# I. Setup virtual machine (VMWare Workstation Pro)

## 1. Install VMWare Workstation Pro and Ubuntu.

### a. Preparation

*Step 1: Download VMWare Workstation Pro.*

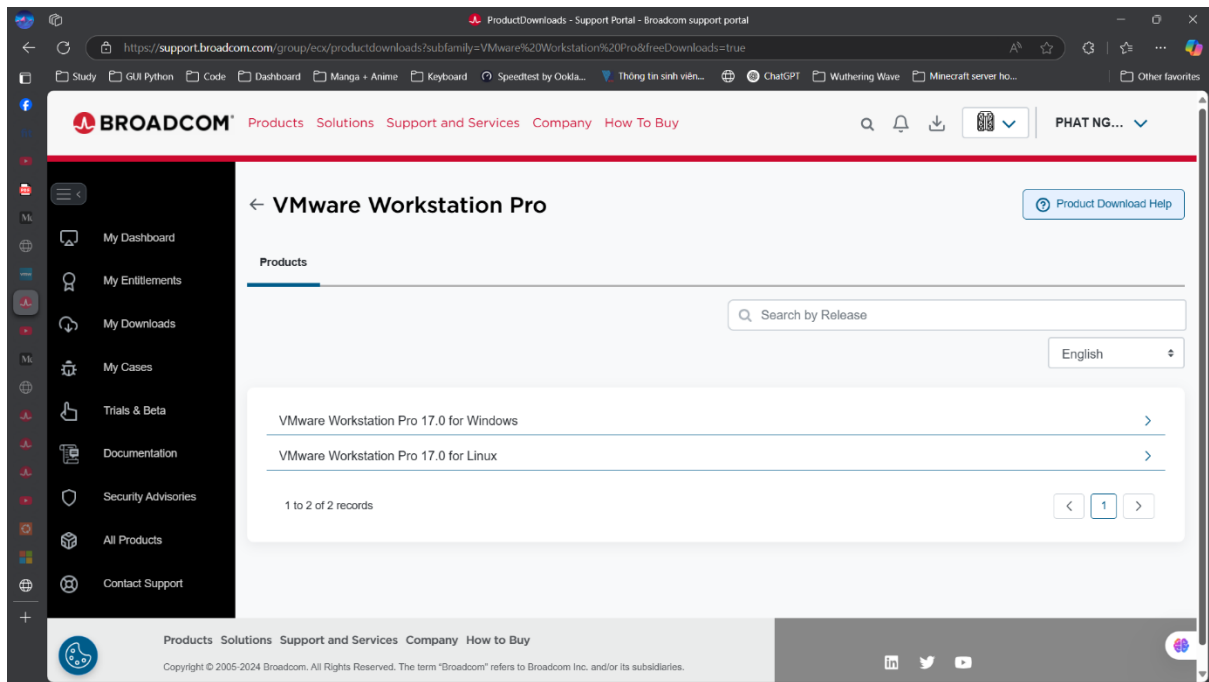- Go to Vmware Workstation Pro website to download the lastest version.



*Figure 1: VMware Workstation Pro official website.*

- **Go to Ubuntu official website to download the suitable (x64) OS**
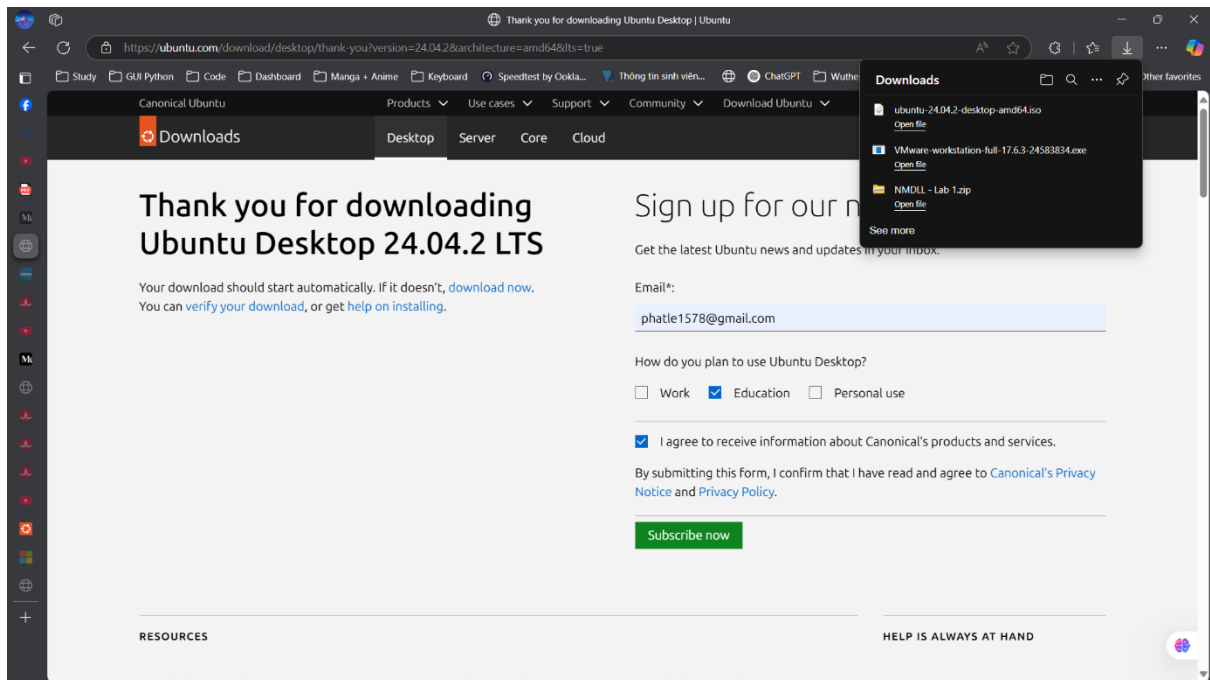


***Figure 2: Ubuntu official download website.***

## b. Setup the VMware Workstation Pro

*Step 1: Run VMware Workstation Pro.*

Choose "*Create a New Virtual Mahchine*" option on the main menu.
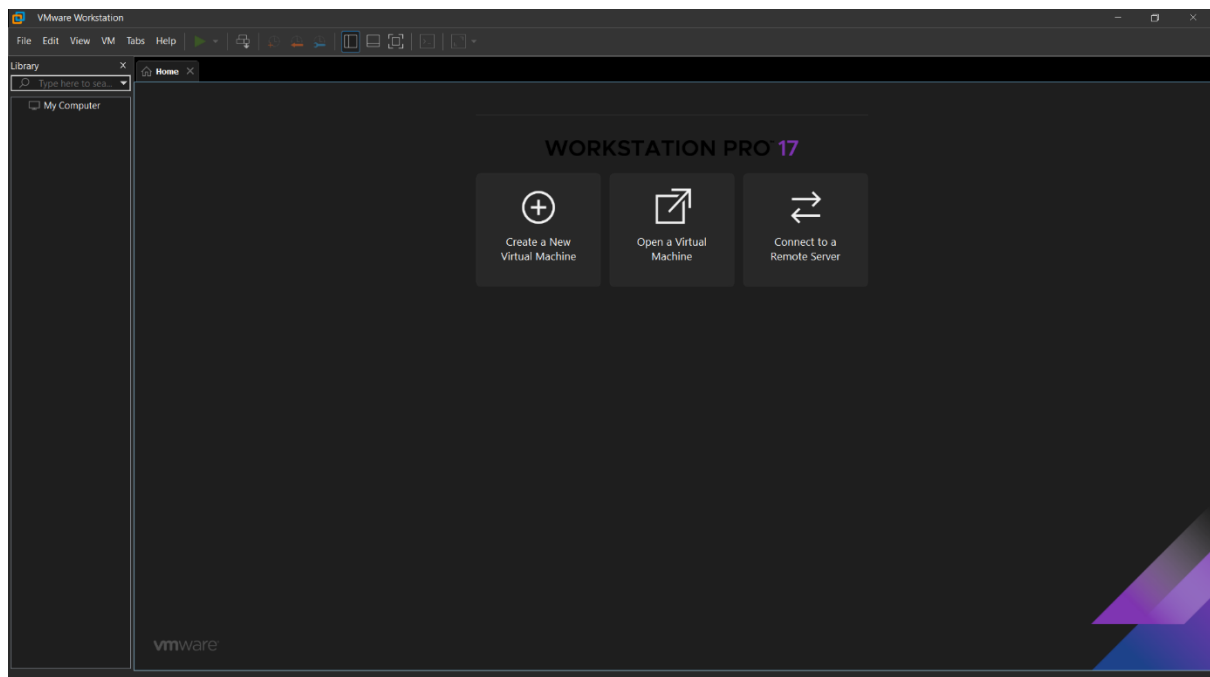


