# Practical No. 2

**AIM**: Install Hadoop, Hadoop Word Count and Mutligram.

| Category | Requirements, Conventions or Software Version Used |
|---|---|
| System | Ubuntu 18.04 |
| Software | Hadoop 2.8.5, Oracle JDK 1.8 |
| Other | Privileged access to your Linux system as root or via the sudo command. |
| Conventions | **#** - requires given linux commands to be executed with root privileges either directly as a root user or by use of sudo command<br>**$** - requires given linux commands to be executed as a regular non-privileged user |

```
zaidubuntu@zaidubuntu-Vostro-15-3568: ~
File Edit View Search Terminal Help
zaidubuntu@zaidubuntu-Vostro-15-3568:~$ adduser hadoop
adduser: Only root may add a user or group to the system.
zaidubuntu@zaidubuntu-Vostro-15-3568:~$ sudo adduser hadoop
[sudo] password for zaidubuntu:
Adding user `hadoop' ...
Adding new group `hadoop' (1002) ...
Adding new user `hadoop' (1001) with group `hadoop' ...
Creating home directory `/home/hadoop' ...
Copying files from `/etc/skel' ...
Enter new UNIX password:
Retype new UNIX password:
passwd: password updated successfully
Changing the user information for hadoop
Enter the new value, or press ENTER for the default
        Full Name []: Hadoop User
        Room Number []:
        Work Phone []:
        Home Phone []:
        Other []:
Is the information correct? [Y/n] Y
zaidubuntu@zaidubuntu-Vostro-15-3568:~$ 
```

# Install and configure the Oracle JDK

Make sure the installed Java is that of Oracle's and not OpenJDK's.

```
zaidubuntu@zaidubuntu-Vostro-15-3568: ~
File Edit View Search Terminal Help
zaidubuntu@zaidubuntu-Vostro-15-3568:~$ java -version
java version "1.8.0_202"
Java(TM) SE Runtime Environment (build 1.8.0_202-b08)
Java HotSpot(TM) 64-Bit Server VM (build 25.202-b08, mixed mode)
zaidubuntu@zaidubuntu-Vostro-15-3568:~$ javac -version
javac 1.8.0_202
zaidubuntu@zaidubuntu-Vostro-15-3568:~$ 
```
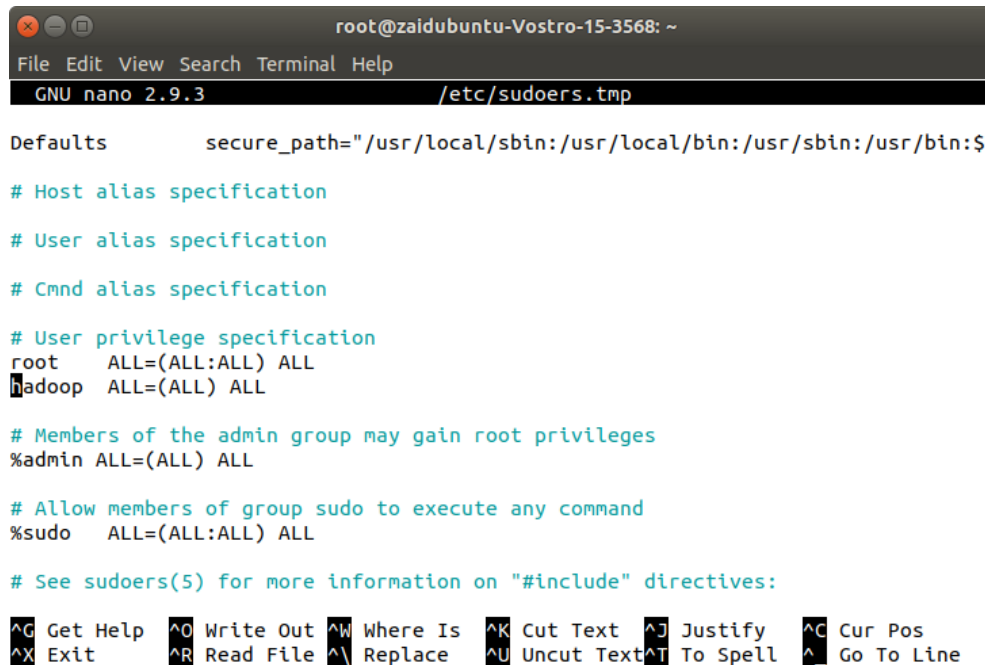
Switch to root user:



