

# MS4S21 - Big Data Engineering and Applications

Moizzah Asif, J418, moizzah.asif@southwales.ac.uk

University of South Wales



# Contents

<b>1 Apache Hadoop multi-node cluster on Ubuntu VMs</b>	<b>3</b>
1.1 VM creation . . . . .	3
1.2 Ubuntu 20.04 installation . . . . .	5
1.3 Setting up Ubuntu for Hadoop-3.2.2 hdfs cluster . . . . .	7
1.4 Java for Hadoop-3.2.2 . . . . .	10
1.5 Download and configure hadoop . . . . .	11
1.5.1 Cloning VM . . . . .	13
1.5.2 Virtual Box network configuration . . . . .	14
1.5.3 Revisiting system files . . . . .	19
1.5.4 Establish Password less SSH access between hadoop nodes . . . . .	21
1.5.5 Configuring hadoop files on all hadoop nodes . . . . .	24
1.5.6 Worker/s node only configuration . . . . .	27
1.6 Starting the hadoop cluster . . . . .	28
<b>2 Mapreduce on Hadoop Cluster</b>	<b>31</b>
2.1 Configuration on hadoop cluster . . . . .	31
2.1.1 Worker/data nodes configuration . . . . .	31
2.1.2 Master/name node configuration . . . . .	31
2.2 Executing a demo mapreduce job on the cluster . . . . .	32
2.3 A word count mapreduce task with python . . . . .	35
<b>3 Twitter App creation for the module</b>	<b>37</b>

# Apache Hadoop multi-node cluster on Ubuntu VMs

This chapter will walk you through step-wise general guidelines of creating a Apache hadoop cluster 3.2.1 at present using Ubuntu latest stable release via VMs.

## 1.1 VM creation

Please note these steps can be followed to create any OS's VM in general, however the Linux Ubuntu 19.10 is used as example here.

1. Open virtual box and click on the new icon.
2. You are about to create a Ubuntu VM, set the values of each field as shown in Fig 1.1, and then press continue. Please note that chose the machine folder based on where you want to save it on your computer. If you are using university computer, then please choose your network drive.

