

hadoop

HADOOP CLUSTER

By: Numaira Zaib

TABLE OF CONTENTS

Introduction.....	3
Advantages of Hadoop.....	3
Methodology.....	4
Hadoop on ubuntu.....	5
Installing Hadoop.....	5
Downloading Hadoop.....	6
Launching Hadoop.....	8
Master & Slave.....	9
MySQL database replication with master & Slave.....	11
Master machine.....	11
Slave machine.....	12
The benefits of master-slave replication.....	13
MapReduce.....	14
Word Count java.....	14
MapReduce with python.....	17
RPC in Hadoop.....	19
Transparency in Hadoop.....	19
Hadoop web interface	20
Conclusion.....	22

INTRODUCTION

Distributed computing is a model in which components of a software system are shared among multiple computers. Even though the components are spread out across multiple computers, they are run as one system. This is done in order to improve efficiency and performance. Hadoop is an open-source software framework for storing data and running applications on clusters of commodity hardware. It provides massive storage for any kind of data, enormous processing power, and the ability to handle virtually limitless concurrent tasks or jobs. The four important libraries of Hadoop are:

1. Hadoop Common – The role of this character is to provide common utilities that can be used across all modules.
2. Hadoop MapReduce - The right hand of our actor, carrying out all the work assigned to it i.e. It does the job scheduling and processing across the cluster. Hadoop is like a data warehousing system so it needs a library like MapReduce to actually process the data.
3. Hadoop Distributed File System (HDFS) – The left hand, maintains all the records i.e. file system management across the cluster.
4. Hadoop YARN – This is the newer and improved version of MapReduce, from version 2.0, and does the same work.

Hadoop is used where there is a large amount of data generated and your business requires insights from that data. The power of Hadoop lies in its framework, as virtually most of the software can be plugged into it and can be used for data visualization. Hadoop is used in big data applications that gather data from disparate data sources in different formats. HDFS is flexible in storing diverse data types.

ADVANTAGES OF HADOOP

Scalable

Hadoop is a highly scalable storage platform because it can store and distribute very large data sets across hundreds of inexpensive servers that operate in parallel. Unlike traditional relational database systems (RDBMS) that can't scale to process large amounts of data.

Cost-effective

Hadoop also offers a cost-effective storage solution for businesses exploding data sets. The problem with traditional relational database management systems is that it is extremely cost prohibitive to scale to such a degree in order to process such massive volumes of data. In an effort to reduce costs, many companies in the past would have had to down-sample data and classify it based on certain assumptions as to which data was the most valuable.

Flexible

Hadoop enables businesses to easily access new data sources and tap into different types of data (both structured and unstructured) to generate value from that data. This means businesses can use Hadoop to derive valuable business insights from data sources such as social media, email conversations.

Fast

Hadoop's unique storage method is based on a distributed file system that basically 'maps' data wherever it is located on a cluster. The tools for data processing are often on the same servers where the data is located, resulting in much faster data processing. If you're dealing with large volumes of unstructured data, Hadoop is able to efficiently process terabytes of data in just minutes, and petabytes in hours.

Resilient to failure

A key advantage of using Hadoop is its fault tolerance. When data is sent to an individual node, that data is also replicated to other nodes in the cluster, which means that in the event of failure, there is another copy available for use.

METHODOLOGY

We opened a virtual machine, then we install java and run other commands and make changes in 5 files to install hadoop. Now hadoop is installed and we can access hadoop directly through the terminal and it is for a single node cluster.

Then we make master and slave from one virtual machine by cloning and making it 'Master' and 'Slave'. Then we find out the addresses of the master and slave and then ping their addresses. Then we install SQL through hadoop. We make changes in the database of slave and when we write command show database in master then it displays the same result as a slave. So we can access the same records and database from the master as well. The connection between master and slave is fine.

Then we write commands for running MapReduce. We used MapReduce 3.2.3 in the project. We implemented MapReduce in java and python in both languages. MapReduce is a software framework and programming model used for processing huge amounts of data. MapReduce program work in two phases, namely, Map and Reduce. Map tasks deal with splitting and mapping of data while Reducing tasks shuffle and reducing the data. Hadoop is capable of running MapReduce programs written in various languages: Java, Ruby, Python, and C++. The programs of Map Reduce in cloud computing are parallel in nature, and thus are very useful for performing large-scale data analysis using multiple machines in the cluster. The input to each phase is key-value pairs. In addition, every programmer needs to specify two functions: map function and reduce function.

HADOOP ON UBUNTU

INSTALLING HADOOP

sudo apt update: This command is used for updating the system.

sudo apt install openjdk-8-jdk -y: This command is used for installing java in the machine.

java -version; javac -version: This command is used for checking java version.

sudo apt install openssh-server openssh-client -y: This command is for opening ssh server. Since hadoop will run in the same machine where my unix box is running so there has to be a connection between hadoop and my unix machine. This ssh server will give me that connection.

sudo adduser hdoop: This command is used for adding a new user 'hdoop'. This user will use hadoop in my system.

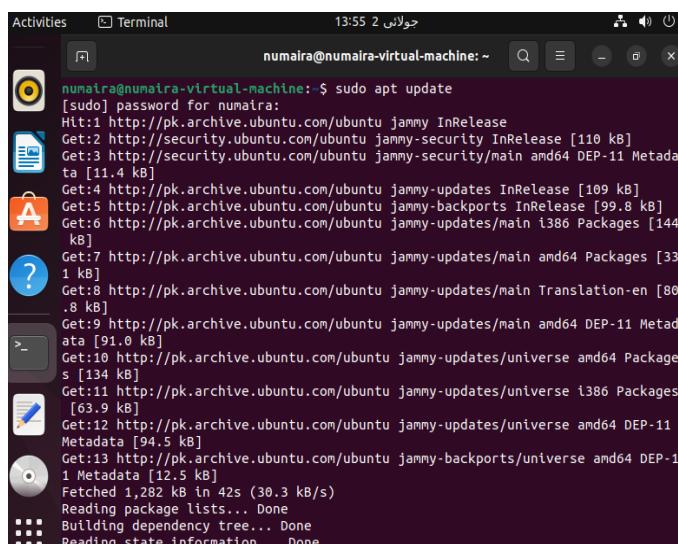
su - hdoop: This command is used for changing the user.

ssh-keygen -t rsa -P "" -f ~/.ssh/id_rsa: This command is used for generating a key which I will use for the connection or the coordination to the local machine to the hadoop.

cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys: This command is used for putting my key in my authorized key list. This is why we are making hdoop user able to make connection with hadoop.

chmod 0600 ~/.ssh/authorized_keys: This command is used for changing the permission of this authorized key file.

ssh localhost: This means the connection is going.



```
Activities Terminal 13:55 2 جولانی 2 numaira@numaira-virtual-machine:~
```

```
numaira@numaira-virtual-machine:~$ sudo apt update
[sudo] password for numaira:
Hit:1 http://pk.archive.ubuntu.com/ubuntu jammy InRelease
Get:2 http://security.ubuntu.com/ubuntu jammy-security InRelease [110 kB]
Get:3 http://security.ubuntu.com/ubuntu jammy-security/main amd64 DEP-11 Metadata [11.4 kB]
Get:4 http://pk.archive.ubuntu.com/ubuntu jammy-updates InRelease [109 kB]
Get:5 http://pk.archive.ubuntu.com/ubuntu jammy-backports InRelease [99.8 kB]
Get:6 http://pk.archive.ubuntu.com/ubuntu jammy-updates/main i386 Packages [144 kB]
Get:7 http://pk.archive.ubuntu.com/ubuntu jammy-updates/main amd64 Packages [331 kB]
Get:8 http://pk.archive.ubuntu.com/ubuntu jammy-updates/main Translation-en [80.8 kB]
Get:9 http://pk.archive.ubuntu.com/ubuntu jammy-updates/main amd64 DEP-11 Metadata [91.0 kB]
Get:10 http://pk.archive.ubuntu.com/ubuntu jammy-updates/universe amd64 Packages [134 kB]
Get:11 http://pk.archive.ubuntu.com/ubuntu jammy-updates/universe i386 Packages [63.9 kB]
Get:12 http://pk.archive.ubuntu.com/ubuntu jammy-updates/universe amd64 DEP-11 Metadata [94.5 kB]
Get:13 http://pk.archive.ubuntu.com/ubuntu jammy-backports/universe amd64 DEP-11 Metadata [12.5 kB]
Fetched 1,282 kB in 42s (30.3 kB/s)
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
```

Numaira Zaib

Activities Terminal 14:23 2 جوادر numaira@numaira-virtual-machine:~

```
Get:12 http://pk.archive.ubuntu.com/ubuntu jammy-updates/universe amd64 DEP-11  
Metadata [94.5 kB]  
Get:13 http://pk.archive.ubuntu.com/ubuntu jammy-backports/universe amd64 DEP-11  
1 Metadata [12.5 kB]  
Fetched 1,282 kB in 42s (30.3 kB/s)  
Reading package lists... Done  
Building dependency tree... Done  
Reading state information... Done  
199 packages can be upgraded. Run 'apt list --upgradable' to see them.  
numaira@numaira-virtual-machine:~$ sudo apt install openjdk-8-jdk -y  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Waiting for cache lock: Could not get lock /var/lib/dpkg/lock-frontend. It is held by process 2984 (unattended-upgr)  
Setting up ncurses-term (0.5-2) ...  
Processing triggers for man-db (2.10.2-1) ...  
Processing triggers for ufw (0.36.1-4build1) ...  
numaira@numaira-virtual-machine:~$ sudo adduser hdoop  
Adding user 'hdoop' ...  
Adding new group 'hdoop' (1001) ...  
Adding new user 'hdoop' (1001) with group 'hdoop' ...  
Creating home directory '/home/hdoop' ...  
Copying files from '/etc/skel' ...  
New password:  
BAD PASSWORD: The password is shorter than 8 characters  
Retype new password:  
Sorry, passwords do not match.  
New password:  
Retype new password:  
passwd: password updated successfully  
Changing the user information for hdoop  
Enter the new value, or press ENTER for the default  
  Full Name []: Room Number []:  
  Work Phone []:  
  Home Phone []:  
  Other []:  
Is the information correct? [Y/n] y  
numaira@numaira-virtual-machine:~$ su - hdoop  
Password:  
hdoop@numaira-virtual-machine:~$
```