And then enter '**visudo**':
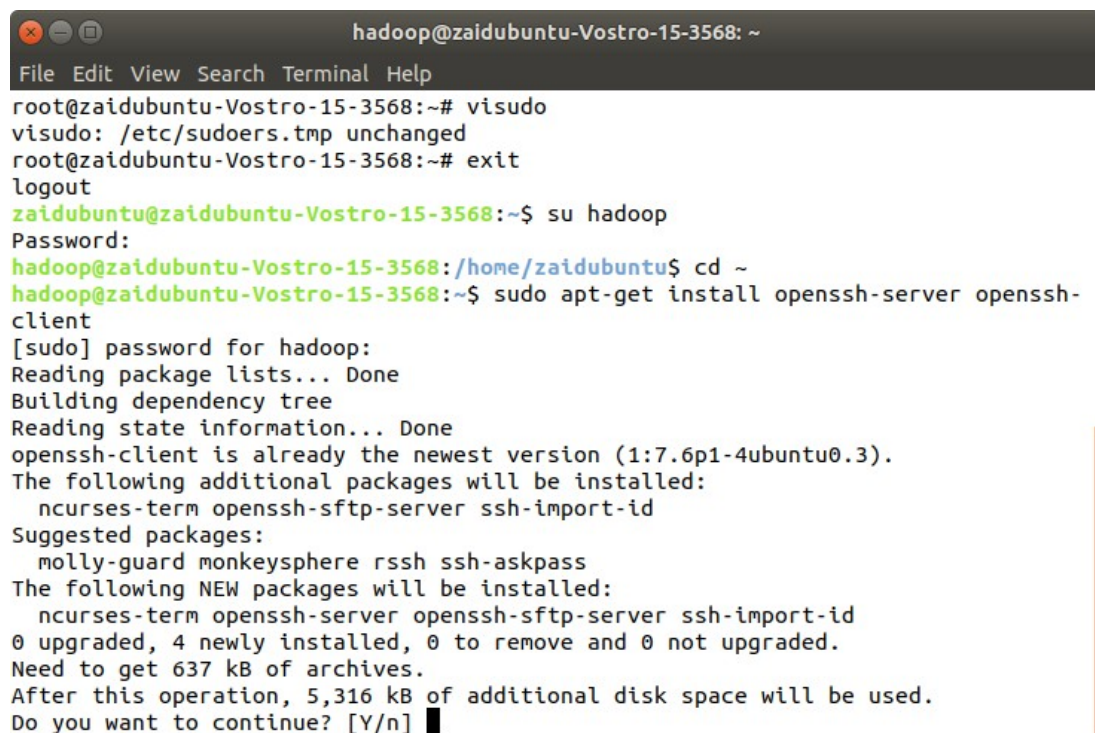
**BEFORE**:



**AFTER ADDING:**hadoop    ALL=(ALL) ALL

Save by pressing CTRL + X and then Save Yes:

Switch to hadoop user and enter:

**This part is required for Accessing and Running MapReduce programs remotely using SSH like in the Lab.**

```
sudo apt-get install openssh-server openssh-client
```
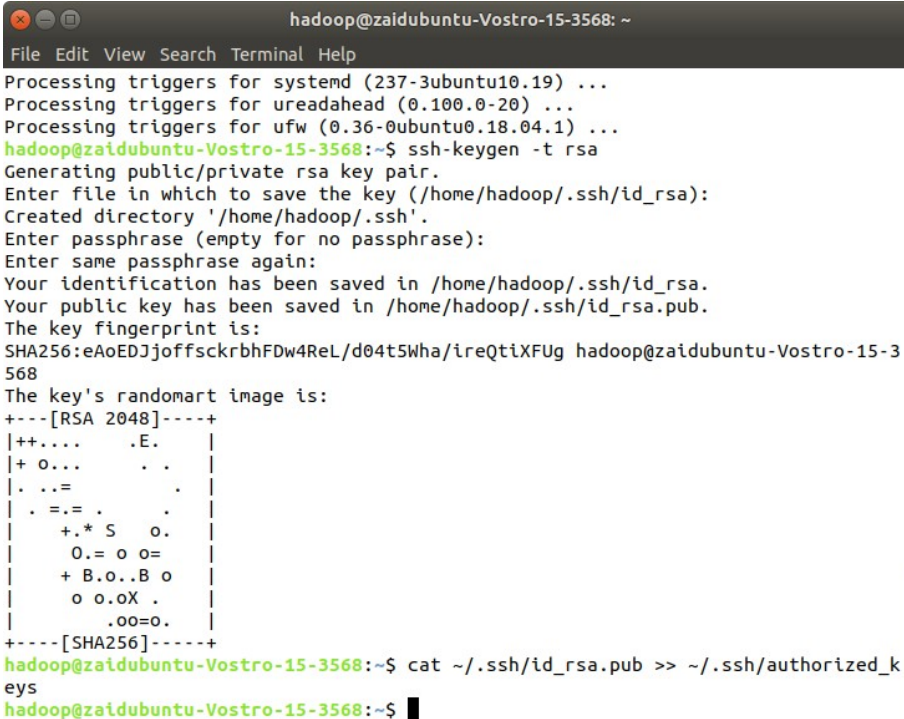
Generate Public and Private Key Pairs with the following command. The terminal will prompt for entering the file name. Press ENTER and proceed. After that copy the public keys form `id_rsa.pub` to `authorized_keys`.
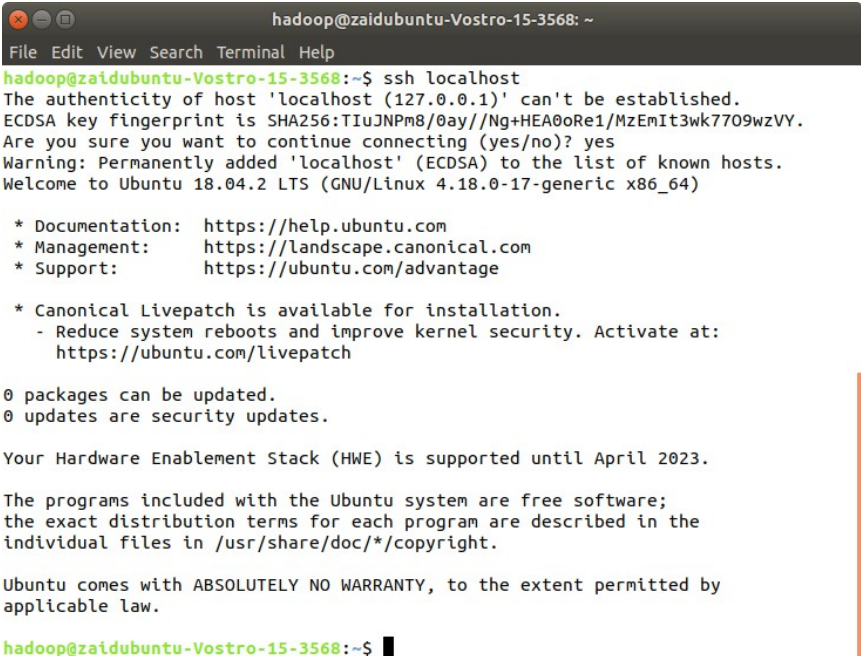
```
$ ssh-keygen -t rsa
$ cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
```



Verify the password-less ssh configuration with the command :

```
$ ssh localhost
```

# Install Hadoop and configure related xml files

Download and extract [Hadoop 2.8.5](#) from Apache official website.

```
# tar -xzvf hadoop-2.8.5.tar.gz
```



