

Polytechnique Montreal  
LOG8415 : Advanced Concepts of Cloud Computing  
Laboratory 2  
MapReduce with Hadoop on AWS

Christian Njon  
Dimitry Kamga  
Michelle Sepkap Sime  
Rui Jie Li

4th November 2022

## 1 Introduction

MapReduce on Hadoop and Spark using AWS is the subject of the second assignment for the LOG8415E course. The objectives of this assignment are to acquire some skills with large data technologies and learn to integrate issues and methods into the MapReduce paradigm. Four primary sections make up this report. First, we will discuss our Word Count application in Hadoop trials. Second, we compare Hadoop's performance to that of Linux. Third, we compare the performance of Spark and Hadoop on AWS. We wrap up by outlining our algorithm and the MapReduce tasks we used to tackle the social network problem. We present our recommendations for connections based on the algorithm.

## 2 Hadoop and Spark

### 2.1 Experiments with Word count Program

Here, we first prepare the lab setting by setting up Hadoop on our computer. We adhered to the assignment's guidelines. Our major goal was to use Hadoop to process a pg4300.txt file. So that the Hadoop Name Node and Data Nodes could share the file, we downloaded it to a local directory and then moved it to the Hadoop Distributed File System (HDFS). The data file pg4300.txt was then moved to the "input" directory we had just created in the Hadoop Distributed File System (HDFS). The wordcount.java program from the Hadoop example directory was then executed. The screen capture of the Hadoop settings on localhost is shown in Figure 1 and Figure 2. The input directory containing the pg4300.txt file is shown in Figure 3.

### 2.2 Performance comparison of Hadoop vs. Linux

In this part, we compared the word frequency computation capabilities of Hadoop with those of a standard PC running Linux.

