# RT4 · Big Data Class Labs

Lab 1 · Batch processing with Hadoop HDFS and Map Reduce

Adam Lahbib · M. Sofiene Barka · Mohamed Rafraf

Mar. 1, 2023 @ INSAT

# Contents

# 1 Introduction

In this lab, we will learn how to use Hadoop to process large amounts of data. We will use the Hadoop Distributed File System (HDFS) to store the data and the MapReduce framework to process it. We will use the Java MapReduce API to write our MapReduce programs.

# 2 Definitions

## 2.1 Hadoop

Hadoop is an open-source software framework for distributed storage and distributed processing of very large data sets on computer clusters built from commodity hardware. All the modules in Hadoop are designed with a fundamental assumption that hardware failures are common and should be automatically handled in software by the framework.

## 2.2 HDFS

HDFS is a distributed file system designed to run on commodity hardware. It has many similarities with existing distributed file systems. However, the differences from other distributed file systems are significant. HDFS is highly fault-tolerant and is designed to be deployed on low-cost hardware. HDFS provides high throughput access to application data and is suitable for applications that have large data sets.

## 2.3 MapReduce

MapReduce is a programming model for processing and generating large data sets with a parallel, distributed algorithm on a cluster. The MapReduce programming model is comprised of the map and reduce operations. The map operation processes a set of input key/value pairs to generate a set of intermediate key/value pairs. The reduce operation merges all intermediate values associated with the same intermediate key to produce a set of output key/value pairs.

## 2.4 Java MapReduce API

The Java MapReduce API is a Java programming interface for writing applications which process vast amounts of data (multi-terabyte data-sets) in-parallel on large clusters (thousands of nodes) of commodity hardware in a reliable, fault-tolerant manner.
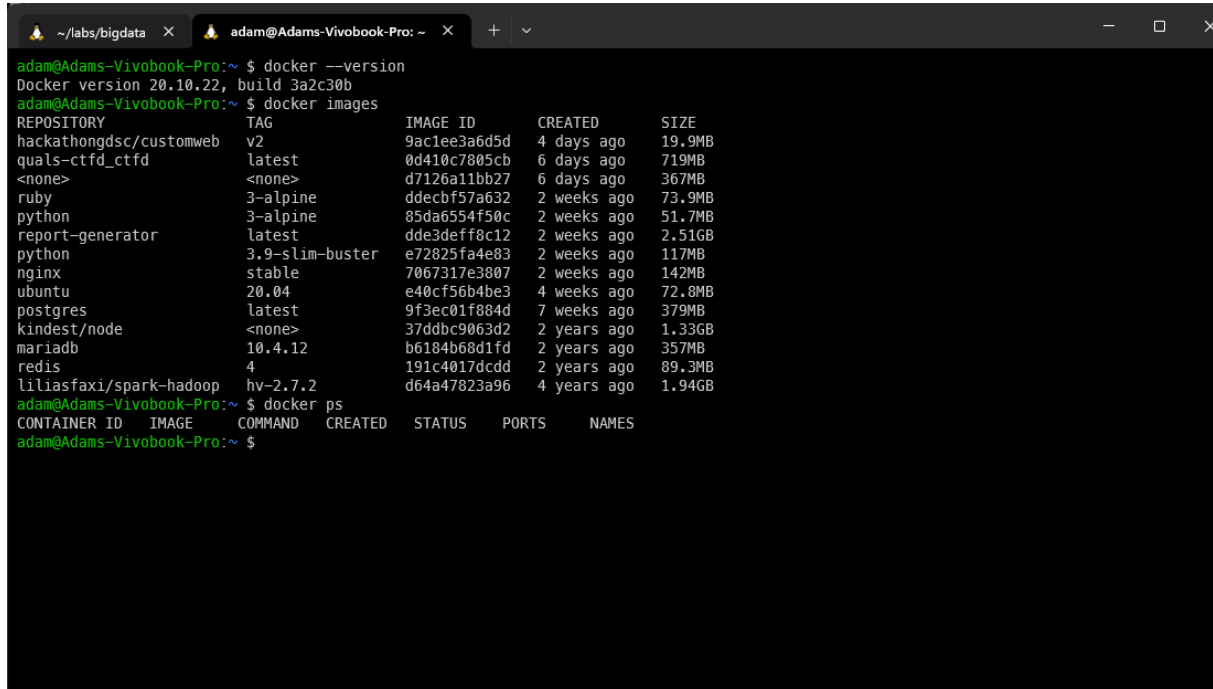
## 2.5  Docker