Edit the `bashrc` for the Hadoop user via setting up the following Hadoop environment variables :

```
export HADOOP_HOME=/home/hadoop/hadoop-2.8.5
export HADOOP_INSTALL=$HADOOP_HOME
export HADOOP_MAPRED_HOME=$HADOOP_HOME
export HADOOP_COMMON_HOME=$HADOOP_HOME
export HADOOP_HDFS_HOME=$HADOOP_HOME
export YARN_HOME=$HADOOP_HOME
export HADOOP_COMMON_LIB_NATIVE_DIR=$HADOOP_HOME/lib/native
export PATH=$PATH:$HADOOP_HOME/sbin:$HADOOP_HOME/bin
export HADOOP_OPTS="-Djava.library.path=$HADOOP_HOME/lib/native"
```





Edit the `hadoop-env.sh` file which is in `/etc/hadoop` inside the Hadoop installation directory and make the following changes and check if you want to change any other configurations.

```
export JAVA_HOME=/home/zaidubuntu/jdk1.8.0_202
```

```
export HADOOP_CONF_DIR=${HADOOP_CONF_DIR:-"/home/hadoop/hadoop-2.8.5/etc/
hadoop"}
```

NOTE: Delete the old `HADOOP_CONF_DIR` and `JAVA_HOME` line



## Configuration Changes in core-site.xml file

Edit the `core-site.xml` with vim or you can use any of the editors. The file is under `/etc/hadoop` inside `hadoop` home directory and add following entries.

```
<configuration>
<property>
<name>fs.defaultFS</name>
<value>hdfs://localhost:9000</value>
</property>
<property>
<name>hadoop.tmp.dir</name>
<value>/home/hadoop/hadooptmpdata</value>
</property>
</configuration>
```

In addition, create the directory under hadoop home folder.

hadoop@zaidubuntu-Vostro-15-3568:~/hadoop-2.8.5/etc/hadoop$ mkdir
/home/hadoop/hadooptmpdata


## Configuration Changes in `hdfs-site.xml` file

Edit the `hdfs-site.xml` which is present under the same location i.e `/etc/hadoop` inside `hadoop` installation directory and create the `Namenode/Datanode` directories under `hadoop` user home directory.

hadoop@zaidubuntu-Vostro-15-3568:~/hadoop-2.8.5/etc/hadoop$ mkdir -p
/home/hadoop/hdfs/namenode

hadoop@zaidubuntu-Vostro-15-3568:~/hadoop-2.8.5/etc/hadoop$ mkdir -p
/home/hadoop/hdfs/datanode

---

```
<configuration>

<property>
<name>dfs.replication</name>
<value>1</value>
<name>dfs.name.dir</name>
<value>file:///home/hadoop/hdfs/namenode</value>
<name>dfs.data.dir</name>
<value>file:///home/hadoop/hdfs/datanode</value>
</property>

</configuration>
```

## Configuration Changes in `mapred-site.xml` file

Copy the `mapred-site.xml` from `mapred-site.xml.template` using `cp` command and then edit the `mapred-site.xml` placed in `/etc/hadoop` under `hadoop` instillation directory with the following changes.

```
$ cp mapred-site.xml.template mapred-site.xml
```

```
<configuration>

<property>
<name>mapreduce.framework.name</name>
<value>yarn</value>
</property>

</configuration>
```

### Configuration Changes in yarn-site.xml file

Edit `yarn-site.xml` with the following entries.

```
<configuration>

<!-- Site specific YARN configuration properties -->
<property>
    <name>yarn.nodemanager.aux-services</name>
    <value>mapreduce_shuffle</value>
</property>

<property>
   <name>yarn.nodemanager.aux-services.mapreduce_shuffle.class</name>
   <value>org.apache.hadoop.mapred.ShuffleHandler</value>
</property>


</configuration>
```



# Starting the Hadoop Cluster

Format the namenode before using it for the first time. As HDFS user run the below command to format the Namenode.

```
$ hdfs namenode -format
```

Once the Namenode has been formatted then start the HDFS using the start-dfs.sh script.

