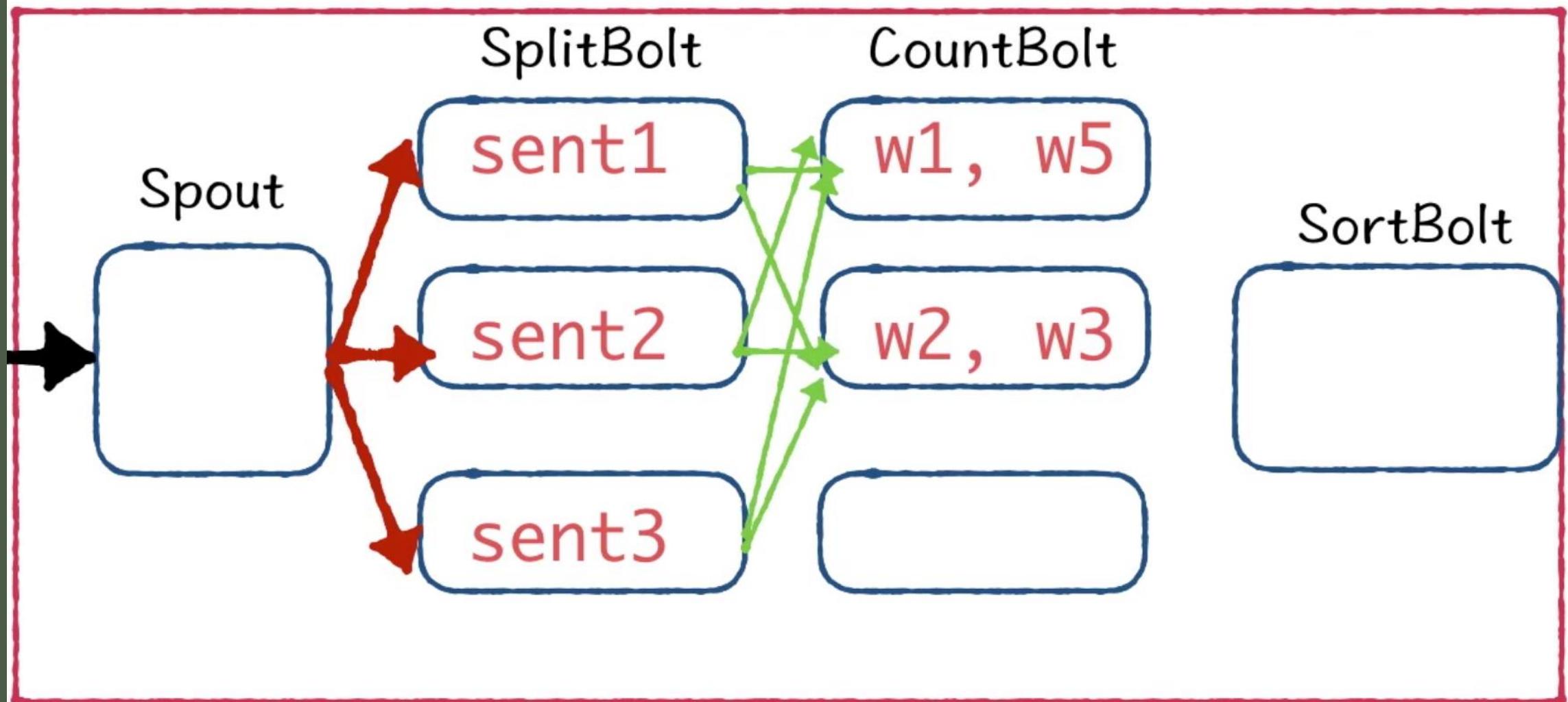
Trending Tweets Topology



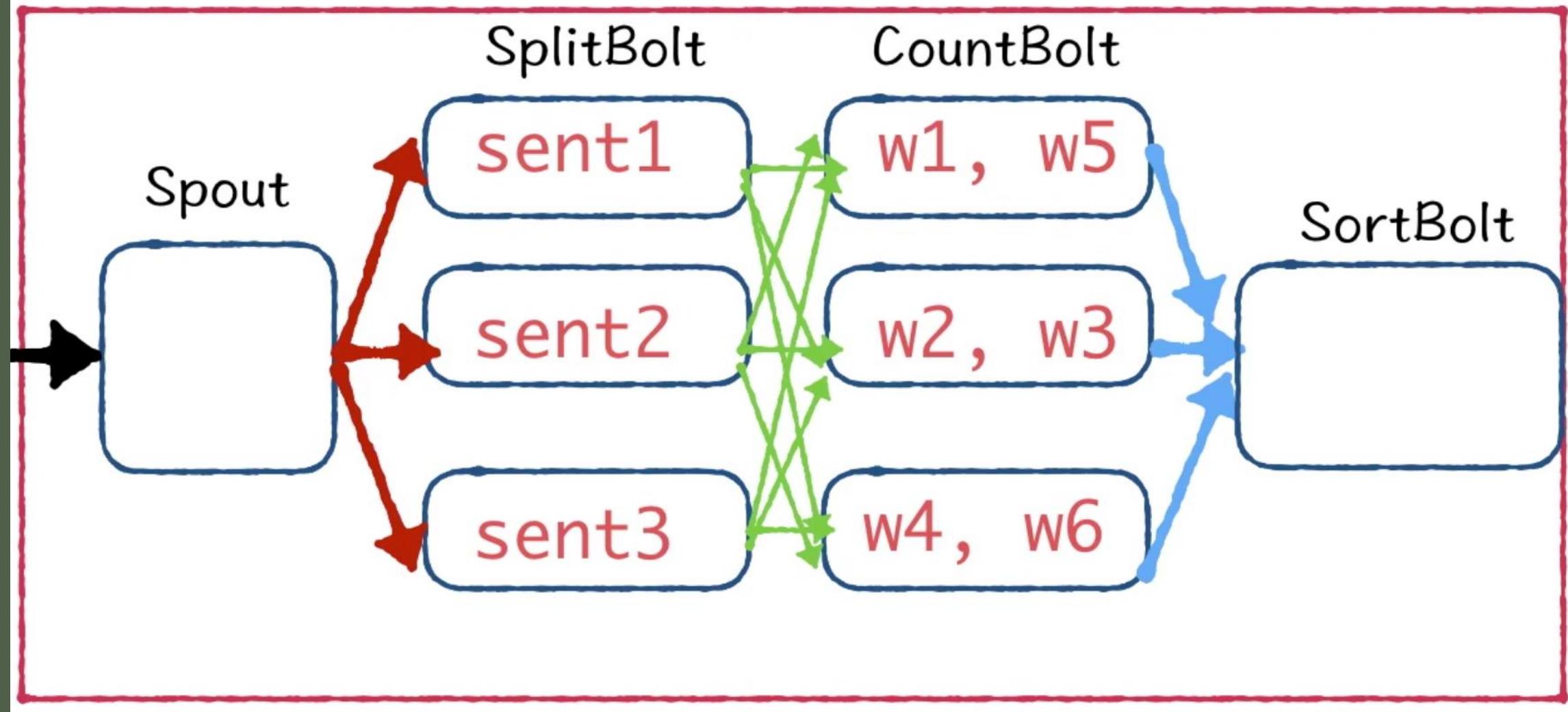
Trending Tweets Topology



Trending Tweets Topology



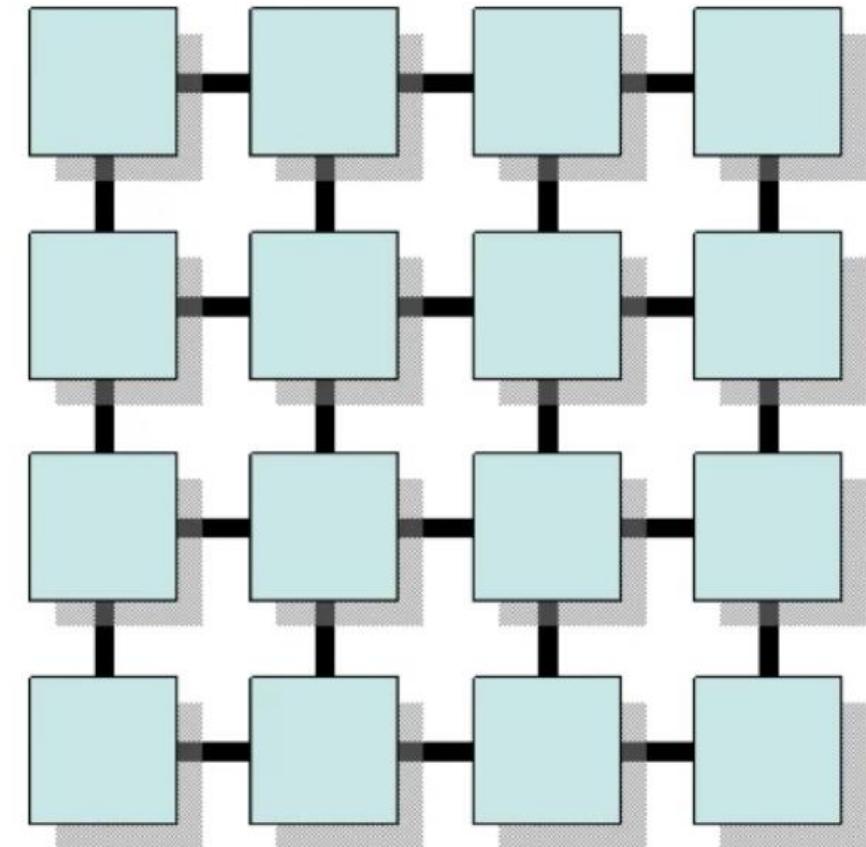
Trending Tweets Topology



Apache Storm

The processes can run
on different machines

Remote Mode



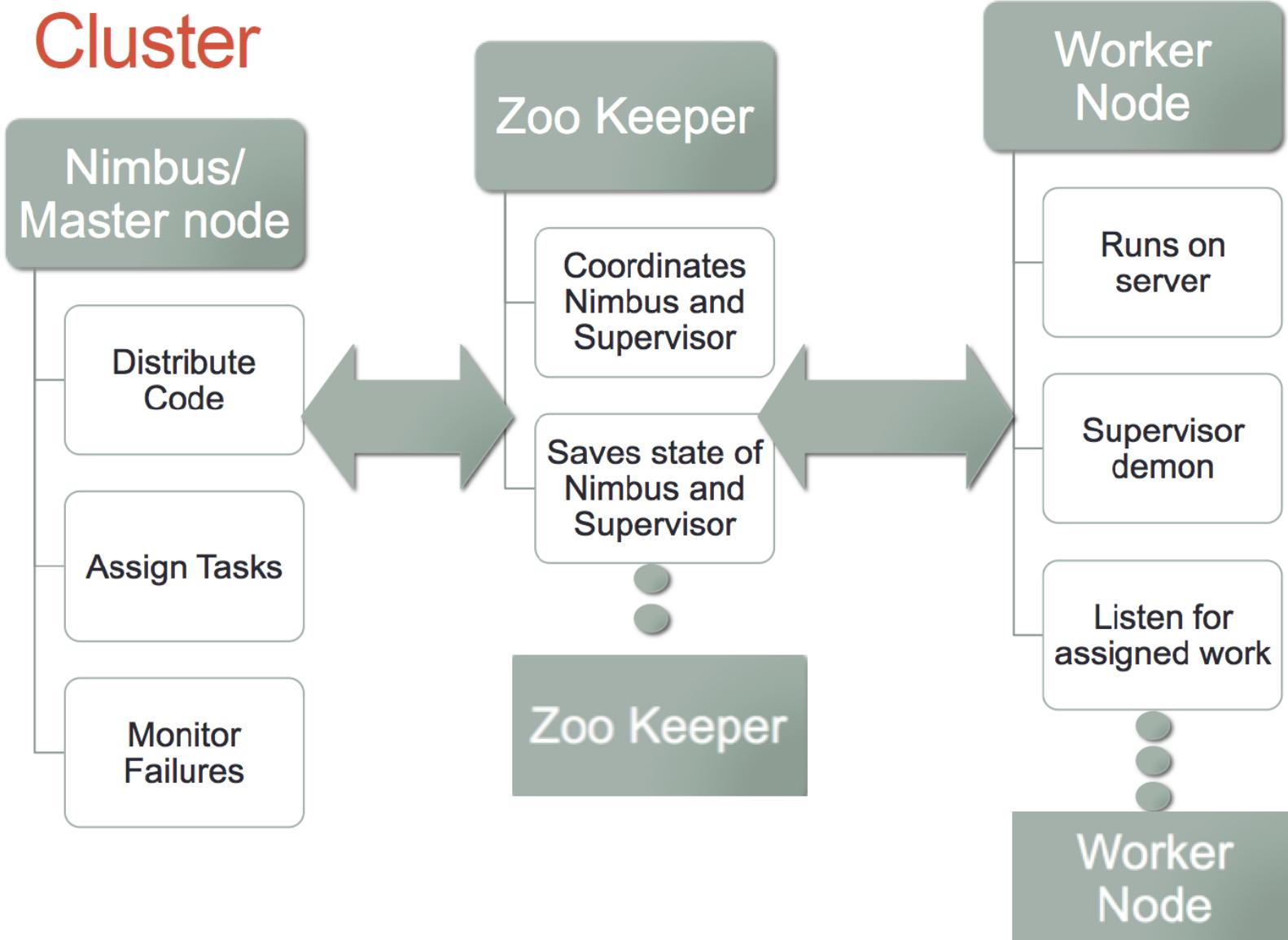
Apache Storm

Nimbus

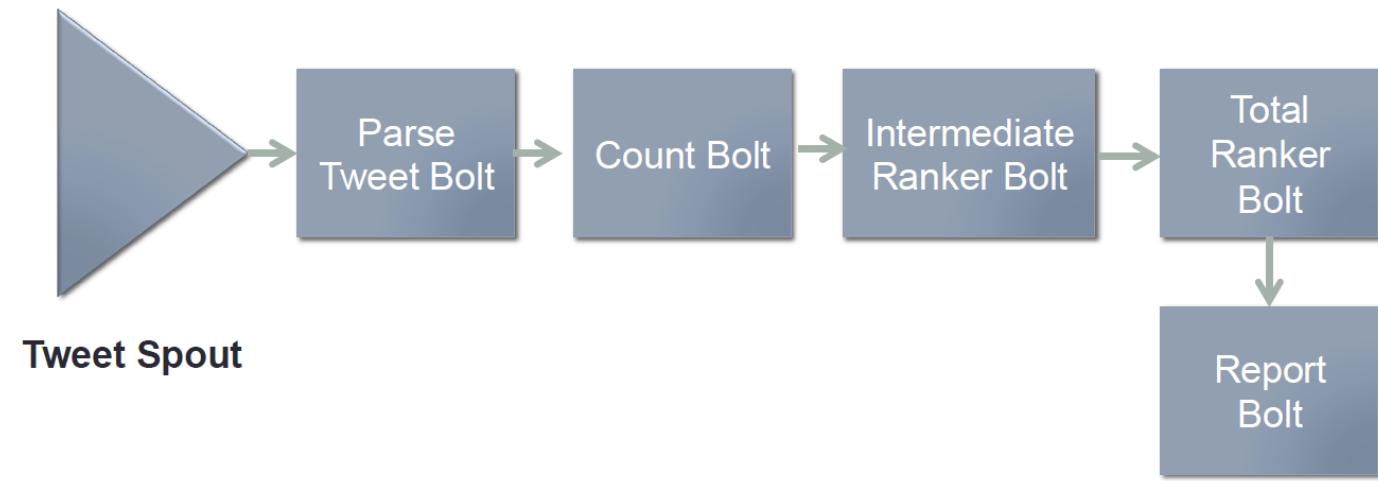
The Nimbus will
start up a bunch of
processes for
Spouts and Bolts

**Storm
Distributed
Processing
Architecture**

Cluster

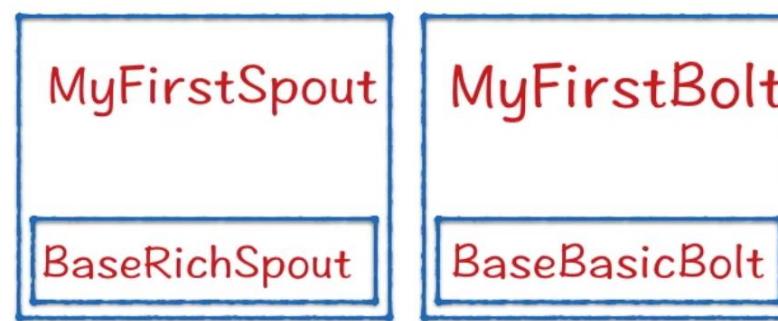


Generic
Flow



Use Case 1

Hello World Topology



TopologyMain

- Build topology
- Submit topology to a cluster

Spout Architecture

D:\Dropbox (VijayMadisetti)\VKM-Files-Private\GT\Cloud-Computing-Fall2017\Storm\SourceCode\SourceCode1-24\src\main\java\examples\ex1_HelloWorld\myFirstSpout.java - Sublime Text
Edit Selection Find View Goto Tools Project Preferences Help

```
1 package examples.ex1_Helloworld;
2
3
4 import org.apache.storm.spout.SpoutOutputCollector;
5 import org.apache.storm.task.TopologyContext;
6 import org.apache.storm.topology.OutputFieldsDeclarer;