Docker is a set of platform as a service (PaaS) products that use OS-level virtualization to deliver software in packages called containers. Containers are isolated from one another and bundle their own software, libraries and configuration files; they can communicate with each other through well-defined channels. All containers are run by a single operating-system kernel and are thus more lightweight than virtual machines. Containers are created from images that specify their precise contents. Images are often created by combining and modifying standard images downloaded from public repositories.

# 3  Environment

Requirements:

- Docker (latest)
- JDK 8
- Maven 4.0.0
- WSL2 (Ubuntu/Debian preferrably) or Linux Distro
- Hadoop image (latest)

## 3.1  WSL2 (Ubuntu) + Docker CLI + Hadoop image (Liliasfaxi's Image in DockerHub)



WSL 2 is a new version of the Windows Subsystem for Linux (WSL) that is based on a real Linux kernel, using a lightweight utility virtual machine (VM) manager called Hyper-V. WSL 2 is a significant improve-

ment over WSL 1, which was based on a Linux compatibility layer called the Windows Subsystem for Linux (WSL) that was built into the Windows 10 kernel. WSL 2 is a significant improvement over WSL 1, which was based on a Linux compatibility layer called the Windows Subsystem for Linux (WSL) that was built into the Windows 10 kernel.

## 3.2  JDK 8 + JDK 11

JDK or Java Development Kit is a software development environment used for developing Java applications and applets. It contains the Java Runtime Environment, an interpreter/loader (java), a compiler (javac), an archiver (jar), a documentation generator (javadoc) and other tools needed in Java development. It is a development environment for building applications, applets, and components using the Java programming language. It includes an interpreter, a compiler, debuggers, and other tools needed to develop applications and applets that run in the Java runtime environment.

I installed SDKMan to manage my Java versions, it's a tool for managing parallel versions of multiple Software Development Kits on most Unix based systems. It provides a convenient Command Line Interface (CLI) and API for installing, switching, removing and listing Candidates. It is written in Groovy and uses Gradle for building. It is available under the Apache 2.0 license.



SDKMan `list java`, can show all the available JDKs

I installed Java JDK 8.0.302-open for this lab and Java JDK 11.0.2-open for myself, I'll keep JDK 11.0.2-open as my default Java version, and switch to JDK 8.0.302-open when I need to run Hadoop using `sdk use java xxx`







