

Exp.1 Downloading and installing Hadoop, Understanding different Hadoop modes, Startup scripts, Configuration files.

AIM:

To Download and install Hadoop, Understanding different Hadoop modes, Startup scripts, Configuration files.

Procedure:

Step 1 : Install Java Development Kit

The default Ubuntu repositories contain Java 8 and Java 11 both. But, Install Java 8 because it only works on this version. Use the following command to install it.

```
$sudo apt update&&sudo apt install openjdk-8-jdk
```

Step 2 : Verify the Java version

Once installed, verify the installed version of Java with the following command:

```
$ java -version
```

Step 3: Install SSH

SSH (Secure Shell) installation is vital for Hadoop as it enables secure communication between nodes in the Hadoop cluster. This ensures data integrity, confidentiality, and allows for efficient distributed processing of data across the cluster.

```
$sudo apt install ssh
```

Step 4 : Create the hadoop user :

All the Hadoop components will run as the user that you create for Apache Hadoop, and the user will also be used for logging in to Hadoop's web interface.

Run the command to create user and set password:

```
$ sudo adduser hadoop
```

Step 5 : Switch user

Switch to the newly created hadoop user:

```
$ su - hadoop
```

Step 6 : Configure SSH

Now configure password-less SSH access for the newly created hadoop user, so didn't enter the key to save file and passphrase. Generate an SSH keypair (generate Public and Private Key Pairs)first

```
$ssh-keygen -t rsa
```

Step 7 : Set permissions :

Next, append the generated public keys from id_rsa.pub to authorized_keys and set proper permission:

```
$ cat ~/.ssh/id_rsa.pub >> ~/.ssh/authorized_keys
```

```
$ chmod 640 ~/.ssh/authorized_keys
```

Step 8 : SSH to the localhost

Next, verify the password less SSH authentication with the following command:

\$ ssh localhost

You will be asked to authenticate hosts by adding RSA keys to known hosts. Type yes and hit Enter to authenticate the localhost:

```

(kali@kali)-[~]
└─$ su hadoop
Password:
(hadoop@kali)-[/home/kali]
└─$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as hadoop in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
localhost: ssh: connect to host localhost port 22: Connection refused
Starting datanodes
localhost: ssh: connect to host localhost port 22: Connection refused
Starting secondary namenodes [kali]
kali: ssh: connect to host kali port 22: Connection refused
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
Starting resourcemanager
Starting nodemanagers
localhost: ssh: connect to host localhost port 22: Connection refused

(hadoop@kali)-[/home/kali]
└─$ ssh localhost
ssh: connect to host localhost port 22: Connection refused

(hadoop@kali)-[/home/kali]
└─$ sudo service restart ssh
[sudo] password for hadoop:
restart: unrecognized service

(hadoop@kali)-[/home/kali]
└─$ sudo service ssh start

(hadoop@kali)-[/home/kali]
└─$ ssh localhost
Linux kali 6.5.0-kali3-amd64 #1 SMP PREEMPT_DYNAMIC Debian 6.5.6-1kali1 (2023-10-09) x86_64
The programs included with the Kali GNU/Linux system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.
Kali GNU/Linux comes with ABSOLUTELY NO WARRANTY, to the extent
permitted by applicable law.
Last login: Wed Aug 28 05:26:53 2024 from ::1

```

Step 9 : Switch user

Again switch to hadoop. So, First, change the user to hadoop with the following command:

\$ su-hadoop

Step 10 : Install hadoop

Next, download the latest version of Hadoop using the wget command:

\$ wget <https://downloads.apache.org/hadoop/common/hadoop-3.3.6/hadoop-3.3.6.tar.gz>

Once downloaded, extract the downloaded file:

\$ tar -xvzf hadoop-3.3.6.tar.gz

Next, rename the extracted directory to hadoop:

empty before retry.

\$ nano ~/.bashrc

Append the below lines to file.

```

export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
export HADOOP_HOME=/home/hadoop/hadoop
export HADOOP_INSTALL=$HADOOP_HOME
export HADOOP_MAPRED_HOME=$HADOOP_HOME
export HADOOP_COMMON_HOME=$HADOOP_HOME
export HADOOP_HDFS_HOME=$HADOOP_HOME
export HADOOP_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"

```

Save and close the file. Then, activate the environment variables with the following command:

s\$ source ~/.bashrc

Next, open the Hadoop environment variable file:

\$ nano \$HADOOP_HOME/etc/hadoop/hadoop-env.sh

Search for the “export JAVA_HOME” and configure it.

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

The screenshot shows the nano text editor interface with the file `/home/hadoop/hadoop/etc/hadoop/hadoop-env.sh` open. The editor displays various configuration comments and settings. The line `JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64` is highlighted. At the bottom, the 'File Name to Write' field shows the full path to the file. The status bar at the very bottom lists keyboard shortcuts for Help, Cancel, DOS Format, Mac Format, Append, Prepend, Backup File, and Browse.

```

File Edit View Search Terminal Help
GNU nano 7.2 /home/hadoop/hadoop/etc/hadoop/hadoop-env.sh *
##
## Precedence rules:
##
## (yarn-env.sh|hdfs-env.sh) > hadoop-env.sh > hard-coded defaults
##
## (YARN_xyz|HDFS_xyz) > HADOOP_xyz > hard-coded defaults
##
# Many of the options here are built from the perspective that users
# may want to provide OVERWRITING values on the command line.
# For example:
#
# JAVA_HOME=/usr/lib/jvm/java-11-openjdk-amd64
#
# Therefore, the vast majority (BUT NOT ALL!) of these defaults
# are configured for substitution and not append. If append
# is preferable, modify this file accordingly.
###
# Generic settings for HADOOP
###
# Technically, the only required environment variable is JAVA_HOME.
# All others are optional. However, the defaults are probably not
# preferred. Many sites configure these options outside of Hadoop,
# such as in /etc/profile.d
#
# The java implementation to use. By default, this environment
# variable is REQUIRED on ALL platforms except OS X!
File Name to Write: /home/hadoop/hadoop/etc/hadoop/hadoop-env.sh
^C Help      ^M-D DOS Format  ^M-A Append      ^M-B Backup File
^C Cancel    ^M-M Mac Format   ^M-P Prepend     ^M-T Browse

```

Save and close the file when you are finished.

Step 11 : Configuring Hadoop :

First, you will need to create the namenode and datanode directories inside the Hadoop user home directory. Run the following command to create both directories:

```
$ cd hadoop/
```

```
$mkdir -p ~/hadoopdata/hdfs/{namenode,datanode}
```

```
$nano $HADOOP_HOME/etc/hadoop/core-site.xml
```

Change the following name as per your system hostname:

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

Save and close the file.

Then, edit the hdfs-site.xml file:

```
$nano $HADOOP_HOME/etc/hadoop/hdfs-site.xml
```

- Change the NameNode and DataNode directory paths as shown below:

```
<configuration>
  <property>
    <name>dfs.replication</name>
    <value>1</value>
  </property>

  <property>
    <name>dfs.namenode.name.dir</name>
    <value>file:///home/hadoop/hadoopdata/hdfs/namenode</value>
  </property>

  <property>
    <name>dfs.datanode.data.dir</name>
    <value>file:///home/hadoop/hadoopdata/hdfs/datanode</value>
  </property>
</configuration>
```

- Then, edit the mapred-site.xml file:

```
$nano $HADOOP_HOME/etc/hadoop/mapred-site.xml
```

- Make the following changes:

```
<configuration>
  <property>
    <name>yarn.app.mapreduce.am.env</name>
    <value>HADOOP_MAPRED_HOME=$HADOOP_HOME/home/hadoop/hadoop/bin/hadoop</value>
  </property>
  <property>
    <name>mapreduce.map.env</name>
    <value>HADOOP_MAPRED_HOME=$HADOOP_HOME/home/hadoop/hadoop/bin/hadoop</value>
  </property>
  <property>
    <name>mapreduce.reduce.env</name>
    <value>HADOOP_MAPRED_HOME=$HADOOP_HOME/home/hadoop/hadoop/bin/hadoop</value>
  </property>
</configuration>
```

- Then, edit the yarn-site.xml file:
\$nano \$HADOOP_HOME/etc/hadoop/yarn-site.xml
- Make the following changes:

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

Save the file and close it .

Step 12 – Start Hadoop Cluster

Before starting the Hadoop cluster. You will need to format the Namenode as a hadoop user.