***Figure 3: VMware Workstation Pro main menu.***

## Step 2: Setup Virtual Machine.

- Choose "*Installer disc image file (iso)*" option and select the Ubuntu file we download previously.
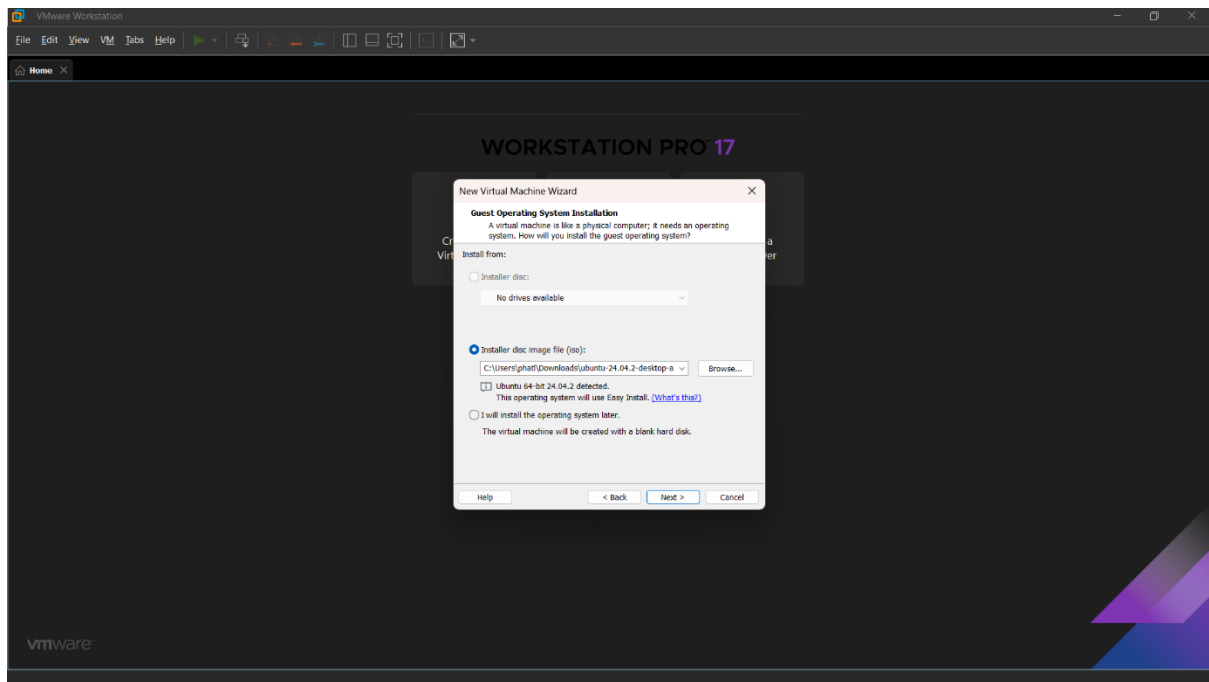


***Figure 4: VMware Workstation Pro virtual machine installer.***
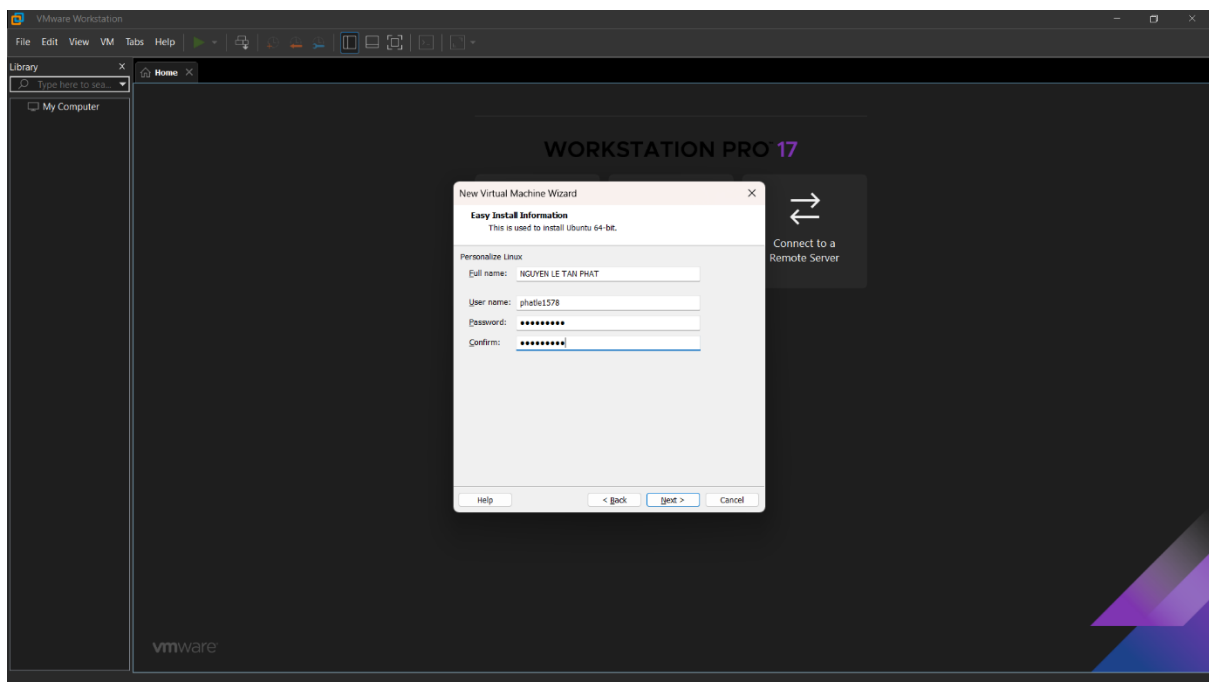
## Step 3: Setup Linux profile.

- Fill in some information.