7 import org.apache.storm.topology.base.BaseRichSpout;
8 import org.apache.storm.tuple.Fields;
9 import org.apache.storm.tuple.Values;
0
1
2 import java.util.Map;
3
4
5 public class myFirstSpout extends BaseRichSpout {
6
7     private SpoutOutputCollector collector;
8
9     private Integer i = 0;
10
11     public void open(Map conf, TopologyContext context,
12                     SpoutOutputCollector collector) {
13
14         this.collector = collector;
15     }
16
17     public void nextTuple() {
```

Big Five Methods

```
26     public void nextTuple() {
27
28         this.collector.emit(new Values(this.i));
29         this.i=this.i+1;
30     }
31
32
33     public void declareOutputFields(OutputFieldsDeclarer declarer) {
34         declarer.declare(new Fields("field"));
35     }
36
37     public void ack(Object msgId) {}
```

- **open** : Provides spout with an environment to execute.
`open(Map conf, TopologyContext context, SpoutOutputCollector collector)`
- **nextTuple** : Emits the generated data.
`nextTuple()`
- **close** : call when a spout is going to close.
`close()`
- **declareOutputFields** : Declares the output schema of the tuple.
`declareOutputFields(OutputFieldsDeclarer declarer)`
- **ack** : Acknowledges that a specific tuple is processed
`ack(Object msgId)`
- **fail** : Specifies that a specific tuple is not processed

Bolt Architecture

D:\Dropbox (VijayMadisetti)\VKM-Files-Private\GT\Cloud-Computing-Fall2017\Storm\SourceCode\SourceCode1-24\src\main\java\examples\ex1_HelloWorld\myFirstBolt.java - Sublime Text

```
1 package examples.ex1_Helloworld;
2
3
4 import org.apache.storm.topology.BasicOutputCollector;
5 import org.apache.storm.topology.OutputFieldsDeclarer;
6 import org.apache.storm.topology.base.BaseBasicBolt;
7 import org.apache.storm.tuple.Fields;
8 import org.apache.storm.tuple.Tuple;
9 import org.apache.storm.tuple.Values;
10
11 public class myFirstBolt extends BaseBasicBolt {
12     public void execute(Tuple input, BasicOutputCollector collector) {
13         collector.emit(new Values(input.getInteger(0)*2));
14     }
15
16     public void declareOutputFields(OutputFieldsDeclarer declarer) {
17         declarer.declare(new Fields("field"));
18     }
19
20     public void cleanup() {}
21 }
```

