# LOG8415
# Advanced Concepts of Cloud Computing
# Lab2: MapReduce with Hadoop and Spark on Microsoft Azure

Marwane Adala
Vahid Majdinasab
Chun-An Bau
{marwane.adala, vahid.majdinasab, chun-an.bau}@polymtl.ca

We introduced how to set up the system, including Hadoop and Spark in section 1, and recorded what should be noticed while performing the WordCount experimentations in different comparison scenarios in sections 2 and 3. In section 4, we described how we managed to solve the social problem to implement a simple "People You Might Know" social network friendship recommendation algorithm. Finally, we showed the procedures of running all the scenarios and benchmarks in section 5. The setup procedures and corresponding scripts could be found in our GitHub repository [4].

## 1 Environment Setup

We set up our testing system environment(s) in the order of **Hadoop**, and **Spark**.

### 1.1 Hadoop Setup [5], [8], [6], [7] and [1]

We used a 20.04 LTS Ubuntu OS (local machine) with Openjdk version "1.8.0_292" and OpenSSH server and client installed on it. Our goal was to configure Hadoop on our system using standalone mode. The first steps were to create a Hadoop user, and add him to the sudo group. Then, we should configure a Password-less SSH using this user.
See the following commands concerning password-less SSH :

```
$ ssh−keygen −t rsa
$ sudo cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
$ sudo chmod 640 ~/.ssh/authorized_keys
$ ssh localhost #to verify if the password−less SSH is functional
```

After downloading, extracting and dealing with permissions of Hadoop version 3.3.1, we made the essential configurations concerning Hadoop and Java environment variables. We also edited core-site.xml, hdfs-site.xml, mapred-site.xml, and yarn-site.xml XML files to specify NameNode URL, define the storing location for node metadata and fsimage files, define MapReduce values and YARN-related settings. See Figures 1, 2, 3, and 4 for the validation of these steps. We refer the reader to Appendix A to see the other aspects of Hadoop configuration in standalone mode.

Figure 1: Hadoop Configuration Validation (1)



Figure 2: Hadoop Configuration Validation (2)



Figure 3: Hadoop Configuration Validation (3)



Figure 4: Hadoop Configuration Validation(4)