## 3.3  Installing Maven

Maven is a build automation tool used primarily for Java projects. Maven addresses two aspects of building software: first, it describes how software is built, and second, it describes its dependencies. Maven is a build automation tool used primarily for Java projects. Maven addresses two aspects of building software: first, it describes how software is built, and second, it describes its dependencies.
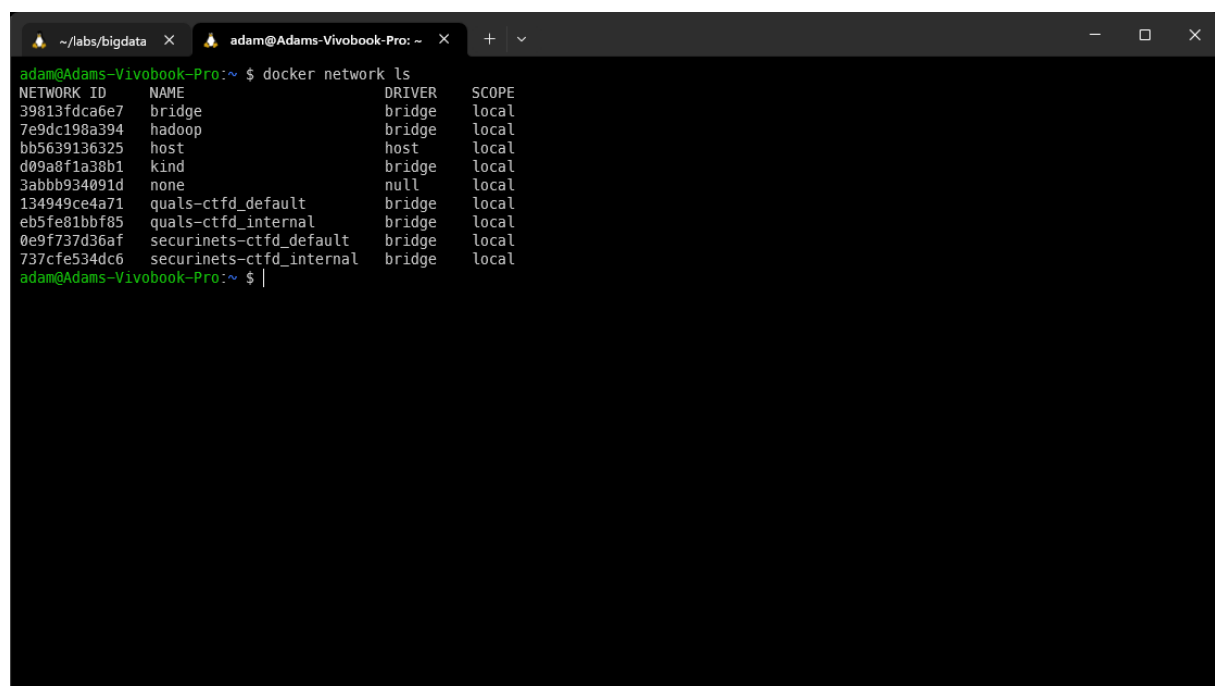
## 3.4  Docker Network

A Docker network is a virtual network that connects containers. It is a layer 2 bridge that connects containers running on the same Docker host. It is also a layer 3 network that connects containers running on different Docker hosts. Docker networks are isolated from each other and from the host by default. Docker networks are defined by a driver. The default driver is bridge. The bridge driver creates a virtual network that connects containers running on the same Docker host. The virtual network is isolated from the host network and other Docker networks by default. The bridge driver is the only driver that supports the default network. The bridge driver supports the creation of user-defined networks.

## 3.5  Master and 2 slaves



## 3.6  Starting Hadoop

## 3.7  First Steps with Hadoop



## 3.8  Hadoop Web UI

The Hadoop Web UI is a web interface that allows you to monitor the status of your Hadoop cluster. It is available at `http://localhost:50070/`. To view job results and progress, go to `http://localhost:8088/`.

# 4  Application

## 4.1  Idea

Total sales per shop is a simple MapReduce application that counts the total sales per shop. The input is a CSV file with the architecure: date, time, shop, product, price, payment method. The output is a CSV file with the architecture: shop, total sales.

## 4.2  Code Modifications

### 4.2.1  `WordCount.java`

```java
package tn.insat.tp1;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.FloatWritable;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {
    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "word count");
        job.setJarByClass(WordCount.class);
        job.setMapperClass(TokenizerMapper.class);
        job.setCombinerClass(IntSumReducer.class);
        job.setReducerClass(IntSumReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(FloatWritable.class);
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}
```

### 4.2.2 `TokenizerMapper.java`

The TokenizerMapper class extends the Mapper class and overrides the map method. The map method takes three parameters: the key, the value, and the Context. The key is the offset of the line in the input file, the value is the line itself, and the Context is the object that allows you to write the output of the map method. The map method is called once for each line in the input file. The map method tokenizes the line into words and writes each word as a key with a value of 1. The output of the map method is written to the Context object.

```java
package tn.insat.tp1;

import org.apache.hadoop.io.FloatWritable;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import java.io.IOException;
import java.util.StringTokenizer;

public class TokenizerMapper
        extends Mapper<Object, Text, Text, FloatWritable>{

    private Text magasin = new Text();
    private FloatWritable cout = new FloatWritable();

    public void map(Object key, Text value, Mapper.Context context) throws IOException,
    ↪    InterruptedException {

        String[] columns = value.toString().split(",");

        if (columns.length == 7) {
            // Extract the 'magasin' and 'cout' values
            magasin.set(columns[3]);
            cout.set(Float.parseFloat(columns[5]));
            System.out.println(magasin);
            System.out.println(cout);
            context.write(magasin, cout);
        }
    }
}
```