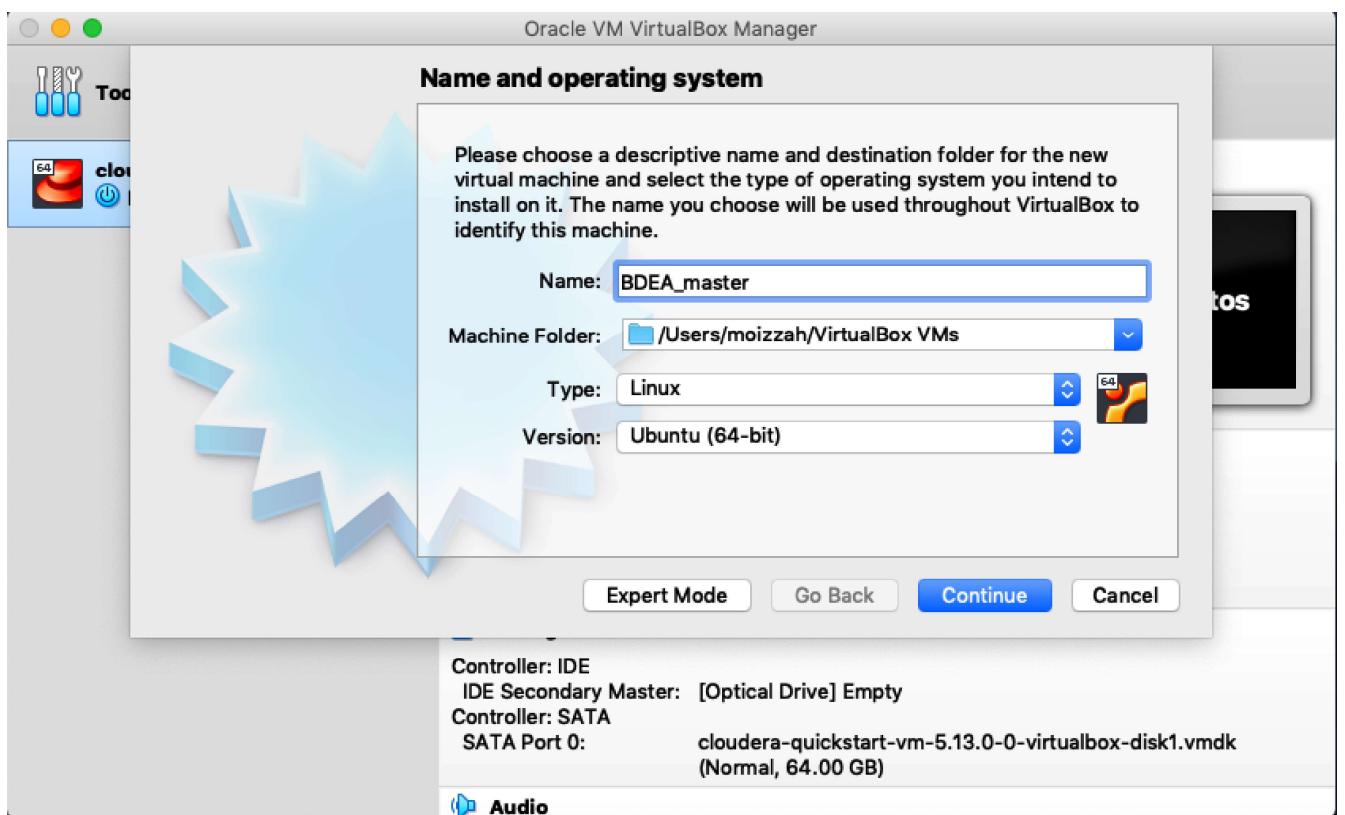


Figure 1.1: Ubuntu VM creation on virtual box - 1

3. Set RAM/memory to 12 GB/ 12288 MBs, and press continue.
4. Select the radio button which enables you to create a virtual hard disk, and then press create.
5. It will be followed by a pop up window where you can specify the hard disk file type as shown in Fig 1.2, please select VHD (virtual Hard Disk) to stay consistent with lectures, and then press continue.