In Figures 5, 6 and 7, we can observe that we can access Hadoop Namenode (using port 9870), the Resource Manage (using port 8088) and browsing HDFS filesystem.
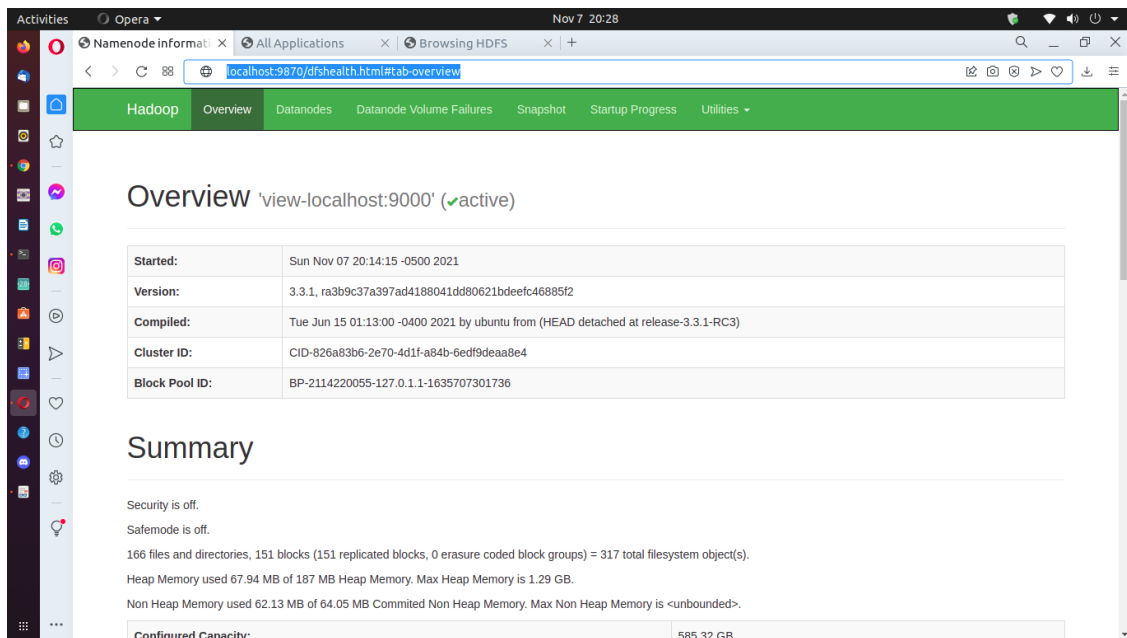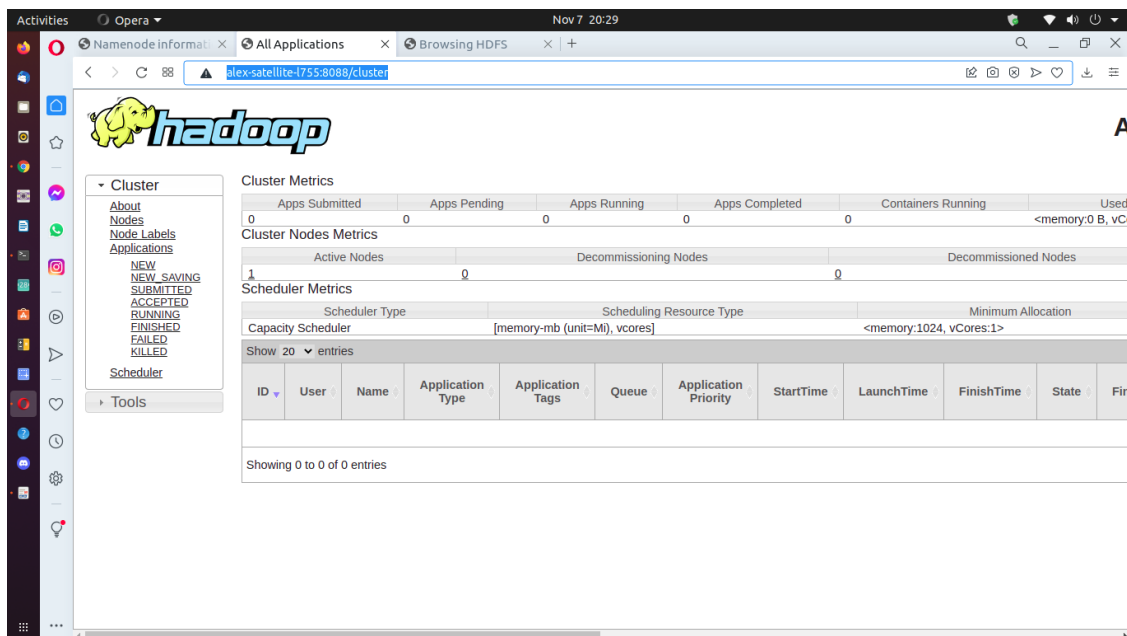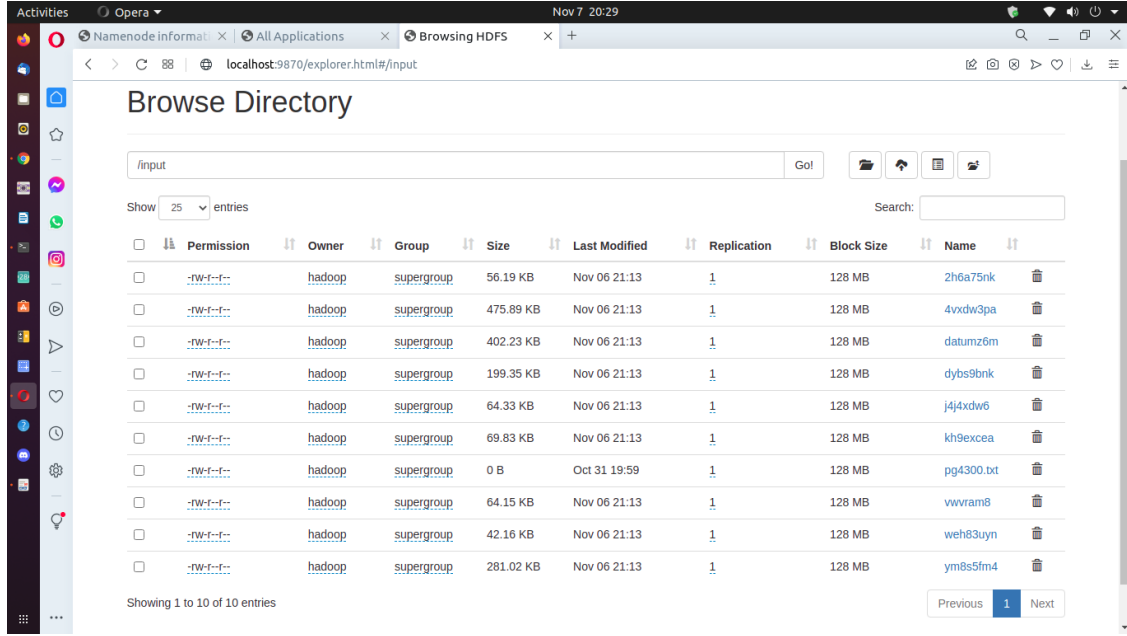
Figure 5: Hadoop Namenode



