# PRACTICAL-1

## AIM:

To install Hadoop framework, configure it and setup a single node cluster. Use web based tools to monitor your Hadoop setup

## IMPLEMENTATION:

The Hadoop framework is written in Java, and its services require a compatible Java Runtime Environment (JRE) and Java Development Kit (JDK).



Apache Hadoop 3.x fully supports Java 8. The OpenJDK 8 package in Ubuntu contains both the runtime environment and development kit.



The OpenJDK or Oracle Java version can affect how elements of a Hadoop ecosystem interact. Hence, we need to be specific.



Install the OpenSSH server and client.

Utilize the adduser command to create a new Hadoop user. The username, in this example, is **hdoop**. You are free the use any username and password you see fit. Switch to the newly created user and enter the corresponding password.



The user now needs to be able to SSH to the localhost without being prompted for a password.



Use the cat command to store the public key as **authorized_keys** in the *ssh* directory

Set the permissions for your user with the chmod command

```
hdoop@parth642001-virtual-machine:~$ cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized
_keys
hdoop@parth642001-virtual-machine:~$ chmod 0600 ~/.ssh/authorized_keys
hdoop@parth642001-virtual-machine:~$ ssh localhost
The authenticity of host 'localhost (127.0.0.1)' can't be established.
ED25519 key fingerprint is SHA256:KwSx/Dngnw2Cxa7TXe4jhuL8olWAbZTZiOHBJLxHXik.
This key is not known by any other names
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'localhost' (ED25519) to the list of known hosts.
Welcome to Ubuntu 22.04 LTS (GNU/Linux 5.15.0-41-generic x86_64)

 * Documentation:  https://help.ubuntu.com
 * Management:     https://landscape.canonical.com
 * Support:        https://ubuntu.com/advantage

132 updates can be applied immediately.
89 of these updates are standard security updates.
To see these additional updates run: apt list --upgradable


The programs included with the Ubuntu system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.

Ubuntu comes with ABSOLUTELY NO WARRANTY, to the extent permitted by
```

Download and extract the Hadoop setup.

```
hdoop@parth642001-virtual-machine:~$ wget https://dlcdn.apache.org/hadoop/commo
n/hadoop-3.3.3/hadoop-3.3.3.tar.gz
--2022-07-21 21:48:48--  https://dlcdn.apache.org/hadoop/common/hadoop-3.3.3/ha
doop-3.3.3.tar.gz
Resolving dlcdn.apache.org (dlcdn.apache.org)... 151.101.2.132, 2a04:4e42::644
Connecting to dlcdn.apache.org (dlcdn.apache.org)|151.101.2.132|:443... connect
ed.
HTTP request sent, awaiting response... 200 OK
Length: 645040598 (615M) [application/x-gzip]
Saving to: 'hadoop-3.3.3.tar.gz'

hadoop-3.3.3.tar.gz 100%[===================>] 615.16M  1.29MB/s    in 3m 11s

2022-07-21 21:51:59 (3.22 MB/s) - 'hadoop-3.3.3.tar.gz' saved [645040598/645040
598]
```

```
hdoop@parth642001-virtual-machine:~$ tar xzf hadoop-3.3.3.tar.gz
hdoop@parth642001-virtual-machine:~$ ls -lrt
total 629936
drwxr-xr-x 10 hdoop hdoop      4096 May  9 23:14 hadoop-3.3.3
-rw-rw-r--  1 hdoop hdoop 645040598 May 11 22:19 hadoop-3.3.3.tar.gz
drwx------  3 hdoop hdoop      4096 Jul 21 21:47 snap
hdoop@parth642001-virtual-machine:~$
```

Hadoop excels when deployed in a fully distributed mode on a large cluster of networked servers. However, if you are new to Hadoop and want to explore basic commands or test applications, you can configure Hadoop on a single node.

This setup, also called pseudo-distributed mode, allows each Hadoop daemon to run as a single Java process. A Hadoop environment is configured by editing a set of configuration files:

- bashrc
- hadoop-env.sh
- core-site.xml
- hdfs-site.xml
- mapred-site-xml

- yarn-site.xml

```
hdoop@parth642001-virtual-machine:~$ sudo nano .bashrc
[sudo] password for hdoop:
```

```
#Hadoop Related Options
export HADOOP_HOME=/home/hdoop/hadoop-3.2.1
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/nativ"
```

```
hdoop@parth642001-virtual-machine:~$ source ~/.bashrc
-bash: export: `HADOOP_OPTS-Djava.library.path=/home/hdoop/hadoop-3.3.3/lib/nat
iv': not a valid identifier
```

```
xid=0 when meet shutdown.
2022-07-21 22:38:19,460 INFO namenode.NameNode: SHUTDOWN_MSG:
/************************************************************
SHUTDOWN_MSG: Shutting down NameNode at parth642001-virtual-machine/127.0.1.1
************************************************************/
```

```
hdoop@parth642001-virtual-machine:~$ ~/hadoop-3.3.3/
-bash: /home/hdoop/hadoop-3.3.3/: Is a directory
hdoop@parth642001-virtual-machine:~$ cd ~/hadoop-3.3.3/sbin
hdoop@parth642001-virtual-machine:~/hadoop-3.3.3/sbin$ ls
distribute-exclude.sh  mr-jobhistory-daemon.sh  start-dfs.sh        stop-balan
cer.sh     workers.sh
FederationStateStore   refresh-namenodes.sh     start-secure-dns.sh stop-dfs.c
md         yarn-daemon.sh
hadoop-daemon.sh       start-all.cmd            start-yarn.cmd      stop-dfs.s
h          yarn-daemons.sh
hadoop-daemons.sh      start-all.sh             start-yarn.sh       stop-secur
e-dns.sh
httpfs.sh              start-balancer.sh        stop-all.cmd        stop-yarn.
cmd
kms.sh                start-dfs.cmd             stop-all.sh         stop-yarn.
sh
hdoop@parth642001-virtual-machine:~/hadoop-3.3.3/sbin$ ./start-dfs.sh
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [parth642001-virtual-machine]
```
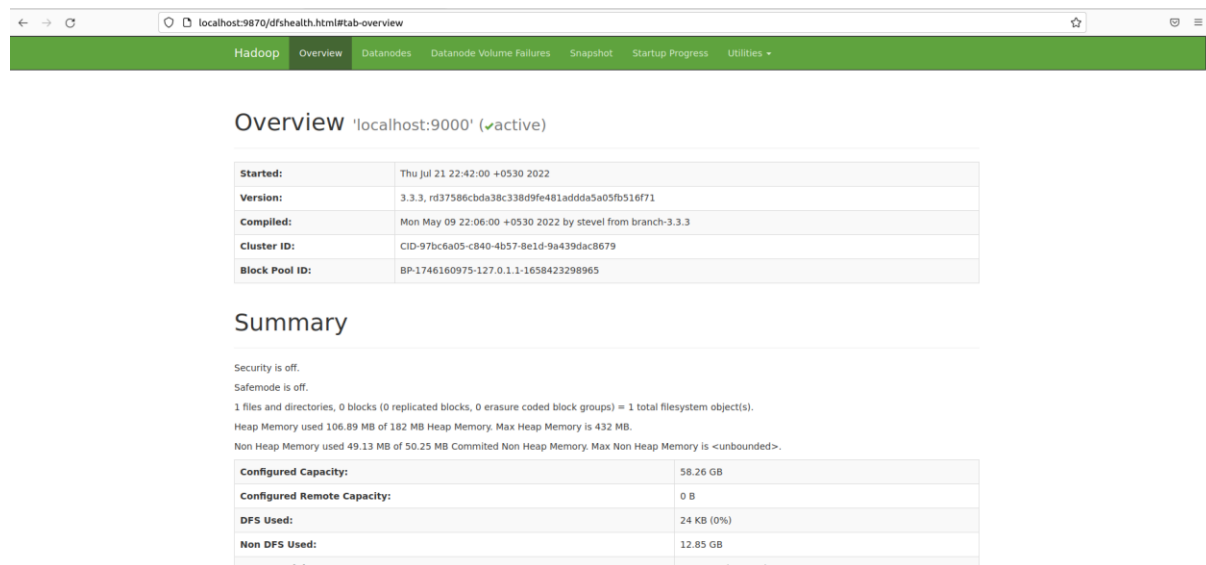
Navigate to the *hadoop-3.2.1/sbin* directory and execute the following commands to start the NameNode and DataNode.

Once the namenode, datanodes, and secondary namenode are up and running, start the YARN resource and nodemanagers

```
hdoop@parth642001-virtual-machine:~/hadoop-3.3.3/sbin$ ./start-dfs.sh
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [parth642001-virtual-machine]
parth642001-virtual-machine: Warning: Permanently added 'parth642001-virtual-ma
chine' (ED25519) to the list of known hosts.
2022-07-21 22:42:11,510 WARN util.NativeCodeLoader: Unable to load native-hadoo
p library for your platform... using builtin-java classes where applicable
hdoop@parth642001-virtual-machine:~/hadoop-3.3.3/sbin$ ./start-yarn.sh
Starting resourcemanager
Starting nodemanagers
hdoop@parth642001-virtual-machine:~/hadoop-3.3.3/sbin$ jps
13920 NameNode
14720 ResourceManager
14833 NodeManager
15177 Jps
14154 DataNode
14459 SecondaryNameNode
hdoop@parth642001-virtual-machine:~/hadoop-3.3.3/sbin$
```
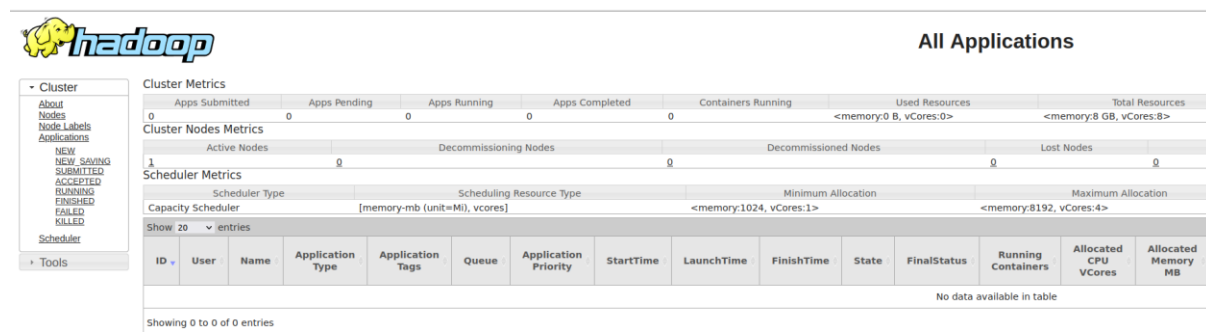
One can also use localhost to access the Hadoop overview.

http://localhost:9870





## CONCLUSION:

By performing this practical, I learnt how to install and configure Hadoop.