Run the following command to format the Hadoop Namenode:

```
$hdfs namenode -format
```

Once the namenode directory is successfully formatted with hdfs file system, you will see the message “Storage directory /home/hadoop/hadoopdata/hdfs/namenode has been successfully formatted “

Then start the Hadoop cluster with the following command.

```
$ start-all.sh
```

```
hadoop@osboxes:~$ start-all.sh
WARNING: Attempting to start all Apache Hadoop daemons as hadoop in 10 seconds.
WARNING: This is not a recommended production deployment configuration.
WARNING: Use CTRL-C to abort.
Starting namenodes on [localhost]
localhost: namenode is running as process 8126. Stop it first and ensure /tmp/hadoop-hadoop-namenode.pid file is empty
before retry.
Starting datanodes
localhost: datanode is running as process 8271. Stop it first and ensure /tmp/hadoop-hadoop-datanode.pid file is empty
before retry.
Starting secondary namenodes [osboxes]
osboxes: secondarynamenode is running as process 8462. Stop it first and ensure /tmp/hadoop-hadoop-secondarynamenode.p
id file is empty before retry.
Starting resource manager
resource manager is running as process 8728. Stop it first and ensure /tmp/hadoop-hadoop-resource manager.pid file is em
pty before retry.
Starting nodemanagers
localhost: nodemanager is running as process 8868. Stop it first and ensure /tmp/hadoop-hadoop-nodemanager.pid file is
empty before retry.
```

You can now check the status of all Hadoop services using the jps command:

\$ jps

```
hadoop@osboxes:~$ jps
10898 Jps
8868 NodeManager
8728 ResourceManager
9212 RunJar
8126 NameNode
8462 SecondaryNameNode
8271 DataNode
```

Step 13 – Access Hadoop Namenode and Resource Manager

- First we need to know our ipaddress, In Ubuntu we need to install net-tools to run ipconfig command,
If you installing net-tools for the first time switch to default user:
\$sudo apt install net-tools
- Then run ifconfig command to know our ip address:

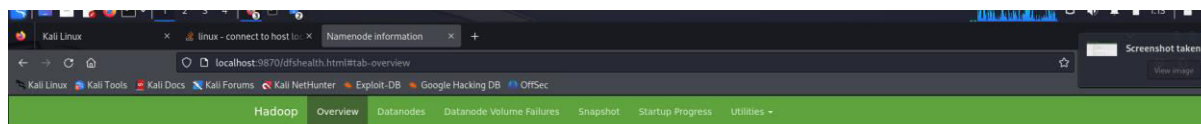
```
hadoop@osboxes:~$ ifconfig
enp0s3: flags=4163<UP,BROADCAST,RUNNING,MULTICAST> mtu 1500
    inet 10.0.2.15 netmask 255.255.255.0 broadcast 10.0.2.255
    inet6 fe80::a00:27ff:fe56:6fa2 prefixlen 64 scopeid 0x20<link>
    ether 08:00:27:56:6f:a2 txqueuelen 1000 (Ethernet)
    RX packets 458730 bytes 616003488 (616.0 MB)
    RX errors 0 dropped 0 overruns 0 frame 0
    TX packets 171067 bytes 18436556 (18.4 MB)
    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 21342 bytes 2170633 (2.1 MB)
    RX errors 0 dropped 0 overruns 0 frame 0
    TX packets 21342 bytes 2170633 (2.1 MB)
    TX errors 0 dropped 0 overruns 0 carrier 0 collisions 0
```

ifconfig

Here my ip address is 192.168.1.6.

- To access the Namenode, open your web browser and visit the URL <http://your-server-ip:9870>.
- You should see the following screen:
<http://192.168.1.6:9870>



Overview 'localhost:9000' (✓active)

Started:	Wed Aug 28 05:28:01 -0400 2024
Version:	3.3.6, r1be78238728da9266a4f88195058f08fd012bf9c
Compiled:	Sun Jun 18 04:22:00 -0400 2023 by ubuntu from (HEAD detached at release-3.3.6-RC1)
Cluster ID:	CID-9ba41419-c275-46b1-aa5e-d4f1e24435af
Block Pool ID:	BP-127595427-127.0.1.1-1724225128501

Summary

Security is off.

Safemode is off.

44 files and directories, 19 blocks (19 replicated blocks, 0 erasure coded block groups) = 63 total filesystem object(s).

Heap Memory used 64.72 MB of 96 MB Heap Memory. Max Heap Memory is 492 MB.

Non Heap Memory used 51.85 MB of 55.44 MB Committed Non Heap Memory. Max Non Heap Memory is <unbounded>.

Configured Capacity:	78.28 GB
Configured Remote Capacity:	0 B
DFS Used:	13.69 MB (0.02%)
Non DFS Used:	17.33 GB

To access Resource Manager, open your web browser and visit the URL <http://your-server-ip:8088>. You should see the following screen:

<http://192.168.16:8088>

Cluster Metrics

Apps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Used Resources	Total Resources
0	0	0	0	0	<memory:0 B, vCores:0>	<memory:8 GB, vCores:8>

Cluster Nodes Metrics

Active Nodes	Decommissioning Nodes	Decommissioned Nodes	Lost Nodes	Unhealthy
1	0	0	0	0

Scheduler Metrics

Scheduler Type	Scheduling Resource Type	Minimum Allocation	Maximum Allocation
Capacity Scheduler	[memory-mb (unit=Mi), vCores]	<memory:1024, vCores:1>	<memory:8192, vCores:4>

Show 20 entries

ID	User	Name	Application Type	Application Tags	Queue	Application Priority	StartTime	LaunchTime	FinishTime	State	FinalStatus	Running Containers	Allocated CPU VCores	Allocated Memory MB	Allocated GPUs
No data available in table															

Showing 0 to 0 of 0 entries

Step 14 – Verify the Hadoop Cluster

At this point, the Hadoop cluster is installed and configured. Next, we will create some directories in the HDFS filesystem to test the Hadoop.

Let's create some directories in the HDFS filesystem using the following command:

```
$ hdfsdfs -mkdir /test1
$ hdfsdfs -mkdir /logs
```

Next, run the following command to list the above directory:

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

You should get the following output:

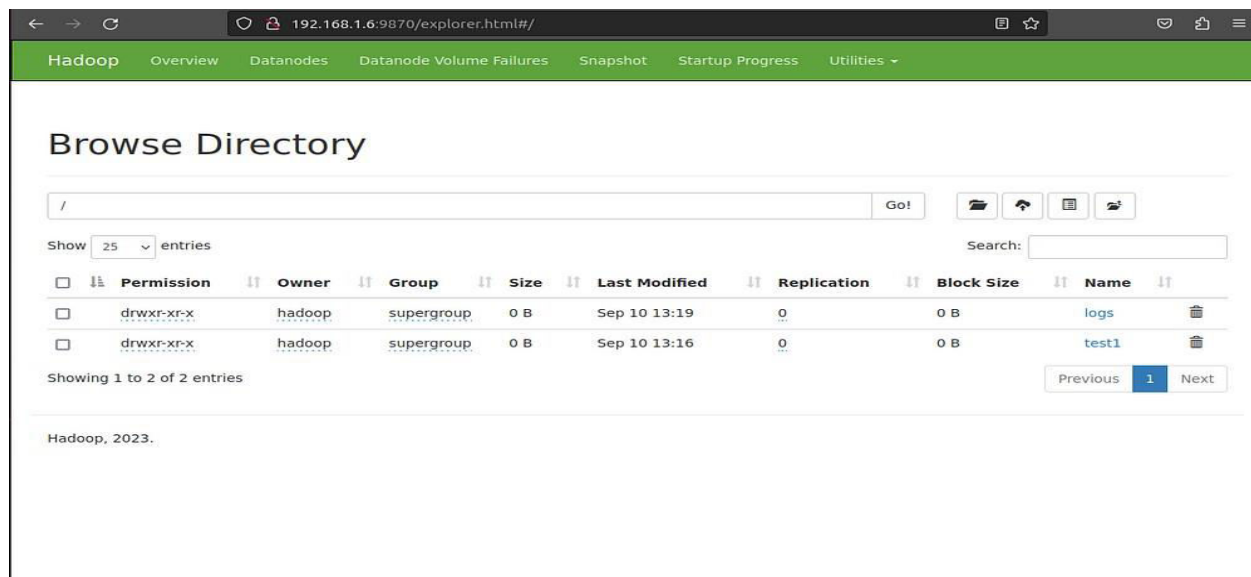
```
hadoop@osboxes:~$ hdfs dfs -ls /
Found 10 items
drwxr-xr-x  - hadoop supergroup      0 2024-09-18 03:06 /hive
drwxr-xr-x  - hadoop supergroup      0 2024-08-31 01:08 /logs
drwxr-xr-x  - hadoop supergroup      0 2024-09-02 03:27 /new_output
drwxr-xr-x  - hadoop supergroup      0 2024-09-04 12:09 /pig_output_data
drwxr-xr-x  - hadoop supergroup      0 2024-09-04 11:46 /piginput
drwxr-xr-x  - hadoop supergroup      0 2024-09-18 05:23 /tmp
drwxr-xr-x  - hadoop supergroup      0 2024-09-04 12:09 /udfs
drwxr-xr-x  - hadoop supergroup      0 2024-09-18 03:10 /user
drwxr-xr-x  - hadoop supergroup      0 2024-09-02 02:25 /weatherdata
drwxr-xr-x  - hadoop supergroup      0 2024-09-01 10:27 /word_count_in_python
hadoop@osboxes:~$
```

Also, put some files to hadoop file system. For the example, putting log files from host machine to hadoop file system.

```
$ hdfs dfs -put /var/log/* /logs/
```

You can also verify the above files and directory in the Hadoop Namenode web interface.