Figure 6: Resource Manager

Figure 7: Browsing HDFS Filesystem

## 1.2 Hadoop and Spark Setup on Azure

Since only one Hadoop node was required, we did not need additional setting related to network traffic. We could run Hadoop by the instructions in section 1.1. We configured Spark on Azure following the online tutorials [10] [9], and the most essential part is shown as below.

```
conf = SparkConf()
sc = SparkContext(conf=conf)
distFile = sc.textFile(input_file)
count = distFile.flatMap(lambda line: line.split(" "))   \
        .map(lambda word: (word, 1))      \
        .reduceByKey(lambda a, b: a + b)     \
        .sortByKey()
count.saveAsTextFile(output_file)
```

The *input_file* and *output_file* in the python script above could be either local storage or HDFS. In our experiments, we retrieved the data sets from HDFS and wrote the result back to HDFS.

## 2 Initial Experiments with WordCount program

In this section, we will go over the various results obtained by running some preliminary tests with Hadoop. After creating an input directory within HDFS, we begin by copying James Joyce's Ulysses book (named here pg4300.txt) to HDFS (see Figures 8, 9 and 10). Later on, we run the WordCount program on this book file in Figures 11, 12, 13 and 14. The Resource Manager interface in Figures 15 and 16 allows us to validate the success of these steps.

4

Figure 8: Making Input Directory in HDFS and Downloading James Joyce's Ulysses book



Figure 9: Exploring Hadoop Related Directories

Figure 10: Check out all Hadoop Daemons And Exploring HDFS files



Figure 11: Performing WordCount Java Program on James Joyce's Ulysses book

```
                    Job Counters
                            Launched map tasks=1
                            Launched reduce tasks=1
                            Other local map tasks=1
                            Total time spent by all maps in occupied slots (ms)=3418
                            Total time spent by all reduces in occupied slots (ms)=3380
                            Total time spent by all map tasks (ms)=3418
                            Total time spent by all reduce tasks (ms)=3380
                            Total vcore-milliseconds taken by all map tasks=3418
                            Total vcore-milliseconds taken by all reduce tasks=3380
                            Total megabyte-milliseconds taken by all map tasks=3500032
                            Total megabyte-milliseconds taken by all reduce tasks=3461120
                    Map-Reduce Framework
                            Map input records=0
                            Map output records=0
                            Map output bytes=0
                            Map output materialized bytes=6
                            Input split bytes=103
                            Combine input records=0
                            Combine output records=0
                            Reduce input groups=0
                            Reduce shuffle bytes=6
                            Reduce input records=0
                            Reduce output records=0
                            Spilled Records=0
                            Shuffled Maps =1
                            Failed Shuffles=0
                            Merged Map outputs=1
                            GC time elapsed (ms)=144
                            CPU time spent (ms)=1390
                            Physical memory (bytes) snapshot=505724928
                            Virtual memory (bytes) snapshot=5144027136
                            Total committed heap usage (bytes)=407896064
                            Peak Map Physical memory (bytes)=297615360
                            Peak Map Virtual memory (bytes)=2567708672
                            Peak Reduce Physical memory (bytes)=208109568
```

