

INDEX

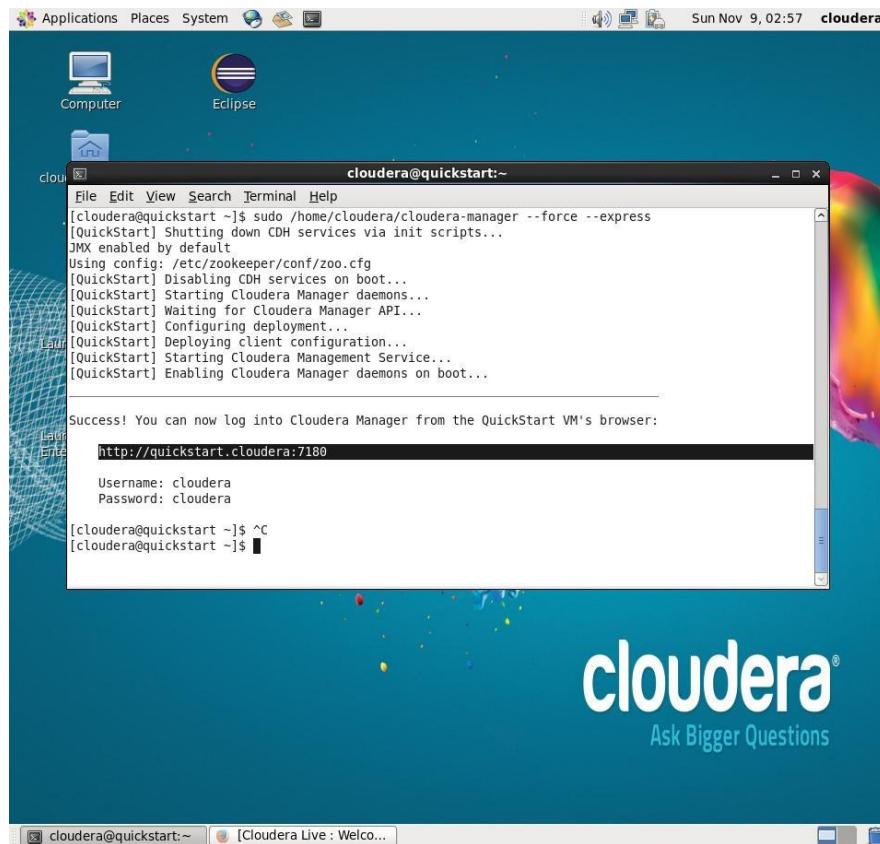
SR NO	NAME OF THE EXPERIMENT	DATE	FACULTY SIGN	CO
1	HDFS: List of Commands (mkdir, touchz, copy from local/put, copy to local/get, move from local, cp, rmr, du, dus, stat)			CO1
2	Map Reduce: <ul style="list-style-type: none"> 1. Write a program in Map Reduce for Word Count operation. 2. Write a program in Map Reduce for Union operation. 3. Write a program in Map Reduce for Intersection operation. 4. Write a program in Map Reduce for Matrix Multiplication. 			CO2
3	MongoDB: <ul style="list-style-type: none"> 1. Installation 2. Sample Database Creation 3. Query the Sample Database using MongoDB querying commands <ul style="list-style-type: none"> a. Create Collection b. Insert Document c. Query Document d. Delete Document e. Indexing 			CO3
4	Hive: <ul style="list-style-type: none"> 1. Hive Data Types 2. Create Database & Table in Hive 3. Hive Partitioning 4. Hive Built-In Operators 5. Hive Built-In Functions 6. Hive Views 			CO4

	7. HiveQL: Select Where, Select OrderBy, Select GroupBy, Select Joins																									
5	Pig: <ol style="list-style-type: none"> 1. Pig Latin Basic 2. Pig Data Types, 3. Download the data 4. Create your Script 5. Save and Execute the Script 6. Pig Operations: Diagnostic Operators, Grouping and Joining, Combining & Splitting, Filtering, Sorting 			CO4																						
6	Spark: <table border="1" data-bbox="278 842 945 1471"> <thead> <tr> <th>Apache Spark Commands in Scala</th> <th>Pair RDD (Key-Value RDD) Operations</th> </tr> </thead> <tbody> <tr><td>Start Spark Shell</td><td>Create Pair RDD</td></tr> <tr><td>Create RDD from Collection</td><td>reduceByKey</td></tr> <tr><td>Read Text File into RDD (from HDFS)</td><td>groupByKey — With Output Viewing</td></tr> <tr><td>Map Transformation</td><td>mapValues</td></tr> <tr><td>Filter Transformation</td><td>sortByKey</td></tr> <tr><td>Reduce Action</td><td>join</td></tr> <tr><td>Collect Action</td><td>cogroup</td></tr> <tr><td>Save RDD to HDFS</td><td>aggregateByKey</td></tr> <tr><td>Cache/Persist RDD</td><td>foldByKey</td></tr> <tr><td></td><td>Read from Local File</td></tr> </tbody> </table>	Apache Spark Commands in Scala	Pair RDD (Key-Value RDD) Operations	Start Spark Shell	Create Pair RDD	Create RDD from Collection	reduceByKey	Read Text File into RDD (from HDFS)	groupByKey — With Output Viewing	Map Transformation	mapValues	Filter Transformation	sortByKey	Reduce Action	join	Collect Action	cogroup	Save RDD to HDFS	aggregateByKey	Cache/Persist RDD	foldByKey		Read from Local File			CO5
Apache Spark Commands in Scala	Pair RDD (Key-Value RDD) Operations																									
Start Spark Shell	Create Pair RDD																									
Create RDD from Collection	reduceByKey																									
Read Text File into RDD (from HDFS)	groupByKey — With Output Viewing																									
Map Transformation	mapValues																									
Filter Transformation	sortByKey																									
Reduce Action	join																									
Collect Action	cogroup																									
Save RDD to HDFS	aggregateByKey																									
Cache/Persist RDD	foldByKey																									
	Read from Local File																									
7	Visualization using Tableau: <p>Tableau: Tool Overview, Importing Data, Analyzing with Charts, Creating Dashboards, working with maps</p>			CO6																						