### 4.2.3 `IntSumReducer.java`

The IntSumReducer class extends the Reducer class and overrides the reduce method. The reduce method takes three parameters: the key, the values, and the Context. The key is the word, the values are the counts, and the Context is the object that allows you to write the output of the reduce method.

The reduce method is called once for each key. The reduce method sums the values for each key and writes the key and the sum as the output. The output of the reduce method is written to the Context object.

```java
package tn.insat.tp1;


import org.apache.hadoop.io.FloatWritable;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;

import java.io.IOException;

public class IntSumReducer
        extends Reducer<Text,FloatWritable,Text,FloatWritable> {

    private FloatWritable result = new FloatWritable();

    public void reduce(Text key, Iterable<FloatWritable> values,
                       Context context
    ) throws IOException, InterruptedException {
        int sum = 0;
        for (FloatWritable val : values) {
            System.out.println("value: "+val.get());
            sum += val.get();
        }
        System.out.println("--> Sum = "+sum);
        result.set(sum);
        context.write(key, result);
    }
}
```

Since prices are rather floats, I used FloatWritable instead of IntWritable.

### 4.2.4 Github repo for source code:
### https://github.com/adamlahbib/BigData_Labs/tree/main/lab1/wordcount

### 4.2.5 `pom.xml` File

This is the pom.xml file, it contains all the dependencies needed for the project, I used the following dependencies:

## 4.3  Local test using VSCode

I will use the following extension to run/debug Java code in VSCode:



I will use this part of the file as my input:

```
428400,2012-02-07,16:17,Memphis,Crafts,409.92,Discover
428401,2012-02-07,16:17,Richmond,Children's Clothing,276.46,MasterCard
428402,2012-02-07,16:17,Honolulu,Crafts,464.81,MasterCard
```

Then thanks to the VSCode extension when I press RUN and add the arguments to the command line, I get the following output:

### 4.3.1 Packages Installation: `mvn package install`



## 4.4 Hadoop test (Dockerized)

For hadoop test, I will use the CSV file as input, instead of the tab delimited file, to make things easier.

```
root@hadoop-master:~# ls
hdfs   purchases.txt  purchases2.txt  run-wordcount.sh  start-hadoop.sh  start-kafka-zookeeper.sh  wordcount-1.jar
root@hadoop-master:~# 
```

Let's remove the old input content as well as the output folder

```
root@hadoop-master:~# hadoop fs -rm -r /user/root/output
23/03/01 18:59:35 INFO fs.TrashPolicyDefault: Namenode trash configuration: Deletion interval = 0 minutes, Emptier interval = 0 minutes.
Deleted /user/root/output
root@hadoop-master:~# hadoop fs -rm -r /user/root/input
23/03/01 19:00:00 INFO fs.TrashPolicyDefault: Namenode trash configuration: Deletion interval = 0 minutes, Emptier interval = 0 minutes.
Deleted /user/root/input
root@hadoop-master:~#
```

Let's now use purchases2.txt

```
root@hadoop-master:~# hadoop fs -mkdir -p input
root@hadoop-master:~# hadoop fs -put purchases2.txt input
root@hadoop-master:~#
```