Figure 12: Performing WordCount Java Program on James Joyce's Ulysses book (2)



```
                            Input split bytes=103
                            Combine input records=0
                            Combine output records=0
                            Reduce input groups=0
                            Reduce shuffle bytes=6
                            Reduce input records=0
                            Reduce output records=0
                            Spilled Records=0
                            Shuffled Maps =1
                            Failed Shuffles=0
                            Merged Map outputs=1
                            GC time elapsed (ms)=144
                            CPU time spent (ms)=1390
                            Physical memory (bytes) snapshot=505724928
                            Virtual memory (bytes) snapshot=5144027136
                            Total committed heap usage (bytes)=407896064
                            Peak Map Physical memory (bytes)=297615360
                            Peak Map Virtual memory (bytes)=2567708672
                            Peak Reduce Physical memory (bytes)=208109568
                            Peak Reduce Virtual memory (bytes)=2576318464
                    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=0
                    File Output Format Counters
                            Bytes Written=0
real    0m37.046s
user    0m7.990s
sys     0m0.440s
hadoop@alex-SATELLITE-L755:/usr/local/hadoop/share/hadoop/mapreduce$
```

Figure 13: Performing WordCount Java Program on James Joyce's Ulysses book (3)

Figure 14: Checking the Output Directory Generated after running WordCount Program



Figure 15: Checking Resource Manager Jobs (1)

Figure 16: Checking Resource Manager Jobs (2)

# 3 Further Experiments with WordCount program

## 3.1 Performance comparison of Hadoop vs. Linux

We ran WordCount Java program example to perform the required tests using these commands, which enabled us to perform 3 experiments on all datasets and measure the user time at each experiment. The Table 1, comparing Linux vs Hadoop, shows the average user time for the nine given datasets. Using the same table and refering to [2] and [3], it is possible to conclude that Linux vastly outperforms Hadoop. In fact, several linux command line tools, such as grep, awk, sort, and uniq, can perform stream processing, so no batching is required and the memory overhead is minimal. Thus, it can be a very simple and rapid approach to pre-process big volumes of data on a local system. Hadoop, on the other hand, is designed to run on lots of machines and process large datasets that would otherwise be impossible to process on a single system.

9

Table 1: Execution time: Hadoop vs. Linux (sec)

|          | Hadoop | Linux |
|----------|--------|-------|
| Dataset1 | 8.257  | 0.254 |
| Dataset2 | 8.011  | 0.041 |
| Dataset3 | 8.070  | 0.088 |
| Dataset4 | 8.234  | 0.213 |
| Dataset5 | 8.585  | 0.041 |
| Dataset6 | 8.161  | 0.15  |
| Dataset7 | 8.063  | 0.04  |
| Dataset8 | 8.033  | 0.033 |
| Dataset9 | 8.437  | 0.035 |

The reader is directed to Figures 17 and 18 for more detailed visualizations of the comparison results.The reader should also refer to appendix B for the commands and steps used in this scenario.



Figure 17: Hadoop vs Linux Comparison. User Time in seconds (1)

Figure 18: Hadoop vs Linux Comparison. User Time in seconds (2)

## 3.2 Performance comparison of Hadoop vs. Spark on Azure

As shown in Table 2, the execution time of Hadoop is significantly longer than Spark's, which matches our expectations. Figures 19 and 20 provide more detailed visualizations of the comparison results. Although both retrieved data from HDFS, Spark stored temporary results in RAM while Hadoop wrote back to HDFS during processing stages, which might be the primary reason for the performance difference.

Table 2: Execution time: Hadoop vs. Spark (sec)

|            | Hadoop | Spark |
|------------|--------|-------|
| Dataset 1  | 11.007 | 0.154 |
| Dataset 2  | 9.781  | 0.153 |
| Dataset 3  | 10.495 | 0.159 |
| Dataset 4  | 10.701 | 0.143 |
| Dataset 5  | 10.761 | 0.164 |
| Dataset 6  | 10.937 | 0.158 |
| Dataset 7  | 10.579 | 0.157 |
| Dataset 8  | 10.513 | 0.158 |
| Dataset 9  | 10.208 | 0.149 |

Figure 19: Hadoop vs Spark Comparison. User Time in seconds (1)



Figure 20: Hadoop vs Spark Comparison. User Time in seconds (2)

# 4 The social network problem

## 4.1 Describe how you have used MapReduce jobs to solve the social network problem.

- We have written two Python scripts: 'mapper.py' and 'reducer.py'.