Go to the web interface, click on the Utilities => Browse the file system. You should see your directories which you have created earlier in the following screen:



Step 15 – Stop Hadoop Cluster

To stop the Hadoop all services, run the following command:

```
$ stop-all.sh
```



```
hadoop@osboxes:~$ stop-all.sh
WARNING: Stopping all Apache Hadoop daemons as hadoop in 10 seconds.
WARNING: Use CTRL-C to abort.
Stopping namenodes on [localhost]
Stopping datanodes
Stopping secondary namenodes [osboxes]
Stopping nodemanagers
Stopping resourcemanager
```

Result:

The step-by-step installation and configuration of Hadoop on Ubuntu linux system have been successfully completed.

EXP 2: Run a basic Word Count Map Reduce program to understand Map Reduce Paradigm.

AIM:

To run a basic Word Count MapReduce program.

Procedure:

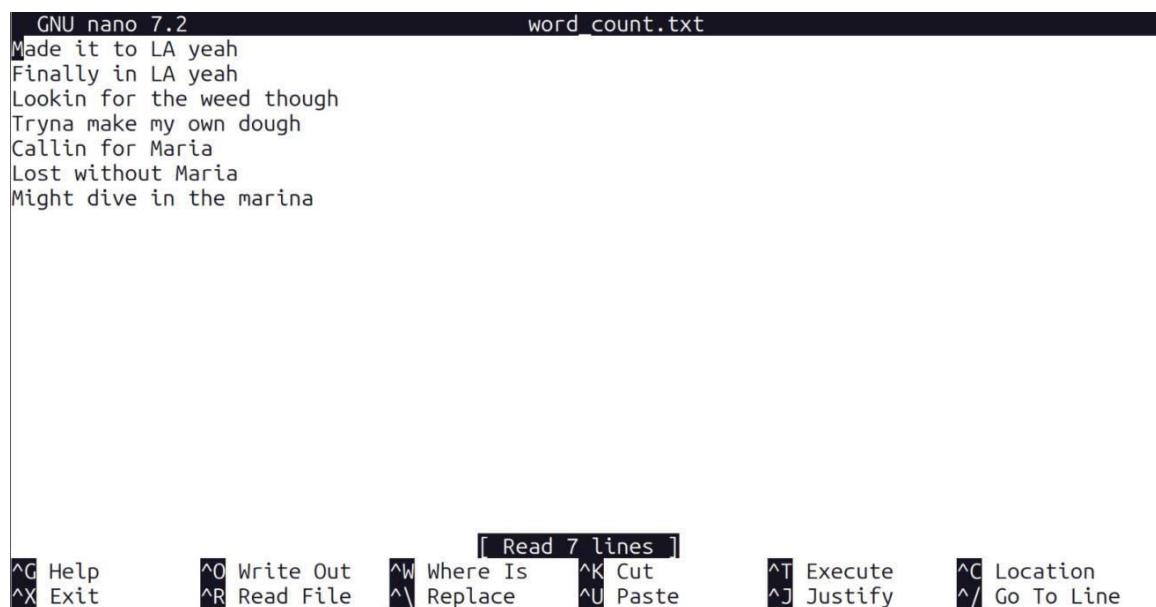
Step 1: Create Data File:

Create a file named "word_count_data.txt" and populate it with text data that you wish to analyse.

Login with your hadoop user.

```
nano word_count.txt
```

Output: Type the below content in word_count.txt



```
GNU nano 7.2 word_count.txt
Made it to LA yeah
Finally in LA yeah
Lookin for the weed though
Tryna make my own dough
Callin for Maria
Lost without Maria
Might dive in the marina

[ Read 7 lines ]
^G Help      ^O Write Out  ^W Where Is   ^K Cut        ^T Execute    ^C Location
^X Exit      ^R Read File  ^\ Replace    ^U Paste      ^J Justify    ^_ Go To Line
```

Step 2: Mapper Logic - mapper.py:

Create a file named "mapper.py" to implement the logic for the mapper. The mapper will read input data from STDIN, split lines into words, and output each word with its count.

```
nano mapper.py
```

```
# Copy and paste the mapper.py code
```

```
#!/usr/bin/env python3
```

```
# import sys because we need to read and write data to STDIN and STDOUT
```

```
#!/usr/bin/python3
```

```
import sys
```

```
for line in sys.stdin:
```

```
    line = line.strip() # remove leading and trailing whitespace
```

```
    words = line.split() # split the line into words
```

```
    for word in words:
```

```
print( '%s\t%s' % (word, 1))
```

Step 3: Reducer Logic - reducer.py:

Create a file named "reducer.py" to implement the logic for the reducer. The reducer will aggregate the occurrences of each word and generate the final output.

```
nano reducer.py
# Copy and paste the reducer.py code
```

reducer.py

```
#!/usr/bin/python3
from operator import itemgetter
import sys
current_word = None
current_count = 0
word = None
for line in sys.stdin:
    line = line.strip()
    word, count = line.split('\t', 1)
    try:
        count = int(count)
    except ValueError:
        continue
    if current_word == word:
        current_count += count
    else:
        if current_word:
            print( '%s\t%s' % (current_word, current_count))
            current_count = count
            current_word = word
        if current_word == word:
            print( '%s\t%s' % (current_word, current_count))
```

Step 4: Prepare Hadoop Environment:

Start the Hadoop daemons and create a directory in HDFS to store your data.

```
start-all.sh
hdfsdfs -mkdir /word_count_in_python
hdfsdfs -copyFromLocal /path/to/word_count.txt/word_count_in_python
```

Step 6: Make Python Files Executable:

Give executable permissions to your mapper.py and reducer.py files.

```
chmod 777 mapper.py reducer.py
```

Step 7: Run Word Count using Hadoop Streaming:

Download the latest hadoop-streaming jar file and place it in a location you can easily access.