*Figure 5: Personalize Linux.*

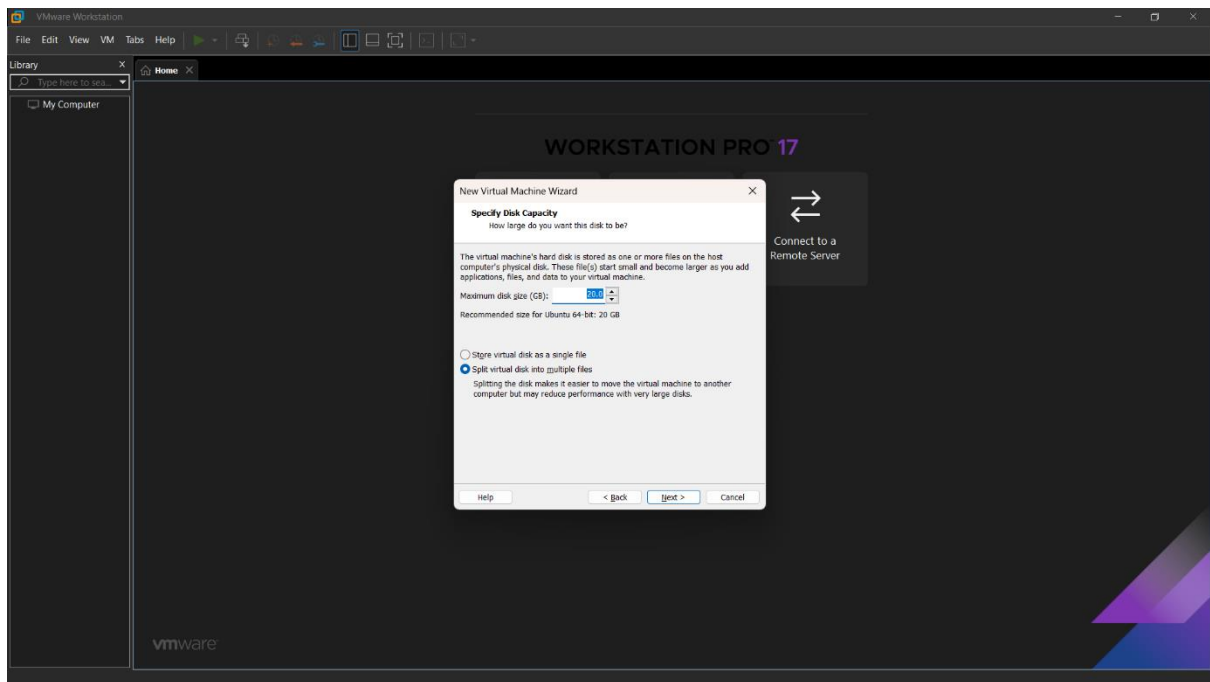*Step 4: Setup virtual machine specs.*

- Keep default and select "*next*"

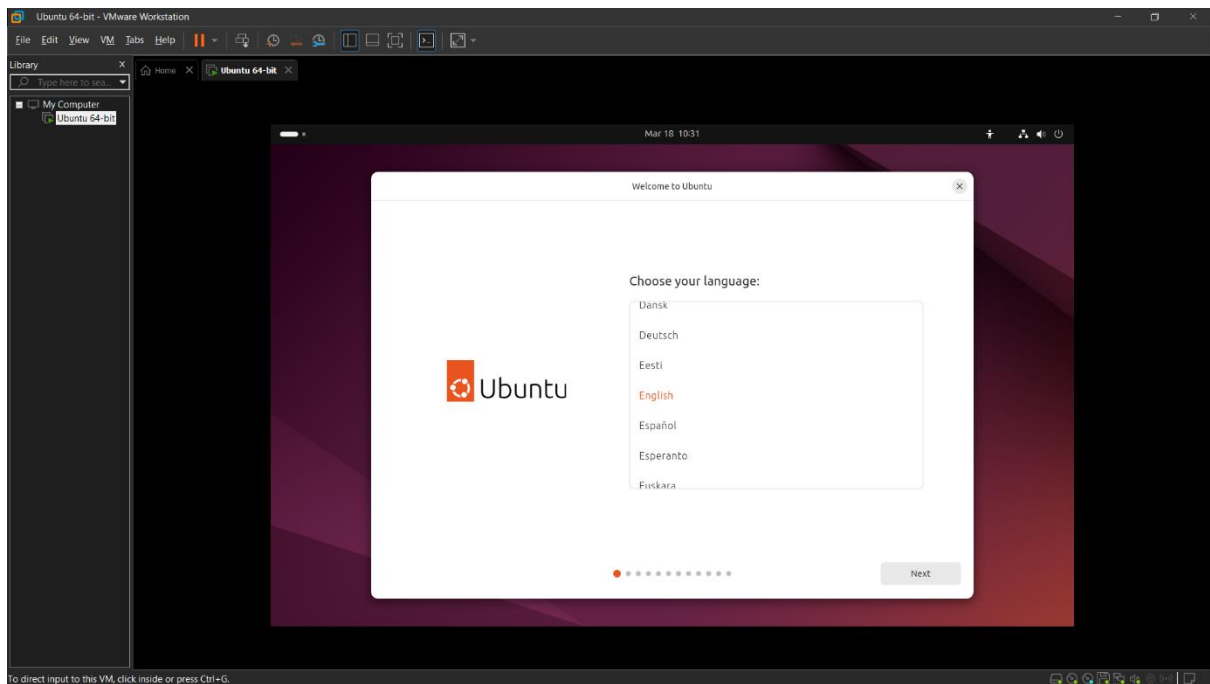

***Figure 6: Virtual machine spec.***



***Figure 7: Ubuntu setup (begin)***

### c. Setup Linux system in the virtual machine.

*Step 1: Setup Ubuntu.*
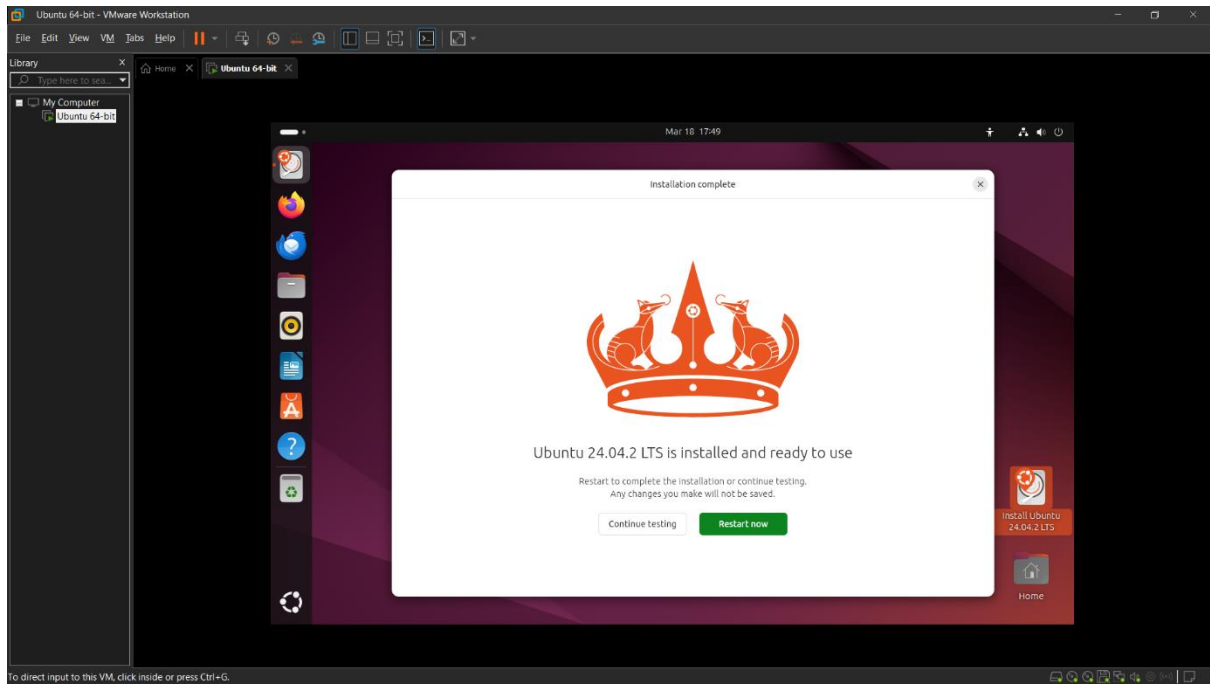
- Choose default settins and just select "*next*"