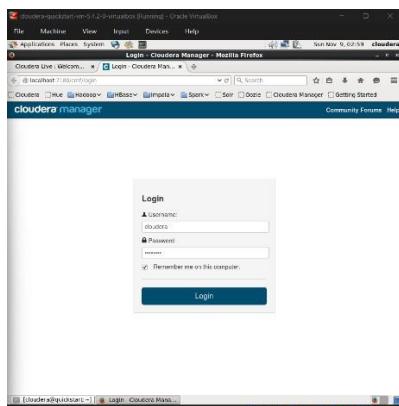
Practical No:1 HDFS

List of Commands (mkdir, touchz, copy from local/put, copy to local/get, move from local, cp, rmr, du, dus, stat)

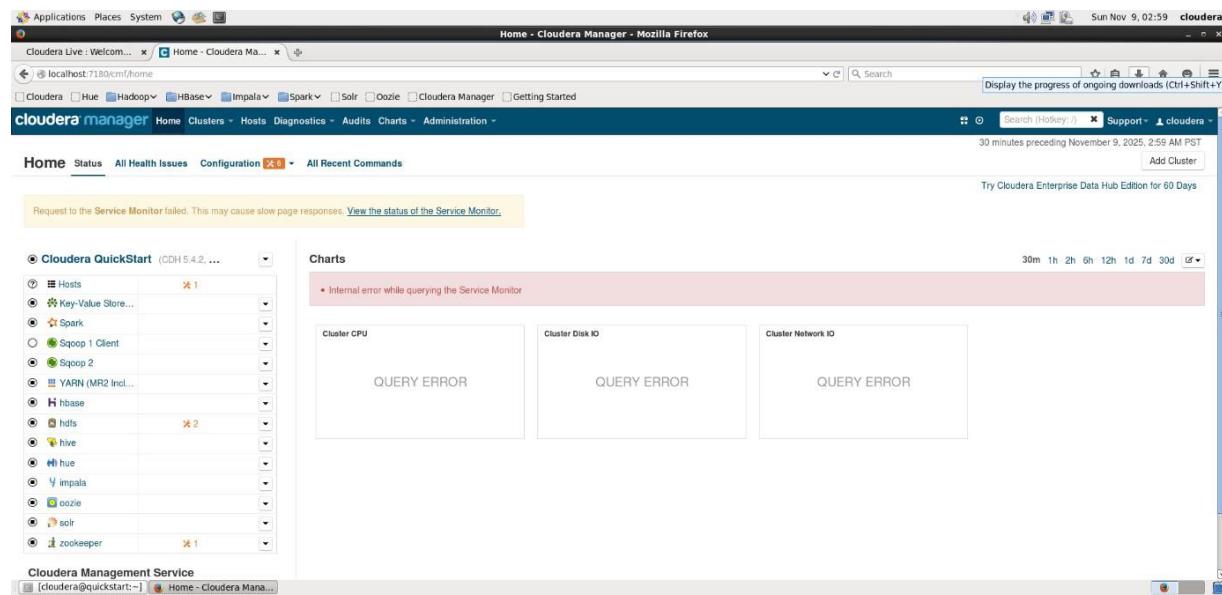
Step 1: Start Cloudera Manager in your Cloudera VM.



Step 2: Enter valid Login credentials – username and password -> will navigate to cloudera manager dashboard.