3/13/2024

- **Prepare** : Provides the bolt with an environment to execute.
prepare(Map conf, TopologyContext context, OutputCollector collector)
- **Execute** : Process a single tuple of input.
execute(Tuple tuple)
- **Cleanup** : When a bolt is going to shutdown.
cleanup()
- **declareOutputFields** : Schema of the tuple.
declareOutputFields(OutputFieldsDeclarer declarer)

Topology of storm

```
1 package examples.ex1_HelloWorld;
2
3
4 import org.apache.storm.Config;
5 import org.apache.storm.LocalCluster;
6 import org.apache.storm.topology.TopologyBuilder;
7
8
9 public class TopologyMain {
10     public static void main(String[] args) throws InterruptedException {
11
12         //Build Topology
13         TopologyBuilder builder = new TopologyBuilder();
14         builder.setSpout("My-First-Spout", new myFirstSpout());
15         builder.setBolt("My-First-Bolt", new myFirstBolt()).shuffleGrouping("My-First-Spout");
16
17         //Configuration
18         Config conf = new Config();
19         conf.setDebug(true);
20
21
22         //Submit Topology to cluster
23         LocalCluster cluster = new LocalCluster();
24         try{
25             cluster.submitTopology("My-First-Topology", conf, builder.createTopology());
26             Thread.sleep(1000);
27
28         finally{
29             cluster.shutdown();}
```

Reading from a File

```
sample.txt *
```

```
1 this
2 is
3 happening
4 on
5 storm
6 a
7 new
8 streaming
9 storm
10 application
11 cool
12 application
```

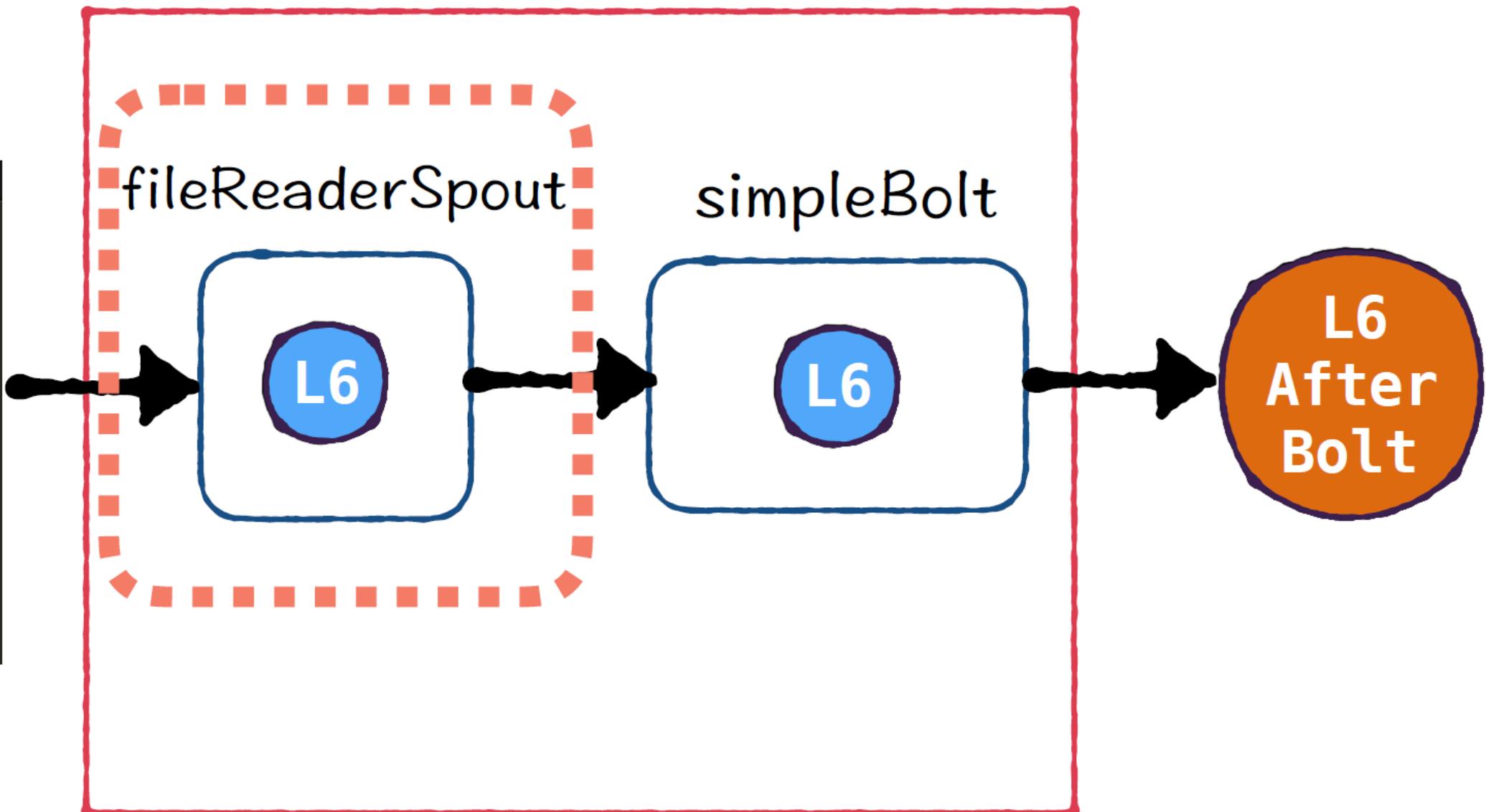


Processed
text

Use Case 2

Reading from a File

```
sample.txt
1 this
2 is
3 happening
4 on
5 storm
6 a
7 new
8 streaming
9 storm
10 application
11 cool
12 application
```



Bolt Architecture

```
1 package examples.ex2_readFromFile;
2
3 import org.apache.storm.task.OutputCollector;
4 import org.apache.storm.topology.BasicOutputCollector;
5 import org.apache.storm.topology.OutputFieldsDeclarer;
6 import org.apache.storm.topology.base.BaseBasicBolt;
7 import org.apache.storm.tuple.Fields;
8 import org.apache.storm.tuple.Tuple;
9 import org.apache.storm.tuple.Values;
10
11 public class simpleBolt extends BaseBasicBolt {
12
13     public void cleanup() {}
14
15
16     public void execute(Tuple input,BasicOutputCollector collector) {
17         collector.emit(new Values(input.getString(0)+" After Bolt"));
18     }
19
20     public void declareOutputFields(OutputFieldsDeclarer declarer) {
21         declarer.declare(new Fields("word"));
22     }
23 }
```

```

1 package examples.ex2_readFromFile;
2
3
4 import org.apache.storm.spout.SpoutOutputCollector;
5 import org.apache.storm.task.TopologyContext;
6 import org.apache.storm.topology.OutputFieldsDeclarer;
7 import org.apache.storm.topology.base.BaseRichSpout;
8 import org.apache.storm.tuple.Fields;
9 import org.apache.storm.tuple.Values;
10
11 import java.io.BufferedReader;
12 import java.io.FileNotFoundException;
13 import java.io.FileReader;
14 import java.util.ArrayList;
15 import java.util.List;
16 import java.util.Map;
17
18 public class fileReaderSpout extends BaseRichSpout {
19
20     private SpoutOutputCollector collector;
21     private boolean completed = false;
22     private FileReader fileReader;
23     private String str;
24     private BufferedReader reader ;
25