Then run the Word Count program using Hadoop Streaming.

```
hadoop jar /path/to/hadoop-streaming-3.3.6.jar \
-input /word_count_in_python/word_count_data.txt \
-output /word_count_in_python/new_output \
-mapper /path/to/mapper.py \
-reducer /path/to/reducer.py
```

```
(hadoop@kali) [~/hadoop/share/hadoop/tools/lib]
$ ls -l hadoop-streaming-3.3.6.jar
-rw-r--r-- 1 hadoop-streaming 141267 Jun 18 2023 hadoop-streaming-3.3.6.jar

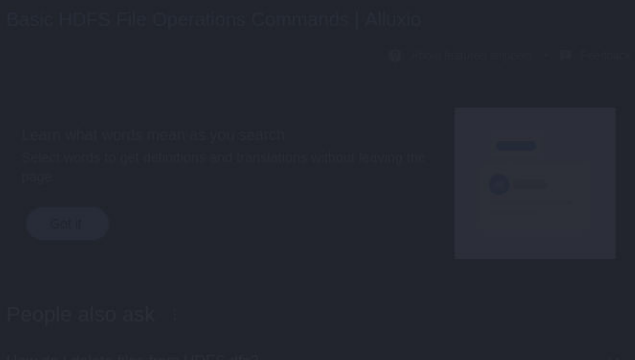
(hadoop@kali) [~/hadoop/share/hadoop/tools/lib]
$ hadoop jar hadoop-streaming-3.3.6.jar -input /word_count_in_python/wordcount.txt -output /word_count_in_python/new_output -mapper mapper.py -reducer reducer.py
2024-08-28 11:56:40,181 INFO impl.MetricsConfig: Loaded properties from hadoop-metrics2.properties
2024-08-28 11:56:40,352 INFO impl.MetricsSystemImpl: Scheduled Metric snapshot period at 10 second(s).
2024-08-28 11:56:40,353 INFO impl.MetricsSystemImpl: JobTracker metrics system started
2024-08-28 11:56:40,366 WARN impl.MetricsSystemImpl: JobTracker metrics system already initialized!
2024-08-28 11:56:40,966 INFO mapred.FileInputFormat: Total input files to process = 1
2024-08-28 11:56:41,605 INFO mapreduce.JobSubmitter: number of splits=1
2024-08-28 11:56:42,164 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_local427539628_0001
2024-08-28 11:56:42,165 INFO mapreduce.JobSubmitter: Executing with tokens: []
2024-08-28 11:56:42,405 INFO mapreduce.Job: The url to track the job: http://localhost:8080/
2024-08-28 11:56:42,418 INFO mapreduce.Job: Running job: job_local427539628_0001
2024-08-28 11:56:42,418 INFO mapred.LocalJobRunner: OutputCommitter set in config null
2024-08-28 11:56:42,421 INFO mapred.LocalJobRunner: OutputCommitter is org.apache.hadoop.mapred.FileOutputCommitter
2024-08-28 11:56:42,429 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2024-08-28 11:56:42,429 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup_temporary folders under output directory:false, ignore cleanup failures: false
2024-08-28 11:56:42,601 INFO mapred.LocalJobRunner: Waiting for map tasks
2024-08-28 11:56:42,612 INFO mapred.LocalJobRunner: Starting task: attempt_local427539628_0001_m_000000_0
2024-08-28 11:56:42,711 INFO output.FileOutputCommitter: File Output Committer Algorithm version is 2
2024-08-28 11:56:42,716 INFO output.FileOutputCommitter: FileOutputCommitter skip cleanup_temporary folders under output directory:false, ignore cleanup failures: false
2024-08-28 11:56:42,750 INFO mapred.Task: Using ResourceCalculatorProcessTree : [ ]
2024-08-28 11:56:42,785 INFO mapred.MapTask: Processing split: hdfs://localhost:9000/word_count_in_python/wordcount.txt:0+150
2024-08-28 11:56:42,845 INFO mapred.MapTask: numReduceTasks: 1
2024-08-28 11:56:43,049 INFO mapred.MapTask: (EQUATOR) 0 kvi 26214396(104857584)
2024-08-28 11:56:43,052 INFO mapred.MapTask: mapreduce.task.io.sort.mb: 100
2024-08-28 11:56:43,052 INFO mapred.MapTask: soft limit at 83886080
2024-08-28 11:56:43,052 INFO mapred.MapTask: bufstart = 0; bufvoid = 104857600
2024-08-28 11:56:43,052 INFO mapred.MapTask: kvstart = 26214396; length = 6553600
2024-08-28 11:56:43,075 INFO mapred.MapTask: Map output collector class = org.apache.hadoop.mapred.MapTask$MapOutputBuffer
2024-08-28 11:56:43,184 INFO streaming.PipeMapRed: PipeMapRed exec [mapper.py]
2024-08-28 11:56:43,386 INFO Configuration.deprecation: mapred.work.output.dir is deprecated. Instead, use mapreduce.task.output.dir
2024-08-28 11:56:43,319 INFO Configuration.deprecation: mapred.local.dir is deprecated. Instead, use mapreduce.cluster.local.dir
2024-08-28 11:56:43,320 INFO Configuration.deprecation: map.input.file is deprecated. Instead, use mapreduce.map.input.file
2024-08-28 11:56:43,321 INFO Configuration.deprecation: map.input.length is deprecated. Instead, use mapreduce.map.input.length
2024-08-28 11:56:43,322 INFO Configuration.deprecation: mapred.job.id is deprecated. Instead, use mapreduce.job.id
2024-08-28 11:56:43,324 INFO Configuration.deprecation: mapred.task.partition is deprecated. Instead, use mapreduce.task.partition
2024-08-28 11:56:43,329 INFO Configuration.deprecation: map.input.start is deprecated. Instead, use mapreduce.map.input.start
```

Step 8: Check Output:

Check the output of the Word Count program in the specified HDFS output directory.

```
hdfs dfs -cat /word_count_in_python/new_output/part-00000
```

```
(hadoop@kali)-[~]
$ hdfs dfs -cat /word_count_in_python/new_output/part-00000
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
Callin 1
Finally 1
LA 1
Lookin 1
Lost 1
Made 1
Marla 1
Might 1
Tryna 1
dive 1
dough 1
for 1
in 1
it 1
make 1
marina 1
my 1
own 1
the 1
though 1
to 1
weed 1
without 1
yeah 1
```

**Result:**

Thus, the program for basic Word Count Map Reduce has been executed successfully.

EXP 3: Map Reduce program to process a weather dataset.**AIM:**

To implement MapReduce program to process a weather dataset.

Procedure:**Step 1: Create Data File:**

Create a file named "word_count_data.txt" and populate it with text data that you wish to analyse.

Login with your hadoop user.

Download the dataset (weather data)**Output:**

dataset - Notepad

File	Edit	Format	View	Help																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																												</
------	------	--------	------	------	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	--	----

Step 2: Mapper Logic - mapper.py:

Create a file named "mapper.py" to implement the logic for the mapper. The mapper will read input data from STDIN, split lines into words, and output each word with its count.

```
nano mapper.py
```

```
# Copy and paste the mapper.py code
```

```
#!/usr/bin/env python
```

```
import sys
```

```
# input comes from STDIN (standard input)
```

```
# the mapper will get daily max temperature and group it by month. so output will be (month,dailymax_temperature)
```



```

for line in sys.stdin:
    # remove leading and trailing whitespace
    line = line.strip()
    # split the line into words
    words = line.split()
    #See the README hosted on the weather website which help us understand how each
position represents a column
    month = line[10:12]
    daily_max = line[38:45]
    daily_max = daily_max.strip()
    # increase counters
    for word in words:
        # write the results to STDOUT (standard output);
        # what we output here will be go through the shuffle process and then
        # be the input for the Reduce step, i.e. the input for reducer.py
        #
        # tab-delimited; month and daily max temperature as output
        print ('%s\t%s' % (month ,daily_max))

```

Step 3: Reducer Logic - reducer.py:

Create a file named "reducer.py" to implement the logic for the reducer. The reducer will aggregate the occurrences of each word and generate the final output.

```

nano reducer.py
# Copy and paste the reducer.py code

```

reducer.py

```

#!/usr/bin/env python

from operator import itemgetter
import sys

#reducer will get the input from stdid which will be a collection of key, value(Key=month ,
value= daily max temperature)
#reducer logic: will get all the daily max temperature for a month and find max temperature
for the month
#shuffle will ensure that key are sorted(month)
current_month = None
current_max = 0
month = None

# input comes from STDIN
for line in sys.stdin:

```

```

# remove leading and trailing whitespace
line = line.strip()
# parse the input we got from mapper.py
month, daily_max = line.split('\t', 1)

# convert daily_max (currently a string) to float
try:
    daily_max = float(daily_max)
except ValueError:
    # daily_max was not a number, so silently
    # ignore/discard this line
    continue

# this IF-switch only works because Hadoop shuffle process sorts map output
# by key (here: month) before it is passed to the reducer
if current_month == month:
    if daily_max > current_max:
        current_max = daily_max
else:
    if current_month:
        # write result to STDOUT
        print('%s\t%s' % (current_month, current_max))
    current_max = daily_max
    current_month = month

# output of the last month
if current_month == month:
    print('%s\t%s' % (current_month, current_max))

```

Step 4: Prepare Hadoop Environment:

Start the Hadoop daemons and create a directory in HDFS to store your data.

```
start-all.sh
```

Step 6: Make Python Files Executable:

Give executable permissions to your mapper.py and reducer.py files.

```
chmod 777 mapper.py reducer.py
```

Step 7: Run the program using Hadoop Streaming:

Download the latest hadoop-streaming jar file and place it in a location you can easily access.