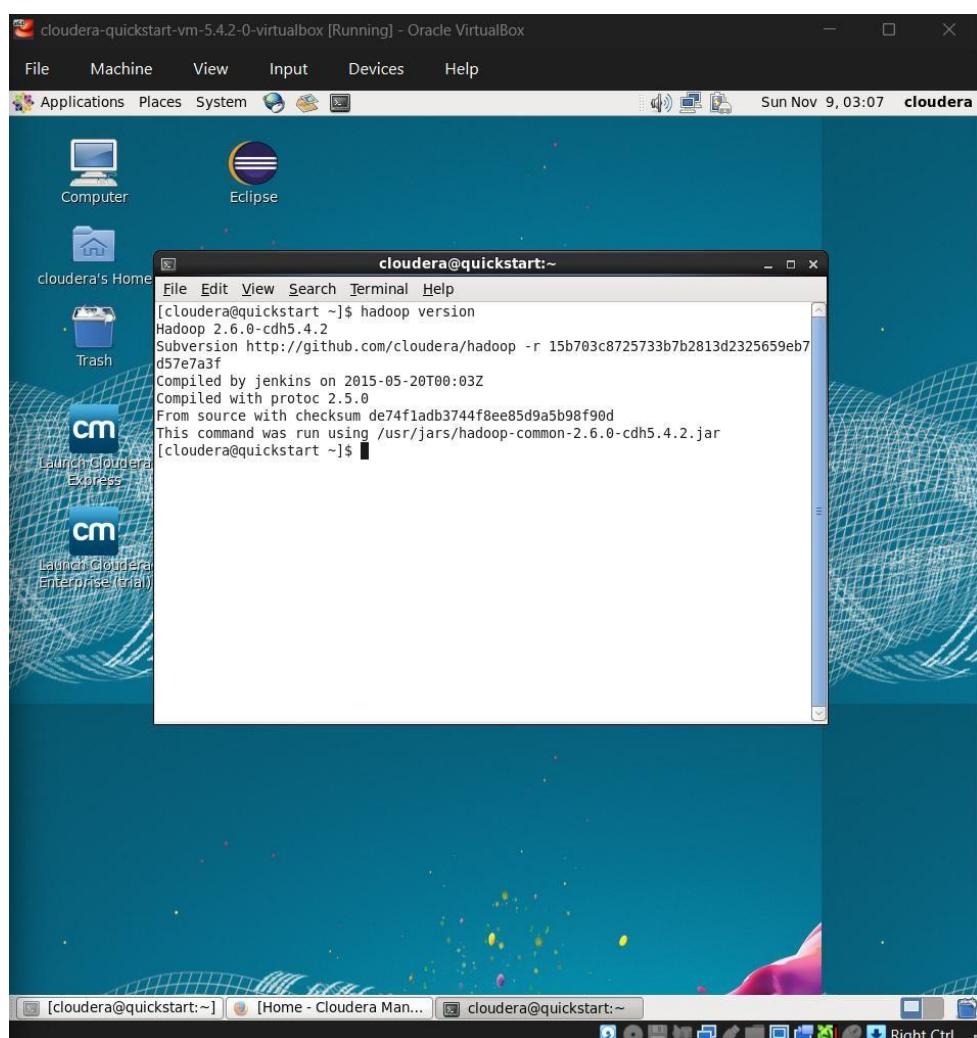


Step 3: Activate HDFS, Map Reduce, Hbase and YARN



Make sure that that should be green tick, in order to perform commands

Step 4: Verify active services by checking Hadoop version.