```

```

27
28     public void open(Map conf, TopologyContext context,
29                      SpoutOutputCollector collector) {
30         try {
31             this.FileReader = new FileReader(conf.get("fileToRead").toString());
32         } catch (FileNotFoundException e) {
33             throw new RuntimeException("Error reading file ["+conf.get("wordFile")+"]");
34         }
35
36         this.collector = collector;
37         this.reader = new BufferedReader(fileReader);
38
39     }
40
41
42     public void nextTuple() {
43
44         if (!completed) {
45             try {
46
47                 this.str = reader.readLine();
48                 if (this.str != null) {
49
50                     this.collector.emit(new Values(str));
51                 } else {
52                     completed = true;
53                     fileReader.close();
54                 }
55
56             } catch (Exception e) {
57                 throw new RuntimeException("Error reading tuple", e);
58             }
59         }
60     }
61
62
63     public void declareOutputFields(OutputFieldsDeclarer declarer) {
64         declarer.declare(new Fields("word"));
65     }
66
67     public void ack(Object msgId) {}
68
69
70     public void fail(Object msgId) {}
71
72

```

Spout Architecture

Storm Topology

```
myFirstBoltJava  myFirstSpoutJava  TopologyMain.java — ex1_HelloWorld  fileReaderSpoutJava  simpleBolt.java  TopologyMain.java — ex2_readFromFile
1 package examples.ex2_readFromFile;
2
3
4 import org.apache.storm.Config;
5 import org.apache.storm.LocalCluster;
6 import org.apache.storm.topology.TopologyBuilder;
7
8
9 public class TopologyMain {
10     public static void main(String[] args) throws InterruptedException {
11
12         //Topology definition
13         TopologyBuilder builder = new TopologyBuilder();
14         builder.setSpout("File-Reader-Spout", new fileReaderSpout());
15         builder.setBolt("Simple-Bolt", new simpleBolt()).shuffleGrouping("File-Reader-Spout");
16
17         //Configuration
18         Config conf = new Config();
19         conf.setDebug(true);
20         conf.put("fileToRead", "/Users/vkm/Desktop/sample.txt");
21
22         LocalCluster cluster = new LocalCluster();
23         try{cluster.submitTopology("File-Reader-Topology", conf, builder.createTopology());
24             Thread.sleep(1000);
25         }finally{
26             cluster.shutdown();
27     }
28 }
```