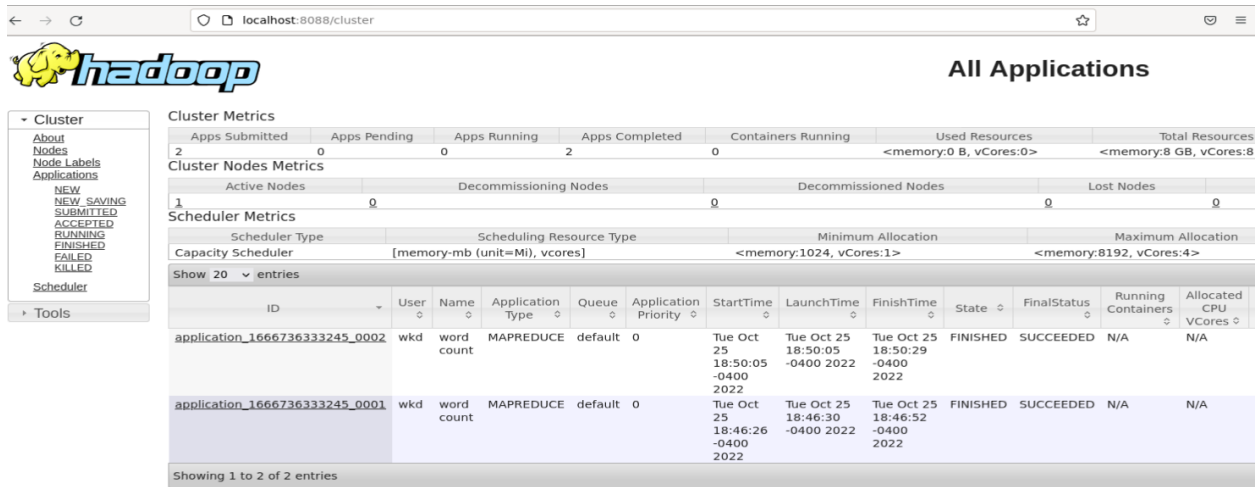


Figure 1: Hadoop Overview GUI - part1

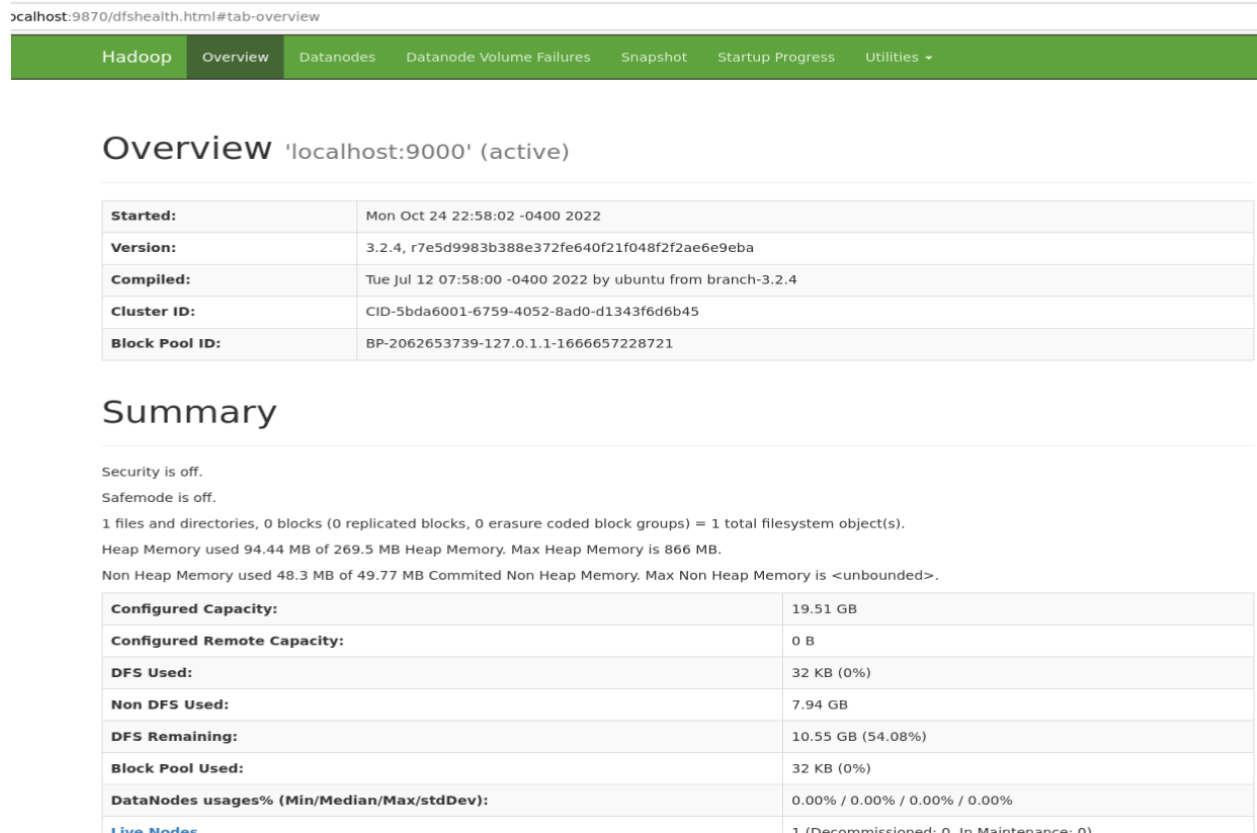


Figure 2: Hadoop Overview GUI - part2

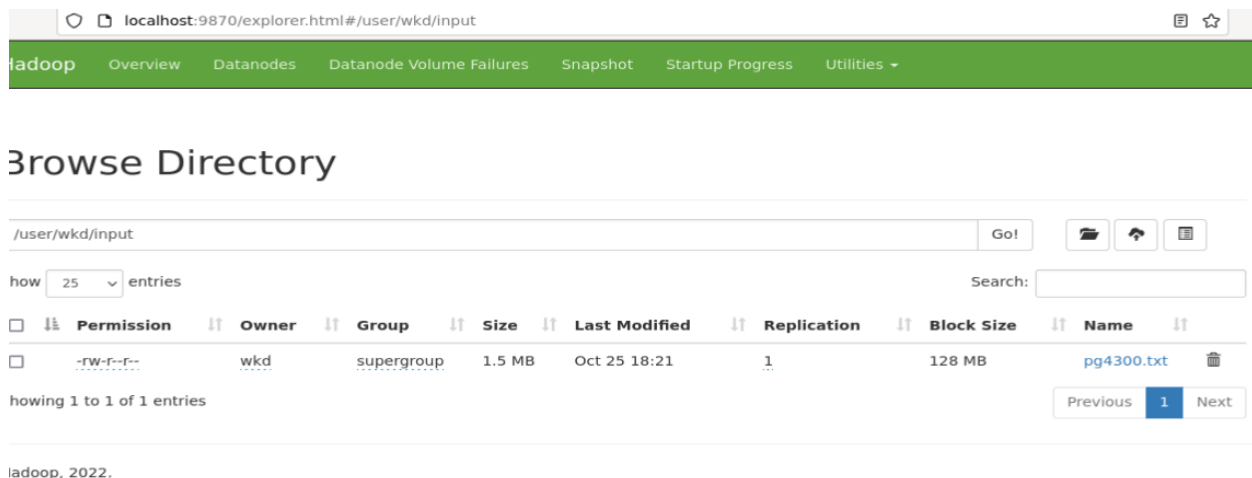


Figure 3: Hadoop directory

First we installed Hadoop and spark binaries. Then we ran the wordcount program with hadoop on a copy of James Joyce’s Ulysses book page 4300 available at [1]. The wordcount program just counts how many times each word appears in a file. We also di the same on an AWS M4.Large instance using the command `cat ./pg4300.txt |tr ' ' '\n' |sort |uniq -c`. Here are the results:

Table 1: Hadoop vs Linux Wordcount

Hadoop	Linux
5.961s	0.170s

As we can see, Linux completes the task more quickly than Hadoop. This is expected because Hadoop is acceptable or suitable for more sophisticated tasks than the one we used.

### 2.3 Performance comparison of Hadoop vs. Spark on AWS

We first set up our infrastructure as follows in order to compare the performances of Hadoop and Spark on AWS. We generated a M4.large linux Ubuntu instance , and we installed Hadoop 3.3.4 and Spark on it. We confirm the installation of all necessary packages. Then, we timed the WordCount program’s execution on both Hadoop and Spark machines three times across the entire dataset. We used 9 text files for this comparison, they can be found in the Datasets folders in Lab2, index shown in Figure 4.

Spark was anticipated to be considerably faster than Hadoop since it makes use of random access memory and this is the case indeed. The results are shown in the Figure 5.

This experiment and the last one has been automated through a bash script described in the section 3.

Dataset	Name	index
<a href="https://tinyurl.com/4vxdw3pa">https://tinyurl.com/4vxdw3pa</a>	buchanj-midwinter-00-t	input1
<a href="https://tinyurl.com/kh9excea">https://tinyurl.com/kh9excea</a>	carman-farhorizons-00-t	input2
<a href="https://tinyurl.com/dybs9bnk">https://tinyurl.com/dybs9bnk</a>	colby-champlain-00-t	input3
<a href="https://tinyurl.com/datumz6m">https://tinyurl.com/datumz6m</a>	cheyneyp-darkbahama-00-t	input4
<a href="https://tinyurl.com/j4j4xdw6">https://tinyurl.com/j4j4xdw6</a>	delamare-bumps-00-t	input5
<a href="https://tinyurl.com/ym8s5fm4">https://tinyurl.com/ym8s5fm4</a>	charlesworth-scene-00-t	input6
<a href="https://tinyurl.com/2h6a75nk">https://tinyurl.com/2h6a75nk</a>	delamare-lucy-00-t	input7
<a href="https://tinyurl.com/vwvram8">https://tinyurl.com/vwvram8</a>	delamare-myfanwy-00-t	input8
<a href="https://tinyurl.com/weh83uyn">https://tinyurl.com/weh83uyn</a>	delamare-penny-00-t	input9