***Figure 8: Ubuntu setup (finish)***

## 2. Setup Hadoop Cluster (Pseudo – Distributed Mode).

### a. Install Hadoop (VMware Workstation Pro).

*Step 1: System preparation.*

- sudo apt install open-jdk-8-jdk: Install OpenJDK 8, a requirement for Hadoop and the most stable version to run Hadoop
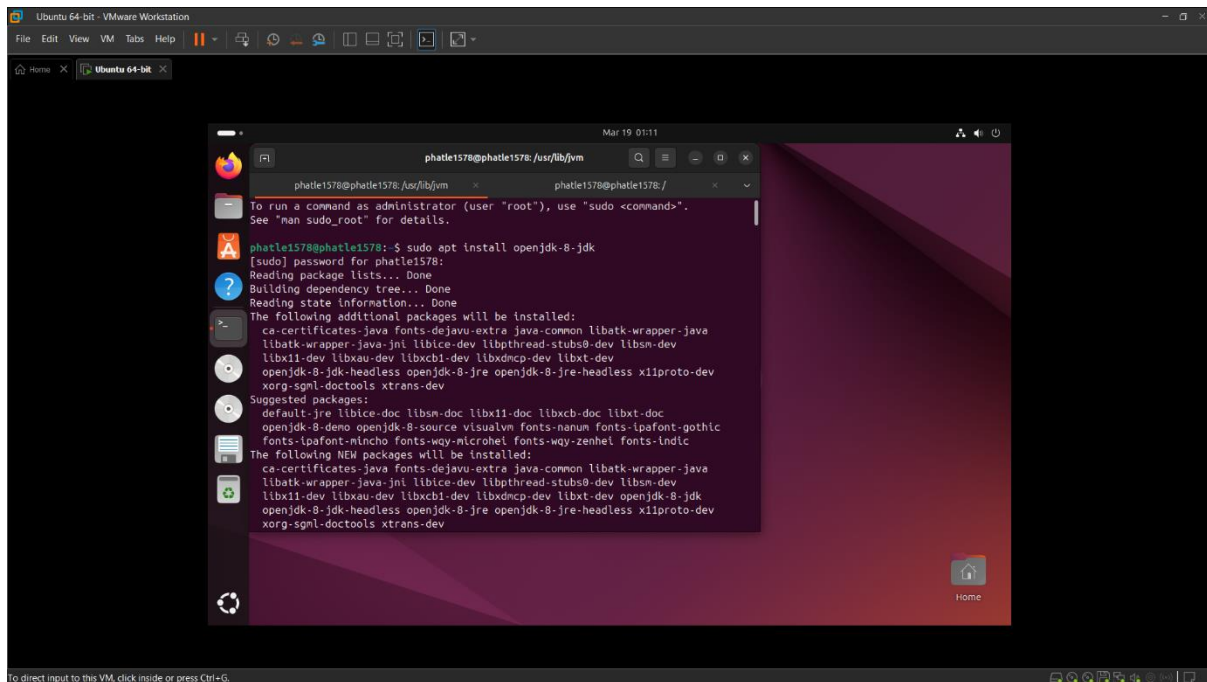


*Figure 9: Install OpenJDK 8*

- sudo nano .bashrc: Configure Hadoop enviroment variables.
- Go to the end of the file and add these line.

```
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
export PATH=$PATH:/usr/lib/jvm/java-8-openjdk-amd64/bin
export HADOOP_HOME=~/hadoop-3.4.1 /
export PATH=$PATH:$HADOOP_HOME/bin
export PATH=$PATH:$HADOOP_HOME/sbin
export HADOOP_MAPRED_HOME=$HADOOP_HOME
export YARN_HOME=$HADOOP_HOME
export HADOOP_CONF_DIR=$HADOOP_HOME/etc/hadoop
export HADOOP_COMMON_LIB_NATIVE_DIR=$HADOOP_HOME/lib/native
export HADOOP_OPTS="-Djava.library.path=$HADOOP_HOME/lib/native"
export HADOOP_STREAMING=$HADOOP_HOME/share/hadoop/tools/lib/hadoop-streaming-3.4.1.jar
export HADOOP_LOG_DIR=$HADOOP_HOME/logs
export PDSH_RCMD_TYPE=ssh
```
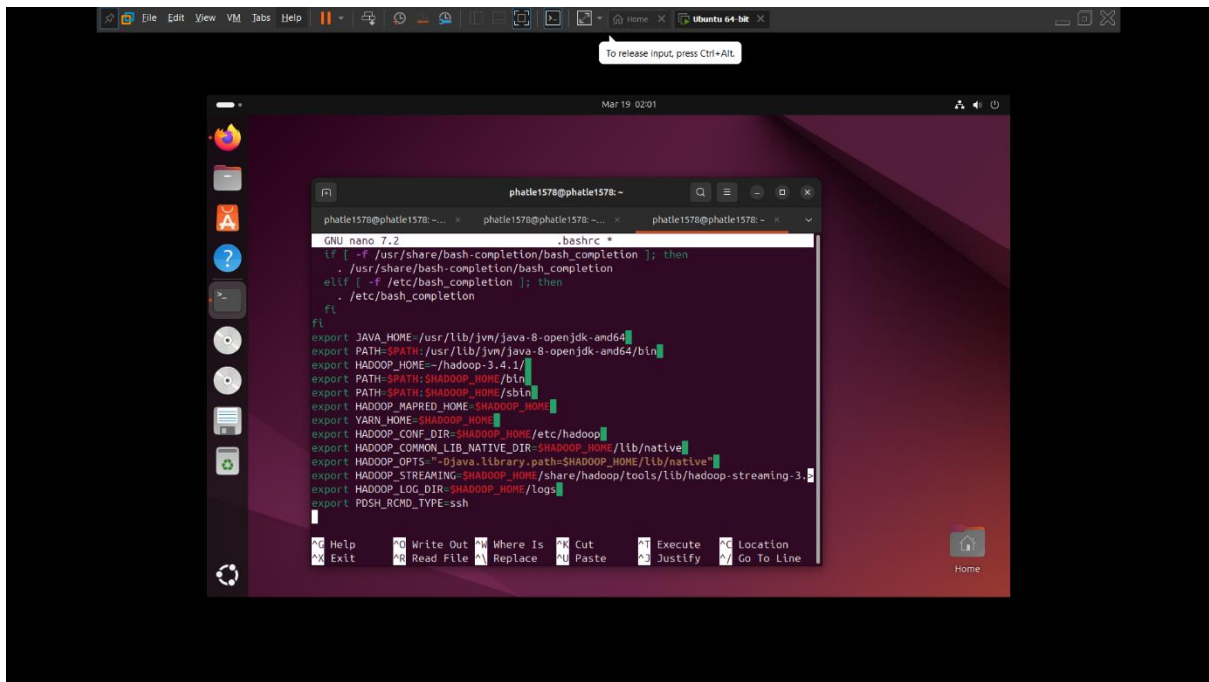
*Figure 10: Configure bashrc file*