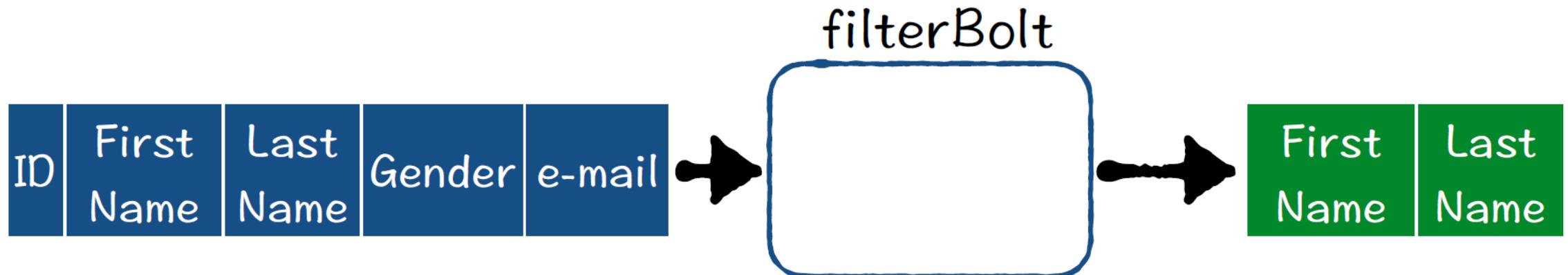
Use Case 3

Fields

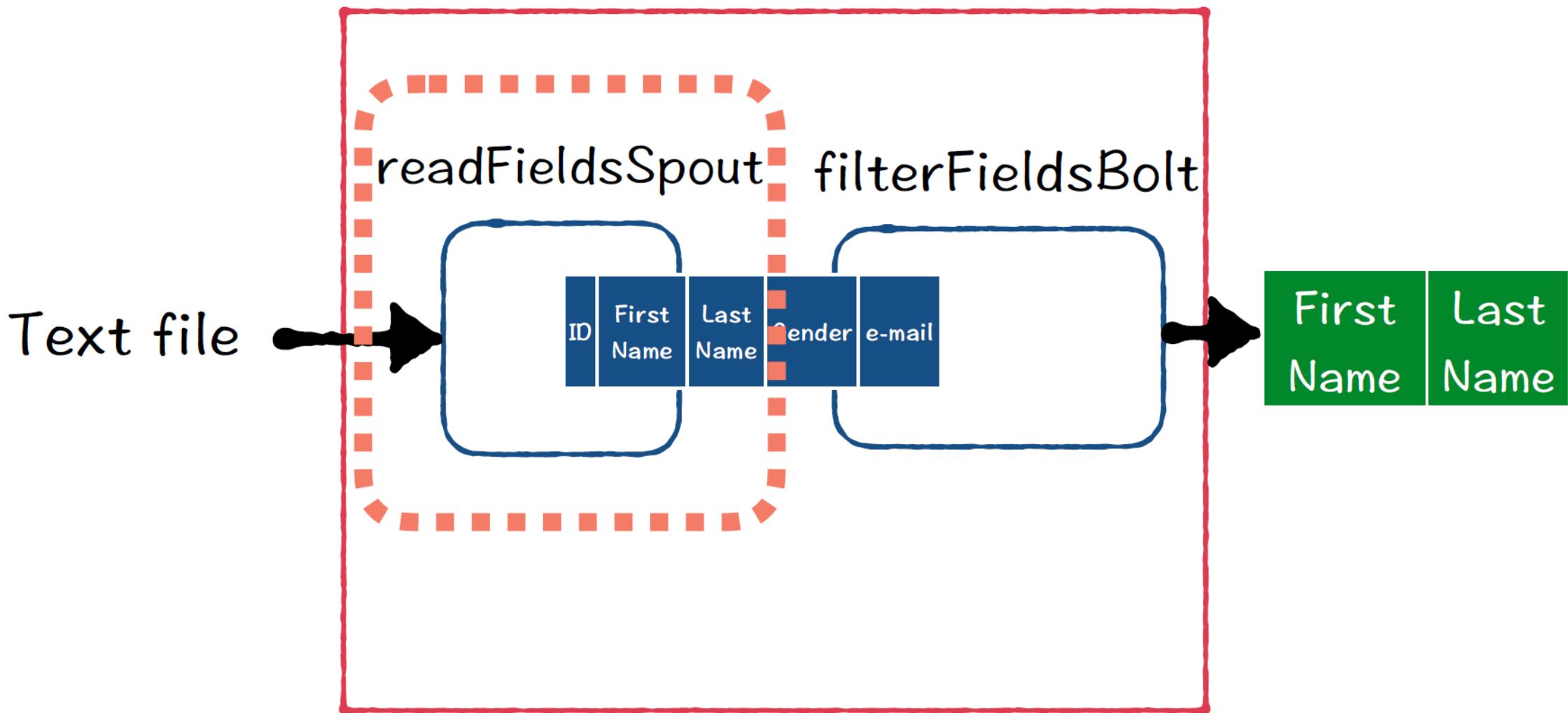
Values

ID	First Name	Last Name	Gender	e-mail
1	Jens	Sue	F	j@gmail

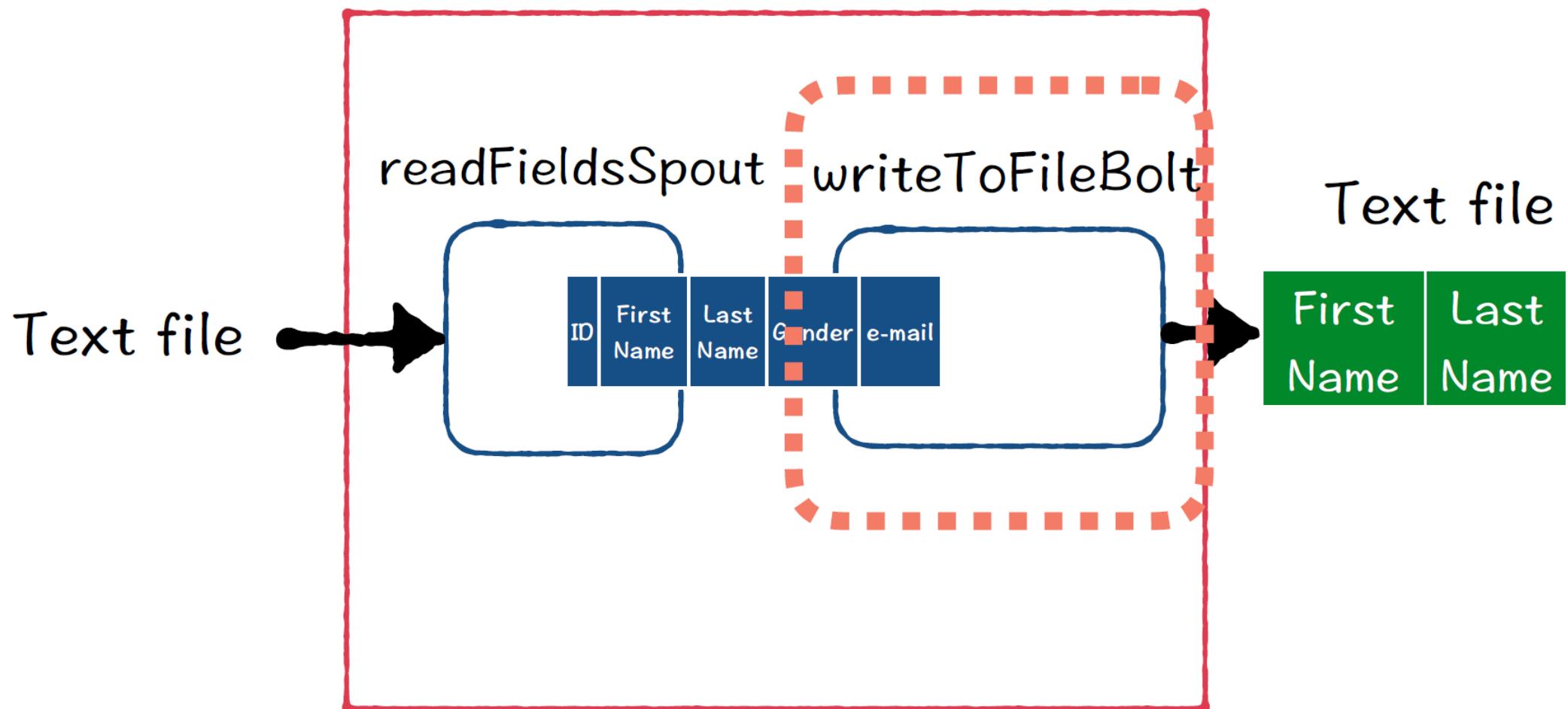
Output fields are fixed for a Storm Component



Extracting Fields



Writing to File



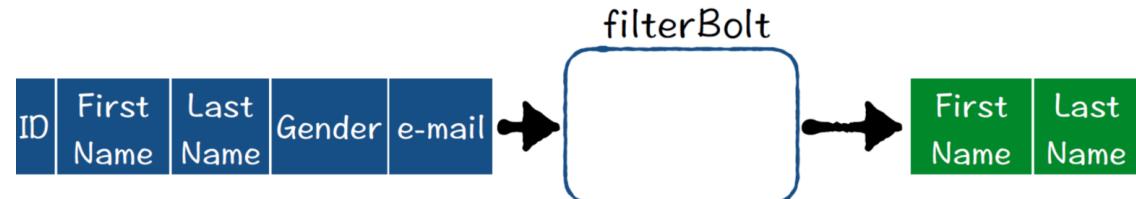
Portion of Storm Spout, `readFieldsSpout`, Generating Fields & Values