Then run the program using Hadoop Streaming.

```
hadoop fs -mkdir -p /weatherdata
```

```
hadoop fs -copyFromLocal /home/sx/Downloads/dataset.txt /weatherdata
```

```
hdfs dfs -ls /weatherdata
```

```
hadoop jar /home/sx/hadoop-3.2.3/share/hadoop/tools/lib/hadoop-streaming-3.2.3.jar \
-input /weatherdata/dataset.txt \
-output /weatherdata/output \
-file "/home/sx/Downloads/mapper.py" \
-mapper "python3 mapper.py" \
-file "/home/sx/Downloads/reducer.py" \
-reducer "python3 reducer.py"
```

```
hdfs dfs -text /weatherdata/output/* > /home/sx/Downloads/outputfile.txt
```

Step 8: Check Output:

Check the output of the program in the specified HDFS output directory.

```
hdfs dfs -text /weatherdata/output/* > /home/sx/Downloads/output/
/part-00000
```

```
drwxr-xr-x - hadoop supergroup 0 2024-08-28 16:58 /test1
drwxr-xr-x - hadoop supergroup 0 2024-09-13 05:21 /tmp
drwxr-xr-x - hadoop supergroup 0 2024-09-11 23:18 /user
drwxr-xr-x - hadoop supergroup 0 2024-09-03 02:03 /weatherdata
drwxr-xr-x - hadoop supergroup 0 2024-08-31 01:03 /word_count_in_python

(hadoop@kali)~]
$ hdfs dfs -ls /weatherdata
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
Found 2 items
drwxr-xr-x - hadoop supergroup 0 2024-09-03 02:03 /weatherdata/output
-rw-r--r-- 1 hadoop supergroup 79205 2024-09-03 01:42 /weatherdata/weatherdata.txt

(hadoop@kali)~]
$ hdfs dfs -cat /weatherdata/output/part-m-00000
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
cat: '/weatherdata/output/part-m-00000': No such file or directory

(hadoop@kali)~]
$ hdfs dfs -cat /weatherdata/output/part-00000
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
01 26.5
02 26.6
03 29.1
04 30.8
05 31.1
06 33.6
07 38.5
08 40.2
09 36.5
10 36.9
11 27.6
12 25.9

(hadoop@kali)~]
```

Result:

Thus, the program for weather dataset using Map Reduce has been executed successfully.

EXP 4: Create UDF in PIG

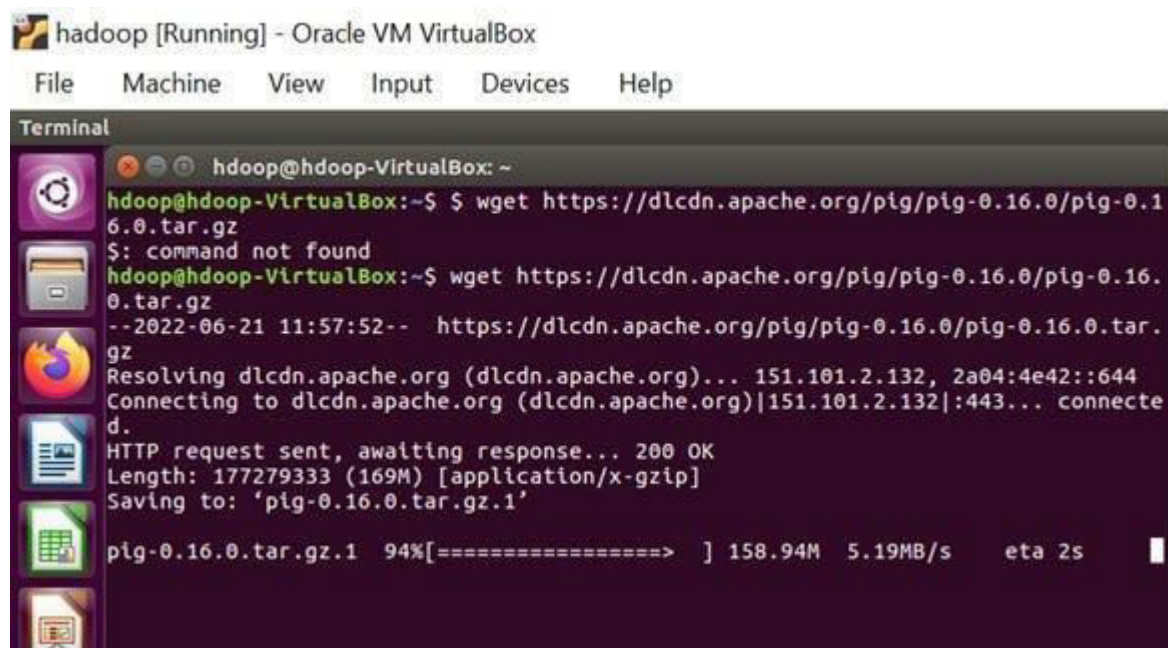
Step-by-step installation of Apache Pig on Hadoop cluster on Ubuntu

Pre-requisite:

- Ubuntu 16.04 or higher version running (I have installed Ubuntu on Oracle VM (Virtual Machine) VirtualBox),
- Run Hadoop on ubuntu (I have installed Hadoop 3.2.1 on Ubuntu 16.04). You may refer to my blog “How to install Hadoop installation” click [here](#) for Hadoop installation).

Pig installation steps

Step 1: Login into Ubuntu



```
hadoop [Running] - Oracle VM VirtualBox
File Machine View Input Devices Help
Terminal
hadoop@hadoop-VirtualBox: ~
hadoop@hadoop-VirtualBox:~$ $ wget https://dlcdn.apache.org/pig/pig-0.16.0/pig-0.16.0.tar.gz
$: command not found
hadoop@hadoop-VirtualBox:~$ wget https://dlcdn.apache.org/pig/pig-0.16.0/pig-0.16.0.tar.gz
--2022-06-21 11:57:52-- https://dlcdn.apache.org/pig/pig-0.16.0/pig-0.16.0.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... connecte
d.
HTTP request sent, awaiting response... 200 OK
Length: 177279333 (169M) [application/x-gzip]
Saving to: 'pig-0.16.0.tar.gz.1'

pig-0.16.0.tar.gz.1 94%[=====] 158.94M 5.19MB/s eta 2s
```

Step 2: Go to <https://pig.apache.org/releases.html> and copy the path of the latest version of pig that you want to install. Run the following comment to download Apache Pig in Ubuntu:

\$ wget <https://dlcdn.apache.org/pig/pig-0.16.0/pig-0.16.0.tar.gz>

Step 3: To untar pig-0.16.0.tar.gz file run the following command:

```
$ tar xvzf pig-0.16.0.tar.gz
```

Step 4: To create a pig folder and move pig-0.16.0 to the pig folder, execute the following command:

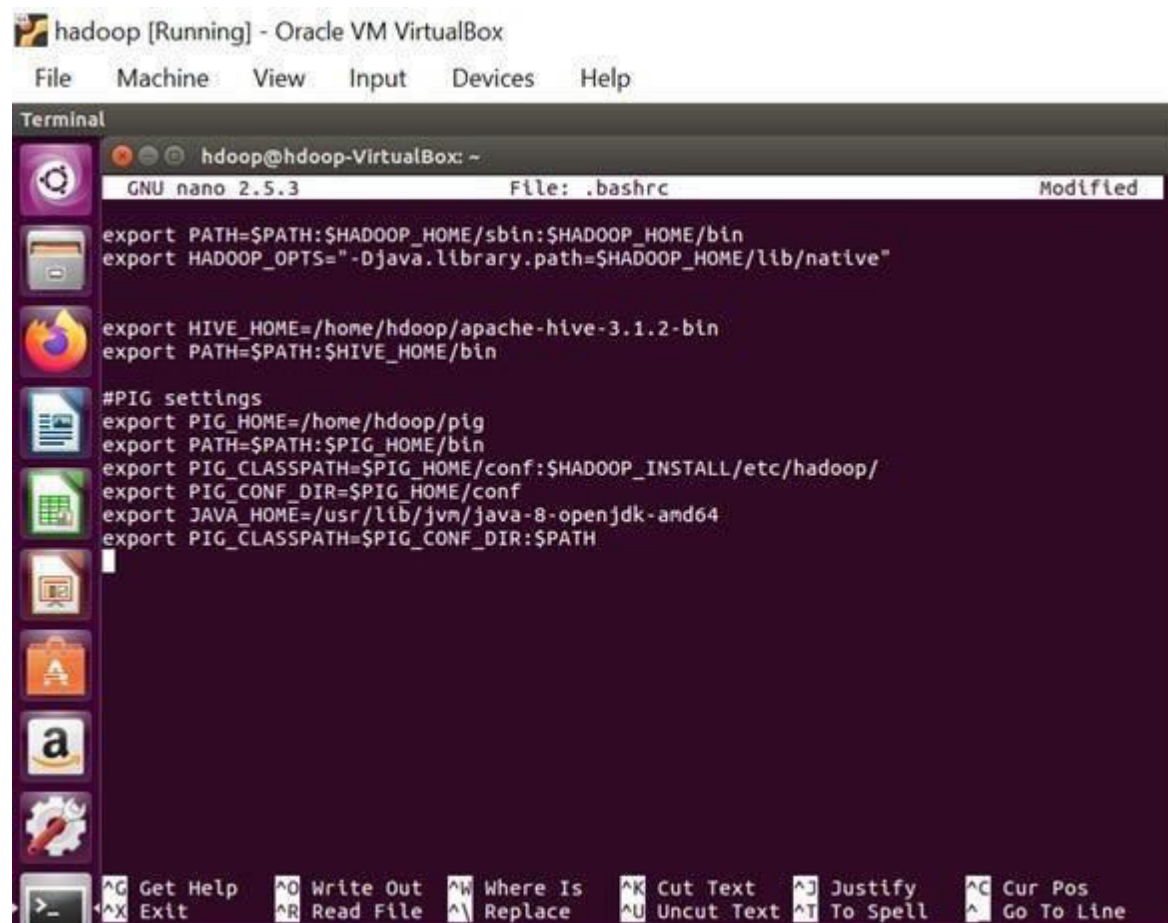
```
$ sudo mv /home/hadoop/pig-0.16.0 /home/hadoop/pig
```