Figure 4: Dataset index

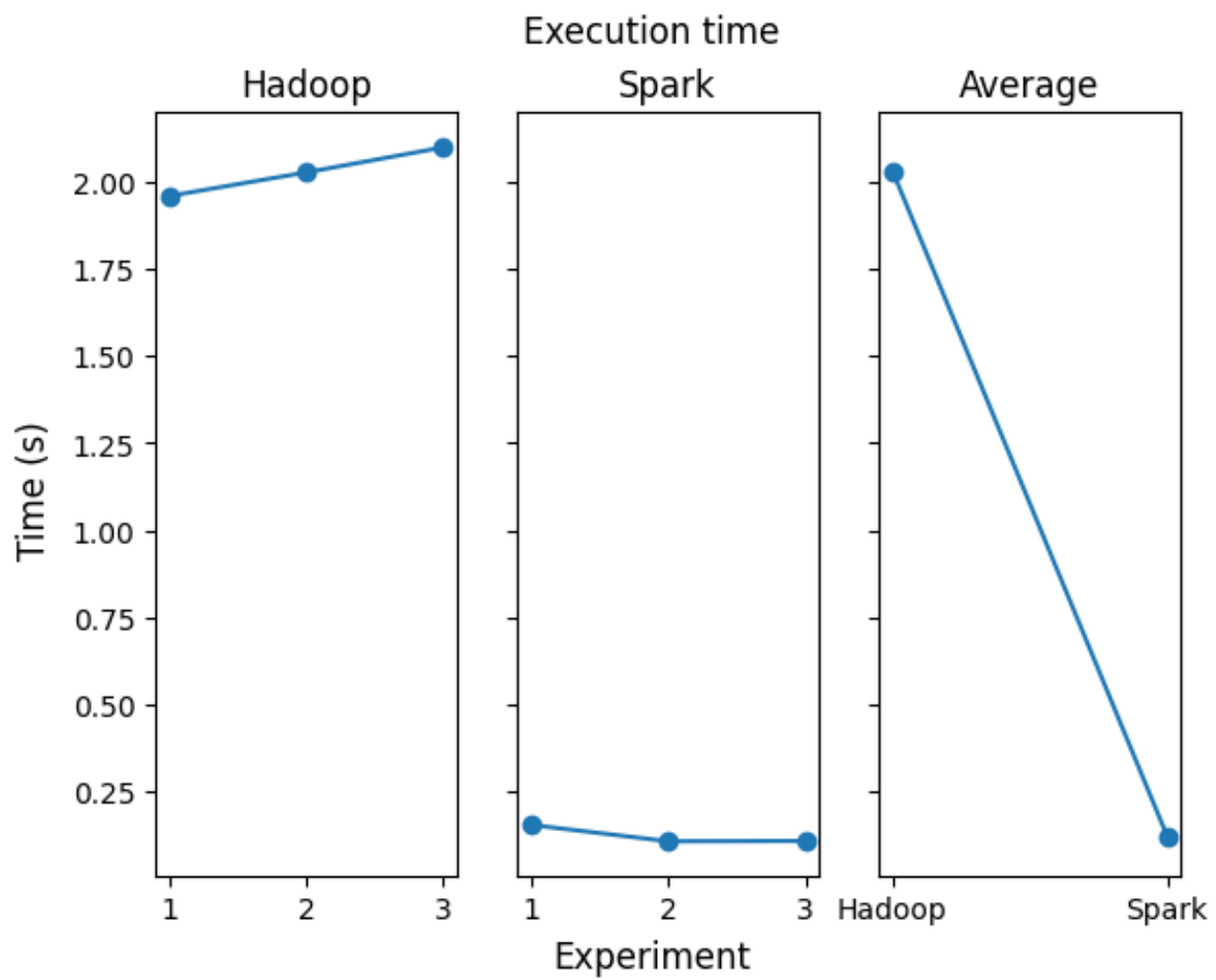


Figure 5: Performance comparison of Hadoop and Spark

### 3 Instructions to run the code

The entry point of the project is the bash script **run.sh** located at the root level of Lab2 folder. What this script does is simply to schedule all steps that need to be done in order to have the performance results. First, it checks whether the necessary credentials (`aws_access_key_id`, `aws_secret_access_key`, `aws_session_token`) and the region config are set. The check proceeds this way :

- Check if the default values have been set by aws cli by means of `configure` command,
- If not, check if they are available among the environment variables,
- If they aren't, get them from user input and export them to make them available for upcoming scripts.

Once credentials and minimum config are set, we create and activate a python virtual environment to install dependencies so that user python environment remains unchanged. Then, we deploy and setup infrastructure, it is composed of one M4.Large instance and a security group to allow SSH access. If the setup fails to complete, we teardown already created infrastructure and exit. During the setup we store SSH private key and public IP address for later use.

Infrastructure step completed, we go to the next one which is executing hadoop and spark programs via SSH (using `paramiko` python library), saving execution time to files (`results.txt` for hadoop vs linux, `hadoop.txt` for hadoop performance and `spark.txt` for spark performance) that are also retrieved by SSH. We save those results in files and plot them. Finally, as soon as we have all we wanted, the infrastructure is destroyed and the virtual environment deactivated.

## 4 Social media problem

### 4.1 Input file

The input file has the following format:

userID of a user, followed by a TAB character, followed by the ids of the friends of the user separated by a comma. For example,

1 2,3

3 1,2

means that user 1 is friends with user 2 and 3, and user 3 is friends with user 2 and 1. Each user and their friends is on a different line (separated by a ENTER character).

### 4.2 Map using Hadoop (Java)

In the `map()` function, Hadoop automatically splits the input file into lines. The function responsible for this is

```
MutualFriends.FriendOfFriendsMapper.map(Object key, Text value, Context context)
```

This function receives a value containing a string in the format of `userID+TAB+list of friends separated by a comma`, for example `1TAB2,3,4`, or if the user didn't add anyone, it will look like `1TAB`. The map functions works by

1. Splitting the value by the TAB character, which will result an an array of strings in the format of [characters before the TAB, characters after the TAB]. If there are no characters following the TAB, the array will have the format of [characters before the TAB]. For example, if the value is 1TAB2,3,4, the array will be ["1", "2,3,4"]; if the value is 1TAB, the array will be ["1"]
2. if the array resulting from splitting by the TAB character have a length of 1, that means the user did not add any friends yet. The key and the sent to the reducer will be
  - (a) key = user ID
  - (b) value = user ID + TAB + null