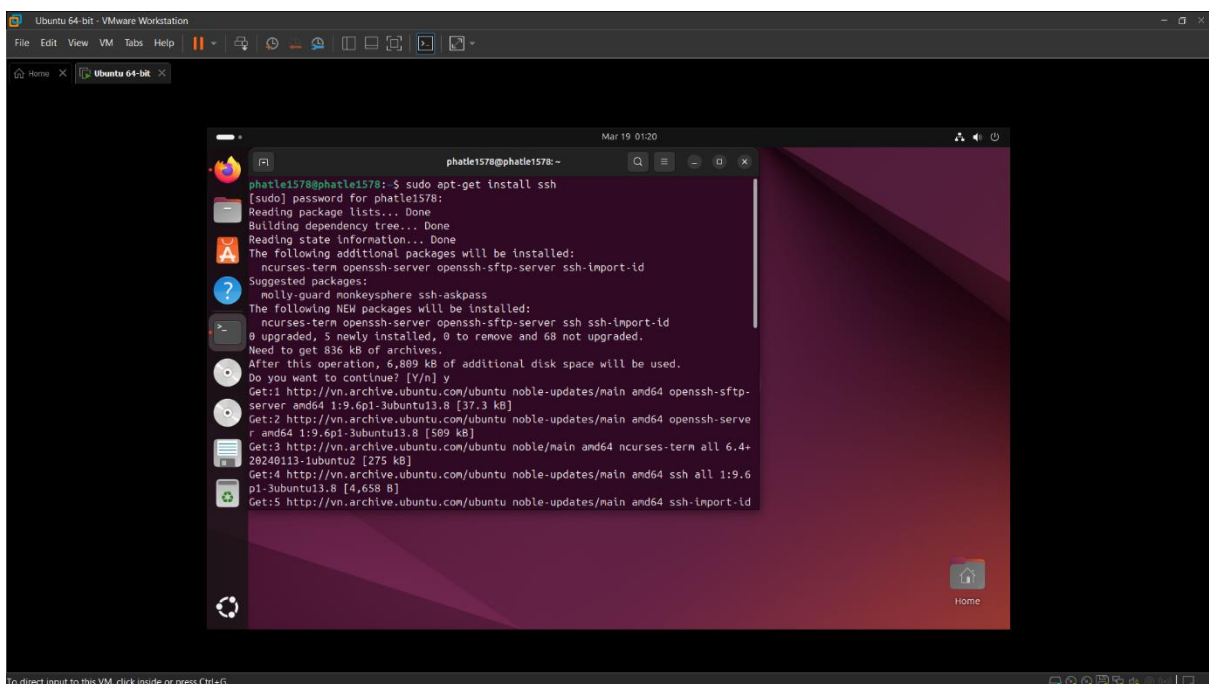- sudo apt-get install ssh: Install ssh.



*Figure 11: Install ssh*

## Step 2: Download Apache Hadoop

- wget https://dlcdn.apache.org/hadoop/common/hadoop-3.4.1/hadoop-3.4.1.tar.gz: download Apache Hadoop lastest version.



*Figure 12: Download Apache Hadoop lastest version*

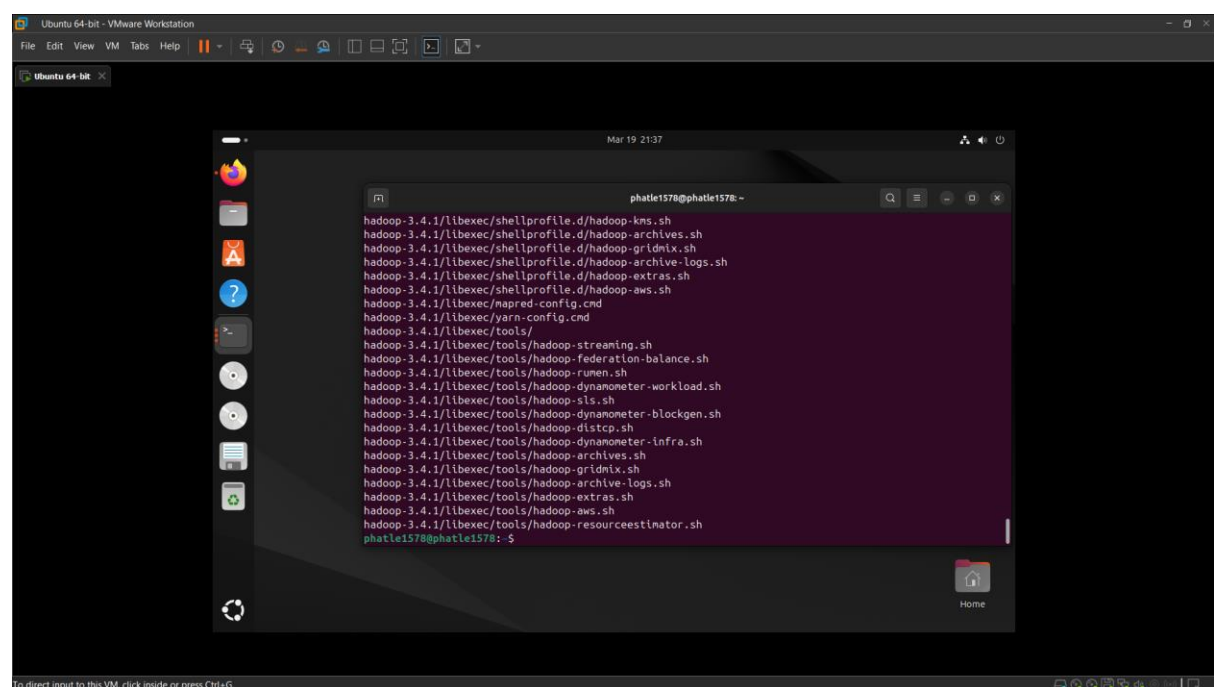- tar -zxvf ~/Downloads/hadoop-3.4.1.tar.gz: Extract the binary Hadoop file we just download.



*Figure 13: Extract the binary Hadoop file.*

*Step 3: Configure java environment variables and other xml files.*

### i. hadoop-env.sh

- cd /etc/hadoop/: go to /etc/hadoop directory
- sudo nano hadoop-env.s: open hadoop-env.sh to configure
- Search for the JAVA_HOME and add this line to set the Java home path:

JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd4



***Figure 14: Configure hadoop-env.sh***
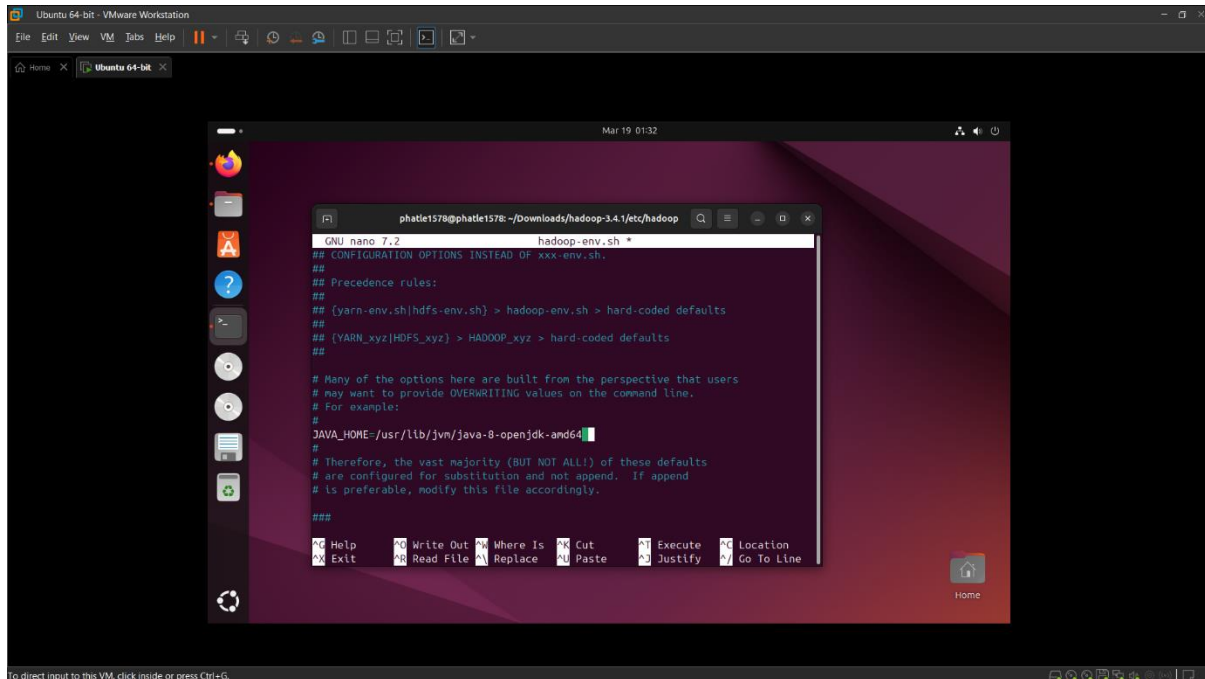
- Ctrl + O: Save
- Ctrl + X: Exit

### ii. core-site.xml

- sudo nano core-site.xml: Open and configure core-site.xml for pseudo-distributed mode
- Navigate to the end of the file and add these lines:

```
<configuration>
   <property>
      <name>fs.defaultFS</name>
      <value>hdfs://localhost:9000</value>
   </property>
</configuration>
```