Everything is ready, let's run the job

```
        File System Counters
                FILE: Number of bytes read=3224
                FILE: Number of bytes written=358056
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=243313966
                HDFS: Number of bytes written=2159
                HDFS: Number of read operations=9
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=2
        Job Counters
                Launched map tasks=2
                Launched reduce tasks=1
                Data-local map tasks=2
                Total time spent by all maps in occupied slots (ms)=63086
                Total time spent by all reduces in occupied slots (ms)=16137
                Total time spent by all map tasks (ms)=63086
                Total time spent by all reduce tasks (ms)=16137
                Total vcore-milliseconds taken by all map tasks=63086
                Total vcore-milliseconds taken by all reduce tasks=16137
                Total megabyte-milliseconds taken by all map tasks=64600064
                Total megabyte-milliseconds taken by all reduce tasks=16524288
        Map-Reduce Framework
                Map input records=4138476
                Map output records=4138476
                Map output bytes=56372650
                Map output materialized bytes=3230
                Input split bytes=242
                Combine input records=4138476
                Combine output records=206
                Reduce input groups=103
                Reduce shuffle bytes=3230
                Reduce input records=206
                Reduce output records=103
                Spilled Records=412
                Shuffled Maps =2
                Failed Shuffles=0
                Merged Map outputs=2
                GC time elapsed (ms)=922
                CPU time spent (ms)=53320
                Physical memory (bytes) snapshot=693579776
                Virtual memory (bytes) snapshot=5885042688
                Total committed heap usage (bytes)=531628032
        Shuffle Errors
                BAD_ID=0
                CONNECTION=0
                IO_ERROR=0
                WRONG_LENGTH=0
                WRONG_MAP=0
                WRONG_REDUCE=0
        File Input Format Counters
                Bytes Read=243313724
        File Output Format Counters
                Bytes Written=2159
root@hadoop-master:~#
```

## 4.5  Monitoring the job

Monitoring the job is done using the Hadoop Web UI, we can see the job running on the Master node:

and Slave 1 for example:

# 5 Homework

## 5.1 Idea

Processing Wireshark PCAP files using Hadoop MapReduce job, to extract the GEO data from PCAP packets using the MaxMind GeoLite2 database.

## 5.2 Architecture

## 5.3 Code

We will, as usual, modify the WordCount (template) example to process the PCAP files.

### 5.3.1 WordCount.java

```java
package tn.insat.tp1;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.conf.Configured;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.input.SequenceFileAsBinaryInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
import org.apache.hadoop.util.Tool;
import org.apache.hadoop.util.ToolRunner;

public class Driver extends Configured implements Tool {

    public int run(String[] args) throws Exception {

        Job job = Job.getInstance(getConf(), "PcapProcessing");
```

```java
        job.setJarByClass(getClass());

        Path in = new Path(args[0]);
        Path out = new Path(args[1]);
        FileInputFormat.setInputPaths(job, in);
        FileOutputFormat.setOutputPath(job, out);

        job.setMapperClass(PcapMapper.class);
        job.setReducerClass(PcapReducer.class);

        job.setInputFormatClass(SequenceFileAsBinaryInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);

        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(IntWritable.class);

        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);

        job.addCacheFile(new Path("GeoLite2-City.mmdb").toUri());

        return job.waitForCompletion(true)?0:1;
    }

    public static void main(String[] args) throws Exception {
        int result = ToolRunner.run(new Configuration(),new Driver(), args);
            System.exit(result);

    }

}
```

### 5.3.2 `TokenizerMapper.java`

```java
package tn.insat.tp1;

import java.io.IOException;
import java.io.InputStream;
import java.net.InetAddress;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.FileSystem;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.BytesWritable;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

import com.maxmind.geoip2.DatabaseReader;
```

```java
import com.maxmind.geoip2.exception.GeoIp2Exception;
import com.maxmind.geoip2.model.CityResponse;

public class PcapMapper extends Mapper<BytesWritable, BytesWritable, Text, IntWritable> {


 private static final String DBNAME="GeoLite2-City.mmdb";

 private Text word = new Text();
 private final static IntWritable one = new IntWritable(1);
 private byte[] binaryValue;
 private byte[] binaryIpAddress = new byte[4];
 private InputStream is;
 private DatabaseReader dbReader;
 private CityResponse dbResponse;

 @Override
protected void setup(Context context) throws IOException, InterruptedException {
    Configuration conf = context.getConfiguration();

    FileSystem fs = FileSystem.getLocal(conf);
    is = fs.open(new Path(DBNAME));
    dbReader = new DatabaseReader.Builder(is).build();
}

@Override
 public void map(BytesWritable key, BytesWritable value, Context contex) throws IOException,
 ↪    InterruptedException {

    // Get Bytes
    binaryValue = value.getBytes();

    // Extract IP address
    binaryIpAddress[0] = binaryValue[27];
    binaryIpAddress[1] = binaryValue[28];
    binaryIpAddress[2] = binaryValue[29];
    binaryIpAddress[3] = binaryValue[30];

    // Find GEO data of the IP address
    try {
        dbResponse = dbReader.city(InetAddress.getByAddress(binaryIpAddress));
        word.set(dbResponse.getCountry().getName()+"-"+dbResponse.getCity().getName());
        contex.write(word, one);
    } catch (GeoIp2Exception e) {
        e.printStackTrace();
    }
 }

@Override
protected void cleanup(Context context) throws IOException, InterruptedException {
    is.close();
}
```

```
}
```