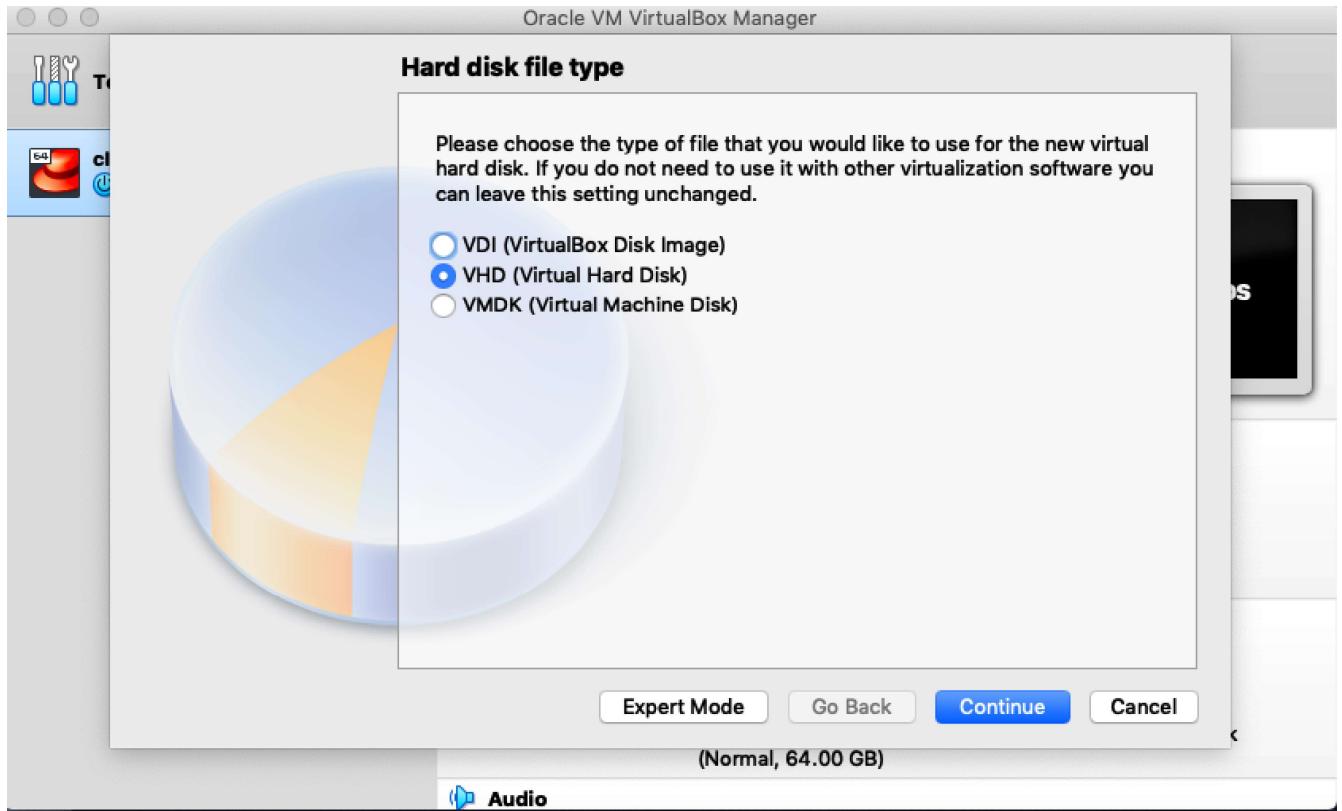


Figure 1.2: Ubuntu VM creation on virtual box - Hard disk file type

6. Select Fixed size radio button for physical storage on hard disk, press continue and then either keep or change the default value of 10 GB. This decision should be made based on your local machine's available hard drive storage, as well as the number of VMs you will have to run at once. In this case consider at least two VMs at a time to simulate the cluster. Press continue and wait for the VM to be created. Please note that you have just configured the physical and memory storage requirement until now. This is similar to bringing home a computer which doesn't have any operating system installed on it.
7. You should be able to see the VM on Virtual Box left pane now. As shown in Fig 1.3 the name of this empty VM should be what you named it in the first step of creation.

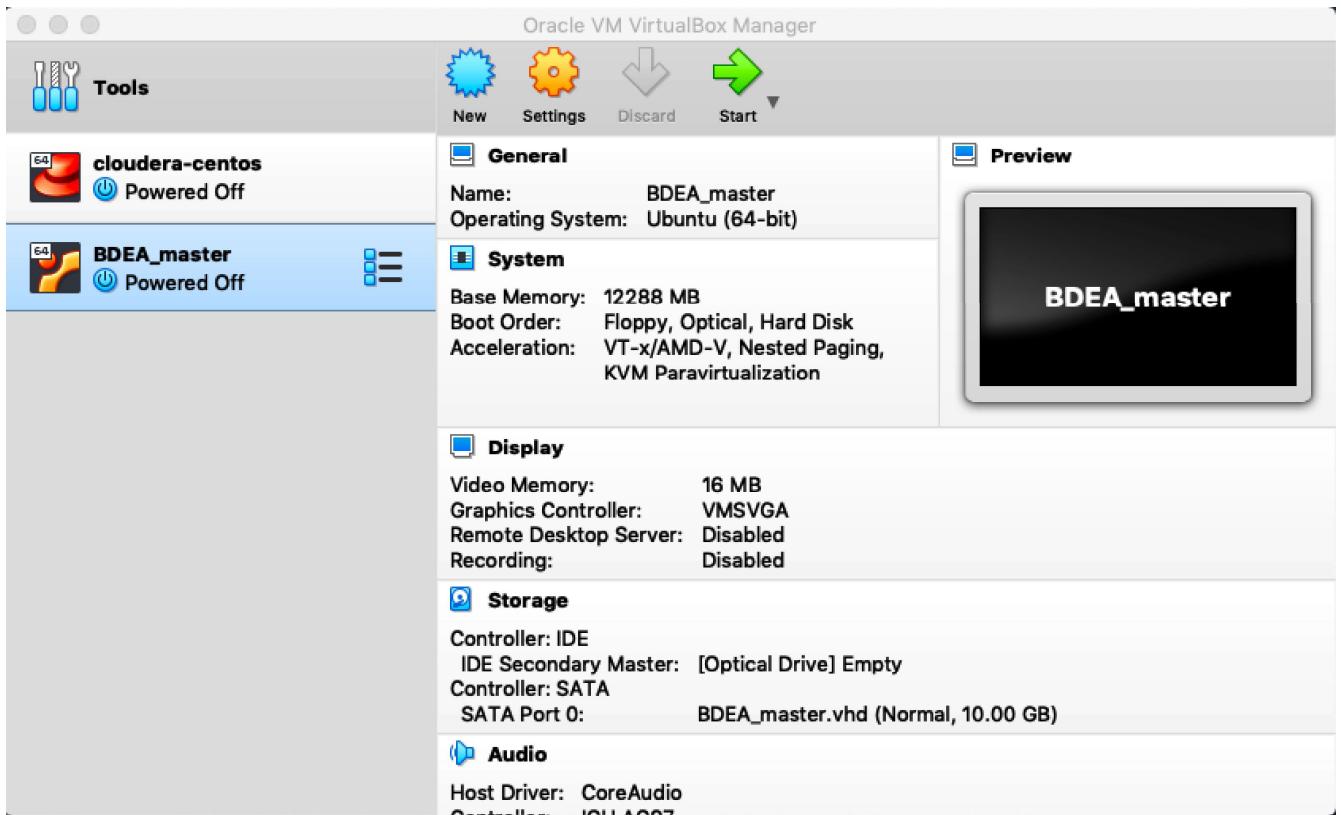


Figure 1.3: Ubuntu VM creation on virtual box - Empty VM

## 1.2 Ubuntu 20.04 installation

You will have to use the .iso image provided in BB Week 1 folder or the one you have downloaded on your computer or saved on your uni network drive.