- Ctrl + O: Save.
- Ctrl + X: Exit

*Figure 15: Configure core-site.xml*

### iii.    hdfs-site.xml

- sudo nano hdfs-site.xml: Open and configure hdfs-site.xml.
- Navigate to the end of the file hdfs-site.xml and add these lines:

```
<configuration>
 <property>
 <name>dfs.replication</name>
 <value>1</value>
 </property>
</configuration>
```
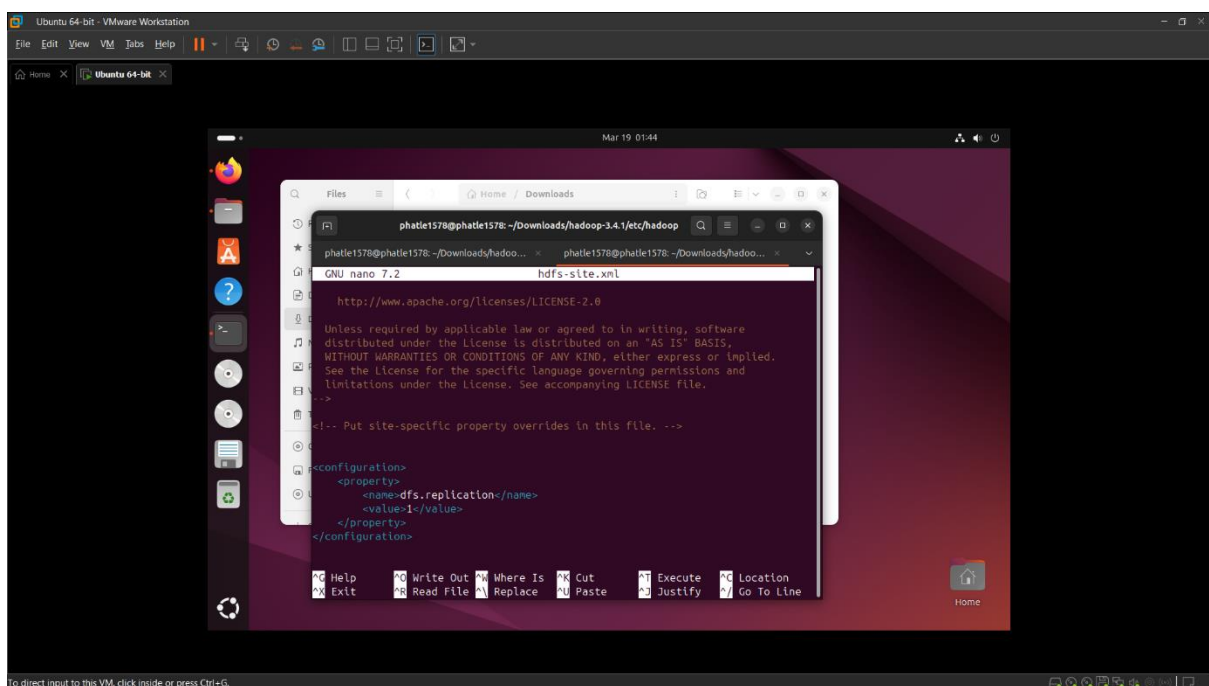
- Ctrl + O: Save.
- Ctrl + X: Exit.



*Figurte 16: Configure hdfs-site.xml*

### iv. mapred-site.xml

- sudo nano mapred-site.xml: Open and configure mapred-site.xml.
- Navigate to the end of the file and add these lines:

```
<configuration>
 <property>
 <name>mapreduce.framework.name</name>  <value>yarn</value>
 </property>
 <property>
 <name>mapreduce.application.classpath</name>

<value>$HADOOP_MAPRED_HOME/share/hadoop/mapreduce/*:$HADOOP_MAPRED_H
OME/share/hadoop/mapreduce/lib/*</value>
 </property>
</configuration>
```

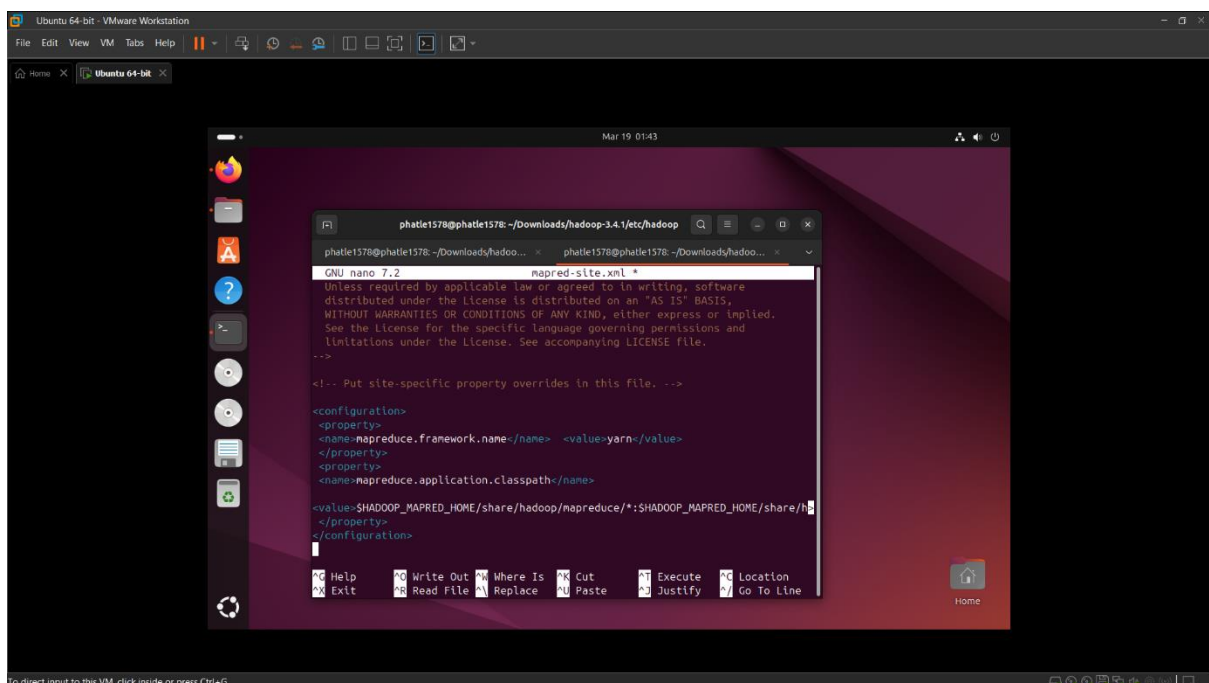- Ctrl + O: Save.
- Ctrl + X: Exit.



