# 1 Overview

Welcome to the Storm tutorial. This tutorial covers the critical skills needed to develop a Storm application in Java.

**Objectives**

Upon completing this tutorial, students will be able to:

- Set up the Storm installation

- Develop a simple Storm application (Word Count)

- Compile Storm applications

- Run Storm applications and examine the output

**Structure**

This guide is designed as a set of step-by-step instructions for hands-on exercises. It teaches you how to operate the MapReduce service in a functional test environment through a series of examples. Exercises should be completed in the order described in the following steps:

1. Download and Setup Storm

2. Develop a simple "Word Count" application and build it

3. Run the application on Storm and examine the output

# 2. Requirements

This tutorial is designed to work on the sample Docker image that we provide. This application will be based on **JDK 8**

# 3. Prepare the Docker image and install Storm

**Step 1:** Start the "default" Docker machine that you created when following the "Tutorial: Docker installation" in week 4, run:

```
1   docker-machine start default
2   docker-machine env
3   # follow the instruction to configure your shell: eval $(...)
```

**Step 2:** Download the Dockerfile and related files for this tutorial, change the current folder, build, and run the docker image, run:

```
1   git clone https://github.com/UIUC-public/tutorial_docker.git
2   cd tutorial_docker
3   docker build -t tutorial_docker .
4   docker run -it tutorial_docker bin/bash
```

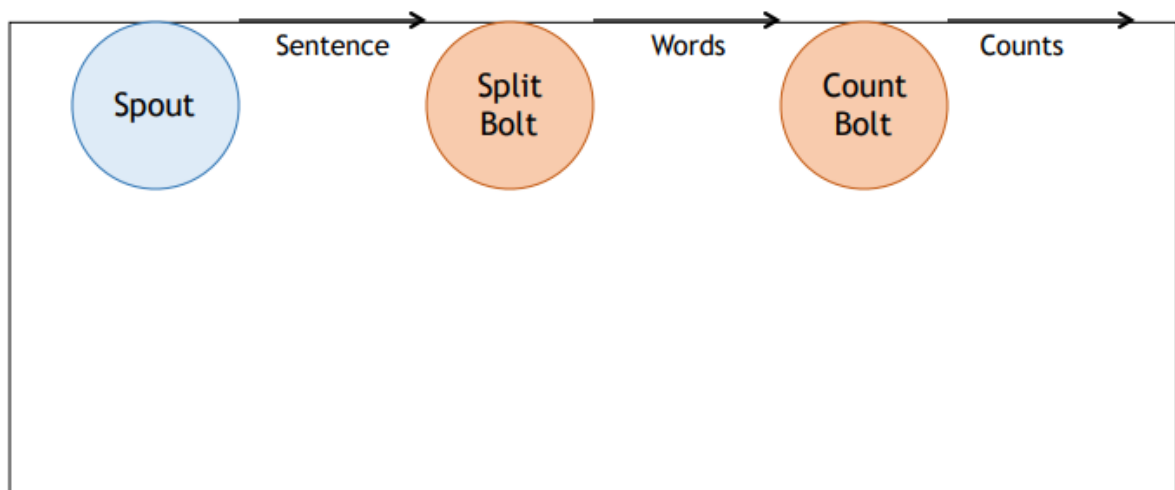**Step 3:** Install storm on the virtual machine and change the PATH system variable:

```
1  curl -s "http://mirrors.advancedhosters.com/apache/storm/apache-storm-1.0.6
     /apache-storm-1.0.6.tar.gz" | tar -xz -C /usr/local/
2  ln -s /usr/local/apache-storm-1.0.6 /usr/local/storm
3  PATH="/usr/local/storm/bin:${PATH}"
```

# 4. Hello World (Word Count)

The **Word Count** application is the canonical example in Strom. In this tutorial, we are going to build all the components needed step-by-step and then build the whole application. In order to have a Word Count application, you need to create a Storm Topology as below:

This topology has three main components:

- **Spout**: In this topology, the Spout is emitting sentences as an output. In this tutorial, we are going to build a Spout that randomly emits a sentence from a predefined set in the Spout.

- **Split Bolt**: This bolt inputs a tuple in the form of a sentence and then emits a tuple for each word in the sentence.

- **Count Bolt**: This bolt keeps track of the count of each word. It does so by inputting tuples in form of words and incrementing the count for each word accordingly. Additionally, it emits the new count for each word. For example, if the tuple being read is the word "apple" and the current count for "apple" is 5 in the bolt, it will increment the internal count to 6 and emit a tuple in form of ("apple", 6).



In the remaining of the tutorial, we are going to explain how to build each of the components mentioned above and how to wire them up to create the whole topology.