1. Start the VM that you have just created in the previous section; click on the folder icon with a green arrow; Click on the add icon in the new pop up window; select the .iso image from the location where you have saved/downloaded it; The file should appear on the pop up windows as shown in Fig 1.4; choose this file and process with installation in the next window that pops up on your screen. The essential options in the process are listed below, please make sure that you have selected them.
  - (a) install ubuntu
  - (b) English UK
  - (c) normal installation
  - (d) erase disk an install ubuntu(it's empty already)

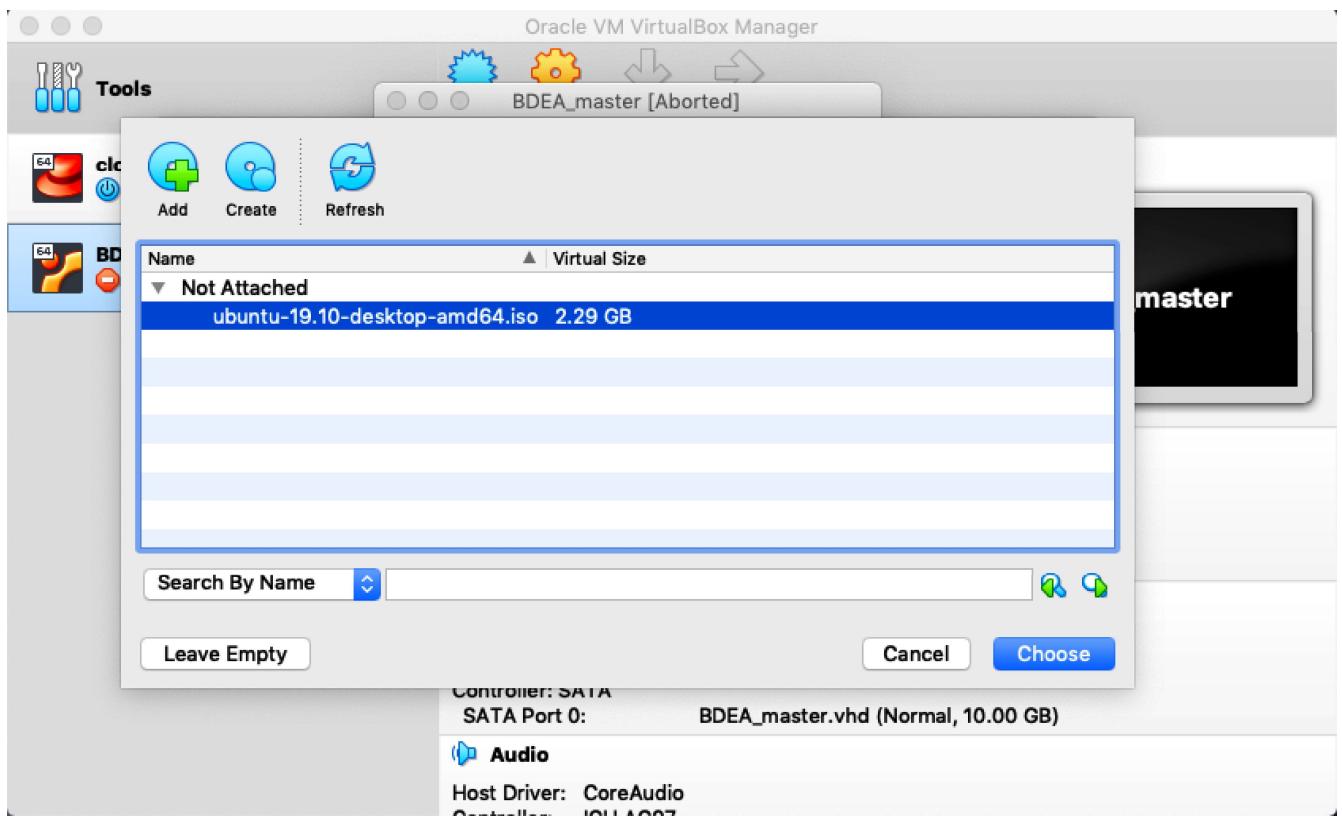


Figure 1.4: Ubuntu VM creation on virtual box - adding iso image 1

2. The screen display shown in Fig 1.5 helps you create a user profile on ubuntu. Please name the user profile intuitively. A user profile BDEA (Big Data Engineering and Applications) will be created for this tutorial's purpose.