Activities Terminal 14:28 2 جوادر numaira@numaira-virtual-machine:~

```
Setting up libice-dev:amd64 (2:1.0.10-1build2) ...  
Setting up libsm-dev:amd64 (2:1.2.3-1build2) ...  
Setting up libxdmcp-dev:amd64 (1:1.1.3-0ubuntu5) ...  
Setting up libxcb1-dev:amd64 (1:1.14-3ubuntu3) ...  
Setting up libx11-dev:amd64 (2:1.7.5-1) ...  
Setting up libxt-dev:amd64 (1:1.2.1-1) ...  
numaira@numaira-virtual-machine:~$ java -version; javac -version  
openjdk version "1.8.0_312"  
OpenJDK Runtime Environment (build 1.8.0_312-8u312-b07-0ubuntu1-b07)  
javac 1.8.0_312  
numaira@numaira-virtual-machine:~$ sudo apt install openssh-server openssh-client  
[sudo] password for numaira:  
Reading package lists... Done  
Building dependency tree... Done  
Reading state information... Done  
openssh-client is already the newest version (1:8.9p1-3).  
openssh-client set to manually installed.  
The following additional packages will be installed:  
  ncurses-term openssh-sftp-server ssh-import-id  
Suggested packages:  
  molly-guard monkeysphere ssh-askpass  
The following NEW packages will be installed:  
  ncurses-term openssh-server openssh-sftp-server ssh-import-id  
0 upgraded, 4 newly installed, 0 to remove and 199 not upgraded.  
Need to get 751 kB of archives.  
After this operation, 6,046 kB of additional disk space will be used.  
Get:1 http://pk.archive.ubuntu.com/ubuntu jammy/main amd64 openssh-sftp-server
```

Activities Terminal 14:40 2 جوادر hdoop@numaira-virtual-machine:~

```
? Password:  
hdoop@numaira-virtual-machine:~$ ssh-keygen -t rsa -P '' -f ~/.ssh/id_rsa  
Generating public/private rsa key pair.  
Created directory '/home/hdoop/.ssh'.  
Your identification has been saved in /home/hdoop/.ssh/id_rsa.  
Your public key has been saved in /home/hdoop/.ssh/id_rsa.pub  
The key fingerprint is:  
SHA256:Loo3nLGrVko4C9vt4ZvpE5bHPxxOGQp+toq/pwHRko hdoop@numaira-virtual-machine  
The key's randomart image is:  
++-[RSA 3072]-++  
| .  
| ..++.. S  
| .+o+o**..  
| .Eo+O*o& .  
| ..*+=X.*  
| oooo+=++#o  
+---[SHA256]---+  
hdoop@numaira-virtual-machine:~$ cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys  
s  
hdoop@numaira-virtual-machine:~$ chmod 6000 ~/.ssh/authorized_keys  
hdoop@numaira-virtual-machine:~$ ssh localhost  
The authenticity of host 'localhost (127.0.0.1)' can't be established.  
ED25519 key fingerprint is SHA256:7Jl6SS7XaUM7lB3+dBXJVLYaq3knAIi6tWPk+UdooB0.  
This key is not known by any other names  
Are you sure you want to continue connecting (yes/no/[fingerprint])? ^[[B
```

DOWNLOADING HADOOP

wget <https://downloads.apache.org/hadoop/common/hadoop-3.2.3/hadoop-3.2.3.tar.gz>

For downloading hadoop from this location.

tar xzf hadoop-3.2.3.tar.gz

For unzipping the file.

Then we edited the 6 important files.

- 1) .bashrc
- 2) hadoop-env.sh
- 3) core-site.xml
- 4) hdfs-site.xml
- 5) mapred-site.xml
- 6) yarn-site.xml

Numaira Zaib

```
Activities Terminal 15:57 2 جولانی 2
hdooop@numaira-virtual-machine:~ $ wget https://downloads.apache.org/hadoop/common/hadoop-3.2.3/hadoop-3.2.3.tar.gz
--2022-07-02 14:47:25-- https://downloads.apache.org/hadoop/common/hadoop-3.2.3/hadoop-3.2.3.tar.gz
Resolving downloads.apache.org (downloads.apache.org)... 135.181.214.104, 88.99.191.2, 2a01:4f8:10a::2, ...
Connecting to downloads.apache.org (downloads.apache.org)|135.181.214.104|:443...
HTTP request sent, awaiting response... 200 OK
Length: 492241961 (469M) [application/x-gzip]
Saving to: 'hadoop-3.2.3.tar.gz'

hadoop-3.2.3.tar.gz 44%[=====] 207.33M -- .KB/s eta
hadoop-3.2.3.tar.gz 100%[=====] 469.44M 110KB/s in 57m 54s
2022-07-02 15:53:21 (138 KB/s) - 'hadoop-3.2.3.tar.gz' saved [492241961/492241961]

hdooop@numaira-virtual-machine:~ $ tar xzf hadoop-3.2.3.tar.gz
hdooop@numaira-virtual-machine:~ $ ls -lrt
total 480720
drwxr-xr-x 9 hdooop hdooop 4096 06:58 20 جولانی.hadoop-3.2.3
-rw-rw-r-- 1 hdooop hdooop 492241961 09:54 28 جولانی.hadoop-3.2.3.tar.gz
drwxr----- 3 hdooop hdooop 4096 14:41 2 snap
hdooop@numaira-virtual-machine:~ $
```

```
For more details see su(1).
numaira@numaira-virtual-machine:~ $ sudo adduser hdooop
adduser: The user 'hdooop' already exists.
numaira@numaira-virtual-machine:~ $ su - hdooop
Password:
To run a command as administrator (user "root"), use "sudo <command>".
See "man sudo_root" for details.

hdooop@numaira-virtual-machine:~ $ sudo nano .bashrc
[sudo] password for hdooop:
hdooop@numaira-virtual-machine:~ $ cat .bashrc
# ~/.bashrc: executed by bash(1) for non-login shells.
# see /usr/share/doc/bash/examples/startup-files (in the package bash-doc)
# for examples

# If not running interactively, don't do anything
Activities Terminal 16:17 2 جولانی 2
hdooop@numaira-virtual-machine:~ $ 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 -q 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
#Hadoop Related Options
export HADOOP_HOME=/home/hdooop/hadoop-3.2.3
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"
hdooop@numaira-virtual-machine:~ $
```

```
hdooop@numaira-virtual-machine:~ $ sudo nano $HADOOP_HOME/etc/hadoop/hdfs-site.xml
[sudo] password for hdooop:
hdooop@numaira-virtual-machine:~ $ Activities Terminal 16:56 2 جولانی 2
GNU nano 6.2 /home/hdooop/hadoop-3.2.3/etc/hadoop/hdfs-site.xml *
http://www.apache.org/licenses/LICENSE-2.0

Unless required by applicable law or agreed to in writing, software
distributed under the License is distributed on an "AS IS" BASIS,
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License. See accompanying LICENSE file.
-->
<!-- Put site-specific property overrides in this file. -->

<configuration>
<property>
<name>dfs.data.dir</name>
<value>/home/hdooop/dfsdata/namenode</value>
</property>
<property>
<name>dfs.data.dir</name>
<value>/home/hdooop/dfsdata/datanode</value>
</property>
<property>
<name>dfs.replication</name>
<value>1</value>
</property>
</configuration>
<!-->
```

```
Activities Terminal 16:06 2 جولانی 2
hdooop@numaira-virtual-machine:~ $ su - numaira
numaira@numaira-virtual-machine:~ $ sudo adduser hdooop sudo
[sudo] password for numaira:
Adding user `hdooop' to group 'sudo' ...
Adding user numaira to group sudo
Done.
numaira@numaira-virtual-machine:~ $ sudo nano .bashrc
Activities Terminal 16:06 2 جولانی 2
GNU nano 6.2 .bashrc *
#!/bin/sh -e
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 -q 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
#Hadoop Related Options
export HADOOP_HOME=/home/hdooop/hadoop-3.2.3
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"
hdooop@numaira-virtual-machine:~ $
```

```
hdooop@numaira-virtual-machine:~ $ source ~/bashrc
-bash: export: 'HADOOP_OPTS-Djava.library.path=/home/hdooop/hadoop-3.2.3/lib/native': not a valid identifier
hdooop@numaira-virtual-machine:~ $ sudo nano $HADOOP_HOME/etc/hadoop/hadoop-env.sh
# It uses the format of (command)_(&subcommand)_USER.
#
# For example, to limit who can execute the namenode command,
# export HDFS_NAMENODE_USER=hdfs
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
hdooop@numaira-virtual-machine:~ $
```