For example, if the value is "1TAB", the output to the reducer will be

- (a) key = 1
  - (b) value = 1TABnull
3. if the array resulting from splitting by the TAB character has a length of more than 1, it means the user added at least some friends. The map function will then create a variable **friendList = user ID + TAB + friends of user ID**. Then, for every user **u** that the user with **user ID** has added (the list can be obtained by splitting the second value of the array by the character ,), the mapper will output to the reducer
    - (a) key = u
    - (b) value = friendList

For example, if the array is ["1", "2,3,4"], the mapper will output

- (a) key = 2
- (b) value = 1TAB2,3,4
- (c) key = 3
- (d) value = 1TAB2,3,4
- (e) key = 4
- (f) value = 1TAB2,3,4

This ensures that the reducer will have every user along with all their friends of friends, and from which friend they come from. For example, if we have

- (a) 1 TAB 2,3,4,5
- (b) 6 TAB 2,3,4,5,7

For the user 2, the reducer will see

- (a) key = 2; value = 1TAB2,3,4,5
- (b) key = 2; value = 6TAB2,3,4,5,7

From this, we see that 2 can reach 2,3,4,5 via 1, and 2 can reach 2,3,4,5,7 via 6.

### 4.3 Reduce with Hadoop (Java)

The reducer collects, for each user, all their friends of friends along with from which friend these come from. For example, if 2 can reach 3 via 1 and 4, then 2 and 3 have two mutual friends, namely 1 and 4. The reduce function basically counts how many times a user ID appears in this lists, and returns the ones that appear most often after removing the ones that the user have already added. If the friend list contains null, the reducer outputs an empty list of suggestions. If not, the following steps will be taken: for example, if we have

- key = 2
- value = 1TAB2,3,4,5

Let `alreadyAdded` be a hashset of hashsets of strings that contains all the users that 2 has already added and `mutualFriendsCount` be a `HashMap` that uses a `HashSet` of strings as a key and an integer as value. The steps executed by the reducer will be the following (same steps for each value corresponding to 2):

1. The string 1TAB2,3,4,5 is splitted by tab to get ["1", "2,3,4,5"], and the value "2,3,4,5" is splitted by , to get the array `potentialFriends` = [2, 3, 4, 5]
2. The reducer will add `HashSet[(2)]` and `HashSet[(2,1)]` in `alreadyAdded` since 2 cannot add 2 and 2 has already added 1.
3. For each user `u` in the list `potentialFriends`, a `HashSet h = HashSet([2,u])` will be created. If `h` is already in `mutualFriendsCount`, it will be added as a key with a value of 1. If it's already added, we will increase its value by 1.

After all the values associated with 2 have been processed, the keys contained in `alreadyAdded` will be removed from `mutualFriendsCount`. The rest will be stored in an array of `Pair` objects, with a `count` attribute that indicates the number of mutual friends and a `uid` attribute that indicates the user ID. The array is then sorted by `Pair.count` (descending order) then by `Pair.uid` (ascending order). The uid of the 10 first `Pair` objects (or the whole array, if less than 10 friends of friends) will be returned. Example:

- key = 2
- value = 1TAB2,3,4,5
- value = 6TAB4,3,7,5
- value = 8TAB6,3,4,5

After processing all the values, `alreadyAdded` will contain

```
HashSet([
  HashSet([2]),
  HashSet([2, 1]),
  HashSet([2, 6]),
  HashSet([2, 8])
])
```

And `mutualFriendsCount` will contain



```
HashMap([
    HashSet([2, 1]): 1,
    HashSet([2, 3]): 3,
    HashSet([2, 4]): 2,
    HashSet([2, 5]): 3,
    HashSet([2, 6]): 2,
    HashSet([2, 7]): 1,
    HashSet([2, 8]): 1,
])
```

After removing the contents of `alreadyAdded` from `mutualFriendsCount`, `mutualFriendsCount` will contain

```
HashMap([
    HashSet([2, 3]): 3,
    HashSet([2, 4]): 2,
    HashSet([2, 5]): 3,
    HashSet([2, 7]): 1,
])
```

And an array containing

```
array = [
    Pair{uid="3", count= 3},
    Pair{uid="4", count= 2},
    Pair{uid="5", count= 3},
    Pair{uid="7", count= 1},
]
```

will be produced, then sorted first according to count, then according to uid (ex. "3" comes before "5" because 3 < 5, even though they both have three mutual friends with 2):

```
array = [
    Pair{uid="3", count= 3},
    Pair{uid="5", count= 3},
    Pair{uid="4", count= 2},
    Pair{uid="7", count= 1},
]
```