Two of the
big Five

```
43    public void nextTuple() {
44
45        if (!completed) {
46            try {
47
48                this.str = reader.readLine();
49                if (this.str != null) {
50
51                    this.collector.emit(new Values(str.split(",")));
52                } else {
53                    completed = true;
54                    fileReader.close();
55                }
56
57            } catch (Exception e) {
58                throw new RuntimeException("Error reading tuple", e);
59            }
60        }
61    }
62
63
64    public void declareOutputFields(OutputFieldsDeclarer declarer) {
65        declarer.declare(new Fields("id","first_name","last_name","gender","email"));
66    }
67}
```

filterBolt for Filtering Fields from Input

```
Selection Find View Goto Tools Project Preferences Help
myFirstBolt.java myFirstSpout.java TopologyMain.java --- ex1_HelloWorld --- fileReaderSpout.java simpleBolt.java TopologyMain.java --- ex2_readFromFile --- filterFieldsBolt.java
1 package examples.ex3_accessingFields;
2
3 import org.apache.storm.task.OutputCollector;
4 import org.apache.storm.topology.BasicOutputCollector;
5 import org.apache.storm.topology.OutputFieldsDeclarer;
6 import org.apache.storm.topology.base.BaseBasicBolt;
7 import org.apache.storm.tuple.Fields;
8 import org.apache.storm.tuple.Tuple;
9 import org.apache.storm.tuple.Values;
10
11 public class filterFieldsBolt extends BaseBasicBolt {
12
13     public void cleanup() {}
14
15     public void execute(Tuple input,BasicOutputCollector collector) {
16
17         String firstName = input.getStringByField("first_name");
18         String lastName = input.getString(2);
19
20             collector.emit(new Values(firstName,lastName));
21     }
22
23     public void declareOutputFields(OutputFieldsDeclarer declarer) {
24         declarer.declare(new Fields("first_name","last_name"));
25     }
26 }
```



Apache Storm Support for Multiple Languages

See Lab 8

- <http://storm.apache.org/releases/current/Multilang-protocol.html>
- <https://docs.microsoft.com/en-us/azure/hdinsight/storm/apache-storm-develop-python-topology>

Defining Bolts in other languages

Bolts can be defined in any language. Bolts written in another language are executed as subprocesses, and Storm communicates with those subprocesses with JSON messages over stdin/stdout. The communication protocol just requires an ~100 line adapter library, and Storm ships with adapter libraries for Ruby, Python, and Fancy.

Here's the definition of the `SplitSentence` bolt from `WordCountTopology`:

```
public static class SplitSentence extends ShellBolt implements IRichBolt {
    public SplitSentence() {
        super("python", "splitsentence.py");
    }

    public void declareOutputFields(OutputFieldsDeclarer declarer) {
        declarer.declare(new Fields("word"));
    }
}
```

`SplitSentence` overrides `ShellBolt` and declares it as running using `python` with the arguments `splitsentence.py`. Here's the implementation of `splitsentence.py`:

```
import storm

class SplitSentenceBolt(storm.BasicBolt):
    def process(self, tup):
        words = tup.values[0].split(" ")
        for word in words:
            storm.emit([word])

SplitSentenceBolt().run()
```

For more information on writing spouts and bolts in other languages, and to learn about how to create topologies in other languages (and avoid the JVM completely), see [Using non-JVM languages with Storm](#).