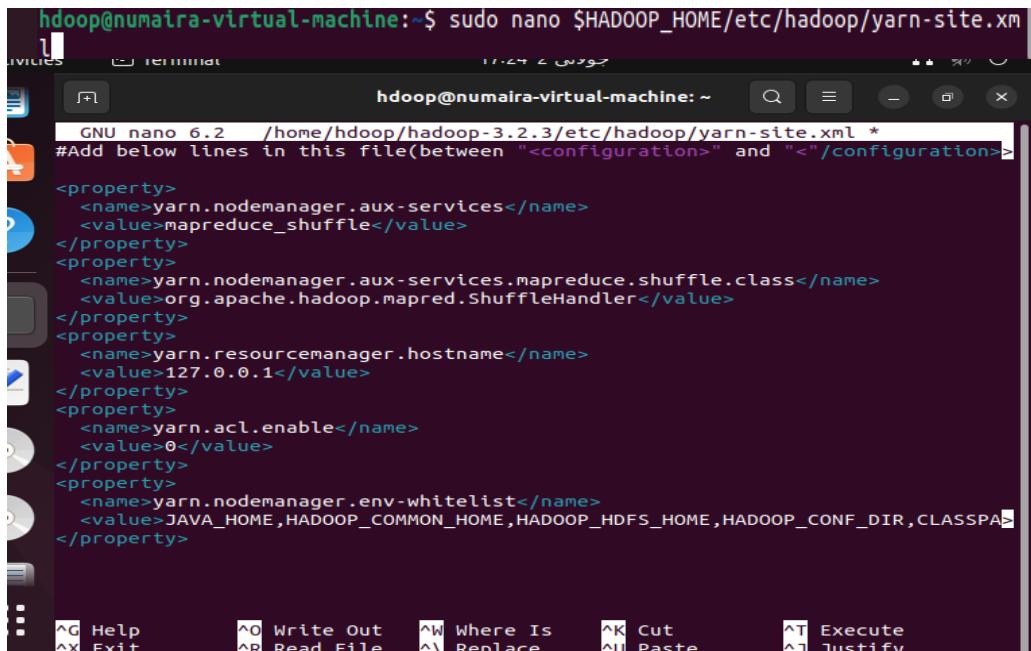
```
hdooop@numaira-virtual-machine:~ $ sudo nano $HADOOP_HOME/etc/hadoop/core-site.xml
Activities Terminal 16:37 2 جولانی 2
GNU nano 6.2 /home/hdooop/hadoop-3.2.3/etc/hadoop/core-site.xml *
http://www.apache.org/licenses/LICENSE-2.0

Unless required by applicable law or agreed to in writing, software
distributed under the License is distributed on an "AS IS" BASIS,
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License. See accompanying LICENSE file.
-->
<!-- Put site-specific property overrides in this file. -->

<configuration>
<property>
<name>hadoop.tmp.dir</name>
<value>/home/hdooop/tmpdata</value>
<description>A base for other temporary directories.</description>
</property>
<property>
<name>fs.default.name</name>
<value>hdfs://localhost:9000</value>
<description>The name of the default file system</description>
</property>
</configuration>
hdooop@numaira-virtual-machine:~ $
```

```
hdooop@numaira-virtual-machine:~ $ sudo nano $HADOOP_HOME/etc/hadoop/mapred-site.xml
Activities Terminal 17:18 2 جولانی 2
GNU nano 6.2 /home/hdooop/hadoop-3.2.3/etc/hadoop/mapred-site.xml *
<!-- Version 1.0 -->
<xsl:stylesheet type="text/xsl" href="configuration.xsl"?>
<!-- Licensed under the Apache License, Version 2.0 (the "License");
you may not use this file except in compliance with the License.
You may obtain a copy of the License at
http://www.apache.org/licenses/LICENSE-2.0

Unless required by applicable law or agreed to in writing, software
distributed under the License is distributed on an "AS IS" BASIS,
WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
See the License for the specific language governing permissions and
limitations under the License. See accompanying LICENSE file.
-->
<configuration>
<property>
<name>mapreduce.framework.name</name>
<value>yarn</value>
</property>
</configuration>
hdooop@numaira-virtual-machine:~ $
```



hadoop@numaira-virtual-machine:~\$ sudo nano \$HADOOP_HOME/etc/hadoop/yarn-site.xml

```
GNU nano 6.2 /home/hadoop/hadoop-3.2.3/etc/hadoop/yarn-site.xml *
#Add below lines in this file(between "<configuration>" and "</configuration>"

<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>
<property>
  <name>yarn.resourcemanager.hostname</name>
  <value>127.0.0.1</value>
</property>
<property>
  <name>yarn.acl.enable</name>
  <value>0</value>
</property>
<property>
  <name>yarn.nodemanager.env-whitelist</name>
  <value>JAVA_HOME,HADOOP_COMMON_HOME,HADOOP_HDFS_HOME,HADOOP_CONF_DIR,CLASSPA</value>
</property>
```

File menu Edit menu View menu Insert menu Search menu Help menu

Help Write Out Where Is Cut Execute

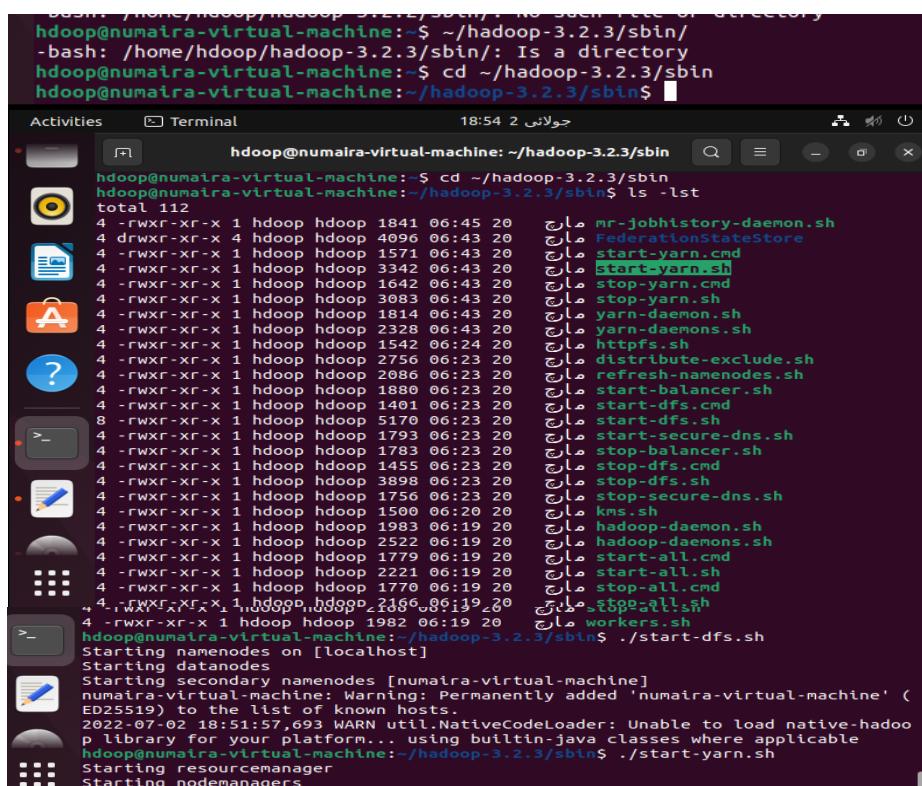
Exit Read File Replace Paste Justify

LAUNCHING HADOOP

hdfs namenode –format

./start-dfs.sh

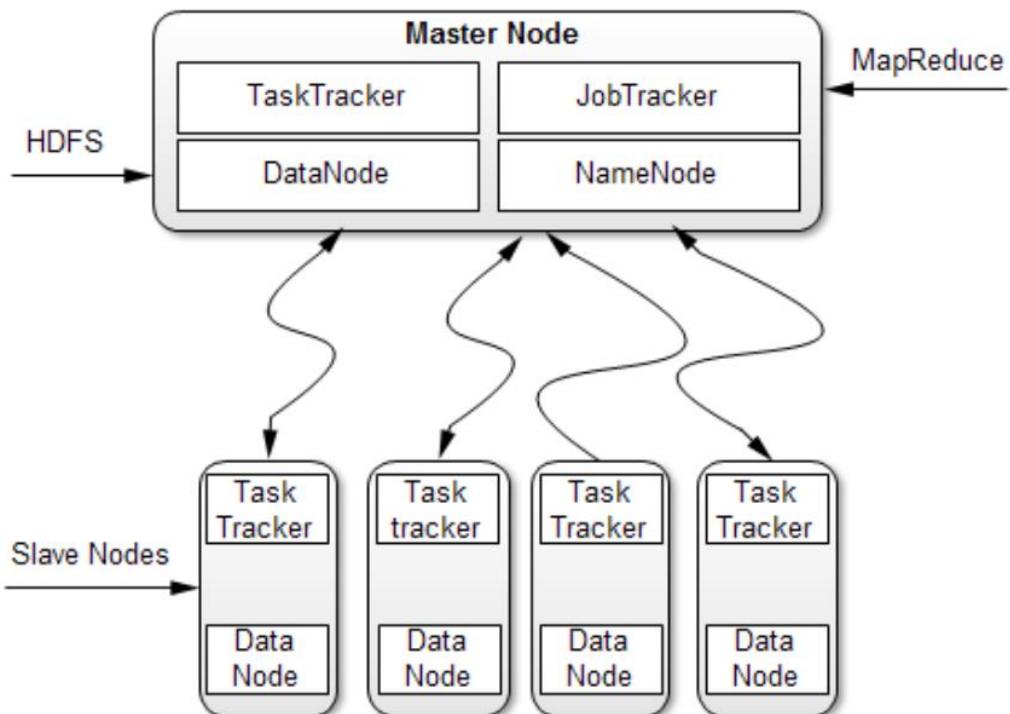
Then we move a local file to hdfs



```
Bash: /home/hadoop/hadoop-3.2.3/sbin/: No such file or directory
hadoop@numaira-virtual-machine:~$ ~/hadoop-3.2.3/sbin/
-bash: /home/hadoop/hadoop-3.2.3/sbin/: Is a directory
hadoop@numaira-virtual-machine:~/hadoop-3.2.3/sbin$ ls -l
total 112
4 -rwxr-xr-x 1 hadoop hadoop 1841 06:45 20  مـ mr-jobhistory-daemon.sh
4 drwxr-xr-x 4 hadoop hadoop 4096 06:43 20  مـ FederationStateStore
4 -rwxr-xr-x 1 hadoop hadoop 1571 06:43 20  مـ start-yarn.cmd
4 -rwxr-xr-x 1 hadoop hadoop 3342 06:43 20  مـ start-yarn.sh
4 -rwxr-xr-x 1 hadoop hadoop 1642 06:43 20  مـ stop-yarn.cmd
4 -rwxr-xr-x 1 hadoop hadoop 3083 06:43 20  مـ stop-yarn.sh
4 -rwxr-xr-x 1 hadoop hadoop 1813 06:43 20  مـ yarn-daemon.sh
4 -rwxr-xr-x 1 hadoop hadoop 2328 06:43 20  مـ yarn-daemons.sh
4 -rwxr-xr-x 1 hadoop hadoop 1542 06:24 20  مـ httpfs.sh
4 -rwxr-xr-x 1 hadoop hadoop 2756 06:23 20  مـ distribute-exclude.sh
4 -rwxr-xr-x 1 hadoop hadoop 2086 06:23 20  مـ refresh-namenodes.sh
4 -rwxr-xr-x 1 hadoop hadoop 1880 06:23 20  مـ start-balancer.sh
4 -rwxr-xr-x 1 hadoop hadoop 1401 06:23 20  مـ start-dfs.cmd
8 -rwxr-xr-x 1 hadoop hadoop 5170 06:23 20  مـ start-dfs.sh
4 -rwxr-xr-x 1 hadoop hadoop 1793 06:23 20  مـ start-secure-dns.sh
4 -rwxr-xr-x 1 hadoop hadoop 1783 06:23 20  مـ stop-balancer.sh
4 -rwxr-xr-x 1 hadoop hadoop 1455 06:23 20  مـ stop-dfs.cmd
4 -rwxr-xr-x 1 hadoop hadoop 3898 06:23 20  مـ stop-dfs.sh
4 -rwxr-xr-x 1 hadoop hadoop 1756 06:23 20  مـ stop-secure-dns.sh
4 -rwxr-xr-x 1 hadoop hadoop 1500 06:20 20  مـ kms.sh
4 -rwxr-xr-x 1 hadoop hadoop 1983 06:19 20  مـ hadoop-daemon.sh
4 -rwxr-xr-x 1 hadoop hadoop 2522 06:19 20  مـ hadoop-daemons.sh
4 -rwxr-xr-x 1 hadoop hadoop 1779 06:19 20  مـ start-all.cmd
4 -rwxr-xr-x 1 hadoop hadoop 2221 06:19 20  مـ start-all.sh
4 -rwxr-xr-x 1 hadoop hadoop 1770 06:19 20  مـ stop-all.cmd
4 -rwxr-xr-x 1 hadoop hadoop 1982 06:19 20  مـ workers.sh
hadoop@numaira-virtual-machine:~/hadoop-3.2.3/sbin$ ./start-dfs.sh
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [numaira-virtual-machine]
numaira-virtual-machine: Warning: Permanently added 'numaira-virtual-machine' (ED25519) to the list of known hosts
2022-07-02 18:51:57,693 WARN util.NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-Java classes where applicable
hadoop@numaira-virtual-machine:~/hadoop-3.2.3/sbin$ ./start-yarn.sh
Starting resourcemanager
Starting nodemanagers
```