The output will be

- key = 2
- value = 3,5,4,7

## 4.4 Result

For the following users, we have:

```
924    439,2409,6995,11860,15416,43748,45881
8941   8943,8944,8940
8942   8939,8940,8943,8944
9019   9022,317,9023
9020   9021,9016,9017,9022,317,9023
```

9021 9020,9016,9017,9022,317,9023  
9022 9019,9020,9021,317,9016,9017,9023  
9990 13134,13478,13877,34299,34485,34642,37941  
9992 9987,9989,35667,9991  
9993 9991,13134,13478,13877,34299,34485,34642,37941

All the users have received either 10 recommendations or as many as possible. For example, one of the result is

0 38737,18591,27383,34211,337,352,1532,12143,12561,17880

we can verify manually that 0 and 38737 have 5 friends in common, and 0 and 18591 have 4 mutual friends in common, and so on.

## 4.5 How to run the program

All the code for running the program is in compile.sh. The only parts that need to be changed manually is the input and the output directory, along with the hadoop path (if necessary). To run compile.sh, just type ./compile.sh in its directory.

However, if it does not work, it can be ran this way (code based on Apache Hadoop word count tutorial, 2022):

```
rm -r path/to/output
export JAVA_HOME="/path/to/java/"
export PATH=${JAVA_HOME}/bin:${PATH}
export HADOOP_CLASSPATH=${JAVA_HOME}/lib/tools.jar
/path/to/hadoop/bin/hadoop com.sun.tools.javac.Main MutualFriends.java Pair.java
jar cf tp2.jar *.class
/path/to/hadoop/bin/hadoop jar tp2.jar MutualFriends /path/input /path/output
```

## 4.6 Results on Azure

We also run this program on Azure (this is a screenshot of the content of the output file):

1. Create a Linux virtual machine with the following setup in the "basic" section (the rest can be left to default settings):

```
name of the virtual machine: mapreduce
availability region: (Canada) Canada central
Select inbound ports: HTTP, SSH, HTTPS
Username: can be anything
key pair name: can be anything
```

Subscription \* ⓘ Azure subscription 1

Resource group \* ⓘ (New) mapreduce\_group  
[Create new](#)

### Instance details

Virtual machine name \* ⓘ mapreduce ✓


Region \* ⓘ (Canada) Canada Central

Availability options ⓘ Availability zone

Availability zone \* ⓘ Zones 1

🔒 You can now select multiple zones. Selecting multiple zones will create one VM per zone. [Learn more](#)

Security type ⓘ Standard

Image \* ⓘ  Ubuntu Server 20.04 LTS - Gen2 (free services eligible)  
[See all images](#) | [Configure VM generation](#)

VM architecture ⓘ  
☐ Arm64  
☒ x64

Run with Azure Spot discount ⓘ ☐

**i** You are in the free trial period. Costs associated with this VM can be covered by any remaining credits on your subscription.  
[Learn more](#) ⓘ

Size \* ⓘ Standard\_B1s - 1 vcpu, 1 GiB memory (\$8.47/month) (free services eligible)  
[See all sizes](#)

Figure 6: Basic configuration

## Create a virtual machine ...

**i** You are in the free trial period. Costs associated with this VM can be covered by any remaining credits on your subscription. [Learn more](#)

Size \* ⓘ

Standard\_B1s - 1 vcpu, 1 GiB memory (\$8.47/month) (free services eligible) ▼

[See all sizes](#)

Administrator account

Authentication type ⓘ

☒ SSH public key

☐ Password

**i** Azure now automatically generates an SSH key pair for you and allows you to store it for future use. It is a fast, simple, and secure way to connect to your virtual machine.

Username \* ⓘ

azure-user-here ✓

SSH public key source

Generate new key pair ▼

Key pair name \*

mapreduce\_key ✓

Inbound port rules

Select which virtual machine network ports are accessible from the public internet. You can specify more limited or granular network access on the Networking tab.

Public inbound ports \* ⓘ

☐ None

☒ Allow selected ports

Select inbound ports \*

HTTP (80), HTTPS (443), SSH (22) ▼

**⚠** This will allow all IP addresses to access your virtual machine. This is only recommended for testing. Use the Advanced controls in the Networking tab to create rules to limit inbound traffic to known IP addresses.

Figure 7: Basic configuration (next section)

2. For disks, select "Standard HDD" for better cost. The rest of the options can be left as default (Networking, Management, Monitoring, Advanced, Tags).