COLLECT

STORE

ETL

PROCESS / ANALYZE

CONSUME

Applications

Web apps



Mobile apps



Data centers



AWS Direct Connect

Logging



AWS Import/Export



Snowball

Messaging



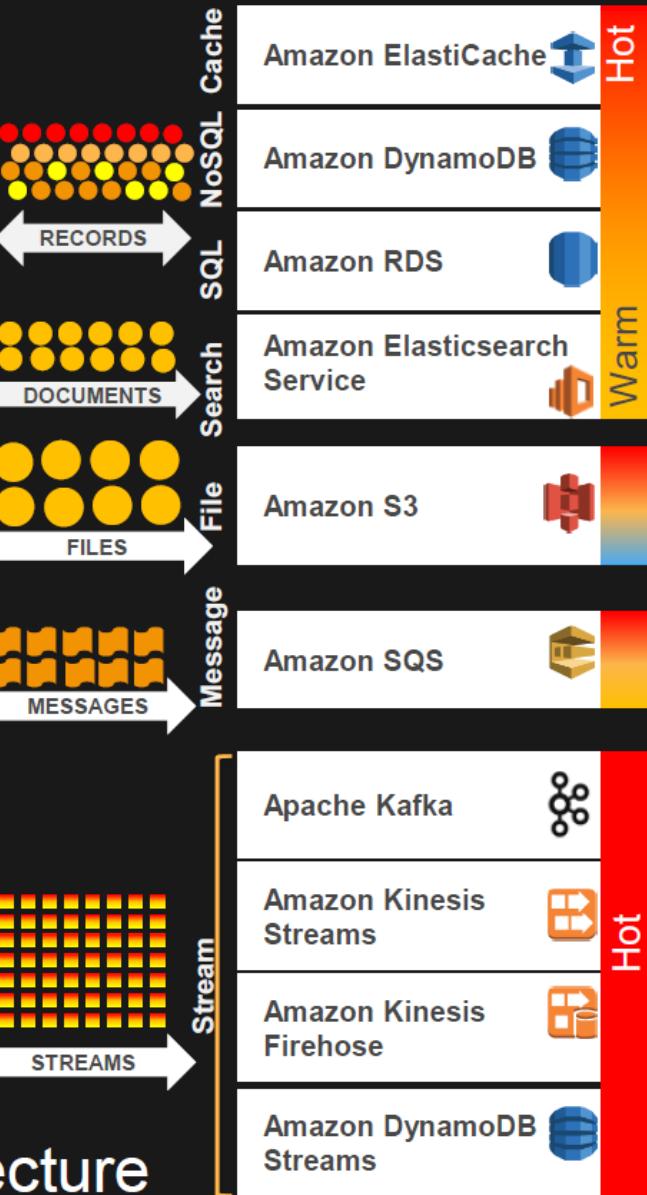
Devices



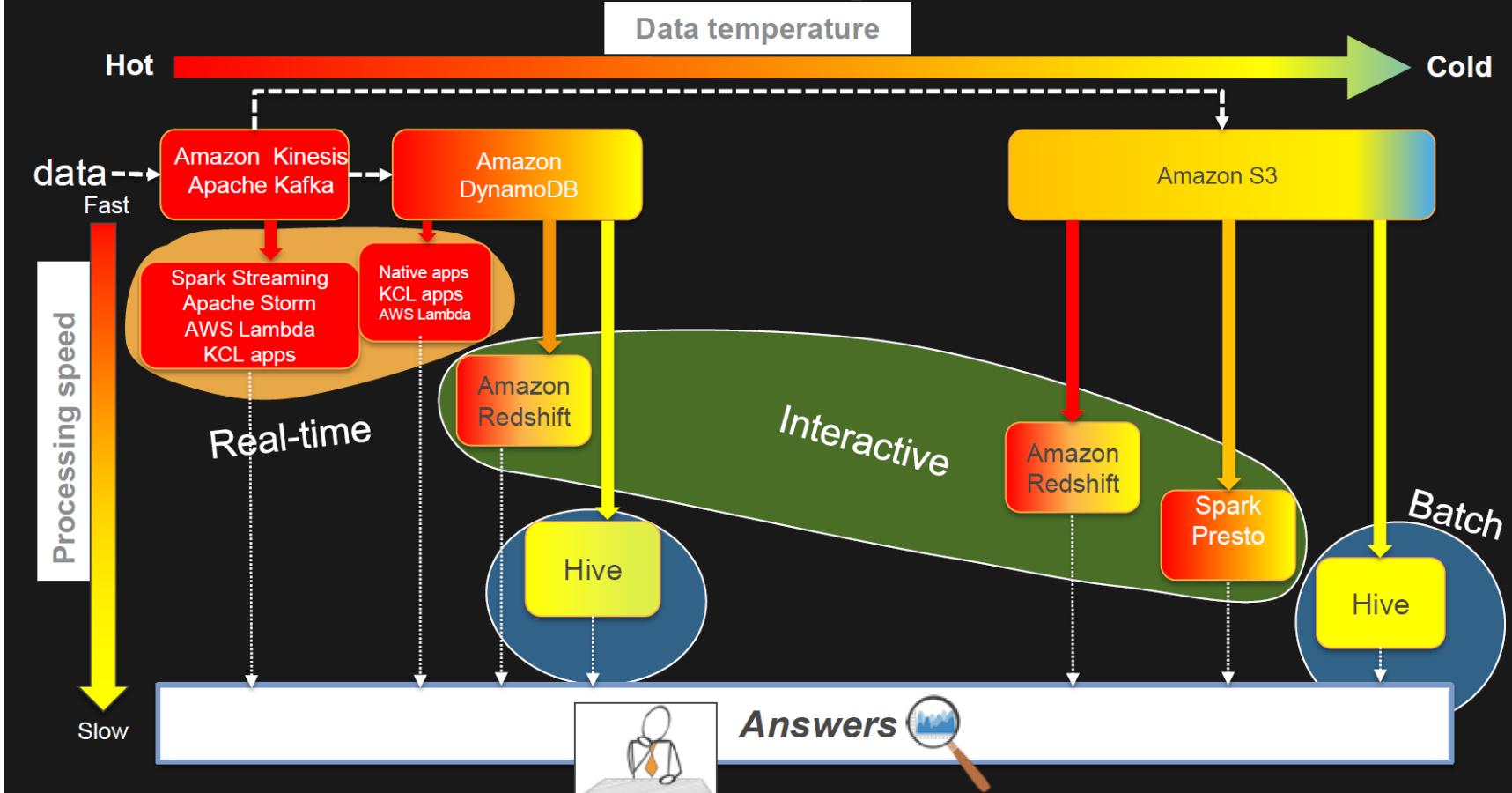
Sensors & IoT platforms



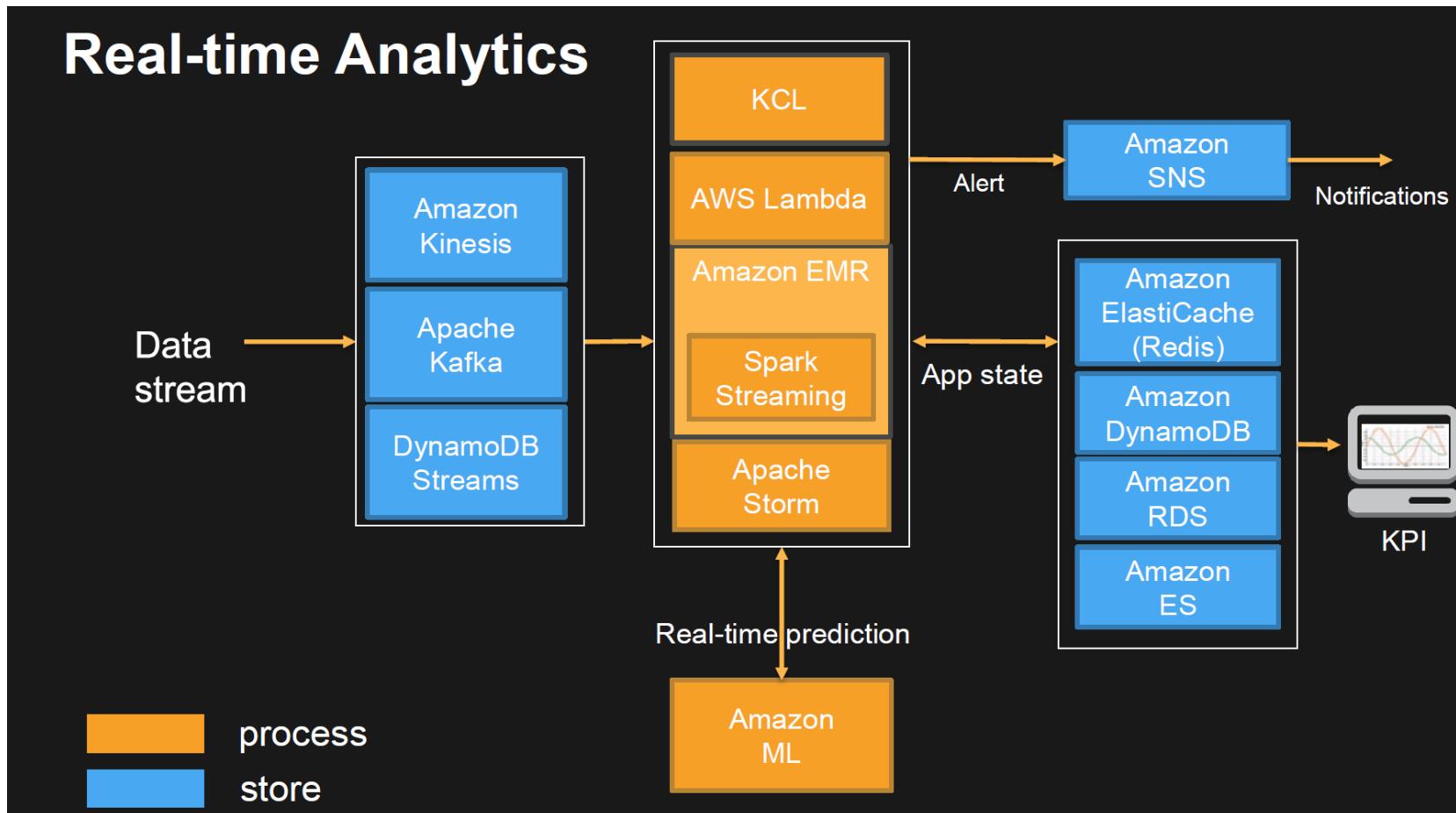
AWS IoT



Data Temperature vs. Processing Speed



Summary - Streaming Implementation Architecture



Summary

- You have looked at the big data stack
- You have look at different patterns for analytics – alpha, beta, ..
- You looked at Storm for stream processing.