MASTER & SLAVE

Hadoop is based on Master/Slave architecture. Master is responsible for NameNode and JobTracker while Slave is responsible for DataNode and TaskTracker. Job Tracker initiates the task and keeps track of each task. TaskTracker manages local data processing and collects its result as per the requirement of application. It also sends a report on the progress of tasks to JobTracker. NameNode and DataNode are responsible for HDFS tasks while JobTracker and TaskTracker are responsible for MapReduce tasks.



SSH MASTER SLAVE

Copy the public key of user master into the slave

```
ssh-copy-id-i.ssh/id_rsa.pub hdoop@slave
```

This command will prompt you for the login password for user hdoop on slave, then copy the public SSH key for you, creating directory and fixing the permissions as necessary.

Connecting from master to master

```
ssh master
```

Connecting from master to slave

```
Ssh slave
```

Numaira Zaib

The image consists of three vertically stacked screenshots of a macOS desktop environment. Each screenshot shows two terminal windows side-by-side.

Top Screenshot: Displays two terminal windows at 12:48. The left window is on the master node (hadoop@master) and the right is on the slave node (hadoop@slave). Both show the output of the 'ifconfig' command, which lists network interfaces (ens33, lo) with their respective statistics and link layer information.

Middle Screenshot: Displays two terminal windows at 13:11. The left window is on the master node (hadoop@master) and the right is on the slave node (hadoop@slave). The master window shows the output of the 'ssh-copy-id -i \$HOME/.ssh/id_rsa.pub hadoop@slave' command, which attempts to copy the public key to the slave host. It includes a warning about skipped keys and a prompt to continue connecting. The slave window shows the output of the 'ping' command to the master host.

Bottom Screenshot: Displays two terminal windows at 13:16. The left window is on the master node (hadoop@master) and the right is on the slave node (hadoop@slave). Both windows show the output of the 'apt update' command, which lists available updates. The master window also shows the output of the 'ssh master' command, which fails due to host key fingerprint mismatch. The slave window shows the output of the 'ssh slave' command, which succeeds.

```
hadoop@master:~$ ifconfig
ens33: flags=4163<UP,BROADCAST,RUNNING,MULTICAST> mtu 1500
    inet 192.168.154.134 netmask 255.255.255.0 broadcast 192.168.154.255
        inet6 fe80::f6ca:6c46:9cc9:203b prefixlen 64 scopeid 0x20<link>
            ether 00:0c:29:70:66:e2 txqueuelen 1000 (Ethernet)
                RX packets 399 bytes 272271 (272.2 KB)
                RX errors 0 dropped 0 overruns 0 frame 0
                TX packets 170 bytes 23816 (23.8 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 171 bytes 16261 (16.2 KB)
                RX errors 0 dropped 0 overruns 0 frame 0
                TX packets 171 bytes 16261 (16.2 KB)
                TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0

hadoop@master:~$ ping 192.168.154.135
PING 192.168.154.135 (192.168.154.135) 56(84) bytes of data.
64 bytes from 192.168.154.135: icmp_seq=1 ttl=64 time=3.96 ms
64 bytes from 192.168.154.135: icmp_seq=2 ttl=64 time=0.915 ms
64 bytes from 192.168.154.135: icmp_seq=3 ttl=64 time=0.672 ms
64 bytes from 192.168.154.135: icmp_seq=4 ttl=64 time=0.746 ms
64 bytes from 192.168.154.135: icmp_seq=5 ttl=64 time=0.604 ms
64 bytes from 192.168.154.135: icmp_seq=6 ttl=64 time=0.554 ms
64 bytes from 192.168.154.135: icmp_seq=7 ttl=64 time=1.31 ms
64 bytes from 192.168.154.135: icmp_seq=8 ttl=64 time=0.903 ms

hadoop@master:~$ ssh-copy-id -i $HOME/.ssh/id_rsa.pub hadoop@slave
/usr/bin/ssh-copy-id: INFO: Source of key(s) to be installed: "/home/hadoop/.ssh/id_rsa.pub"
The authenticity of host 'slave (192.168.154.135)' can't be established.
ED25519 key fingerprint is SHA256:7J16SS7XauUM7iB3+BXJVLYaq3knAi6iWPk+UdooB0.
This host key is known by the following other names/addresses:
    ~/.ssh/known_hosts:1: [hashed name]
    ~/.ssh/known_hosts:4: [hashed name]
    ~/.ssh/known_hosts:5: [hashed name]
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
/usr/bin/ssh-copy-id: INFO: attempting to log in with the new key(s), to filter
out any that are already installed
/usr/bin/ssh-copy-id: WARNING: All keys were skipped because they already exist
on the remote system.
(if you think this is a mistake, you may want to use -f option)

hadoop@master:~$ 
```

```
hadoop@slave:~$ ifconfig
ens33: flags=4163<UP,BROADCAST,RUNNING,MULTICAST> mtu 1500
    inet 192.168.154.135 netmask 255.255.255.0 broadcast 192.168.154.255
        inet6 fe80::3669:982b:db5f:33ee prefixlen 64 scopeid 0x20<link>
            ether 00:0c:29:59:f7:44 txqueuelen 1000 (Ethernet)
                RX packets 407 bytes 272338 (272.3 KB)
                RX errors 0 dropped 0 overruns 0 frame 0
                TX packets 175 bytes 24444 (24.4 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 178 bytes 15771 (15.7 KB)
                RX errors 0 dropped 0 overruns 0 frame 0
                TX packets 178 bytes 15771 (15.7 KB)
                TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0

hadoop@slave:~$ ping 192.168.154.134
PING 192.168.154.134 (192.168.154.134) 56(84) bytes of data.
64 bytes from 192.168.154.134: icmp_seq=1 ttl=64 time=1.80 ms
64 bytes from 192.168.154.134: icmp_seq=2 ttl=64 time=1.02 ms
64 bytes from 192.168.154.134: icmp_seq=3 ttl=64 time=1.03 ms
...
-- 192.168.154.134 ping statistics --
3 packets transmitted, 3 received, 0% packet loss, time 2002ms
rtt min/avg/max/ndev = 1.023/1.282/1.798/0.364 ms
hadoop@slave:~$ 
```

```
hadoop@master:~$ /usr/bin/ssh-copy-id -i $HOME/.ssh/id_rsa.pub hadoop@slave
/usr/bin/ssh-copy-id: WARNING: All keys were skipped because they already exist
on the remote system.
(if you think this is a mistake, you may want to use -f option)

hadoop@master:~$ ssh master
Welcome to Ubuntu 22.04 LTS (GNU/Linux 5.15.0-25-generic x86_64)

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

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

Last login: Sat Jul  2 20:20:22 2022 from 192.168.154.134
hadoop@master:~$ ssh slave
Welcome to Ubuntu 22.04 LTS (GNU/Linux 5.15.0-25-generic x86_64)

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

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

Last login: Sat Jul  2 20:21:49 2022 from 192.168.154.135
hadoop@slave:~$ 
```

```
hadoop@master:~$ apt update
* Management: https://landscape.canonical.com
* Support: https://ubuntu.com/advantage

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

Last login: Sat Jul  2 20:21:49 2022 from 192.168.154.135
hadoop@master:~$ ssh master
The authenticity of host 'master (192.168.154.134)' can't be established.
ED25519 key fingerprint is SHA256:7J16SS7XauUM7iB3+BXJVLYaq3knAi6iWPk+UdooB0.
This host key is known by the following other names/addresses:
    ~/.ssh/known_hosts:1: [hashed name]
    ~/.ssh/known_hosts:4: [hashed name]
    ~/.ssh/known_hosts:5: [hashed name]
Are you sure you want to continue connecting (yes/no/[fingerprint])? yes
Warning: Permanently added 'master' (ED25519) to the list of known hosts.
Welcome to Ubuntu 22.04 LTS (GNU/Linux 5.15.0-25-generic x86_64)

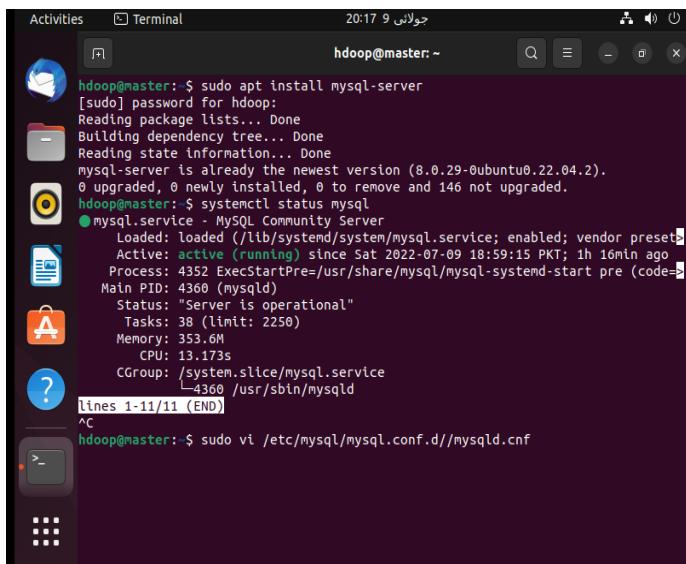
 * Documentation: https://help.ubuntu.com
 * Management: https://landscape.canonical.com
 * Support: https://ubuntu.com/advantage

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

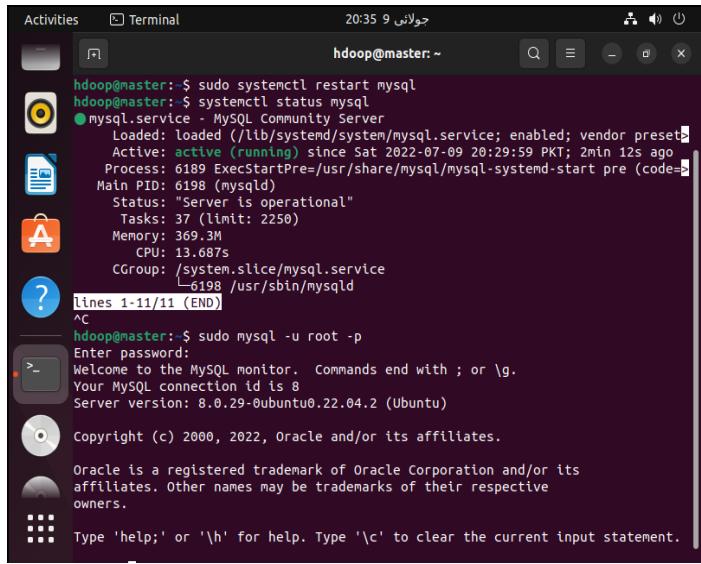
Last login: Sun Jul  3 13:12:23 2022 from 192.168.154.134
hadoop@master:~$ 
```