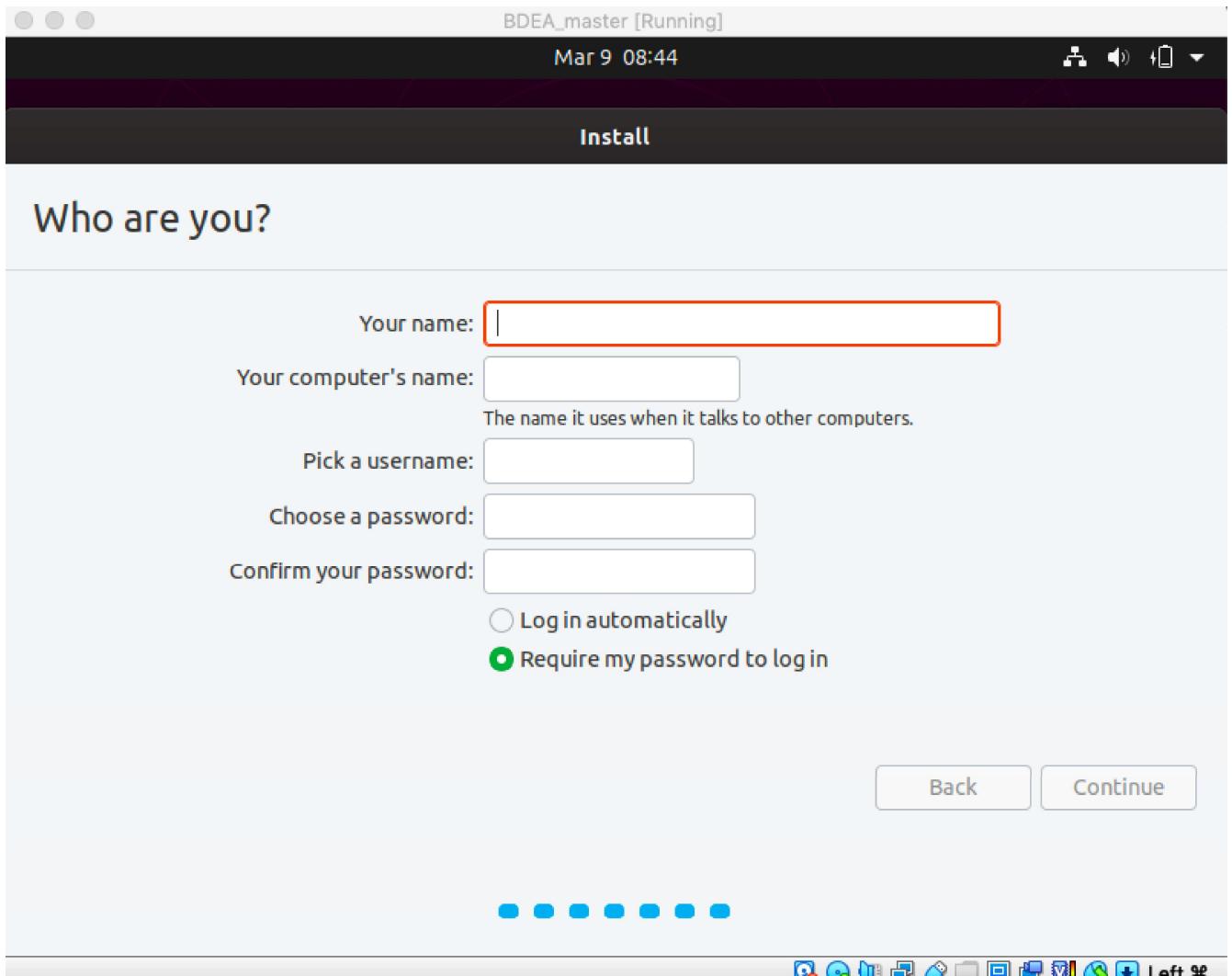


Figure 1.5: Ubuntu user creation

### 1.3 Setting up Ubuntu for Hadoop-3.2.2 hdfs cluster

Once you arrive at your Ubuntu's desktop, open Terminal so that you may set up the machine for hadoop-3.2.2 installation. Follow the guidelines listed below after opening terminal.

*For more information on linux terminal please use the forum in blackboard module.*