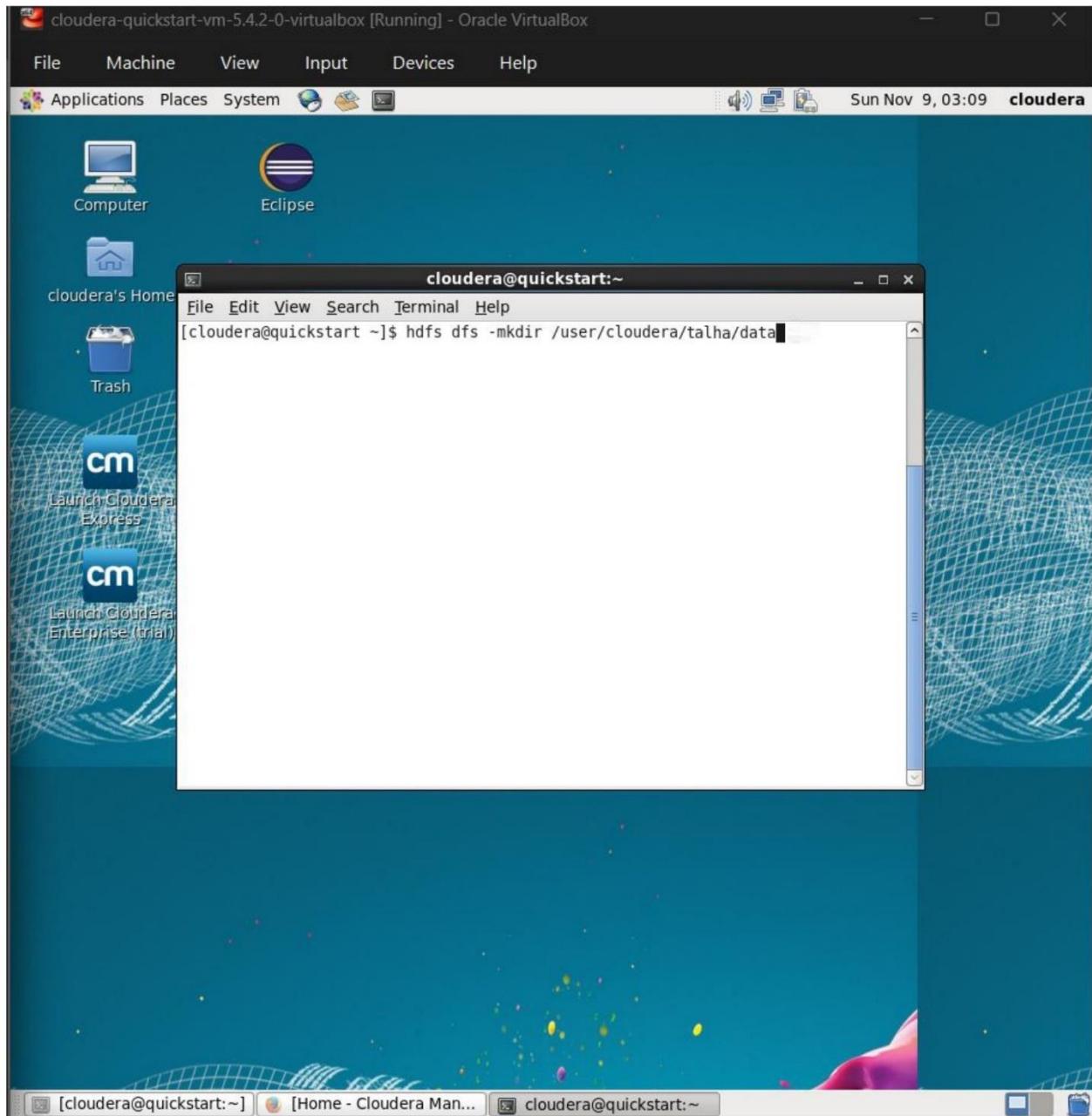


List of Commands to Perform in the terminal.

1. mkdir

Creates a new directory in HDFS.

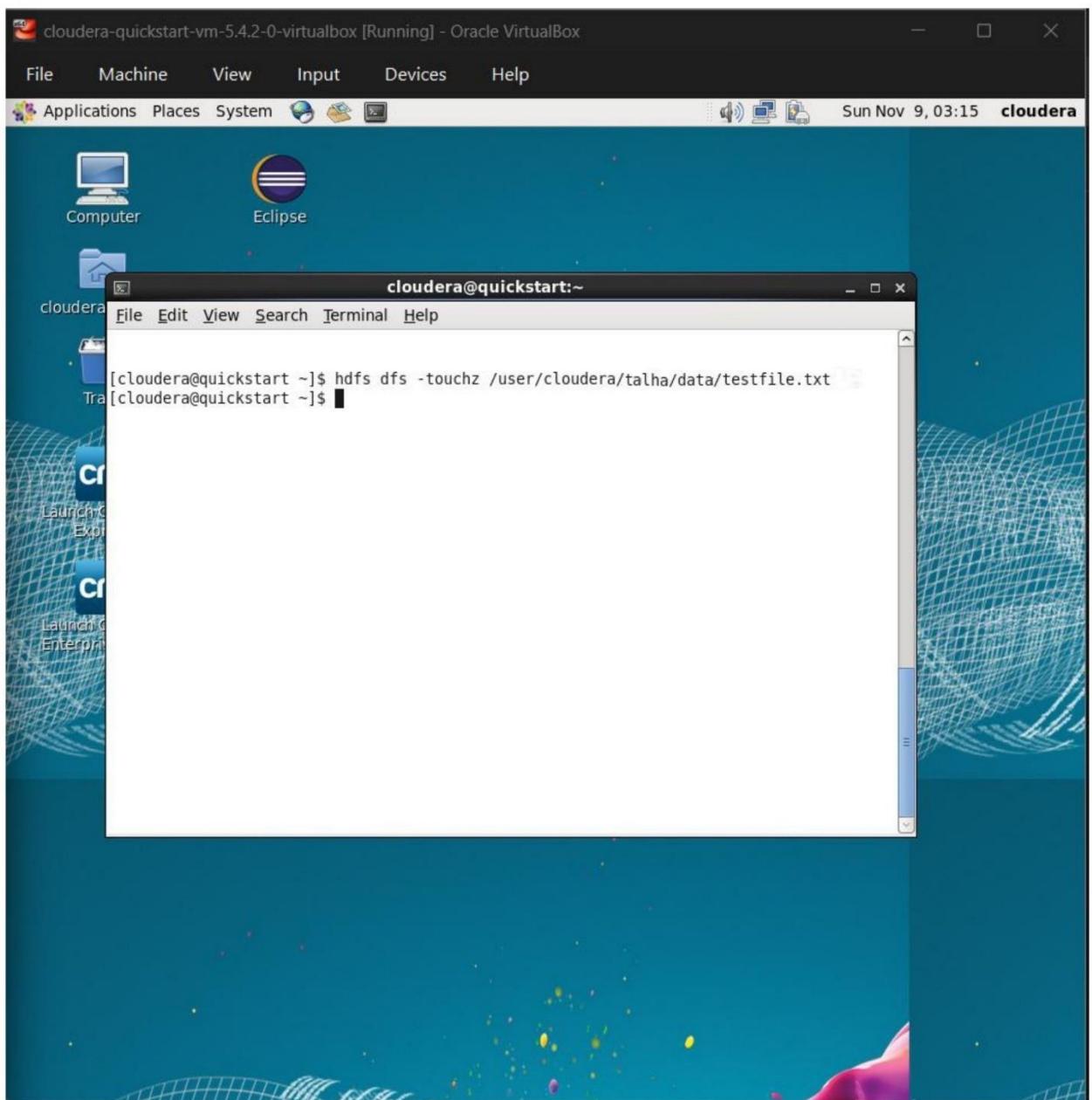
```
hdfs dfs -mkdir /user/cloudera/talha/data
```



2. touchz

Creates an empty file in HDFS.