We know that in a pcap file, the Source IP address is stored in the 27th to 30th bytes of the packet, so we extract it and use it to find the GEO data.

### 5.3.3 `IntSumReducer.java`

```java
package tn.insat.tp1;

import java.io.IOException;
import java.util.Iterator;


import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;



public class PcapReducer extends Reducer<Text, IntWritable, Text, IntWritable> {

 private IntWritable totalWordCount = new IntWritable();

 @Override
 public void reduce(Text key, Iterable<IntWritable> values, Context context) throws
 ↪   IOException, InterruptedException {

  int wordCount = 0;

  Iterator<IntWritable> it=values.iterator();
  while (it.hasNext()) {
   wordCount += it.next().get();
  }

  // Write keys with count > 1
  if(wordCount > 1){
      totalWordCount.set(wordCount);
      context.write(key, totalWordCount);
  }
 }
}
```

You can view the full code at: https://github.com/adamlahbib/BigData_Labs/tree/main/homework

I will use a dummy PCAP file from Wireshark samples as input.

Link: https://tranalyzer.com/download/data/faf-exercise.pcap *I will rename it to file.pcap simply.*

We will need to convert the pcap file to a sequence file using the following tool: https://github.com/marouni/pcap2seq

I will also need this DB file to get the GEO data of the IP addresses. https://github.com/P3TERX/GeoLite.mmdb/raw/download/GeoLite2-City.mmdb it weights ~ 60 MB and I wget it to my home directory.