MySQL DATABASE REPLICATION WITH MASTER & SLAVE

MASTER MACHINE



```
hadoop@master:~$ sudo apt install mysql-server
[sudo] password for hadoop:
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
mysql-server is already the newest version (8.0.29-0ubuntu0.22.04.2).
0 upgraded, 0 newly installed, 0 to remove and 146 not upgraded.
hadoop@master:~$ systemctl status mysql
● mysql.service - MySQL Community Server
   Loaded: loaded (/lib/systemd/system/mysql.service; enabled; vendor preset: enabled)
   Active: active (running) since Sat 2022-07-09 18:59:15 PKT; 1h 16min ago
     Process: 6189 ExecStartPre=/usr/share/mysql/mysql-systemd-start pre (code=0)
      Main PID: 6198 (mysqld)
        Status: "Server is operational"
       Tasks: 37 (limit: 2250)
      Memory: 369.3M
         CPU: 13.687s
      CGroup: /system.slice/mysql.service
              └─ 6198 /usr/sbin/mysqld
lines 1-11/11 (END)
^C
hadoop@master:~$ sudo vi /etc/mysql/mysql.conf.d/mysqld.cnf
```

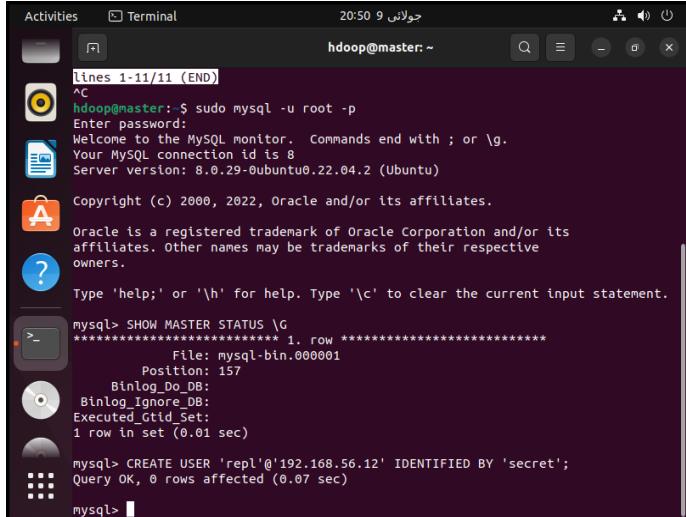


```
hadoop@master:~$ sudo systemctl restart mysql
hadoop@master:~$ systemctl status mysql
● mysql.service - MySQL Community Server
   Loaded: loaded (/lib/systemd/system/mysql.service; enabled; vendor preset: enabled)
   Active: active (running) since Sat 2022-07-09 20:29:59 PKT; 2min 12s ago
     Process: 6189 ExecStartPre=/usr/share/mysql/mysql-systemd-start pre (code=0)
      Main PID: 6198 (mysqld)
        Status: "Server is operational"
       Tasks: 37 (limit: 2250)
      Memory: 369.3M
         CPU: 13.687s
      CGroup: /system.slice/mysql.service
              └─ 6198 /usr/sbin/mysqld
lines 1-11/11 (END)
^C
hadoop@master:~$ sudo mysql -u root -p
Enter password:
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 8
Server version: 8.0.29-0ubuntu0.22.04.2 (Ubuntu)

Copyright (c) 2000, 2022, Oracle and/or its affiliates.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.
```



```
hadoop@master:~$ sudo mysql -u root -p
Enter password:
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 8
Server version: 8.0.29-0ubuntu0.22.04.2 (Ubuntu)

Copyright (c) 2000, 2022, Oracle and/or its affiliates.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> SHOW MASTER STATUS \G
***** 1. row *****
    File: mysql-bin.000001
    Position: 157
    Binlog_Do_DB:
    Binlog_Ignore_DB:
    Executed_Gtid_Set:
1 row in set (0.01 sec)

mysql> CREATE USER 'repl'@'192.168.56.12' IDENTIFIED BY 'secret';
Query OK, 0 rows affected (0.07 sec)

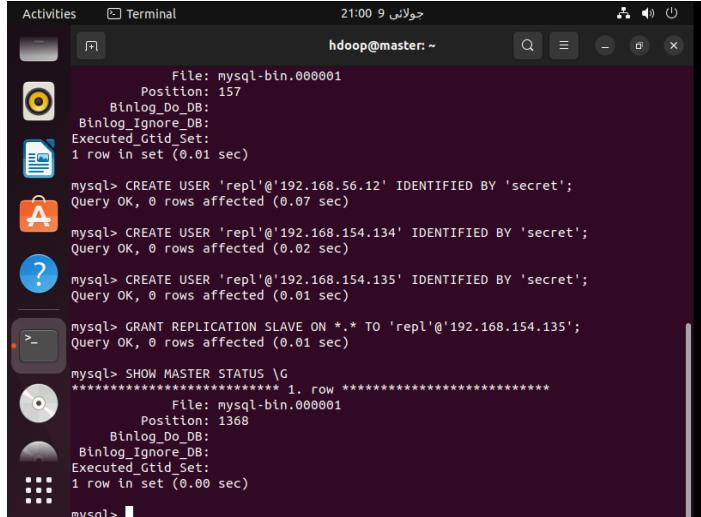
mysql> CREATE USER 'repl'@'192.168.154.134' IDENTIFIED BY 'secret';
Query OK, 0 rows affected (0.02 sec)

mysql> CREATE USER 'repl'@'192.168.154.135' IDENTIFIED BY 'secret';
Query OK, 0 rows affected (0.01 sec)

mysql> GRANT REPLICATION SLAVE ON *.* TO 'repl'@'192.168.154.135';
Query OK, 0 rows affected (0.01 sec)

mysql> SHOW MASTER STATUS \G
***** 1. row *****
    File: mysql-bin.000001
    Position: 1368
    Binlog_Do_DB:
    Binlog_Ignore_DB:
    Executed_Gtid_Set:
1 row in set (0.00 sec)

mysql>
```



```
File: mysql-bin.000001
Position: 157
Binlog_Do_DB:
Binlog_Ignore_DB:
Executed_Gtid_Set:
1 row in set (0.01 sec)

mysql> CREATE USER 'repl'@'192.168.56.12' IDENTIFIED BY 'secret';
Query OK, 0 rows affected (0.07 sec)

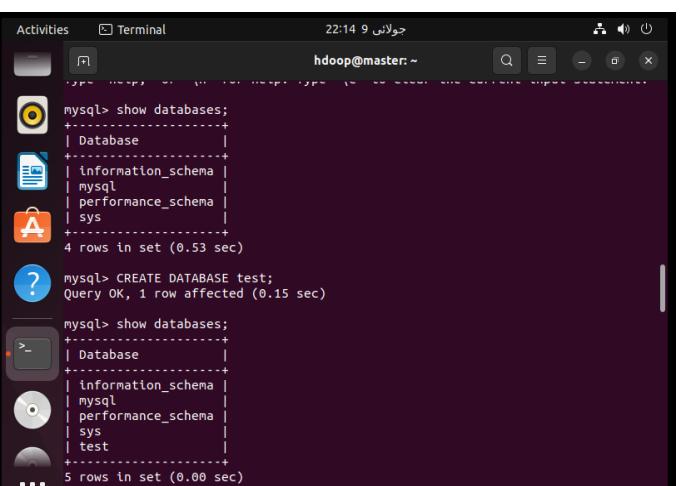
mysql> CREATE USER 'repl'@'192.168.154.134' IDENTIFIED BY 'secret';
Query OK, 0 rows affected (0.02 sec)

mysql> CREATE USER 'repl'@'192.168.154.135' IDENTIFIED BY 'secret';
Query OK, 0 rows affected (0.01 sec)

mysql> GRANT REPLICATION SLAVE ON *.* TO 'repl'@'192.168.154.135';
Query OK, 0 rows affected (0.01 sec)

mysql> SHOW MASTER STATUS \G
***** 1. row *****
    File: mysql-bin.000001
    Position: 1368
    Binlog_Do_DB:
    Binlog_Ignore_DB:
    Executed_Gtid_Set:
1 row in set (0.00 sec)

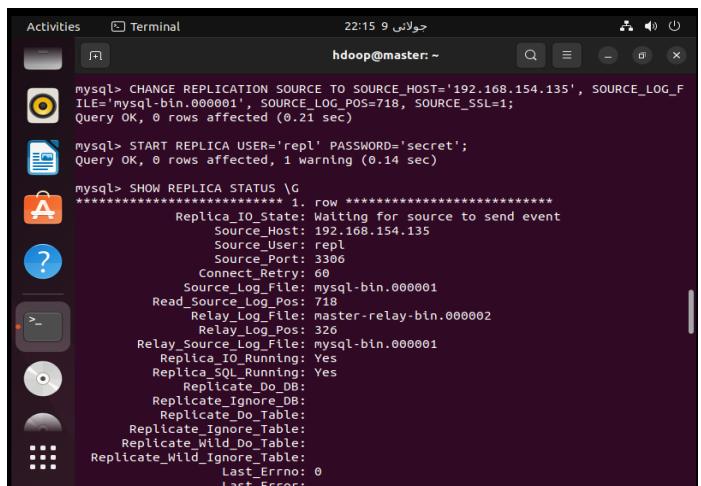
mysql>
```



```
mysql> show databases;
+-----+
| Database |
+-----+
| information_schema |
| mysql |
| performance_schema |
| sys |
+-----+
4 rows in set (0.53 sec)

mysql> CREATE DATABASE test;
Query OK, 1 row affected (0.15 sec)

mysql> show databases;
+-----+
| Database |
+-----+
| information_schema |
| mysql |
| performance_schema |
| sys |
| test |
+-----+
5 rows in set (0.00 sec)
```



```
mysql> CHANGE REPLICATION SOURCE TO SOURCE_HOST='192.168.154.135', SOURCE_LOG_FILE='mysql-bin.000001', SOURCE_LOG_POS=718, SOURCE_SSL=1;
Query OK, 0 rows affected (0.21 sec)

mysql> START REPLICA USER='repl' PASSWORD='secret';
Query OK, 0 rows affected, 1 warning (0.14 sec)

mysql> SHOW REPLICA STATUS \G
***** 1. row *****
Replica_IO_State: Waiting for source to send event
Source_Host: 192.168.154.135
Source_User: repl
Source_Port: 3306
Connect_Retry: 60
Source_Log_File: mysql-bin.000001
Read_Source_Log_Pos: 718
Relay_Log_File: master-relay-bin.000002
Relay_Log_Pos: 326
Relay_Source_Log_File: mysql-bin.000001
Replica_IO_Running: Yes
Replica_SQL_Running: Yes
Replicate_Do_DB:
Replicate_Ignore_DB:
Replicate_Do_Table:
Replicate_Ignore_Table:
Replicate_Wild_Do_Table:
Replicate_Wild_Ignore_Table:
Last_Error: 
Last_Error:
```

```

hadoop@master:~$ corresponds to your MySQL server version for the right syntax to use near 'name varchar(255))' at line 1
mysql>
mysql> CREATE TABLE students (id int, name varchar(255));
Query OK, 0 rows affected (0.22 sec)

mysql> show databases;
+-----+
| Database |
+-----+
| information_schema |
| mysql |
| performance_schema |
| sys |
| test |
| test1 |
+-----+
6 rows in set (0.00 sec)

mysql> show tables;
+-----+
| Tables_in_test1 |
+-----+
| students |
+-----+
1 row in set (0.01 sec)

mysql>