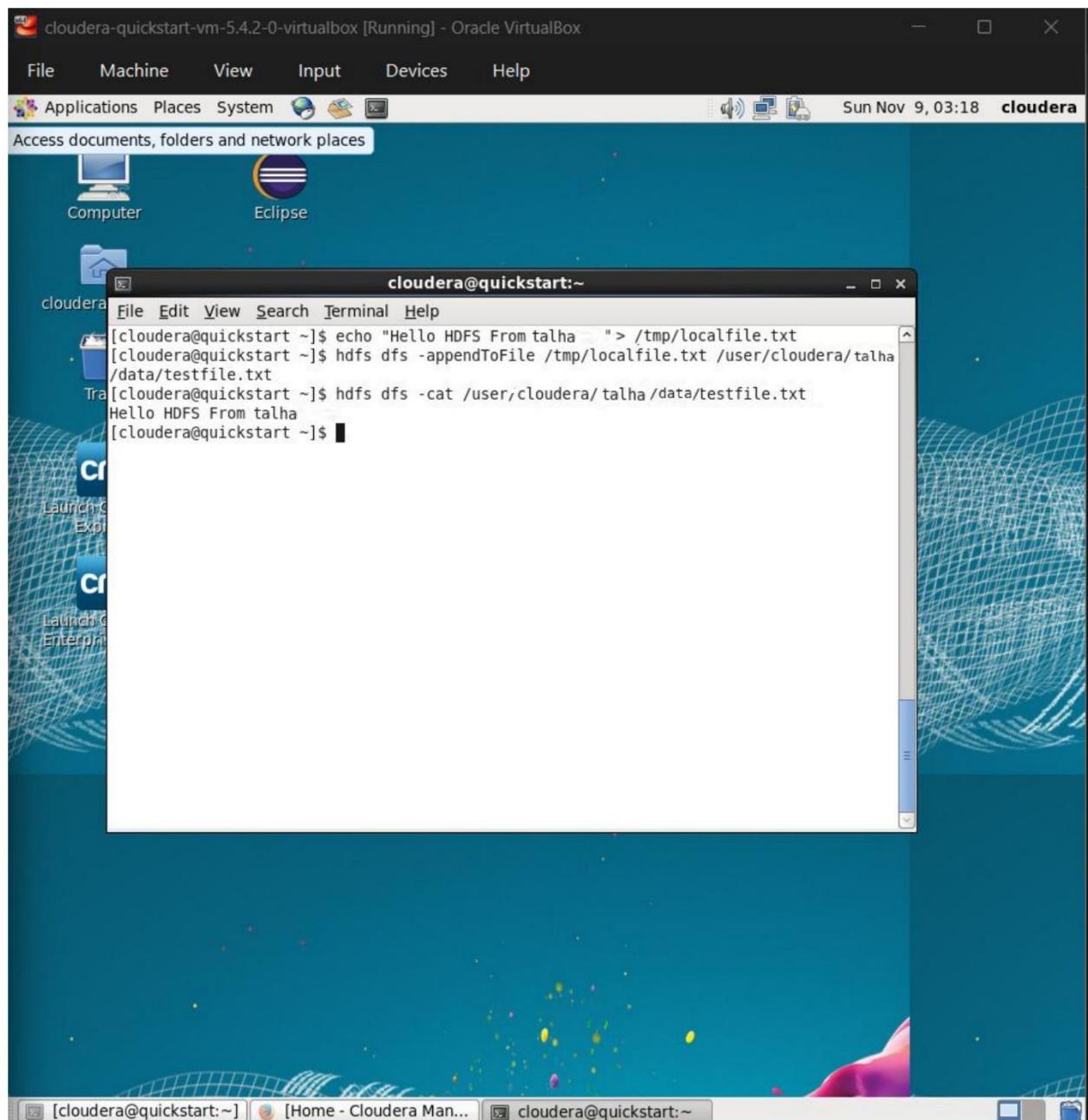
```
hdfs dfs -touchz /user/cloudera/talha/data/testfile.txt
```



3. appendToFile

Appends local file content to an HDFS file.