**step 1:** create a temporary folder for this tutorial and enter

```
1  mkdir storm-tutorial
2  cd storm-tutorial
```

**step 2:** Create a "pom.xml" file by using the following command, which will be used by maven in order to build and package the project

(option 1)

```
1  vim pom.xml
2  #then type or copy/paste the XML code at the end of this tutorial in the editor;
     save and quit
```

(option 2) Download the pom.xml from the github:

```
1  wget -c https://github.com/UIUC-public/storm-tutorial/raw/master/pom.xml
```

**step 3:** Create a "src" directory and go into that directory by using the command below. This directory is where we will place all the java classes. In the following steps, we will create each component and then create the whole topology

```
1   mkdir src
2   cd src
3
```

**step 4: [Spout]**

(option 1) Create a new Java source code file "RandomSentenceSpout.java" using the following command:

```
1   vim RandomSentenceSpout.java
2   #then type or copy/paste the RandomSentenceSpout java code at the end of this
      tutorial in the editor; save and quit
```

(option 2) Download the RandomSentenceSpout.java from the github:

```
1   wget -c https://github.com/UIUC-public/storm-tutorial/raw/master/src
      /RandomSentenceSpout.java
2
```

As a result, this spout randomly chooses one of the sentences in "String[] sentences", and emits it as a tuple.

**step 5: [Split Bolt]**

(option 1) Create a new Java source code file "SplitSentenceBolt.java" using the following command:

```
1   vim SplitSentenceBolt.java
2   #then type or copy/paste the SplitSentenceBolt java code at the end of this
      tutorial in the editor; save and quit
```

(option 2) Download the SplitSentenceBolt.java from the github:

```
1   wget -c https://github.com/UIUC-public/storm-tutorial/raw/master/src
      /SplitSentenceBolt.java
2
```

As a result, this Bolt reads a sentence and splits it into words, then it emits a tuple for each word.

**step 6: [Count Bolt]**

(option 1) Create a new Java source code file "WordCountBolt.java" using the following command:

```
1   vim WordCountBolt.java
2   #then type or copy/paste the WordCountBolt java code at the end of this tutorial
      in the editor; save and quit
```

(option 2) Download the WordCountBolt.java from the github:

```
1   wget -c https://github.com/UIUC-public/storm-tutorial/raw/master/src
      /WordCountBolt.java
2
```

This bolt keeps an in-memory hash map of words to their count value. Upon receipt of a tuple it increments the count for that word and emits the new count as a tuple in the format of ("word", count).

**step 7:** Word Count Topology

Now that we have all the components, we have to build the Storm Topology.

(option 1) Create a new Java source code file "WordCountTopology.java" using the following command:

```
1    vim WordCountTopology.java
2    #then type or copy/paste the WordCountTopology java code at the end of this
       tutorial in the editor; save and quit
```

(option 2) Download the WordCountTopology.java from the github:

```
1    wget -c https://github.com/UIUC-public/storm-tutorial/raw/master/src
       /WordCountTopology.java
```

This class creates a "ToplogyBuilder" instance which wires up all the components and submits the topology. It first adds the 5 instance of the spout. Then it connects 8 instances of the SplitBolt to the spout, and 12 instances of the CountBolt to the SplitBolt.

After building the topology, it submits the topology, runs the topology for 10 seconds, and then stops the application.

**step 8:** Build and Submitting the Job

8.1: Go to the root directory of the project using the following command:

```
1    cd ..
```

8.2: Compile the source code using the following command:

```
1    mvn clean package
2    # this will create a folder called "target" and build a fat jar under that
       folder
```

8.3: Go to the target folder using the following command:

```
1    cd target
```

8.4: Run the application and save its output using the following command:

```
1    storm jar storm-example-0.0.1-SNAPSHOT.jar WordCountTopology > wordcount.txt
```