```



```

Activities Terminal 22:23 چلچلہ
hadoop@master:~$ | test1 |
+-----+
6 rows in set (0.00 sec)

mysql> show tables;
+-----+
| Tables_in_test1 |
+-----+
| students |
+-----+
1 row in set (0.01 sec)

mysql> show databases;
+-----+
| Database |
+-----+
| ebadullah |
| information_schema |
| mahnorsaghir |
| mysql |
| numairazaib |
| performance_schema |
| sys |
| test |
| test1 |
+-----+
9 rows in set (0.01 sec)

mysql>

```

SLAVE MACHINE

```

Activities Terminal 21:10 چلچلہ
hadoop@slave:~$ sudo apt install mysql-server
[sudo] password for hadoop:
Reading package lists... Done
Building dependency tree... Done
Reading state information... Done
The following additional packages will be installed:
  libaio1 libcgifast-perl libcgipm-perl libevent-core-2.1-7
  libevent-pthreads-2.1-7 libcgibin libfcgi-perl libfcgi0ldbl
  libhtml-template-perl libmecab2 libprotobuf-lite23 mecab-ipadic
  mecab-ipadic-utf8 mecab-utils mysql-client-8.0 mysql-common
  mysql-server-8.0 mysql-server-core-8.0
Suggested packages:
  libipcs-sharedcache-perl mailx tinyca
The following NEW packages will be installed:
  libaio1 libcgifast-perl libcgipm-perl libevent-core-2.1-7
  libevent-pthreads-2.1-7 libcgibin libfcgi-perl libfcgi0ldbl
  libhtml-template-perl libmecab2 libprotobuf-lite23 mecab-ipadic
  mecab-ipadic-utf8 mecab-utils mysql-client-8.0 mysql-common mysql-server
  mysql-server-8.0 mysql-server-core-8.0
0 upgraded, 19 newly installed, 0 to remove and 141 not upgraded.
Need to get 25.8 MB of archives.
Do you want to continue? [Y/n] y
Get:1 http://pk.archive.ubuntu.com/ubuntu jammy/main amd64 mysql-common all 5.8
+1.0.8 [7,212 B]
Get:2 http://pk.archive.ubuntu.com/ubuntu jammy-updates/main amd64 mysql-client
-8.0 amd64 8.0.29-0ubuntu0.22.04.2 [22.7 kB]
Get:3 http://pk.archive.ubuntu.com/ubuntu jammy/main amd64 libaio1 amd64 0.3.11
-2.13build1 [7,176 B]

Activities Terminal 21:11 چلچلہ
hadoop@slave:~$ emitting matrix : 100% |#####
done!
update-alternatives: using /var/lib/mecab/dic/ipadic-utf8 to provide /var/lib/mecab/dic/debian (mecab-dictionary) in auto mode
Setting up mysql-server-8.0 (8.0.29-0ubuntu0.22.04.2) ...
update-alternatives: using /etc/mysql/mysql.cnf to provide /etc/mysql/my.cnf (my.cnf) in auto mode
Renaming removed key_buffer and myisam-recover options (if present)
mysqld will log errors to /var/log/mysql/error.log
mysqld is running as pid 4616
Created symlink /etc/systemd/system/multi-user.target.wants/mysql.service → /lib/systemd/system/mysql.service.
Setting up mysql-server (8.0.29-0ubuntu0.22.04.2) ...
Processing triggers for man-db (2.10.2-1) ...
Processing triggers for mysql (2.35-0ubuntu3) ...
hadoop@slave:~$ systemctl status mysql
● mysql.service - MySQL Community Server
   Loaded: loaded (/lib/systemd/system/mysql.service; enabled; vendor preset: enabled)
   Active: active (running) since Sat 2022-07-09 21:08:18 PKT; 2min 11s ago
     Process: 4836 ExecStartPre=/usr/share/mysql/mysql-systemd-start pre (code=exec)
   Main PID: 4836 (mysqld)
      Status: "Server is operational"
        Tasks: 37 (limit: 2250)
       Memory: 353.3M
          CPU: 6.495s
         CGroup: /system.slice/mysql.service
                  └─4836 /usr/sbin/mysqld

```

```

Activities Terminal 21:30 چلچلہ
hadoop@slave:~$ sudo vim -r /etc/mysql/mysql.conf.d/mysqld.cnf
[sudo] password for hadoop:
Sorry, try again.
[sudo] password for hadoop:
hadoop@slave:~$ sudo vim -r /etc/mysql/mysql.conf.d/mysqld.cnf
hadoop@slave:~$ sudo vim -r /etc/mysql/mysql.conf.d/mysqld.cnf
hadoop@slave:~$ sudo systemctl restart mysql
hadoop@slave:~$ sudo mysql -u root -p
Enter password:
ERROR 1045 (28000): Access denied for user 'root'@'localhost' (using password: NO)
hadoop@slave:~$ 
hadoop@slave:~$ 
hadoop@slave:~$ sudo mysql -u root -p
Enter password:
Welcome to the MySQL monitor. Commands end with ; or \g.
Your MySQL connection id is 9
Server version: 8.0.29-0ubuntu0.22.04.2 (Ubuntu)

Copyright (c) 2000, 2022, Oracle and/or its affiliates.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql>

```

```

Activities Terminal 21:40 چلچلہ
hadoop@slave:~$ mysql> CHANGE REPLICATION SOURCE TO SOURCE_HOST='192.168.154.134', SOURCE_LOGFILE='mysql-bin.000001', SOURCE_LOG_POS=1368, SOURCE_SSL=1;
Query OK, 0 rows affected (0.10 sec)

mysql> START REPLICA USER='rep1' PASSWORD='secret';
Query OK, 0 rows affected, 1 warning (0.07 sec)

mysql> SHOW REPLICATION STATUS \G
***** 1. row *****
Replica_IO_Status: Connecting to source
Source_Host: 192.168.154.134
Source_User: rep1
Source_Port: 3306
Connect_Retry: 60
Source_Log_File: mysql-bin.000001
Read_Source_Log_Pos: 1368
Relay_Log_File: slave-relay-bin.000001
Relay_Log_Pos: 4
Relay_Source_Log_File: mysql-bin.000001
Replica_IO_Running: Connecting
Replica_SQL_Running: Yes
  Replicate_Do_DB:
  Replicate_Ignore_DB:
  Replicate_Do_Table:
  Replicate_Ignore_Table:
  Replicate_Wild_Do_Table:
  Replicate_Wild_Ignore_Table:
Last_Error: 0
Last_Error:

```

```
Source_Retry_Count: 86400
Source_Bind:
Last_IO_Error_Timestamp:
Last_SQL_Error_Timestamp:
Source_SSL_Crl:
Source_SSL_Crlpath:
Retrieved_Gtid_Set:
Executed_Gtid_Set:
Auto_Position: 0
Replicate_Rewrite_DB:
Channel_Name:
Source_TLS_Version:
Source_public_key_path:
Get_source_public_key: 0
Network_Namespace:
1 row in set (0.01 sec)

mysql> show databases;
+-----+
| Database |
+-----+
| information_schema |
| mysql |
| performance_schema |
| sys |
+-----+
4 rows in set (0.19 sec)

mysql>
```



```
mysql> CREATE DATABASE ebad_ullah;
ERROR 1064 (42000): You have an error in your SQL syntax; check the manual that corresponds to your MySQL server version for the right syntax to use near 'ull ah' at line 1
mysql> CREATE DATABASE numairazaiib;
Query OK, 1 row affected (0.03 sec)
mysql> CREATE DATABASE mahoorsaghir;
Query OK, 1 row affected (0.01 sec)
mysql> show databases;
+-----+
| Database |
+-----+
| ebadullah |
| information_schema |
| mahoorsaghir |
| mysql |
| numairazaiib |
| performance_schema |
| sys |
| test1 |
+-----+
8 rows in set (0.01 sec)

mysql>
```

The benefits of master-slave replication

- 1. Data is more secure:** Made data redundancy, No data loss due to single server downtime
- 2. Improved performance:** One master, many followers, Different users read from different databases, Performance improvement
- 3. Better scalability:** When the flow rate increases, you can easily add from the server, which does not affect the use of the system
- 4. Load balancing:** One master and many slaves share the host task, Load balancing is done.



MAPREDUCE

MapReduce is a programming paradigm that enables massive scalability across hundreds or thousands of servers in a Hadoop cluster.

WORDCOUNT PROGRAM (JAVA)

WordCount is a simple application that counts the number of occurrences of each word in a given input set.

❖ **sudo su - hadoop**

We can't run the HDFS command directly. First of all we have to run all the demons and we have to ensure that all the demons of Hadoop is up and running. Then type ls and you can find out that there's a script start-all.sh. This will start all the demons of Hadoop. Then we check once by typing jps for ensuring that all demons of Hadoop is running. Now we can work with the HDFS file system commands.

❖ **hdfs dfs – help**

It shows all the existing command inside the Hadoop.

❖ **hdfs fsck /**

For checking the health of file system.

❖ **hdfs dfs - ls ~/**

List all the files in a particular path.

❖ **hdfs dfs - ls ~/f1**

For checking what is there inside this particular directory.

❖ **hdfs dfs - mkdir - p ~/test**

For creating a folder test if it is present then it will leave it as it is and it is not present it will create a test folder as we used a prefix p.

❖ **hdfs dfs - ls ~/**

For checking whether the folder is created or not.

❖ **hdfs dfs - rm - r ~/test/**

This command is for recursively remove all the files first and then remove the directory.

❖ **hdfs dfs - touchz ~sample.txt**

This command is used for creating an empty file inside the HDFS.

❖ **hdfs dfs - du - s ~sample.txt**

This command is used to check the size of the file.

❖ **hdfs dfs - appendToFile - ~sample.txt**

This command is used for adding content to the file by using terminal.

❖ **hdfs dfs - cat ~sample.txt**

This command is used to see the contents of the file.

❖ **vim host.txt**

This is the second way of creating a file inside the HDFS

❖ **hdfs dfs - copyFromLocal host.txt ~/**

This command is used for moving a particular file from NFS TO HDFS.

❖ **vim readme.md**

For creating a new file.

❖ **hdfs dfs - put readme.md ~/**

This command is used for recent version of hadoop to move file from NFS to HDFS.

❖ **hdfs dfs - copyToLocal ~/readme.md ~/sample/**

This command is used to move file from HDFS to NFS.

❖ **hdfs dfs - get ~/host.txt ~/sample/**

This is the other way of getting the files from the HDFS to the local file system.

❖ **hdfs dfs - mv ~/host.txt ~/test/**

mv command is going to help us move the contents inside the HDFS File System itself,

Numaira Zaib

The image displays a Linux desktop environment with several open windows, likely a Unity interface. There are four main terminal windows and two file explorer windows.

- Top Left Terminal:** Shows the execution of `./stop-all.sh` and `./start-all.sh` commands on the Hadoop master node.
- Top Right Terminal:** Shows the execution of `hdfs fsck`, which ends with "FSCK ended at Sun Jul 03 20:11:18 PKT 2022 in 88 milliseconds". It also shows the creation of a sample file and its contents being read back.
- Middle Left Terminal:** Shows the creation of a directory `sample` and copying files from the local machine to HDFS using `hdfs dfs -copyToLocal`.
- Middle Right Terminal:** Shows the general command-line syntax for Hadoop commands, including options for files, jars, archives, and specific command-line arguments.
- Bottom Left Terminal:** Shows navigating through the Hadoop share directory to the `mapreduce` folder.
- Bottom Right Terminal:** Shows running the `hadoop jar` command to execute the `mapreduce` examples.
- File Explorers:** Two Nautilus file explorers are visible, showing the file structures on the local machine and in HDFS.
- Bottom Window:** Shows the results of a word count job, displaying the counts for various names.

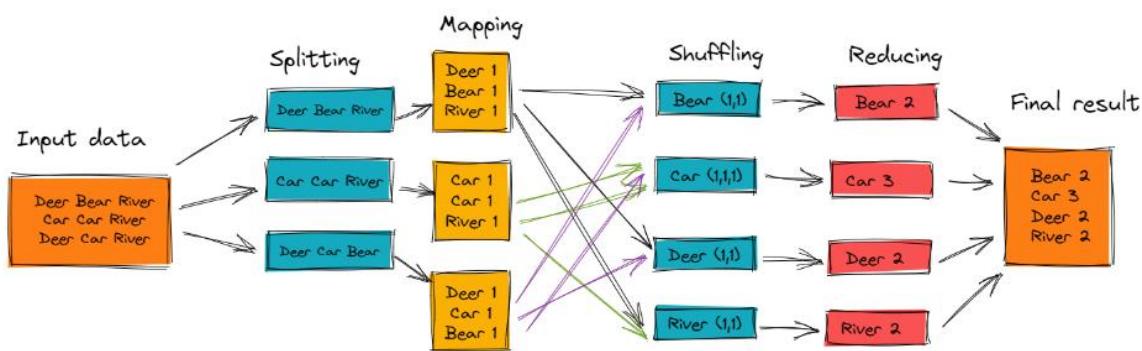
MAPREDUCE WITH PYTHON

MapReduce With Python (WORDCOUNT PROGRAM)

Mapper: It is the first phase of MapReduce application which is used to process a series of key-value pairs.

Shuffler/Combiner: It is the second phase of MapReduce application used to sort, group and shuffle the output coming from the Mapper function. This process of moving output from the mappers to the reducers is known as shuffling.

Reducer: It is the third phase of MapReduce application, it aggregates the values for each unique key and produces zero or more output key-value pairs.



Hadoop Streaming

The mapper and reducer are both executables that read input, line by line, from the standard input (stdin), and write output to the standard output (stdout)

- Mapper
 - Reads the input from stdin line by line.
 - Splits the word separated by tab, space or comma etc.
 - Pairs key-values.
- Reducer
 - Reads the result of mapper from stdout
 - Sums the occurrences of each word
 - Writes the result to stdout

WORD COUNT PROGRAM

mapper.py

First the path is provided then we used sys library. Reading the line from stdin. If there are spaces between the lines that will be removed using strip function then we have only

consecutive lines. From the lines, we will be splitting the words one by one and then for each word we will be assigning the value “1” so this is called peering.

reducer.py

We used sys library, then we initialized 2 variables prev_word and prev_count. Then reading the output of the mapper line by line so it is stripping again. One by one they are all read it from the standard input and they will be assigned to these 2 variables “word” and “count”. Then converting this count by typecasting to integer. Then we put conditions to check if prev_word is equal to word so just increment its counter or else if there is any prev_word, then just print it. Then again initializing the prev_count to the count and prev_word to the word. Then we used one last if statement for checking whether the prev_word is equal to word and if it is, then also print the last word.

Testing locally

To test the python programs locally before running them as a Map-Reduce job.

It is highly recommended to test all programs locally before running them across a Hadoop cluster.

Once the mapper and reducer programs are executing successfully against tests, they can be run as a MapReduce application using the Hadoop streaming utility.