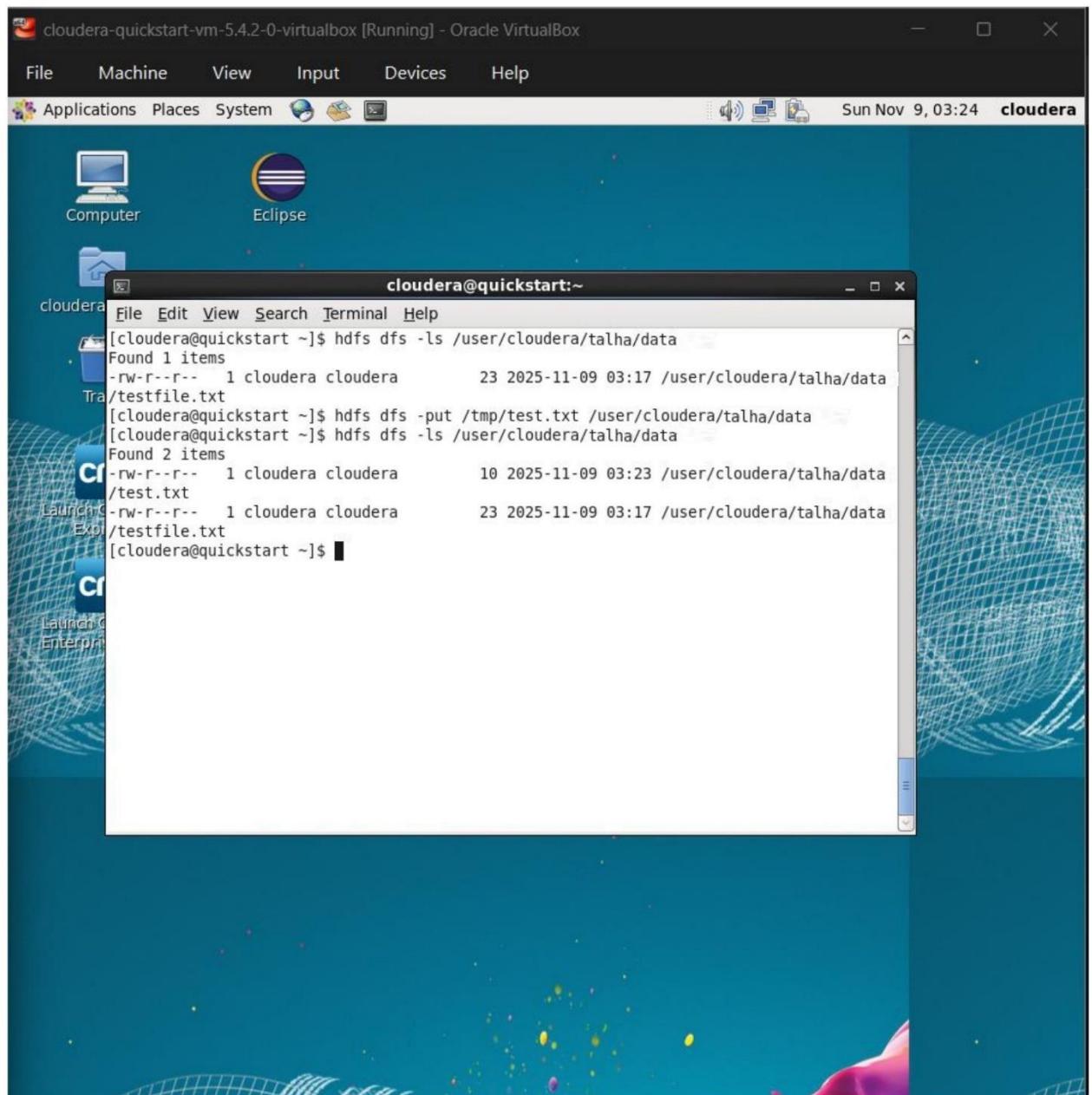
```
hdfs dfs -appendToFile /tmp/localfile.txt /user/cloudera/talha/data/testfile.txt
```



4. put

Copies file from local system to HDFS.

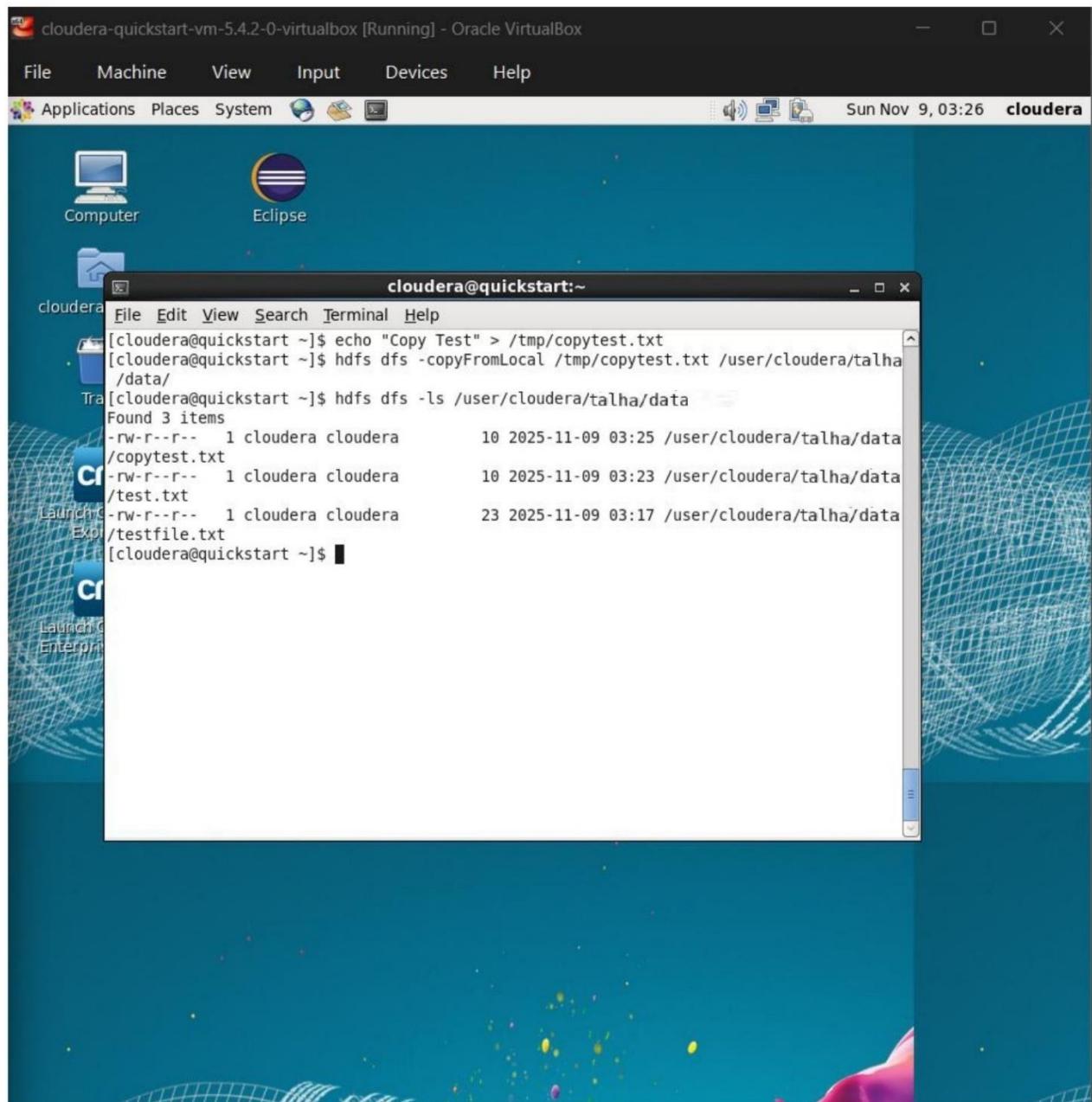
```
hdfs dfs -put /tmp/test.txt /user/cloudera/talha/data/
```