Step 5: Now open the .bashrc file to edit the path and variables/settings for pig. Run the following command:

```
$ sudo nano .bashrc
```

Add the below given to .bashrc file at the end and save the file.

```
#PIG settings
export PIG_HOME=/home/hadoop/pig
export PATH=$PATH:$PIG_HOME/bin
export PIG_CLASSPATH=$PIG_HOME/conf:$HADOOP_INSTALL/etc/hadoop/
export PIG_CONF_DIR=$PIG_HOME/conf
export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
export PIG_CLASSPATH=$PIG_CONF_DIR:$PATH
#PIG setting ends
```



Step 6: Run the following command to make the changes effective in the .bashrc file:

```
$ source .bashrc
```

Step 7: To start all Hadoop daemons, navigate to the `hadoop-3.2.1/sbin` folder and run the following commands:

```
$ ./start-dfs.sh $ ./start-yarn$ jps
```

```
hadoop@hadoop-VirtualBox:~$ cd hadoop-3.2.1/sbin
hadoop@hadoop-VirtualBox:~/hadoop-3.2.1/sbin$ ./start-dfs.sh
Starting namenodes on [localhost]
Starting datanodes
Starting secondary namenodes [hadoop-VirtualBox]
hadoop@hadoop-VirtualBox:~/hadoop-3.2.1/sbin$ ./start-yarn.sh
Starting resourcemanager
Starting nodemanagers
hadoop@hadoop-VirtualBox:~/hadoop-3.2.1/sbin$ jps
4817 DataNode
5298 ResourceManager
5000 SecondaryNameNode
5450 NodeManager
4683 NameNode
5982 Jps
hadoop@hadoop-VirtualBox:~/hadoop-3.2.1/sbin$
```

Step 8: Now you can launch pig by executing the following command:

```
$ pig
```

```
hadoop@kali: ~/dalab/expt4
File Actions Edit View Help
Success!
Job Stats (time in seconds):
JobId  Maps  Reduces  MaxMapTime  MinMapTime  AvgMapTime  MedianMapTime  MaxReduceTime  MinReduceTime  AvgReduceTime  MedianReduceTime  Alias  Feature Outputs
job_local1214476415_0001  1  0  n/a  n/a  n/a  0  0  data,uppercased_data  MAP_ONLY  hdfs:///piginput/pig_output_data,

Input(s):
Successfully read 4 records (42400458 bytes) from: "hdfs:///piginput/sample.txt"

Output(s):
Successfully stored 4 records (42400292 bytes) in: "hdfs:///piginput/pig_output_data"

Counters:
Total records written : 4
Total bytes written : 42400292
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local1214476415_0001

2024-09-12 11:41:19.749 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2024-09-12 11:41:19.755 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2024-09-12 11:41:19.764 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2024-09-12 11:41:19.799 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2024-09-12 11:41:19.838 [main] INFO org.apache.pig.Main - Pig script completed in 29 seconds and 719 milliseconds (29719 ms)

(hadoop@kali)~/dalab/expt4
$ hdfs dfs -ls /piginput
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
Found 3 items
drwxr-xr-x - hadoop supergroup 0 2024-09-12 11:41 /piginput/pig_output_data
-rw-r--r-- 1 hadoop supergroup 27 2024-09-04 23:25 /piginput/sample.txt
-rw-r--r-- 1 hadoop supergroup 166 2024-09-12 04:17 /piginput/uppercase_udf.py

(hadoop@kali)~/dalab/expt4
$ hdfs dfs -cat /piginput/pig_output_data
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
cat: /piginput/pig_output_data: Is a directory

(hadoop@kali)~/dalab/expt4
$ hdfs dfs -cat /piginput/pig_output_data/part-m-00000
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
1,JOHN
2,JANE
3,JOE
4,EMMA

(hadoop@kali)~/dalab/expt4
```

Step 9: Now you are in pig and can perform your desired tasks on pig. You can come out of the pig by the quit command:

```
> quit;
```


CREATE USER DEFINED FUNCTION(UDF)

Aim : To create User Define Function in Apache Pig and execute it on map reduce.

Procedure:

Create a sample text file

```
hadoop@Ubuntu:~/Documents$ nano sample.txt
```

Paste the below content to sample.txt

1,John

2,Jane

3,Joe

4,Emma

```
hadoop@Ubuntu:~/Documents$ hadoop fs -put sample.txt /home/hadoop/piginput/
```

Create PIG File

```
hadoop@Ubuntu:~/Documents$ nano demo_pig.pig
```

paste the below the content to demo_pig.pig

-- Load the data from HDFS

```
data = LOAD '/home/hadoop/piginput/sample.txt' USING PigStorage(',') AS (id:int>
```

-- Dump the data to check if it was loaded correctly

```
DUMP data;
```

Run the above file

```
hadoop@Ubuntu:~/Documents$ pig demo_pig.pig
```

2024-08-07 12:13:08,791 [main] INFO

org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil

- Total input paths to process : 1

(1,John)

(2,Jane)

(3,Joe)

(4,Emma)

Create udf file an save as uppercase_udf.py

uppercase_udf.py

```
def uppercase(text):
```

```
    return text.upper()
```

```
if __name__ == "__main__":
```

```
    import sys
```

```
    for line in sys.stdin:
```

```
        line = line.strip()
```

```
        result = uppercase(line)
```

```
        print(result)
```

Create the udfs folder on hadoop

```
hadoop@Ubuntu:~/Documents$ hadoop fs -mkdir /home/hadoop/udfs
```

put the upppercase_udf.py in to the abv folder

```
hadoop@Ubuntu:~/Documents$ hdfs dfs -put uppercase_udf.py /home/hadoop/udfs/
```

```
hadoop@Ubuntu:~/Documents$ nano udf_example.pig
```

copy and paste the below content on udf_example.pig

-- Register the Python UDF script

```
REGISTER 'hdfs:///home/hadoop/udfs/uppercase_udf.py' USING jython AS udf;
```

```
-- Load some data
data = LOAD 'hdfs:///home/hadoop/sample.txt' AS (text:chararray);

-- Use the Python UDF
uppercased_data = FOREACH data GENERATE udf.uppercase(text) AS uppercase_text;

-- Store the result
STORE uppercased_data INTO 'hdfs:///home/hadoop/pig_output_data';
```

place sample.txt file on hadoop

```
hadoop@Ubuntu:~/Documents$ hadoop fs -put sample.txt /home/hadoop/
```

To Run the pig file

```
hadoop@Ubuntu:~/Documents$ pig -f udf_example.pig
```

finally u get

Success!

Job Stats (time in seconds):

```
JobId Maps Reduces MaxMapTimeMinMapTime AvgMapTime MedianMapTime
MaxReduceTime MinReduceTime AvgReduceTime MedianReductime
Alias Feature Outputs
```

```
job_local1786848041_0001 1 0 n/a n/a n/a n/a 00 0 0
```

```
data,uppercased_data MAP_ONLY hdfs:///home/hadoop/pig_output_data,
```

Input(s):

Successfully read 4 records (42778068 bytes) from: "hdfs:///home/hadoop/sample.txt"

Output(s):

Successfully stored 4 records (42777870 bytes) in: "hdfs:///home/hadoop/pig_output_data"

Counters:

Total records written : 4

Total bytes written : 42777870

Spillable Memory Manager spill count : 0

Total bags proactively spilled: 0

Total records proactively spilled: 0

Job DAG:

job_local1786848041_0001

2024-08-07 13:33:04,631 [main] WARN
org.apache.hadoop.metrics2.impl.MetricsSystemImpl -

JobTracker metrics system already initialized!

2024-08-07 13:33:04,639 [main] WARN
org.apache.hadoop.metrics2.impl.MetricsSystemImpl -

JobTracker metrics system already initialized!