1. upgrade and update apt with sudo privileges (APT - advanced packaging tool) - `sudo apt update/upgrade`
2. Install openssh server - `sudo apt install openssh-server`; verify the status now: `sudo systemctl status ssh`. You should be able to get an output similar to Fig 1.6
3. check hostname - `sudo hostname`; rename to hadoop-master - `sudo hostname hadoop-master`; verify successful renaming as shown in Fig 1.6.

```

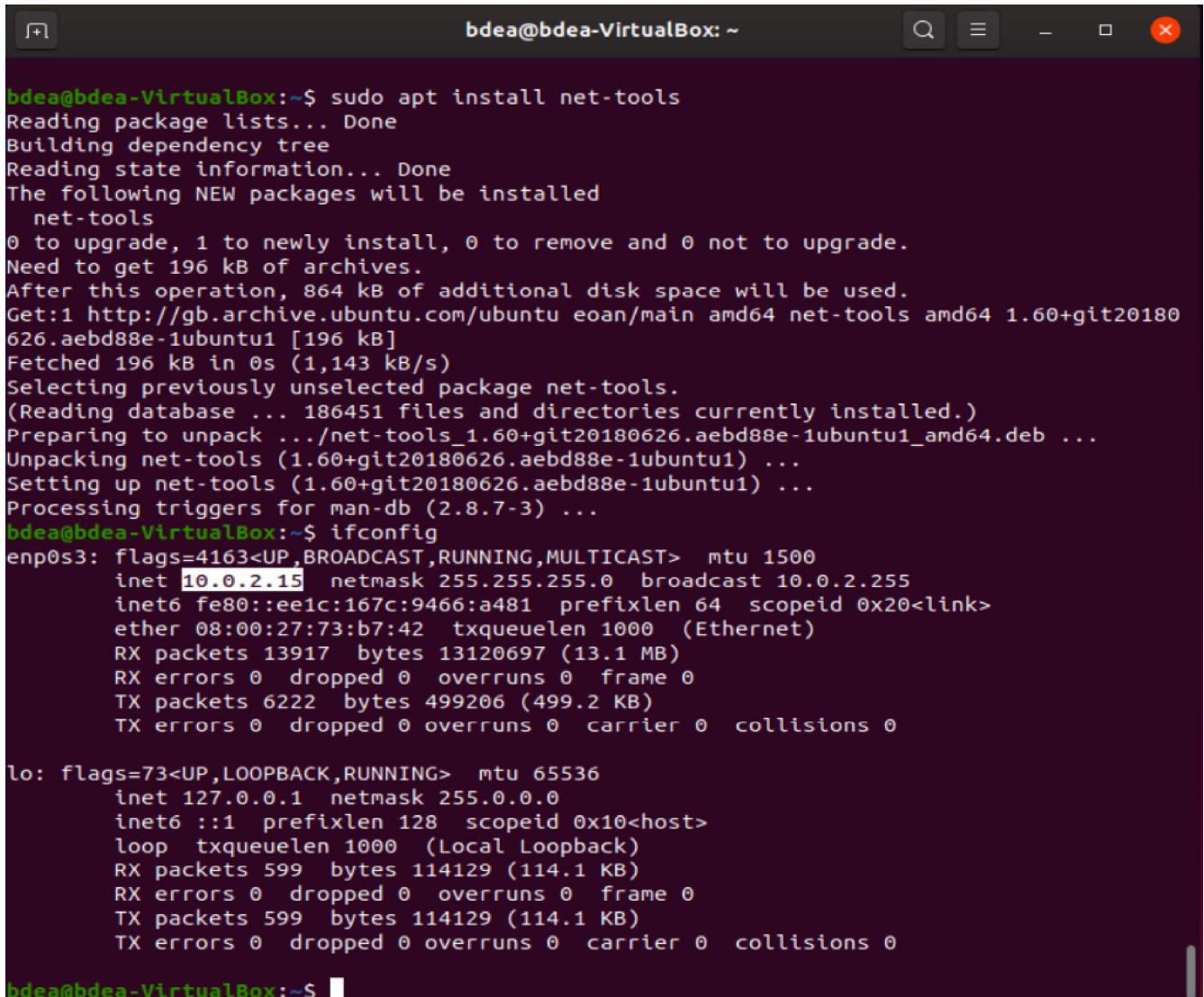
bdea@bdea-VirtualBox:~$ sudo systemctl status ssh
● ssh.service - OpenBSD Secure Shell server
  Loaded: loaded (/lib/systemd/system/ssh.service; enabled; vendor preset: enabled)
  Active: active (running) since Mon 2020-03-09 09:05:48 GMT; 1min 21s ago
    Docs: man:sshd(8)
          man:sshd_config(5)
 Main PID: 27955 (sshd)
   Tasks: 1 (limit: 4915)
  Memory: 1.1M
   CGroup: /system.slice/ssh.service
           └─27955 /usr/sbin/sshd -D

Mar 09 09:05:48 bdea-VirtualBox systemd[1]: Starting OpenBSD Secure Shell server...
Mar 09 09:05:48 bdea-VirtualBox sshd[27955]: Server listening on 0.0.0.0 port 22.
Mar 09 09:05:48 bdea-VirtualBox sshd[27955]: Server listening on :: port 22.
Mar 09 09:05:48 bdea-VirtualBox systemd[1]: Started OpenBSD Secure Shell server.
bdea@bdea-VirtualBox:~$ sudo hostname
bdea-VirtualBox
bdea@bdea-VirtualBox:~$ sudo hostname hadoop-master
bdea@bdea-VirtualBox:~$ sudo hostname
hadoop-master
bdea@bdea-VirtualBox:~$ █