You should see an output similar to the one below in wordcount.txt:

```
7506 [Thread-28-split-executor[6 6]] INFO  o.a.s.d.executor - TRANSFERING tuple [dest: 2 tuple: source: split:6, stream: default, id: {}, [keeps]]
7506 [Thread-28-split-executor[6 6]] INFO  o.a.s.d.task - Emitting: split default [the]
7506 [Thread-28-split-executor[6 6]] INFO  o.a.s.d.executor - TRANSFERING tuple [dest: 3 tuple: source: split:6, stream: default, id: {}, [the]]
7506 [Thread-28-split-executor[6 6]] INFO  o.a.s.d.task - Emitting: split default [doctor]
7507 [Thread-28-split-executor[6 6]] INFO  o.a.s.d.executor - TRANSFERING tuple [dest: 2 tuple: source: split:6, stream: default, id: {}, [doctor]]
7507 [Thread-28-split-executor[6 6]] INFO  o.a.s.d.task - Emitting: split default [away]
7507 [Thread-28-split-executor[6 6]] INFO  o.a.s.d.executor - TRANSFERING tuple [dest: 2 tuple: source: split:6, stream: default, id: {}, [away]]
7507 [Thread-28-split-executor[6 6]] INFO  o.a.s.d.executor - BOLT ack TASK: 6 TIME: -1 TUPLE: source: spout:9, stream: default, id: {}, [an apple a day keeps the doctor away]
7507 [Thread-28-split-executor[6 6]] INFO  o.a.s.d.executor - Execute done TUPLE source: spout:9, stream: default, id: {}, [an apple a day keeps the doctor away] TASK: 6 DELTA: -1
7510 [Thread-20-count-executor[2 2]] INFO  o.a.s.d.executor - Processing received message FOR 2 TUPLE: source: split:6, stream: default, id: {}, [apple]
7510 [Thread-20-count-executor[2 2]] INFO  o.a.s.d.task - Emitting: count default [apple, 1]
7510 [Thread-20-count-executor[2 2]] INFO  o.a.s.d.executor - BOLT ack TASK: 2 TIME: -1 TUPLE: source: split:6, stream: default, id: {}, [apple]
7510 [Thread-20-count-executor[2 2]] INFO  o.a.s.d.executor - Execute done TUPLE source: split:6, stream: default, id: {}, [apple] TASK: 2 DELTA: -1
7510 [Thread-20-count-executor[2 2]] INFO  o.a.s.d.executor - Processing received message FOR 2 TUPLE: source: split:6, stream: default, id: {}, [keeps]
7510 [Thread-20-count-executor[2 2]] INFO  o.a.s.d.task - Emitting: count default [keeps, 1]
7511 [Thread-20-count-executor[2 2]] INFO  o.a.s.d.executor - BOLT ack TASK: 2 TIME: -1 TUPLE: source: split:6, stream: default, id: {}, [keeps]
7511 [Thread-20-count-executor[2 2]] INFO  o.a.s.d.executor - Execute done TUPLE source: split:6, stream: default, id: {}, [keeps] TASK: 2 DELTA: -1
```

**To learn more about developing Storm applications in Java, visit:**

**https://github.com/apache/storm/tree/master/examples/storm-starter**

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

**pom.xml**

```xml
1   <project xmlns="http://maven.apache.org/POM/4.0.0" xmlns:xsi="http://www.w3.org
    /2001/XMLSchema-instance"
2   xsi:schemaLocation="http://maven.apache.org/POM/4.0.0 http://maven.apache.org
    /xsd/maven-4.0.0.xsd">
3   <modelVersion>4.0.0</modelVersion>
4
5   <groupId>storm.example</groupId>
6   <artifactId>storm-example</artifactId>
7   <version>0.0.1-SNAPSHOT</version>
8   <packaging>jar</packaging>
9
10  <properties>
11    <project.build.sourceEncoding>UTF-8</project.build.sourceEncoding>
12  </properties>
13
14  <repositories>
15    <repository>
16      <id>clojars.org</id>
17      <url>http://clojars.org/repo</url>
18    </repository>
19  </repositories>
20
21  <dependencies>
22    <!-- https://mvnrepository.com/artifact/org.apache.storm/storm -->
23    <dependency>
24      <groupId>org.apache.storm</groupId>
25      <artifactId>storm-core</artifactId>
26      <version>1.0.5</version>
27      <!-- keep storm out of the jar-with-dependencies -->
28      <scope>provided</scope>
29    </dependency>
30  </dependencies>
31
32  <build>
33    <sourceDirectory>src</sourceDirectory>
34    <plugins>
35      <plugin>
36        <groupId>org.apache.maven.plugins</groupId>
37        <artifactId>maven-compiler-plugin</artifactId>
38        <configuration>
39          <source>1.6</source>
40          <target>1.6</target>
41        </configuration>
42      </plugin>
43
44    </plugins>
45  </build>
46 </project>
```