*Figure 17: Configure mapred-stie.xml*

### v. yarn-site.xml

- sudo nano yarn-site.xml: Open and configure yarn-site.xml.
- Navigate to the end of the file and add these lines:

```
<configuration>
 <property>
 <name>yarn.nodemanager.aux-services</name>
 <value>mapreduce_shuffle</value>
 </property>
 <property>
 <name>yarn.nodemanager.env-whitelist</name>
```

```
<value>JAVA_HOME,HADOOP_COMMON_HOME,HADOOP_HDFS_HOME,HADOOP_
CONF_DIR,CLASSPATH_PREP
END_DISTCACHE,HADOOP_YARN_HOME,HADOOP_MAPRED_HOME</value>
 </property>
</configuration>
```
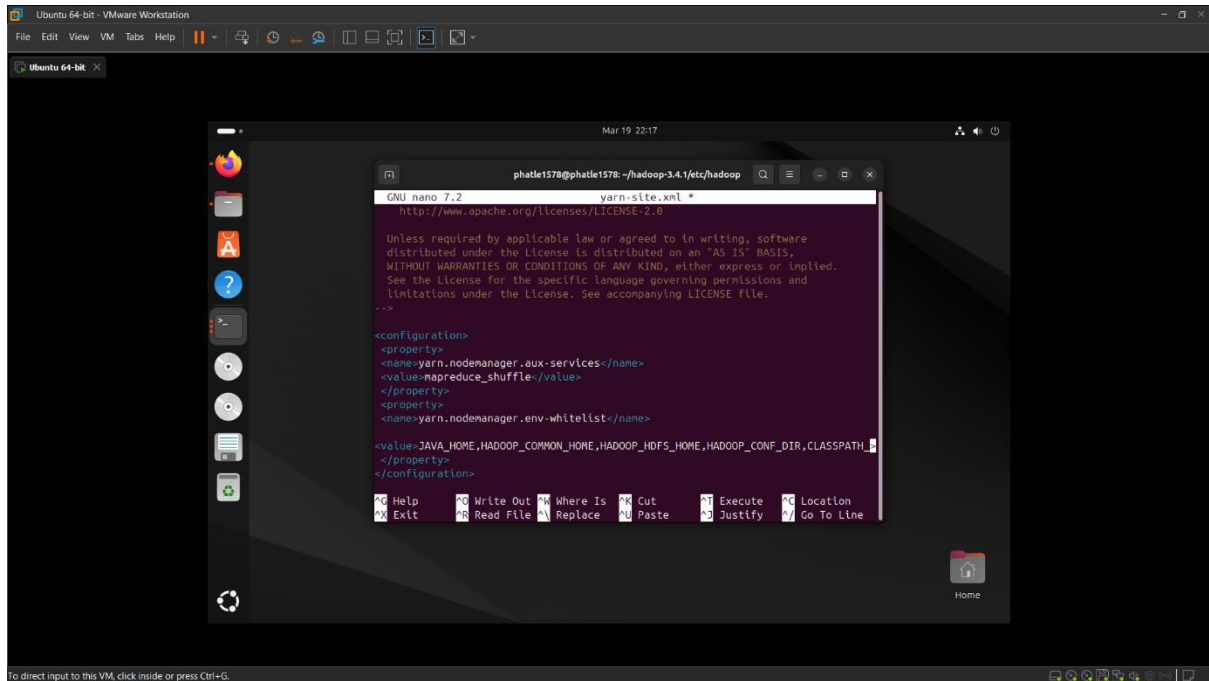
- Ctrl + O: Save.
- Ctrl + X: Exit.



*Figure 18: Configure yarn-site.xml*

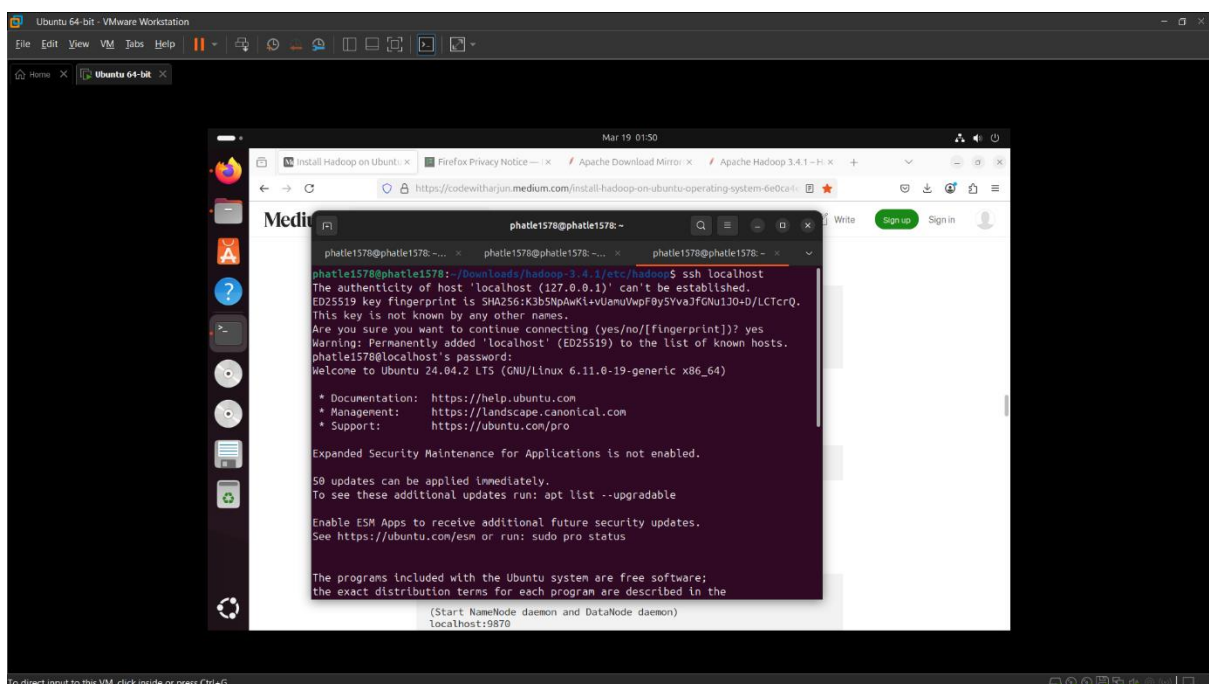### Step 4: Setup ssh

- ssh localhost: Start the ssh localhost.



*Figure 19: Start ssh localhost*

- ssh-keygen -t rsa -P '' -f ~/.ssh/id_rsa: Generate private key at ~/.ssh/id_rsa and public key at ~/.ssh/id_rsa.pub
- cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys: add public key we just create previously to the authorized _keys list.



*Figure 20: Generate ssh key.*

- chmod 0600 ~/.ssh/authorized_keys:  Set the access permission to the key.  Only the owner can read  and write.



*Figure 21: Set key permission.*

- export PDSH_RCMD_TYPE=sshFormat: delete all the current metadata of HDFS and restart the system.



*Figure 22: Restart the distributed system.*

## b. Start Hadoop Cluster.

- start-all.sh: Start all hadoop services.



*Figure 23: Start all hadoop services*

- hadoop fs -mkdir /hcmus: Create hcmus folder in hadoop



*Figure 24: Create hcmus folder in hadoop.*

- sudo adduser khtn_22120262: Create new user with name khtn_22120262



*Figure 24: Create new user with name khtn_22120262*

- hadoop fs -mkdir /hcmus/22120262: Create a subfolder 22120262 in hcmus folder.

- Download the require file from courses.fit.hcmus.edu.vn
- Right click in the folder Downloads to open in terminal



*Figure 26: Downloads require file.*

- unzip NMDLL\ -\ Lab\ 1.zip - ~\: Unzip the file we just download into the /home/phatle1578 directory..



*Figure 27: Unzip the file.*
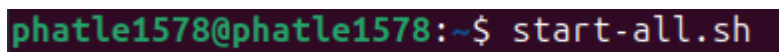
- cp ~/NMDLL\ -\ Lab\ 1/hadoop-test.jar ~/: Move the hadoop-test.jar file out to /home/phatle1578 directory for more convenient.