- mapper.py: Reads the source text file line by line. For each line in the source file which is in the format (UserID)(tab)(FriendsID), mapper outputs: (UserID)(tab)(FriendID)('friend'). For each FriendID in UserID's friends list, it generates a combination in this format (FriendID)(tab)(OtherFriendID)('not friend').

- reducer.py: Reads the outputs of mapper.py line by line. Since we set UserID as key, all UserID's data goes to the same reducer. Each reducer, counts the and relations of the users that it reads.

## 4.2 Describe your algorithm to tackle the social network problem.

- Our idea is to use UserID's as keys and FriendID information as values. For example, if we have user 0 which is friends with users 2, 3, 4, and 5, the output of the mapper will be:

  – 0 2 friend
  – 0 3 friend
  – 0 4 friend
  – 0 5 friend
  – 2 3 not friend
  – 2 4 not friend
  – 2 5 not friend
  – 3 4 not friend
  – 3 5 not friend
  – 4 5 not friend

  Since friendships have been defined as both ways (if 0 is friends with 1, 1 is friends with 0), we sort the users by numerical order before generating outputs. This means that we will be able to have consistent keys for each user and send all that user's information to the same reducer. The reducers read the outputs of mappers. Since we are using the UserID as key and we are sorting by numerical value, our reducers have the entire information for that UserID.

- Each reducer contains a dictionary keyed by UserIDs. Each UserID's dictionary is a nested dictionary that contains the number of times each User has seen other users. If a relation of 'friend' has been seen between users the relation flag will be set to True and this user will not be counted for recommendation. If the 'not friend' is seen for two users, we will consider them as recommendation candidates. Finally, we filter the user's dictionary by 'not friend' flag to see which users should be recommended as friends.

- We tried multiple times with ideas inspired from the WordCount example. We first tried setting the key as UserID and FriendID. For the same example as above, the output would be:

  - (0 2) friend
  - (0 3) friend
  - (0 4) friend
  - (0 5) friend
  - (2 3) not friend
  - ....

  Like our final solution, these keys would be sorted numerically. However, even though this algorithm worked on a single machine, on a distributed system, it meant that not all users' information would go to the same reducer and thus reducers would have incomplete information.

- The other solution that we tried was the same as our current solution but without sorting the friend lists numerically when we were generating mapper's outputs. This solution counted duplicates for the same keys and was not able to find correct recomendations.

## 4.3  Presents your recommendations of connection for the users with following user IDs: 924, 8941, 8942, 9019, 9020, 9021, 9022, 9990, 9992, 9993

For each user, each recommendation contains two numbers. The first is the other UserID and the second is the number of mutual friends they have. Our results are presented below:

- User 924: (6995, 1) - (439, 1) - (2409, 1) - (11860, 1) - (15416, 1) - (43748, 1) - (45881, 1)

- User 8941: (8943, 2) - (8944, 2) - (8940, 1)

- User 8942: (8939, 3) - (8940, 1) - (8943, 1) - (8944, 1)

- User 9019: (9022, 2) - (317, 1) - (9023, 1)

- User 9020: (9021, 3) - (9016, 2) - (9017, 2) - (9022, 2) - (317, 1) - (9023, 1)

- User 9021: (9020, 3) - (9016, 2) - (9017, 2) - (9022, 2) - (317, 1) - (9023, 1)

- User 9022: (9019, 2) - (9020, 2) - (9021, 2) - (317, 1) - (9016, 1) - (9017, 1) - (9023, 1)

- User 9990: (13134, 1) - (13478, 1) - (13877, 1) - (34299, 1) - (34485, 1) - (34642, 1) - (37941, 1)

- User 9992: (9987, 4) - (9989, 4) - (35667, 3) - (9991, 2)