**RandomSentenceSpout.java**

```
 1   import org.apache.storm.spout.SpoutOutputCollector;
 2   import org.apache.storm.task.TopologyContext;
 3   import org.apache.storm.topology.OutputFieldsDeclarer;
 4   import org.apache.storm.topology.base.BaseRichSpout;
 5   import org.apache.storm.tuple.Fields;
 6   import org.apache.storm.tuple.Values;
 7   import org.apache.storm.utils.Utils;
 8
 9   import java.util.Map;
10   import java.util.Random;
11
12   public class RandomSentenceSpout extends BaseRichSpout {
13     SpoutOutputCollector _collector;
14     Random _rand;
15
16
17     @Override
18     public void open(Map conf, TopologyContext context, SpoutOutputCollector
         collector) {
19       _collector = collector;
20       _rand = new Random();
21     }
22
23     @Override
24     public void nextTuple() {
25       Utils.sleep(100);
26       String[] sentences = new String[]{ "the cow jumped over the moon", "an apple
           a day keeps the doctor away",
27           "four score and seven years ago", "snow white and the seven dwarfs", "i
               am at two with nature" };
28       String sentence = sentences[_rand.nextInt(sentences.length)];
29       _collector.emit(new Values(sentence));
30     }
31
32     @Override
33     public void ack(Object id) {
34     }
35
36     @Override
37     public void fail(Object id) {
38     }
39
40     @Override
41     public void declareOutputFields(OutputFieldsDeclarer declarer) {
42       declarer.declare(new Fields("word"));
43     }
44
45   }
```

**SplitSentenceBolt.java**

```
 1   import org.apache.storm.topology.BasicOutputCollector;
 2   import org.apache.storm.topology.OutputFieldsDeclarer;
 3   import org.apache.storm.topology.base.BaseBasicBolt;
 4   import org.apache.storm.tuple.Fields;
 5   import org.apache.storm.tuple.Tuple;
 6   import org.apache.storm.tuple.Values;
 7
 8   public class SplitSentenceBolt extends BaseBasicBolt {
 9
10   @Override
11     public void execute(Tuple tuple, BasicOutputCollector collector) {
12         String sentence = tuple.getString(0);
13         String[]words=sentence.split("[\\s~`!@#$%^&*(-)+=_:;'\",.<>?/\\\\0-9"
             +"\\]\\[\\}\\{]+");
14
15         for(String word:words){
16             collector.emit(new Values(word));
17         }
18     }
19     @Override
20     public void declareOutputFields(OutputFieldsDeclarer declarer) {
21       declarer.declare(new Fields("word"));
22     }
23   }
```

**WordCountBolt.java**

```
1   import org.apache.storm.topology.BasicOutputCollector;
2   import org.apache.storm.topology.OutputFieldsDeclarer;
3   import org.apache.storm.topology.base.BaseBasicBolt;
4   import org.apache.storm.tuple.Fields;
5   import org.apache.storm.tuple.Tuple;
6   import org.apache.storm.tuple.Values;
7
8   import java.util.HashMap;
9   import java.util.Map;
10
11  public class WordCountBolt extends BaseBasicBolt {
12      Map<String, Integer> counts = new HashMap<String, Integer>();
13
14      @Override
15      public void execute(Tuple tuple, BasicOutputCollector collector) {
16        String word = tuple.getString(0);
17        Integer count = counts.get(word);
18        if (count == null)
19          count = 0;
20        count++;
21        counts.put(word, count);
22        collector.emit(new Values(word, count));
23      }
24
25      @Override
26      public void declareOutputFields(OutputFieldsDeclarer declarer) {
27        declarer.declare(new Fields("word", "count"));
28      }
29  }
```

**WordCountTopology.java**

```
1   import org.apache.storm.Config;
2   import org.apache.storm.LocalCluster;
3   import org.apache.storm.StormSubmitter;
4   import org.apache.storm.topology.BasicOutputCollector;
5   import org.apache.storm.topology.OutputFieldsDeclarer;
6   import org.apache.storm.topology.TopologyBuilder;
7   import org.apache.storm.topology.base.BaseBasicBolt;
8   import org.apache.storm.tuple.Fields;
9   import org.apache.storm.tuple.Tuple;
10  import org.apache.storm.tuple.Values;
11
12  /**
13   * This topology demonstrates Storm's stream groupings and multilang
       capabilities.
14   */
15  public class WordCountTopology {
16
17
18    public static void main(String[] args) throws Exception {
19
20      TopologyBuilder builder = new TopologyBuilder();
21
22      builder.setSpout("spout", new RandomSentenceSpout(), 5);
23
24      builder.setBolt("split", new SplitSentenceBolt(), 8).shuffleGrouping("spout"
          );
25      builder.setBolt("count", new WordCountBolt(), 12).fieldsGrouping("split",
          new Fields("word"));
26
27      Config conf = new Config();
28      conf.setDebug(true);
29
30
31      if (args != null && args.length > 0) {
32        conf.setNumWorkers(3);
33
34        StormSubmitter.submitTopology(args[0], conf, builder.createTopology());
35      }
36      else {
37        conf.setMaxTaskParallelism(3);
38
39        LocalCluster cluster = new LocalCluster();
40        cluster.submitTopology("word-count", conf, builder.createTopology());
41
42        Thread.sleep(10000);
43
44        cluster.shutdown();
45      }
46    }
47  }
```