2024-08-07 13:33:04,644 [main] WARN
org.apache.hadoop.metrics2.impl.MetricsSystemImpl -

JobTracker metrics system already initialized!

2024-08-07 13:33:04,667 [main] INFO

org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher -
Success!

Note :

If any error check jython package is installed and check the path specified on the above steps are give correctly

To check the output file is created

hadoop@Ubuntu:~/Documents\$ hdfs dfs -ls /home/hadoop/pig_output_data

Found 2 items

If you need to examine the files in the output folder, use:

To view the output

hadoop@Ubuntu:~/Documents\$ hdfs dfs -cat /home/hadoop/pig_output_data/part-m-00000

```
File Actions Edit View Help
JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxReduceTime MinReduceTime AvgReduceTime MedianReducetime
job_local1214476415_0001 1 0 n/a n/a n/a 0 0 0 data,uppercased_data MAP_ONLY hdfs

Input(s):
Successfully read 4 records (42400458 bytes) from: "hdfs:///piginput/sample.txt"

Output(s):
Successfully stored 4 records (42400292 bytes) in: "hdfs:///piginput/pig_output_data"

Counters:
Total records written : 4
Total bytes written : 42400292
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0

Job DAG:
job_local1214476415_0001

2024-09-12 11:41:19,749 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2024-09-12 11:41:19,755 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2024-09-12 11:41:19,764 [main] WARN org.apache.hadoop.metrics2.impl.MetricsSystemImpl - JobTracker metrics system already initialized!
2024-09-12 11:41:19,799 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2024-09-12 11:41:19,838 [main] INFO org.apache.pig.Main - Pig script completed in 29 seconds and 719 milliseconds (29719 ms)

(hadoop@kali)~[~/dalab/expt4]
$ hdfs dfs -ls /piginput
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
Found 3 items
drwxr-xr-x - hadoop supergroup 0 2024-09-12 11:41 /piginput/pig_output_data
-rw-r--r-- 1 hadoop supergroup 27 2024-09-04 23:25 /piginput/sample.txt
-rw-r--r-- 1 hadoop supergroup 166 2024-09-12 04:17 /piginput/uppercase_udf.py

(hadoop@kali)~[~/dalab/expt4]
$ hdfs dfs -cat /piginput/pig_output_data
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
cat: '/piginput/pig_output_data': Is a directory

(hadoop@kali)~[~/dalab/expt4]
$ hdfs dfs -cat /piginput/pig_output_data/part-m-00000
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
1,JOHN
2,JANE
3,JOE
4,EMMA

(hadoop@kali)~[~/dalab/expt4]
$
::1 ff02::1 ff02::2 ip6-allnodes ip6-allrouters ip6-localhost ip6-loopback kali localhost
(hadoop@kali)~[~/dalab/expt4]
$
```

Result:

Thus the program is executed successfully

Exp5: Installation of Hive on Ubuntu

Aim:

To Download and install Hive, Understanding Startup scripts, Configuration files.

Procedure:

Step 1: Download and extract it

Download the Apache hive and extract it use tar, the commands given below:

```
$wgethttps://downloads.apache.org/hive/hive-3.1.2/apache-hive-3.1.2-bin.tar.gz
```

```
$ tar -xvf apache-hive-3.1.2-bin.tar.gz
```

Step 2: Place different configuration properties in Apache Hive

In this step, we are going to do two things

- Placing Hive Home path in bashrc file

```
$nano .bashrc
```

And append the below lines in it

```
export HIVE_HOME=/home/hadoop/apache-hive-3.1.2-bin
export PATH=$PATH:$HIVE_HOME/bin
export HADOOP_USER_CLASSPATH_FIRST=true
```

- Exporting **Hadoop path in Hive-config.sh** (To communicate with the Hadoop ecosystem we are defining Hadoop Home path in hive config field) **Open the hive-config.sh as shown in below**

```
$cd apache-hive-3.1.2-bin/bin
```

```
$cp hive-env.sh.template hive-env.sh
```

```
$nano hive-env.sh
```

Append the below commands on it

```
export HADOOP_HOME=/home/Hadoop/Hadoop
export HIVE_CONF_DIR=/home/Hadoop/apache-hive-3.1.2/conf
```

```
# Set HADOOP_HOME to point to a specific hadoop install directory
# HADOOP_HOME=${bin}/../.. /hadoop
export HADOOP_HOME=/home/hadoop/hadoop

# Hive Configuration Directory can be controlled by:
# export HIVE_CONF_DIR=
export HIVE_CONF_DIR=/home/hadoop/apache-hive-3.1.2-bin/conf
# Folder containing extra libraries required for hive compilation/execution can be controlled by:
```

Step 3: Install mysql

- Install mysql in Ubuntu by running this command:

```
$sudo apt update
```

```
$sudo apt install mysql-server
```

- Alter username and password for MySQL by running below commands:

```
$sudomysql
```

Opens command line interface for MySQL and run the below SQL queries to change username and set password

```
mysql> SELECT user, host, plugin FROM mysql.user WHERE user = 'root';
```



```
mysql> ALTER USER 'root'@'localhost' IDENTIFIED WITH 'mysql_native_password' BY  
'your_new_password';  
mysql> FLUSH PRIVILEGES;
```

Step 4: Config hive-site.xml

Config the hive-site.xml by appending this xml code and change the username and password according to your MySQL.

```
$cd apache-hive-3.1.2-bin/bin
```

```
$cp hive-default.xml.template hive-site.xml
```

```
$nano hive-site.xml
```

Append these lines into it

Replace root as your username of MySQL

Replace your_new_password as with your password of MySQL

```
<configuration>
```

```
<property>
```

```
  <name>javax.jdo.option.ConnectionURL</name>
```

```
  <value>jdbc:mysql://localhost/metastore?createDatabaseIfNotExist=true</value>
```

```
</property>
```

```
<property>
```

```
  <name>javax.jdo.option.ConnectionDriverName</name>
```

```
  <value>com.mysql.cj.jdbc.Driver</value>
```

```
</property>
```

```
<property>
```

```
  <name>javax.jdo.option.ConnectionUserName</name>
```

```
  <value>root</value>
```

```
</property>
```

```
<property>
<name>javax.jdo.option.ConnectionPassword</name>
<value>your_new_password</value>
</property>
```

```
<property>
<name>datanucleus.autoCreateSchema</name>
<value>true</value>
</property>
```

```
<property>
<name>datanucleus.fixedDatastore</name>
<value>true</value>
</property>
```

```
<property>
<name>datanucleus.autoCreateTables</name>
<value>True</value>
</property>
```

```
</configuration>
```

Step 5: Setup MySQL java connector:

First, you'll need to download the *MySQL Connector/J*, which is the *JDBC* driver for *MySQL*. You can download it from the below link

https://drive.google.com/file/d/1QFhB7Kvc7a4LzDRe6GcmZva1yAxKz-/view?usp=drive_link

Copy the downloaded *MySQL Connector/J* JAR file to the Hive library directory. By default, the Hive library directory is usually located at */path/to/apache-hive-3.1.2/lib/* on Ubuntu. Use the following command to copy the JAR file:

```
$sudo cp /path/to/mysql-connector-java-8.0.15.jar /path/to/apache-hive-3.1.2/lib/
```

Replace */path/to/* with the actual path to the JAR file.

Step 6: Initialize the Hive Metastore Schema:

Run the following command to initialize the Hive metastore schema:

```
$$HIVE_HOME/bin/schematool -initSchema -dbTypemysql
```

Step 7: Start hive:

```
nadoop@osboxes:~$ wget https://archive.apache.org/dist/hive/hive-3.1.2/apache-hive-3.1.2-bin.tar.gz
--2024-09-18 04:57:39-- https://archive.apache.org/dist/hive/hive-3.1.2/apache-hive-3.1.2-bin.tar.gz
Resolving archive.apache.org (archive.apache.org)... 65.108.204.189, 2a01:4f9:1a:a084::2
Connecting to archive.apache.org (archive.apache.org)[65.108.204.189]:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 278813748 (266M) [application/x-gzip]
Saving to: 'apache-hive-3.1.2-bin.tar.gz'

apache-hive-3.1.2-b 55%[=====>] 146.97M ---KB/s in 14m 57s

2024-09-18 05:12:38 (168 KB/s) - Read error at byte 15411881/278813748 (Connection reset by peer). Retrying.

--2024-09-18 05:12:39-- (try: 2) https://archive.apache.org/dist/hive/hive-3.1.2/apache-hive-3.1.2-bin.tar.gz
Connecting to archive.apache.org (archive.apache.org)[65.108.204.189]:443... connected.
HTTP request sent, awaiting response... 206 Partial Content
Length: 278813748 (266M), 124701867 (119M) remaining [application/x-gzip]
Saving to: 'apache-hive-3.1.2-bin.tar.gz'

apache-hive-3.1.2-b 100%[++++++>] 265.90M 548KB/s in 3m 9s

2024-09-18 05:15:49 (644 KB/s) - 'apache-hive-3.1.2-bin.tar.gz' saved [278813748/278813748]
```

You can test Hive by running the Hive shell: Copy code hive You should be able to run Hive queries, and metadata will be stored in your MySQL database.