- User 9993: (9991, 5) - (13134, 1) - (13478, 1) - (13877, 1) - (34299, 1) - (34485, 1) - (34642, 1) - (37941, 1)

# 5 Summary of results and instructions to run your code

This section deals with steps to follow in order to reproduce our work (**we uploaded everything at our GitHub repository**[4]). First, we need to follow the instructions in the environment setup section 1 to install Hadoop and Spark adequately. Next, we should start the Apache Hadoop Cluster using the next commands:

```
$ start-dfs.sh
$ start-yarn.sh
$ jps # to verify all the running components
```

After that, Spark script can retrieve data from the input folder of HDFS and run the test.

For the first scenario which concerns comparing Hadoop to Linux, we should download our Github repo [4]. Then, we copy this script and the 9 datasets to the directory **/usr/local/hadoop/share /hadoop/mapreduce**. After that, we simply run the former script three times and get the average user time. We can see the results in this Excel file. At this stage, we can easily perform the comparison between Hadoop and Linux. For the second scenario which concerns comparing Hadoop to Spark, this file contains the comparison results, which shows that Spark has a much better performance. As for the social network scenario, before running the scripts, we needed to install Hadoop on our Azure VM [5]. After that, since we have used python instead of Java, we have used the instructions from this source [11]. In short, after installing Hadoop, first put the source file in Hadoop's HDFS directory. Then, we leverage Hadoop's streaming API to pass files between 'mapper.py' and 'reducer.py' by STDIN and STDOUT.

```
azureuser@ubuntu:/usr/local/hadoop
$ bin/hadoop jar contrib/streaming/hadoop-*streaming*.jar \
-mapper /home/azureuser/mapper.py \
-reducer /home/azureuser/reducer.py \
-input /user/azureuser/recom_example/cc_tp2.txt \
-output /user/azureuser/recome_output
```

# References

[1] *BASIC HDFS FILE OPERATIONS COMMANDS*. URL: https://www.alluxio.io/learn/hdfs/basic-file-operations-commands/.

[2] *Command-line Tools can be 235x Faster than your Hadoop Cluster*. URL: https://adamdrake.com/command-line-tools-can-be-235x-faster-than-your-hadoop-cluster.html.

[3] *Command-line Tools vs Hadoop*. URL: https://www.reddit.com/r/programming/comments/8ljjzm/commandline_tools_can_be_235x_faster_than_your/.

[4] *Github Repo for Lab2*. URL: https://github.com/marsup13/CCF_Lab2_Fall2021.

[5] *Hadoop single node cluster setup*. URL: https://hadoop.apache.org/docs/stable/hadoop-project-dist/hadoop-common/SingleCluster.html.

[6] *How to Install and Configure Hadoop on Ubuntu 20.04*. URL: https://tecadmin.net/install-hadoop-on-ubuntu-20-04/.

[7]     *How to Install Hadoop on Ubuntu 18.04 or 20.04.* URL: https://phoenixnap.com/kb/install-hadoop-ubuntu.

[8]     *Install and Configure Apache Hadoop on Ubuntu 20.04.* URL: https://www.vultr.com/docs/install-and-configure-apache-hadoop-on-ubuntu-20-04.

[9]     *RDD APIs.* URL: http://spark.apache.org/examples.html.

[10]    *RDD programming guide.* URL: https://spark.apache.org/docs/latest/rdd-programming-guide.html#resilient-distributed-datasets-rdds.

[11]    *Writing An Hadoop MapReduce Program In Python.* URL: https://www.michael-noll.com/tutorials/writing-an-hadoop-mapreduce-program-in-python/.

# A Appendix: Hadoop Configuration



Figure 21: Testing Hadoop (1)



Figure 22: Testing Hadoop (2)

Figure 23: Testing Hadoop (3)

# B  Appendix: Hadoop vs Linux Comparison



Figure 24: Looking through all datasets

Figure 25: Measuring the execution time of dataset 1



Figure 26: Measuring the execution time of dataset 1 (2)

Figure 27: Deleting output after each experiment



Figure 28: Using Linux command line to perform word frequencies of the dataset

Figure 29: Using Linux command line to perform word frequencies of the dataset (2)