The figure consists of three vertically stacked screenshots of a Linux desktop environment, likely Ubuntu, showing terminal windows. Each terminal window has a dark purple background and a light purple header bar with the text "Activities" and "Terminal". The time at the top of each window is 21:28, 21:29, and 21:30 respectively. The first terminal window shows the command "hadoop@master: \$ ls" followed by a list of files and directories including "hadoop", "mapper.py", "reducer.py", "input.txt", and "output.txt". The second terminal window shows the command "hadoop@master: \$ cat input.txt |python mapper.py" followed by the output of the mapper program, which lists names and the number 1. The third terminal window shows the command "hadoop@master: \$ cat input.txt |python mapper.py |sort |python reducer.py" followed by the final output of the reducer program, which shows the words and their counts.

```
hadoop@master: ~
hadoop@master: ~
hadoop@master: ~
```

```
hadoop@master: $ ls
desktop hadoop-3.2.3 mapper.py readme.md Templates
dfsdta hadoop-3.2.3.tar.gz Music reducer.py tmpdata
Documents host.txt Pictures sample Videos
Downloads input.txt Public snap

hadoop@master: $ cat input.txt
Muhammad Ebad Ullah Khan
Numaira Zaib
Mahnoor saghir
Ebad
Ebad
Numaira
Khan
Zaib
Muhammad
Nemo
Ebad
Numaira
Hadoop
Hadoop
Ebad
Reducemap
WordCount
Ebunk
mahnoor
saghir
mahnoor
saghir
project work

hadoop@master: ~
hadoop@master: ~
hadoop@master: ~
```

```
hadoop@master: $ cat input.txt |python mapper.py
Muhammad 1
Ebad 1
Ullah 1
Khan 1
Numaira 1
Zaib 1
Mahnoor 1
saghir 1
Ebad 1
Ebad 1
Numaira 1
Khan 1
Zaib 1
Muhammad 1
Nemo 1
Ebad 1
Numaira 1
Hadoop 1
Hadoop 1
Ebad 1
Reducemap 1
WordCount 1
Ebunk 1
mahnoor 1
saghir 1
mahnoor 1
saghir 1
project 1

hadoop@master: ~
hadoop@master: ~
hadoop@master: ~
```

```
hadoop@master: $ cat input.txt |python mapper.py |sort |python reducer.py
Hadoop 1
Ebad 1
Reducemap 1
WordCount 1
Ebunk 1
Mahnoor 1
saghir 1
mahnoor 1
saghir 1
project 1
work 1
hadoop@master: ~
hadoop@master: ~
hadoop@master: ~
```

RPC IN HADOOP

It allows one computer program to remotely call another computer's subprogram without having to care about the details of the underlying network communication, which is transparent to us. Therefore, it is often used in distributed network communications.

Hadoop's inter-process interaction is carried out through RPC, such as the direct connection between Namenode and Datanode, and between Jobtracker and Tasktracker.

Therefore, it can be said that the operation of Hadoop is based on RPC . Hadoop RPC is an important part of Hadoop, which provides object call functions in a distributed environment. The source code is in org.apache.hadoop.ipc. HBase has almost completely copied the source code of this part, but has changed the configuration items.

RPC mechanism in Hadoop

Like other RPC frameworks, Hadoop RPC is divided into four parts:

Serialization layer: The information communicated between Client and Server uses the serialization class provided in Hadoop or the custom Writable type

Function call layer: Hadoop RPC realizes function call through dynamic proxy and java reflection

Network transmission layer: Hadoop RPC uses a socket mechanism based on TCP/IP

Server-side framework layer: RPC Server uses java NIO and an event-driven I/O model to improve the concurrent processing capabilities of RPC Server

Hadoop RPC is widely used throughout Hadoop, and the communication between Client, DataNode, and NameNode depends on it. For example: When we usually operate HDFS, we use the FileSystem class, which has a DFSClient object inside, which is responsible for dealing with the NameNode. At runtime, DFSClient creates a NameNode agent locally, and then operates the agent. The agent will call the method of the NameNode remotely through the network and return the value.

TRANSPARENCY IN HADOOP

HDFS implements transparent, end-to-end encryption. Once configured, data read from and written to special HDFS directories is transparently encrypted and decrypted without requiring changes to user application code. This encryption is also end-to-end, which means the data can only be encrypted and decrypted by the client. HDFS never stores or has access to unencrypted data or unencrypted data encryption keys. This satisfies two typical requirements for encryption: at-rest encryption (meaning data on persistent media, such as a disk) as well as in-transit encryption (e.g. when data is travelling over the network).

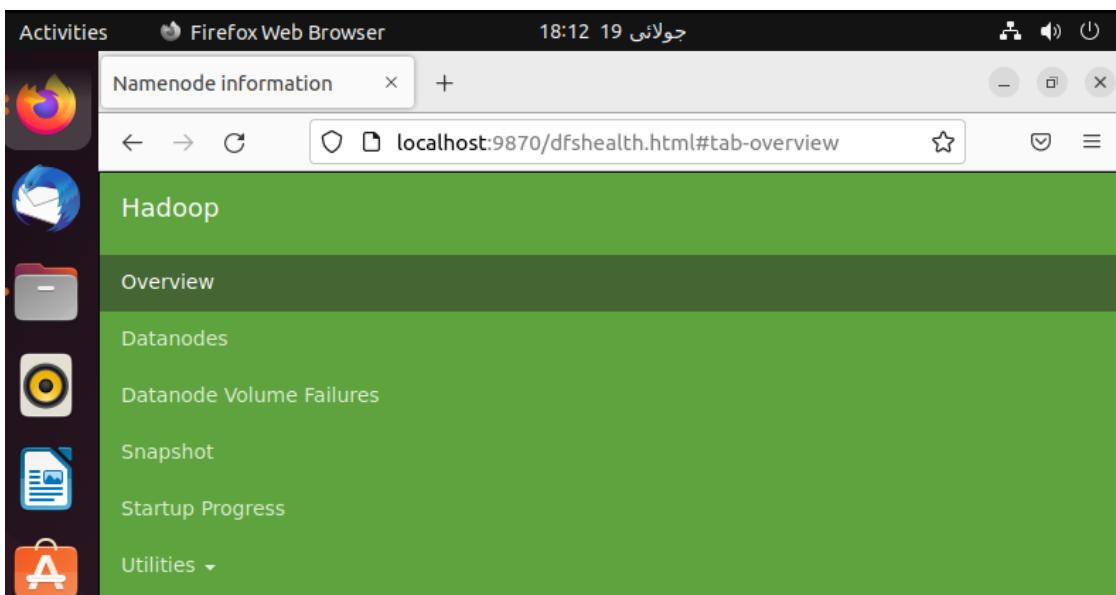
Encryption can be done at different layers in a traditional data management software/hardware stack. Choosing to encrypt at a given layer comes with different advantages and disadvantages.

- **Application-level encryption.** This is the most secure and most flexible approach. The application has ultimate control over what is encrypted and can precisely reflect the requirements of the user. However, writing applications to do this is hard. This is also not an option for customers of existing applications that do not support encryption.
- **Database-level encryption.** Similar to application-level encryption in terms of its properties. Most database vendors offer some form of encryption. However, there can be performance issues. One example is that indexes cannot be encrypted.
- **Filesystem-level encryption.** This option offers high performance, application transparency, and is typically easy to deploy. However, it is unable to model some application-level policies. For instance, multi-tenant applications might want to encrypt based on the end user. A database might want different encryption settings for each column stored within a single file.
- **Disk-level encryption.** Easy to deploy and high performance, but also quite inflexible. Only really protects against physical theft.

HDFS-level encryption fits between database-level and filesystem-level encryption in this stack. This has a lot of positive effects. HDFS encryption is able to provide good performance and existing Hadoop applications are able to run transparently on encrypted data. HDFS also has more context than traditional filesystems when it comes to making policy decisions.

HDFS-level encryption also prevents attacks at the filesystem-level and below (so-called “OS-level attacks”). The operating system and disk only interact with encrypted bytes, since the data is already encrypted by HDFS.

HADOOP WEB INTERFACE



Activities Firefox Web Browser چولانی 19 18:13

Namenode information +

localhost:9870/dfshealth.html#tab-overview ☆ ☰ ☱ ☳

Overview 'localhost:9000' (active)

Started:	Tue Jul 19 17:59:15 +0500 2022
Version:	3.2.3, rabe5358143720085498613d399be3bbf01e0f131
Compiled:	Sun Mar 20 06:18:00 +0500 2022 by ubuntu from branch-3.2.3
Cluster ID:	CID-bd2951b9-d84c-4d58-8258-b76b297101a2
Block Pool ID:	BP-644048595-192.168.154.134-1658235522900

Activities Firefox Web Browser چولانی 19 18:14

Namenode information +

localhost:9870/dfshealth.html#tab-overview ☆ ☰ ☱ ☳

Summary

Security is off.

Safemode is off.

1 files and directories, 0 blocks (0 replicated blocks, 0 erasure coded block groups) = 1 total filesystem object(s).

Heap Memory used 50.49 MB of 144 MB Heap Memory. Max Heap Memory is 432 MB.

Non Heap Memory used 46.54 MB of 47.7 MB Committed Non Heap Memory. Max Non Heap Memory is <unbounded>.

Configured Capacity:	0 B
Configured Remote Capacity:	0 B
DFS Used:	0 B (100%)
Non DFS Used:	0 B
DFS Remaining:	0 B (0%)

Activities Firefox Web Browser چولانی 19 18:15

Namenode information +

localhost:9870/dfshealth.html#tab-overview ☆ ☰ ☱ ☳

Block Pool Used:	0 B (100%)
DataNodes usages% (Min/Median/Max/stdDev):	0.00% / 0.00% / 0.00% / 0.00%
Live Nodes	0 (Decommissioned: 0, In Maintenance: 0)
Dead Nodes	0 (Decommissioned: 0, In Maintenance: 0)
Decommissioning Nodes	0
Entering Maintenance Nodes	0
Total Datanode Volume Failures	0 (0 B)
Number of Under-Replicated Blocks	0
Number of Blocks Pending Deletion (including replicas)	0
Block Deletion Start Time	Tue Jul 19 17:59:15 +0500 2022
Last Checkpoint Time	Tue Jul 19 17:58:43 +0500 2022

The image contains three screenshots of a Linux desktop environment, likely Ubuntu, showing Hadoop NameNode status and log files.

- Screenshot 1: NameNode Journal Status**
Shows the current transaction ID (1) and a table of Journal Managers. One entry is for a FileJournalManager with root at /home/hadoop/tmpdata/dfs/name.
- Screenshot 2: NameNode Storage**
Shows the NameNode Storage configuration with one storage directory at /home/hadoop/tmpdata/dfs/name, which is an IMAGE_AND_EDITS type and active.
- Screenshot 3: DFS Storage Types**
Shows the DFS Storage Types table with one entry for Hadoop, dated 2022.
- Screenshot 4: Directory: /logs/**
Shows a file listing for the /logs directory. The table has columns for Name, Last Modified, and Size. The logs are named in a specific pattern related to Hadoop components like datanode and numaira-virtual-machine.

CONCLUSION

Hadoop is playing a significant role in today's life and can be fitted with any domain as per requirements we can adopt this technology in our organization as per our business requirements. Hadoop has been a very effective solution for companies dealing with data in petabytes. It has solved many problems in the industry related to huge data management and distributed systems.