## 5.4  Results

To conduct a test on the cluster, I had to put the DB file `GeoLite2-City.mmdb` in HDFS, the `pcap2seq.jar` file along the **fat jar** `wordocunt-1.jar` (meaning a jar packed with dependencies) in the hadoop root, and the `file.seq` file in the root of the HDFS (according to the tool!).

```
root@hadoop-master:~# ls
file.pcap   pcap2seq-1.2.jar  purchases2.txt   start-hadoop.sh            wordcount-1.jar
hdfs        purchases.txt     run-wordcount.sh start-kafka-zookeeper.sh
root@hadoop-master:~# hadoop jar pcap2seq-1.2.jar file.pcap file.seq org.apache.hadoop.io.compress.BZip2Codec
PCAP FILE FORMAT : SWAPPED
23/03/01 23:25:37 WARN bzip2.Bzip2Factory: Failed to load/initialize native-bzip2 library system-native, will use pure-Java version
23/03/01 23:25:37 INFO compress.CodecPool: Got brand-new compressor [.bz2]
Converting pcap file to Hadoop sequence file ...
Converted 5902 packets.
Read a total of 4993414 bytes.
root@hadoop-master:~#
```

```
root@hadoop-master:~# hadoop jar wordcount-1.jar tn.insat.tp1.WordCount file.seq output
23/03/02 00:15:29 INFO client.RMProxy: Connecting to ResourceManager at hadoop-master/172.24.0.2:8032
23/03/02 00:15:31 INFO input.FileInputFormat: Total input paths to process : 1
23/03/02 00:15:31 INFO mapreduce.JobSubmitter: number of splits:1
23/03/02 00:15:31 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1677686996211_0022
23/03/02 00:15:31 INFO impl.YarnClientImpl: Submitted application application_1677686996211_0022
23/03/02 00:15:31 INFO mapreduce.Job: The url to track the job: http://hadoop-master:8088/proxy/application_1677686996211_0022/
23/03/02 00:15:31 INFO mapreduce.Job: Running job: job_1677686996211_0022
23/03/02 00:15:38 INFO mapreduce.Job: Job job_1677686996211_0022 running in uber mode : false
23/03/02 00:15:38 INFO mapreduce.Job:  map 0% reduce 0%
23/03/02 00:15:59 INFO mapreduce.Job: Task Id : attempt_1677686996211_0022_m_000000_0, Status : FAILED
Error: Java heap space
23/03/02 00:16:08 INFO mapreduce.Job: Task Id : attempt_1677686996211_0022_m_000000_1, Status : FAILED
Error: Java heap space
23/03/02 00:16:19 INFO mapreduce.Job:  map 100% reduce 0%
23/03/02 00:16:21 INFO mapreduce.Job: Task Id : attempt_1677686996211_0022_m_000000_2, Status : FAILED
Error: Java heap space
23/03/02 00:16:22 INFO mapreduce.Job:  map 0% reduce 0%
23/03/02 00:16:27 INFO mapreduce.Job:  map 100% reduce 100%
23/03/02 00:16:27 INFO mapreduce.Job: Job job_1677686996211_0022 failed with state FAILED due to: Task failed task_1677686996211_0022_m_000000
Job failed as tasks failed. failedMaps:1 failedReduces:0

23/03/02 00:16:27 INFO mapreduce.Job: Counters: 16
        Job Counters
                Failed map tasks=4
                Killed reduce tasks=1
                Launched map tasks=4
                Other local map tasks=3
                Data-local map tasks=1
                Total time spent by all maps in occupied slots (ms)=35038
                Total time spent by all reduces in occupied slots (ms)=0
                Total time spent by all map tasks (ms)=35038
                Total time spent by all reduce tasks (ms)=0
                Total vcore-milliseconds taken by all map tasks=35038
                Total vcore-milliseconds taken by all reduce tasks=0
                Total megabyte-milliseconds taken by all map tasks=35878912
                Total megabyte-milliseconds taken by all reduce tasks=0
        Map-Reduce Framework
                CPU time spent (ms)=0
                Physical memory (bytes) snapshot=0
                Virtual memory (bytes) snapshot=0
root@hadoop-master:~#
```

Unfortunately, as you can see, I got the Java Heap Space error, I tried to increase the memory of the JVM but it didn't work.

Let's try a seq file locally.

```
~/l/b/labs/homework on main !9 ?2 ❯ /usr/bin/env /home/adam/.sdkman/candidates/java/8.0.282.j9-adpt/bin/java -cp /tmp/cp_
f2nu587fs1zcd0enk0zzoz1bn.jar tn.insat.tp1.WordCount src/main/resources/input/file.seq src/main/resources/output
log4j:WARN No appenders could be found for logger (org.apache.hadoop.util.Shell).
log4j:WARN Please initialize the log4j system properly.
log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more info.
```

```
 1    1 Canada-Thornhill
 2    2 United States-San Marcos
 3    3 United States-Wilmington
 4    4 United States-Palo Alto
 5    5 United States-Clifton
 6    6 United States-Ashburn
 7    7 United States-North Bergen
 8    8 United States-Hampton
 9    9 United States-Meredosia
10   10 United States-Seattle
11   11 United States-Hampton
12   12 United States-Mountain View  High anonymity  2886 kB/s
13   13 United States-Houston
14   14 United States-Ashburn
15   15 United States-Ashburn
```

# 6 Conclusion

In this lab, I learned how to install Hadoop on a cluster of 3 nodes, and how to run a simple MapReduce job using Hadoop. I also learned how to use the Hadoop Web UI to monitor the job.

---

Lab By: Mrs. Rabaa Youssef.