[Basics](#)
[Disks](#)
[Networking](#)
[Management](#)
[Monitoring](#)
[Advanced](#)
[Tags](#)
[Review + create](#)

Azure VMs have one operating system disk and a temporary disk for short-term storage. You can attach additional data disks. The size of the VM determines the type of storage you can use and the number of data disks allowed. [Learn more](#)

### VM disk encryption

Azure disk storage encryption automatically encrypts your data stored on Azure managed disks (OS and data disks) at rest by default when persisting it to the cloud.

Encryption at host ☐

Encryption at host is not registered for the selected subscription. [Learn more about enabling this feature](#)

### OS disk

OS disk type \* ☐ Standard HDD (locally-redundant storage) ☐  
 If performance is critical for your workloads, choose Premium SSD disks for lower latency, higher IOPS and bandwidth, and bursting. [Learn more](#)

Delete with VM ☒

Key management ☐ Platform-managed key ☐

Enable Ultra Disk compatibility ☐

### Data disks for mapreduce

You can add and configure additional data disks for your virtual machine or attach existing disks. This VM also comes with a temporary disk.

LUN	Name	Size (GiB)	Disk type	Host caching	Delete with VM <input type="checkbox"/>
<a href="#">Create and attach a new disk</a> <a href="#">Attach an existing disk</a>					

☒ Advanced

[Review + create](#)
[< Previous](#)
[Next : Networking >](#)

Figure 8: Disk configuration

- After the machine is up and running, an SSH key pair is generated on a local linux machine, and the public key of the virtual machine (on Azure) is updated using the "reset password" functionality

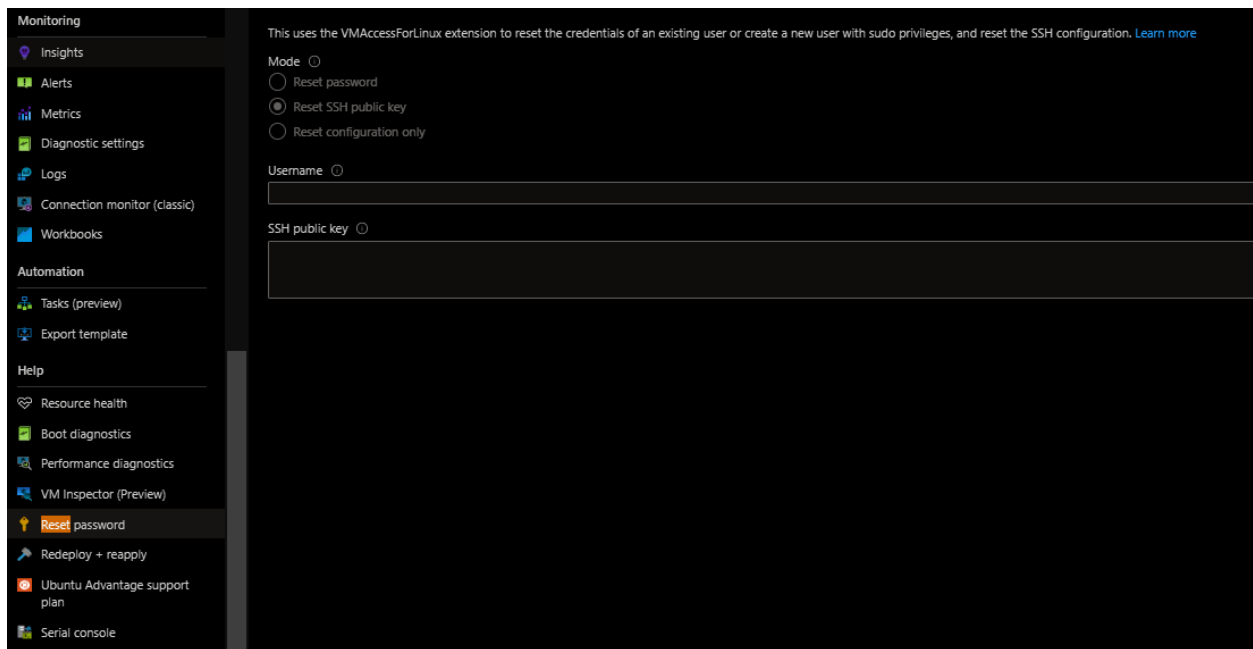


Figure 9: Changing public key (public key not shown for security reasons)