5. copyFromLocal

Alternate command to copy file from local to HDFS.

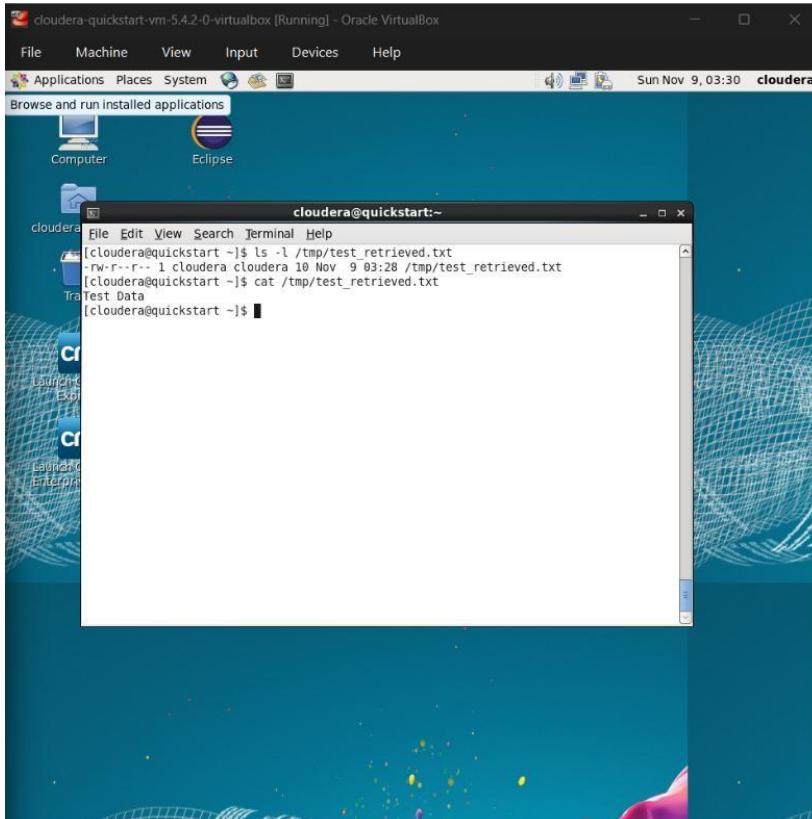
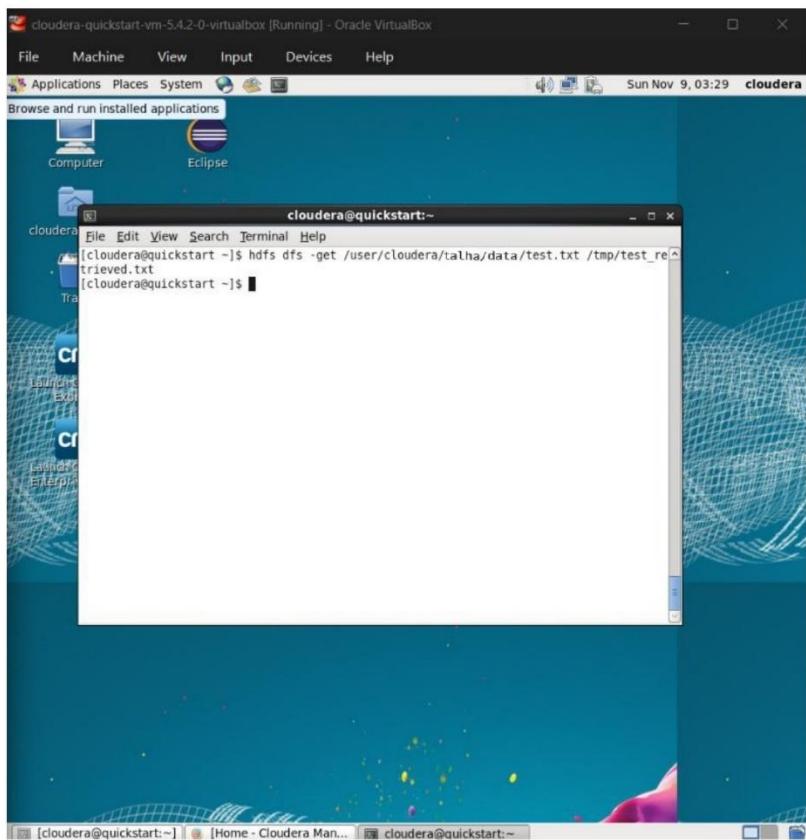
```
hdfs dfs -copyFromLocal /tmp/copytest.txt /user/cloudera/talha/data/
```