\$hive

```

(hadoop@kali):~$ cd hive
(hadoop@kali):~/hive$ cd bin
(hadoop@kali):~/hive/bin$ hive
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
Picked up _JAVA_OPTIONS: -Dawt.useSystemAAFontSettings=on -Dswing.aatext=true
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/home/hadoop/hive/lib/log4j-slf4j-impl-2.10.0.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/home/hadoop/hadoop/share/hadoop/common/lib/slf4j-reload4j-1.7.36.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.apache.logging.slf4j.Log4jLoggerFactory]
Hive Session ID = bd58a5e7-c2a5-4cda-b0e4-df27859d8bc

Logging initialized using configuration in jar:file:/home/hadoop/hive/lib/hive-common-3.1.2.jar!/hive-log4j2.properties Async: true
Hive Session ID = d4d0de0-467d-422c-9ef9-7fb273c9b0da
Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
hive> CREATE DATABASE financials;
OK
Time taken: 0.991 seconds
hive> USE financials;
OK
Time taken: 0.033 seconds
hive> CREATE TABLE financial_table(id INT,name STRING);
OK
Time taken: 1.002 seconds
hive> INSERT INTO TABLE financial_table VALUES (1, 'Alice'), (2, 'Bob'), (3, 'Charlie');
MismatchedTokenException(24:374)
    at org.antlr.runtime.BaseRecognizer.recoverFromMismatchedToken(BaseRecognizer.java:617)
    at org.antlr.runtime.BaseRecognizer.match(BaseRecognizer.java:115)
    at org.apache.hadoop.hive.ql.parse.HiveParser_IdentifierParser.expressionsInParenthesis(HiveParser_IdentifierParser.java:2376)
    at org.apache.hadoop.hive.ql.parse.HiveParser.expressionsInParenthesis(HiveParser.java:45260)
    at org.apache.hadoop.hive.ql.parse.HiveParser.FromClauseParser.valueRowConstructor(HiveParser_FromClauseParser.java:6214)
    at org.apache.hadoop.hive.ql.parse.HiveParser.FromClauseParser.valuesTableConstructor(HiveParser_FromClauseParser.java:6131)
    at org.apache.hadoop.hive.ql.parse.HiveParser.FromClauseParser.valuesClause(HiveParser_FromClauseParser.java:6045)
    at org.apache.hadoop.hive.ql.parse.HiveParser.valuesClause(HiveParser.java:45342)
    at org.apache.hadoop.hive.ql.parse.HiveParser.regularBody(HiveParser.java:39614)
    at org.apache.hadoop.hive.ql.parse.HiveParser.queryStatementExpressionBody(HiveParser.java:38900)
    at org.apache.hadoop.hive.ql.parse.HiveParser.queryStatementExpression(HiveParser.java:38788)
    at org.apache.hadoop.hive.ql.parse.HiveParser.execStatement(HiveParser.java:2396)
    at org.apache.hadoop.hive.ql.parse.HiveParser.statement(HiveParser.java:1420)
    at org.apache.hadoop.hive.ql.parse.ParseDriver.parse(ParseDriver.java:220)
    at org.apache.hadoop.hive.ql.parse.ParseUtils.parse(ParseUtils.java:74)

```

Result:

Thus, the Apache Hive installation is completed successfully on Ubuntu.

Exp5a: Design and test various schema models to optimize data storage and retrieval Using Hive.

Aim:

To Design and test various schema models to optimize data storage and retrieval Using Hbase.

Procedure:

Step 1: Start Hive

Open a terminal and start Hive by running:

```
$hive
```

Step 2: Create a Database

Create a new database in Hive:

```
hive>CREATE DATABASE financials;
```

```
hive> CREATE DATABASE financials;
```

```
OK
```

```
Time taken: 0.063 seconds
```

Step 3: Use the Database:

Switch to the newly created database:

```
hive>use financials;
```

```
hive> use financials;
```

```
OK
```

```
Time taken: 0.066 seconds
```

Step 4: Create a Table:

Create a simple table in your database:

```
hive>CREATE TABLE finance_table( id INT, name STRING );
```

```
hive> CREATE TABLE finance_table (  
    > id INT,  
    > name STRING  
    > );
```

```
OK
```

```
Time taken: 0.768 seconds
```

Step 5: Load Sample Data:

You can insert sample data into the table:

```
hive>INSERT INTO finance_tableVALUES (1, 'Alice'), (2, 'Bob'), (3, 'Charlie');
```

```

hive> INSERT INTO finance_table VALUES
      > (1, 'Alice'),
      > (2, 'Bob'),
      > (3, 'Charlie');
Query ID = hadoop_20231028192937_fdebeb4e-abf7-4bad-a248-ac908246e3c1
Total jobs = 3
Launching Job 1 out of 3
Number of reduce tasks determined at compile time: 1
In order to change the average load for a reducer (in bytes):
  set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
  set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
  set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2023-10-28 19:29:41,158 Stage-1 map = 0%,  reduce = 0%

```

Step 6: Query Your Data

Use SQL-like queries to retrieve data from your table:

```
hive> CREATE VIEW myview AS SELECT name, id FROM finance_table;
```

Step 7: View the data:

To see the data in the view, you would need to query the view

```
hive> SELECT * FROM myview;
```

```

hive> SELECT * FROM myview;
OK
Alice    1
Bob      2
Charlie  3
Time taken: 0.238 seconds, Fetched: 3 row(s)

```

Step 8: Describe a Table:

You can describe the structure of a table using the DESCRIBE command:

```
hive> DESCRIBE finance_table;
```

```

hive> DESCRIBE finance_table;
OK
id                int
name              string
Time taken: 0.081 seconds, Fetched: 2 row(s)

```

Step 9: Alter a Table:

You can alter the table structure by adding a new column:

```
hive> ALTER TABLE finance_table ADD COLUMNS (age INT);
```

```

hive> ALTER TABLE finance_table ADD COLUMNS (age INT);
OK
Time taken: 0.165 seconds

```

Step 10: Quit Hive:

To exit the Hive CLI, simply type:

```
hive> quit;
```

>quit;

```

set hive.exec.reducers.bytes.per.reducer=<number>
In order to limit the maximum number of reducers:
set hive.exec.reducers.max=<number>
In order to set a constant number of reducers:
set mapreduce.job.reduces=<number>
Job running in-process (local Hadoop)
2024-09-14 03:51:34,630 Stage-1 map = 0%, reduce = 0%
2024-09-14 03:51:39,686 Stage-1 map = 100%, reduce = 100%
Ended Job = job_local584158451_0001
Stage-4 is selected by condition resolver.
Stage-3 is filtered out by condition resolver.
Stage-5 is filtered out by condition resolver.
Moving data to directory hdfs://localhost:9000/user/hive/warehouse/financials.db/financial_table/.hive-staging_hive_
2024-09-14_03-51-25_798_4467520597677782042-1/-ext-10000
Loading data to table financials.financial_table
MapReduce Jobs Launched:
Stage-Stage-1: HDFS Read: 0 HDFS Write: 212 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
Time taken: 15.156 seconds
hive> CREATE VIEW myview AS SELECT name, id FROM finance_table;
FAILED: SemanticException [Error 10001]: Line 1:43 Table not found 'finance_table'
hive> CREATE VIEW myview AS SELECT name, id FROM finance_table;
FAILED: SemanticException [Error 10001]: Line 1:43 Table not found 'finance_table'
hive> CREATE VIEW myview AS SELECT name, id from financial_table;
OK
Time taken: 0.233 seconds
hive> SELECT * FROM myview;
OK
Alice 1
Bob 2
Charlie 3
Time taken: 0.235 seconds, Fetched: 3 row(s)
hive> DESCRIBE finance_table;
FAILED: SemanticException [Error 10001]: Table not found finance_table
hive> DESCRIBE financial_table;
OK
id int
name string
Time taken: 0.053 seconds, Fetched: 2 row(s)
hive> ALTER TABLE financial_table ADD COLUMNS (age INT);
OK
Time taken: 0.208 seconds
hive> █

```

Result:

Thus, the usage of various commands in Hive has been successfully completed.