To start the YARN services you need to execute the yarn start script i.e. start-yarn.sh

To verify all the Hadoop services/daemons are started successfully you can use the jps command.



Now we can check the current Hadoop version you can use below command :

$ hadoop version

or

```
$ hdfs version
```

# HDFS Command Line Interface

To access the HDFS and create some directories top of DFS you can use HDFS CLI.

```
$ hdfs dfs -mkdir /test
```

```
$ hdfs dfs -ls /
```

# Overview of Hadoop Cluster :

# INSTALLATION OF HADOOP ON MULTIPLE MACHINES

One (1) Name Node : 192.168.1.1 (hadoopmaster)

Three (3) Data Nodes : 192.168.1.2 (hadoopslave1),  192.168.1.3 (hadoopslave2)  192.168.4 (hadoop slave3)

After installation of SINGLE NODE HADOOP CLUSTER. You are going to CLONE that ubuntu image and named it hadoopmaster.

Open terminal and run "ifconfig" command to see the IPv4 Address.

If it is IPv6 then you have to disable the IPv6 address. Here is a link of tutorial of how to disable IPv6 address - http://www.michael-noll.com/tutorials/running-hadoop-on-ubuntu-linux-single-node-cluster/#disabling-ipv6.

Now change the host file -->> "$ sudo gedit /etc/hosts".

Add the following lines :

hadoopmaster  192.168.1.1

hadoopslave1  192.168.1.2

hadoopslave2  192.168.1.3

hadoopslave3  192.168.1.4

Change the hostname --->> "$ sudo gedit /etc/hostname'.

Add : hadoopmaster.

Go to hadoop director and do changes in its files :     "$ cd /usr/local/hadoop/etc/hadoop".

Edit core site xml file ------>> "$ sudo gedit core-site.xml".

replace localhost as hadoopmaster.

Edit hdfs site xml file ------>> $ sudo gedit hdfs-site.xml.

replace value 1 as 3 (represents no of datanode).

Edit yarn site xml file ------>> $ sudo gedit yarn-site.xml.

Add these properties files inside Configuration tag :

```
<property>

 <name>yarn.resourcemanager.resource-tracker.address</name>

 <value>hadoopmaster:8025</value>
```

&lt;property&gt;

&lt;property&gt;

&lt;name&gt;yarn.resourcemanager.scheduler.address&lt;/name&gt;

&lt;value&gt;hadoopmaster:8030&lt;/value&gt;

&lt;property&gt;

&lt;property&gt;

&lt;name&gt;yarn.resourcemanager.address&lt;/name&gt;

&lt;value&gt;hadoopmaster:8050&lt;/value&gt;

&lt;/property&gt;

Edit couple of things in yarn site xml file -----------------&gt;&gt;     $ sudo gedit yarn-site.xml.

replace mapreduce.framework.name as mapred.job.tracker

replace yarn as hadoopmaster:54311

SHUT DOWN the hadoopmaster Ubuntu Image.

--------------------------------------------------------------------------------------------------------------------------------------------------------------------

## master node setup complete

--------------------------------------------------------------------------------------------------------------------------------------------------------------------

Clone hadoopmaster Node as hadoopslave1, hadoopslave2, hadoopslave3.

Change hadoop master host : $ sudo gedit /usr/local/hadoop/etc/hadoop/master.

replace localhost to hadoopmaster.

Change hadoop slaves : $ sudo gedit /usr/local/hadoop/etc/hadoop/slave.

replace localhost to hadoopslave1 \n hadoopslave2    \n   hadoopslave3.

Change hdfs site xml file --------&gt;&gt;     $ sudo gedit /usr/local/hadoop/etc/hadoop/hdfs-site.xml.

remove dfs.datanode.data.dir property section.

--------------------------------------------------------------------------------------------------------------------------------------------------------------------

Initial Network setup -

In a virtual machine IDE.

select hadoopmaster ubuntu image and go to its settings.