```

Figure 1.6: openssh-server and hostname verification

4. install gedit text editor for smooth editing - *sudo apt install gedit*. If you are working on remote hosts, this editor may not be available and you may have to work with vim. (A healthy discussion on vim had been conducted during lectures. for more info, please create a thread in the forum, or use one if it is already there.)
5. update the hosts file by adding hadoop-master as a host.
  - (a) Install ifconfig to find the ip address of your machine via terminal - *sudo apt install net-tools*. Fig 1.7 shows the output and script of installing net-tools and listing the internet configurations of your machine.
  - (b) Find the ip address of your vm - *ifconfig*. Fig 1.7 shows the ip address highlighted manually.



```

bdea@bdea-VirtualBox:~$ sudo apt install net-tools
Reading package lists... Done
Building dependency tree
Reading state information... Done
The following NEW packages will be installed
  net-tools
0 to upgrade, 1 to newly install, 0 to remove and 0 not to upgrade.
Need to get 196 kB of archives.
After this operation, 864 kB of additional disk space will be used.
Get:1 http://gb.archive.ubuntu.com/ubuntu eoan/main amd64 net-tools amd64 1.60+git20180626.aebd88e-1ubuntu1 [196 kB]
Fetched 196 kB in 0s (1,143 kB/s)
Selecting previously unselected package net-tools.
(Reading database ... 186451 files and directories currently installed.)
Preparing to unpack .../net-tools_1.60+git20180626.aebd88e-1ubuntu1_amd64.deb ...
Unpacking net-tools (1.60+git20180626.aebd88e-1ubuntu1) ...
Setting up net-tools (1.60+git20180626.aebd88e-1ubuntu1) ...
Processing triggers for man-db (2.8.7-3) ...
bdea@bdea-VirtualBox:~$ ifconfig
enp0s3: flags=4163<UP,BROADCAST,RUNNING,MULTICAST>  mtu 1500
      inet 10.0.2.15  netmask 255.255.255.0  broadcast 10.0.2.255
        inet6 fe80::ee1c:167c:9466:a481  prefixlen 64  scopeid 0x20<link>
          ether 08:00:27:73:b7:42  txqueuelen 1000  (Ethernet)
            RX packets 13917  bytes 13120697 (13.1 MB)
            RX errors 0  dropped 0  overruns 0  frame 0
            TX packets 6222  bytes 499206 (499.2 KB)
            TX errors 0  dropped 0  overruns 0  carrier 0  collisions 0

lo: flags=73<UP,LOOPBACK,RUNNING>  mtu 65536
      inet 127.0.0.1  netmask 255.0.0.0
        inet6 ::1  prefixlen 128  scopeid 0x10<host>
          loop  txqueuelen 1000  (Local Loopback)
            RX packets 599  bytes 114129 (114.1 KB)
            RX errors 0  dropped 0  overruns 0  frame 0
            TX packets 599  bytes 114129 (114.1 KB)
            TX errors 0  dropped 0  overruns 0  carrier 0  collisions 0
bdea@bdea-VirtualBox:~$ 

```

Figure 1.7: net-tools installation and ifconfig output

- (c) open the hosts file for editing, to add hadoop-master as a host and map it to its ip address  
*- sudo gedit hosts.* The hosts file configuration shown in Fig 1.8 shows the correct format of adding a new hosts; in this example the new host is hadoop-master. Please note the highlighted line refers to a host which does not exist any more, please remove this line to avoid any carried forward errors. You may as well comment this line to keep as a reminder for future references.

```

Open  +  *hosts
Save  -  ×
127.0.0.1      localhost
127.0.1.1      bdea-VirtualBox
10.0.2.15      hadoop-master

# The following lines are desirable for IPv6 capable hosts
::1      ip6-localhost ip6-loopback
fe00::0 ip6-localnet
ff00::0 ip6-mcastprefix
ff02::1 ip6-allnodes
ff02::2 ip6-allrouters

```

Figure 1.8: interim configuration of hosts file

6. Download, install and configure the latest stable Hadoop-3.2.2 compatible version of Java. Please follow the instructions in next section.

## 1.4 Java for Hadoop-3.2.2

Hadoop-3.2.2 is not compatible with Java 11 for development purposes. A hdfs cluster could be successfully simulated with Java 11, however, some of the big data services on hadoop cluster can not work with Java 11, such as: mapreduce framework on hadoop.

Therefore, we will work with Java 8, as it the latest stable and compatible version of Java with hadoop-3.2.2. You can follow the next set of instructions for Java 11 as well.