*Figure 28: Move the hadoop-test.jar file to home directory.*

- ls -l ~/: Verify it has been copy to the home directory

```
phatle1578@phatle1578:~$ ls -l ~/
total 62384
drwxr-xr-x  2 phatle1578 phatle1578     4096 Mar 19 01:00 Desktop
drwxr-xr-x  2 phatle1578 phatle1578     4096 Mar 19 01:00 Documents
drwxr-xr-x  2 phatle1578 phatle1578     4096 Mar 19 10:47 Downloads
drwxr-xr-x 11 phatle1578 phatle1578     4096 Mar 19 02:10 hadoop-3.4.1
-rw-rw-r--  1 phatle1578 phatle1578 63828099 Mar 19 17:16 hadoop-test.jar
drwxrwxr-x  3 phatle1578 phatle1578     4096 Mar 19 10:49 __MACOSX
drwxr-xr-x  2 phatle1578 phatle1578     4096 Mar 19 01:00 Music
drwxr-xr-x  2 phatle1578 phatle1578     4096 Mar 19 10:52 NMDLL-Lab1
drwxr-xr-x  2 phatle1578 phatle1578     4096 Mar 19 01:00 Pictures
drwxr-xr-x  2 phatle1578 phatle1578     4096 Mar 19 01:00 Public
drwx------  4 phatle1578 phatle1578     4096 Mar 19 01:15 snap
drwxr-xr-x  2 phatle1578 phatle1578     4096 Mar 19 01:00 Templates
drwxr-xr-x  2 phatle1578 phatle1578     4096 Mar 19 01:00 Videos
```
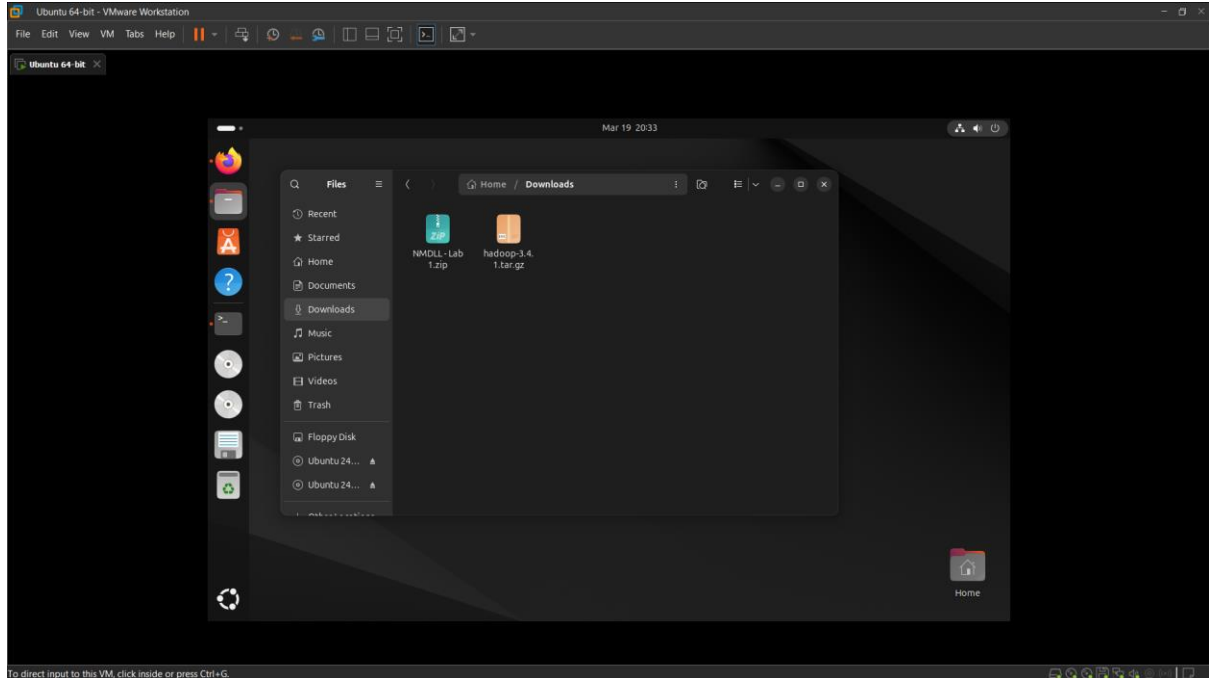
*Figure 29: Verify the file is exist in home directory.*

- hdfs dfs -put ~/hadoop-test.jar /hcmus/22120262: Move the hadoop-test.jar into the path /hcmus/22120262 in hdfs.

```
phatle1578@phatle1578:~$ hdfs dfs -put ~/hadoop-test.jar /hcmus/22120262
```

*Figure 30: Move the file hadoop-test.jar into /hcmus/22120262 in hdfs.*

- hdfs dfs -ls /hcmus/22120262: Verify it has been move to /hcmus/22120262 in hdfs

```
phatle1578@phatle1578:~$ hdfs dfs -ls /hcmus/22120262
Found 1 items
-rw-r--r--   1 phatle1578 supergroup   63828099 2025-03-19 17:17 /hcmus/22120
262/hadoop-test.jar
```

*Figure 31: Verify the file is exist in /hcmus/22120262 directory.*

- hdfs dfs -chown khtn_22120262 /hcmus/22120262: Set ownership of /hcmus/22120262/ to khtn_22120262 user.

```
phatle1578@phatle1578:~$ hdfs dfs -chown khtn_22120262 /hcmus/22120262
```

*Figure 32: Set ownership.*

- hdfs dfs -chmod 744 /hcmus/22120262/hadoop-test.jar: Set the file permissions to 744

```
phatle1578@phatle1578:~$ hdfs dfs -chmod 744 /hcmus/22120262/hadoop-test.jar
```

*Figure 33: Set the file permissions to 744.*

- java -jar ~/hadoop-test.jar 9000 /hcmus/22120262: Execute the hadoop-test.jar file



```
phatle1578@phatle1578:~$ java -jar ~/hadoop-test.jar 9000 /hcmus/22120262
Trying to read /hcmus/22120262
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.
Found hdfs://localhost:9000/hcmus/22120262/hadoop-test.jar
Your student ID: 22120262 (ensure it matches your student ID)
The first method to get MAC address is failed: Could not get network interface
Trying the alternative method
The first method to get MAC address is failed: Could not get network interface
Trying the alternative method
File written at /home/phatle1578/22120262_verification.txt
```

*Figure 35: Execute the file.*



```
MAC=00-0C-29-BB-66-1B
179f4cfb620bf7dcfa33f2f310c104889d3f04da846b5cfb97c98c5e95698100
```

*Figure 36: Verification file*

# II.    Word Count

## 1. Mappper.

The script reads input from *sys.stdin*, extracts words that contain only character in the alphabet using *re.findall(r'[a-zA-Z]+', line)*, checks if the first letter (lowercased) is in *key_set*, and if so, prints the letter and the word separated by a tab (*\t*).

```python
1   import re
2   import sys
3
4   key_set = {'a', 'f', 'j', 'g', 'h', 'c', 'm', 'u', 's'}
5
6   for line in sys.stdin:
7       line = line.strip()
8       words = re.findall(r'[a-zA-Z]+', line)
9       for word in words:
10          if word[0].lower() in key_set:
11              print(f"{word[0].lower()}\t{word}")
12
```

*Figure 37: mapper.py*

## 2. Reducer.

The script reads input from *sys.stdin*, processes key-value pairs (tab-separated), counts occurrences of each key, and stores results in output. It then sorts the output based on key_list order using a dictionary lookup (key_dict) and prints the sorted results.

```python
import sys

current_key = None
count = 0
output = []
key_list = ['a', 'f', 'j', 'g', 'h', 'c', 'm', 'u', 's']

for line in sys.stdin:
    line = line.strip()
    if not line:
        continue
    key, item = line.split('\t')

    if key != current_key:
        if current_key is not None:
            output.append(f"{current_key}\t{count}")
        current_key = key
        count = 0

    count+=1

if current_key is not None:
    output.append(f"{current_key}\t{count}")

key_dict = {char: idx for idx, char in enumerate(key_list)}
def sortKey(string):
    return key_dict.get(string[0], float('inf'))

output.sort(key=sortKey)
for ln in output:
    print(ln)
```

*Figure 38: reducer.py*

# 3. Result.

Result is the txt file with the TSV-formatted.

```
a          32921
f          18793
j          4530
g          16002
h          20911
c          42817
m          27239
u          24301
s          59567
```

*Figure 38: result.txt*