Go to network

choose attached to option as "internal network".

Give name : "hadoop multinode network".

Go to its advanced settings.

"Allow all" --- promiscuous mode.

Do this for all the slave machines as well.

-----------------------------------------------------------------------------------------------------------------------------------------------------------------

Open all 3 slave nodes and run ------->>  $ sudo gedit /etc/hostname.

replace hadoopmaster to hadoopslave1, hadoopslav2, hadoopslave3 respectively to all the three slave virtual machines.

Reboot all the slave nodes/machines.

On hadoopmaster node run below command to remove all the hadoop data:

"Remove hadoop data ------>>  $ sudo rm -rf /usr/local/hadoop/hadoop_data."

On hadoopmaster node === >>    $ sudo mkdir -p /usr/local/hadoop/hadoop_data/hdfs/namenode

 Run this command ------->>  $ sudo chown -R username:username /usr/local/hadoop

On all hadoopslave nodes ===>>  Run following commands ===>>

$ sudo rm -rf /usr/local/hadoop/hadoop_data

$ sudo mkdir -p /usr/local/hadoop/hadoop_data/hdfs/datanode.

$ sudo chown -R username:username /usr/local/hadoop.

Change hdfs site xml file for all slave nodes ----->>$ sudo gedit /usr/local/hadoop/etc/hadoop/hdfs-site.xml.

remove dfs.namenode.data.dir property section.

On hadoopmaster node,

Run the command ---->>$ sudo ssh-copy-id -i ~/.ssh/id_dsa.pub username@hadoopmaster.

If you get error then solution of your problem is :

OpenSSH is not installed. For installation : sudo apt-get install openssh-client.

OR you will get this error "permission denied for root@localhost for ssh connection" .

Solution of 2nd problem is : http://askubuntu.com/questions/497895/permission-denied-for-rootlocalhost-for-ssh-connection.

Next problem might be your Internal network is not setup.


-------------------------Internal Network Setup between all the 4 virtual machine-------------------------------------------------


Click on top right WIFI or Internet icon.

Go to edit connection.

click on add/edit connection

Give connection name : "master connection"

Go to IPv4 settings.

Change method from automatic to manual.

Enter IP Address like for master node : 192.168.1.1

Enter net mask address : 255.255.255.0

save it and do the above steps for all the nodes.

On hadoopmaster machine -

Run following commands -

$ sudo ssh-copy-id -i ~/.ssh/id_dsa.pub chaalpritam@hadoopmaster

$ sudo ssh-copy-id -i ~/.ssh/id_dsa.pub chaalpritam@hadoopslave1

$ sudo ssh-copy-id -i ~/.ssh/id_dsa.pub chaalpritam@hadoopslave2

$ sudo ssh-copy-id -i ~/.ssh/id_dsa.pub chaalpritam@hadoopslave3

Now we can access the machines using SSH -

$ sudo ssh hadoopmaster

$ exit

$ sudo ssh hadoopslave1

$ exit

$ sudo ssh hadoopslave2

$ exit

$ sudo ssh hadoopslave3

$ exit

If you able to do this you are accessing all the nodes using SSH.

Next, Format namenode and start hadoop -

$ hadoop namenode -format

$ start-all.sh

$ jps (check in all 3 datanodes)

http://hadoopmaster:8088/

http://hadoopmaster:50070/

http://hadoopmaster:50090/

http://hadoopmaster:50075/

## Conclusion:

We started installing Hadoop on a pseudo node/single node on our laptops. Then we tried installing Hadoop on distributed machine(s). We learned that it is very important checking the versions of Java and that OpenJDK should not be used.

There are also the problem that Output directory should not exist before running and that even a program fails from executing it creates the output directory before terminating and throwing an exception.

To be on the safe side don't use folders having spaces in them and don't use names too similar to existing keywords and program names.

The Word count program is the classic Hello World program and does mostly executes without giving any errors.

For Multigrams, different strategies could be used. One being (Generalized) where we could use StringTokenizer to append string in runtime in an optimized manner.