1. Install Java 8 JDK - *sudo apt install openjdk-8-jdk*

The Java Development Kit (JDK) is a software development environment used for developing Java applications and applets. It includes the Java Runtime Environment (JRE), an interpreter/loader (Java), a compiler (javac), an archiver (jar), a documentation generator (Javadoc) and other tools needed in Java development

2. Configure bash file to set user variable JAVA\_HOME. Setting this variable will allow the OS to locate Java directory whenever Java is called by any process/application. Either enter the absolute to path to bash file (it is located in your home directory) or open it from the home directory as follows - *sudo gedit .bashrc*

The value of JAVA\_HOME should be set as - */usr/lib/jvm/java-8-openjdk-amd64/jre*. Please refer to Fig for exact syntax of variable initialisation in bash file.

```

.bashrc
/home/bdea
alias egrep='egrep --color=auto'
fi

# colored GCC warnings and errors
#export GCC_COLORS='error=01;31:warning=01;35:note=01;36:caret=01;32:locus=01:quote=01'

# some more ls aliases
alias ll='ls -alF'
alias la='ls -A'
alias l='ls -CF'

# Add an "alert" alias for long running commands.  Use like so:
# sleep 10; alert
alias alert='notify-send --urgency=low -i "$( [ $? = 0 ] && echo terminal || echo error)" "$
(history|tail -n1|sed -e '\''s/^\\s*[0-9]\\+\\s*/;s/[;&]\\s*alert$//'\''")'

# Alias definitions.
# You may want to put all your additions into a separate file like
# ~/.bash_aliases, instead of adding them here directly.
# See /usr/share/doc/bash-doc/examples in the bash-doc package.

if [ -f ~/.bash_aliases ]; then
  . ~/.bash_aliases
fi

# enable programmable completion features (you don't need to enable
# this, if it's already enabled in /etc/bash.bashrc and /etc/profile
# sources /etc/bash.bashrc).
if ! shopt -oq posix; then
  if [ -f /usr/share/bash-completion/bash_completion ]; then
    . /usr/share/bash-completion/bash_completion
  elif [ -f /etc/bash_completion ]; then
    . /etc/bash_completion
  fi
fi

export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64/jre

```

sh ▾ Tab Width: 8 ▾ Ln 1, Col 1 ▾ INS

Figure 1.9: Initialising JAVA\_HOME in bash file

## 1.5 Download and configure hadoop

Unlike Java, hadoop is not installed. Hadoop's directories are downloaded and configuration files are modified to ensure smooth execution and run of cluster.

1. Download hadoop-3.2.2 from apache hadoop website -  
`wgethttps://archive.apache.org/dist/hadoop/common/hadoop-3.2.2/hadoop-3.2.2.tar.gz`
2. unzip and delete the tar file was follows:
  - `tar -xvf hadoop-3.2.2.tar.gz`
  - `rm hadoop-3.2.2.tar.gz`
3. move hadoop to the usr local directory while renaming it to hadoop - `sudo mv hadoop-3.2.2 /usr/local/hadoop`
4. Add hadoop bin and sbin directories to the PATH variable using bash file. Please follow the syntax in Fig 1.10 to concatentae the paths to PATH variable.

```

*.bashrc
/home/bdea
Save
Open ▾
+ ×

alias grep='grep --color=auto'
alias fgrep='fgrep --color=auto'
alias egrep='egrep --color=auto'
fi

# colored GCC warnings and errors
#export
GCC_COLORS='error=01;31:warning=01;35:note=01;36:caret=01;32:locus=01:quote=01;33:help=01;36'

# some more ls aliases
alias ll='ls -alF'
alias la='ls -A'
alias l='ls -CF'

# Add an "alert" alias for long running commands.  Use like so:
#   sleep 10; alert
alias alert='notify-send --urgency=low -i "$( [ $? = 0 ] && echo terminal
|| echo error)" "$(history|tail -n1|sed -e '\''s/^\\s*[0-9]\\+\s*//;s/[;&]
\\s*alert$/'\''")"

# Alias definitions.
# You may want to put all your additions into a separate file like
# ~/.bash_aliases, instead of adding them here directly.
# See /usr/share/doc/bash-doc/examples in the bash-doc package.

if [ -f ~/.bash_aliases ]; then
    . ~/.bash_aliases
fi

# enable programmable completion features (you don't need to enable
# this, if it's already enabled in /etc/bash.bashrc and /etc/profile
# sources /etc/bash.bashrc).
if ! shopt -oq posix; then
    if [ -f /usr/share/bash-completion/bash_completion ]; then
        . /usr/share/bash-completion/bash_completion
    elif [ -f /etc/bash_completion ]; then
        . /etc/bash_completion
    fi
fi

export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64/jre
export PATH=$PATH:/usr/local/hadoop/bin:/usr/local/hadoop/sbin
export CONF=/usr/local/hadoop/etc/hadoop

```

Figure 1.10: Bash file with Hadoop bin and sbin added to path

5. create another environment variable using the bash file, so that accessing the hadoop config directory becomes easier. Note the creation of CONF variable in Fig 1.10

Don't forget to source the bash file whenever you update it.

Before we start editing the hadoop configuration files, it's is essential that we clone the master VM to create a worker VM. Cloning will help save time and effort at this stage.