- the following commands are used to connect to the VM:

```
# on a local machine (step 3)
chmod 400 path/to/private_key.pem
ssh -i path/to/private_key.pem <azure user name>@<public IP of VM>
```

```
# on the Azure Virtual machine
sudo apt-get update
sudo apt-get install openjdk-17-jre-headless
export JAVA_HOME="/path/to/java"
sudo wget <link to hadoop version>
sudo tar -xvzf path/to/hadoop.tar.gz
```

```
# on the local machine
sudo scp -i path/to/private/key tp2.jar \
    <azure user name>@<public IP of VM>:/home/azureuser
sudo scp -i path/to/compile.sh \
    <azure user name>@<public IP of VM>:/home/azureuser
./compile.sh
```

```
rulika@mapreduce: ~
File Edit View Search Terminal Help
9950 2618,4875,9959,9962,9969,9971,9976,11818,22868,24983
9950 2882,6543,9953,9955,9958,9975,18459,18499,12844,12848
996 935,936,937,938,939,948,941,942,943,944
9968 4875,9962,9969,9971,9976,982,1113,2106,2231,2244
9961 12289,28682,13273,7184,13287,13289,28693,778,13274,13281
9962 14138,13888,13891,13911,13966,14843,14844,14851,14873,14894
9963 2822,9788,14264,17401,52,134,457,575,596,680
9964 4839,4858,9969,18584,28924,39171,348,540,1238,1513
9965 348,2196,2244,4717,4991,5895,5254,5273,7229,7839
9966 4848,562,579,2244,2882,3937,5033,9969,18468,12187
9967 24811,24819,2244,3388,4782,9958,18488,12438,14383,14339
9968 4991,2196,9972,9973,339,439,583,1888,1188,1118
9969 7479,9927,9971,13854,14271,14284,14293,14318,14562,23178
997 935,936,937,938,939,948,941,942,943,944
9978 3931,4234,7852,15843,18523,22147,34825,3782,4197,4282
9971 9969,9953,9958,9968,9962,9976,19984,338,2196,2244
9972 2196,4991,9968,22517,2244,4742,5957,5968,7259,7482
9973 3937,548,2981,8707,9778,9778,9958,9968,14881,18532
9974 5937,28,216,985,1886,2892,2196,2244,3456,3887
9975 14318,622,8154,9954,9955,9959,9969,18789,18881,18814
9976 14882,18532,7697,12229,16631,3412,39719,39728,1423,3782
9977 9958,186,622,1196,1488,2196,2244,2535,2828,2823
9978 588,4389,4991,9924,8815,8871,9415,18279,18498,11383
9979 8457,9978,21388,25866,354,455,577,2892,2899,2196
998 19,935,936,937,938,939,948,941,942,943
9988 52,351,4081,9869,7611,7651,8768,9951,9962,11612
9981 38489,44295,127,336,364,365,772,438,546,615
9982 31278,14454,17999,38819,48618,48896,37,123,237,299
9983 4589,17199,1426,1973,11757,1387,17188,1421,17193,23993
9984 4408,4589,1382,1387,1388,1421,1772,1973,1982,7284
9985 18411,522,889,428,9774,11191,16185,28698,22147,33481
9986 522,579,18532,4839,4841,21935,38868,562,621,2882
9987 9992,13134,13478,13877,34299,34485,34642,37941
9988 35667
9989 9992,13134,13478,13877,34299,34485,34642,37941
999 935,936,937,938,939,948,941,942,943,944
9998 13134,13478,13877,34299,34485,34642,37941
9991 9993,9994,9992,13134,13478,13877,34299,34485,34642,37941
9992 9987,9989,35667,9991
9993 9991,13134,13478,13877,34299,34485,34642,37941
9994 9991,13134,13478,13877,34299,34485,34642,37941
9995 18807,37356,36693,37188,18116,18123,37135,36665,36679,48748
9996 36679,36854,18828,18135,18888,18888,18826,18872,18896,44858
9997 18818,18864,18886,18385,9998,9999,18814,18827,18859,18863
9998 25256,27964,23176,23197,23676,7749,27571,27586,27598,27617
9999 36764,44132,18858,44888,36765,36889,18855,44868,44876,18888
rulika@mapreduce: ~$ ls
compile.sh  hadoop-3.1.4  input  output  tp2.jar
rulika@mapreduce: ~$ ls output
_SUCCESS  part-r-88888
rulika@mapreduce: ~$
```

Figure 10: Output of MapReduce from Microsoft Azure

## 5 References

- [1] Gutenberg textual data. <http://www.gutenberg.org/cache/epub/4300/pg4300.txt>.
- [2] Github repo. <https://github.com/MichelleSS1/Lab8415>
- [3] Apache Software Foundation. (2022, July 29). *Mapreduce tutorial*. Apache Hadoop. Retrieved November 4, 2022, from <https://hadoop.apache.org/docs/stable/hadoop-mapreduce-client/hadoop-mapreduce-client-core/MapReduceTutorial.html>