6. get

Copies file from HDFS to local system.

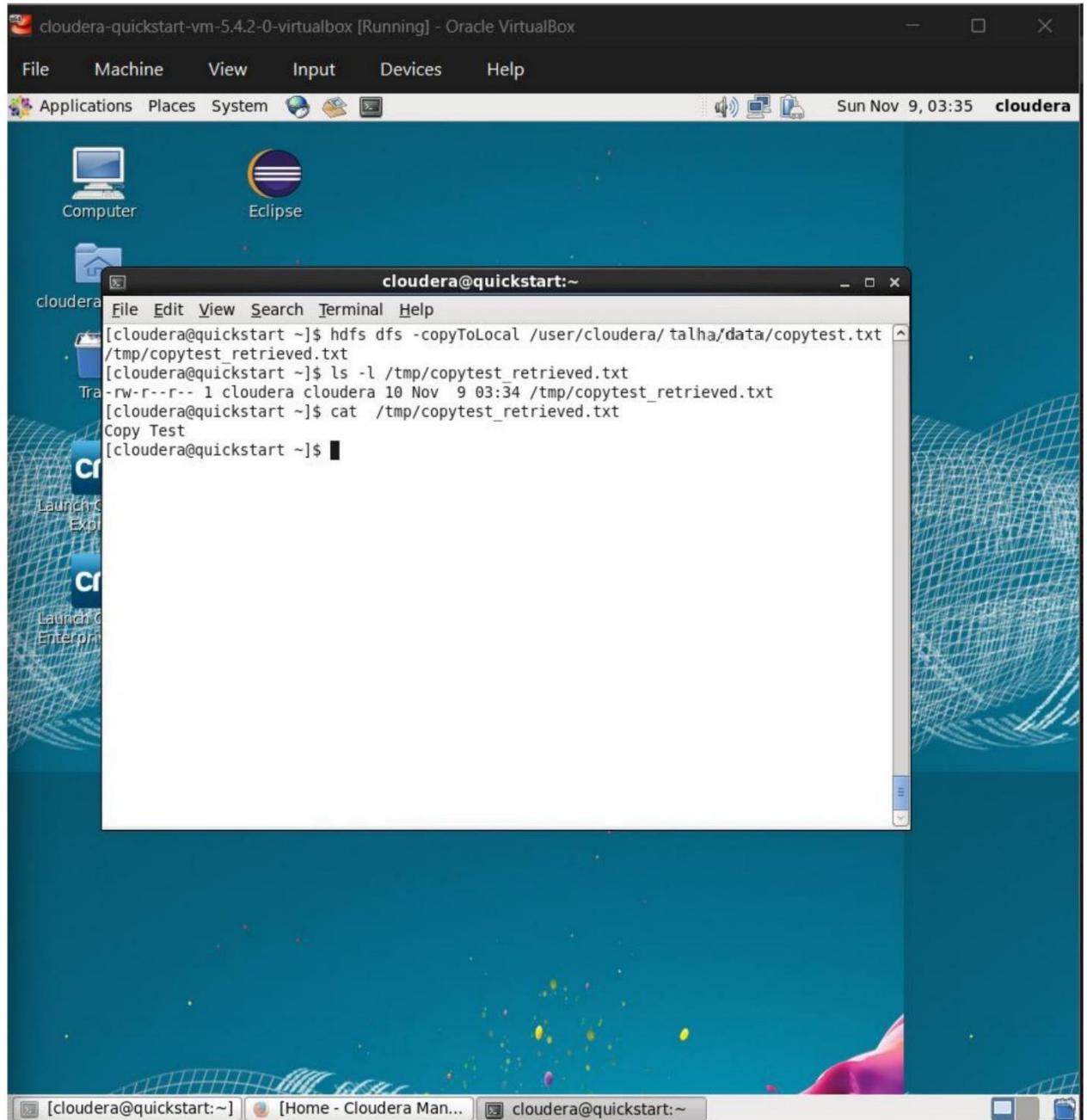
```
hdfs dfs -get /user/cloudera/talha/data/test.txt /tmp/test_retrieved.txt
```



7. copyToLocal

Alternate command to copy file from HDFS to local.

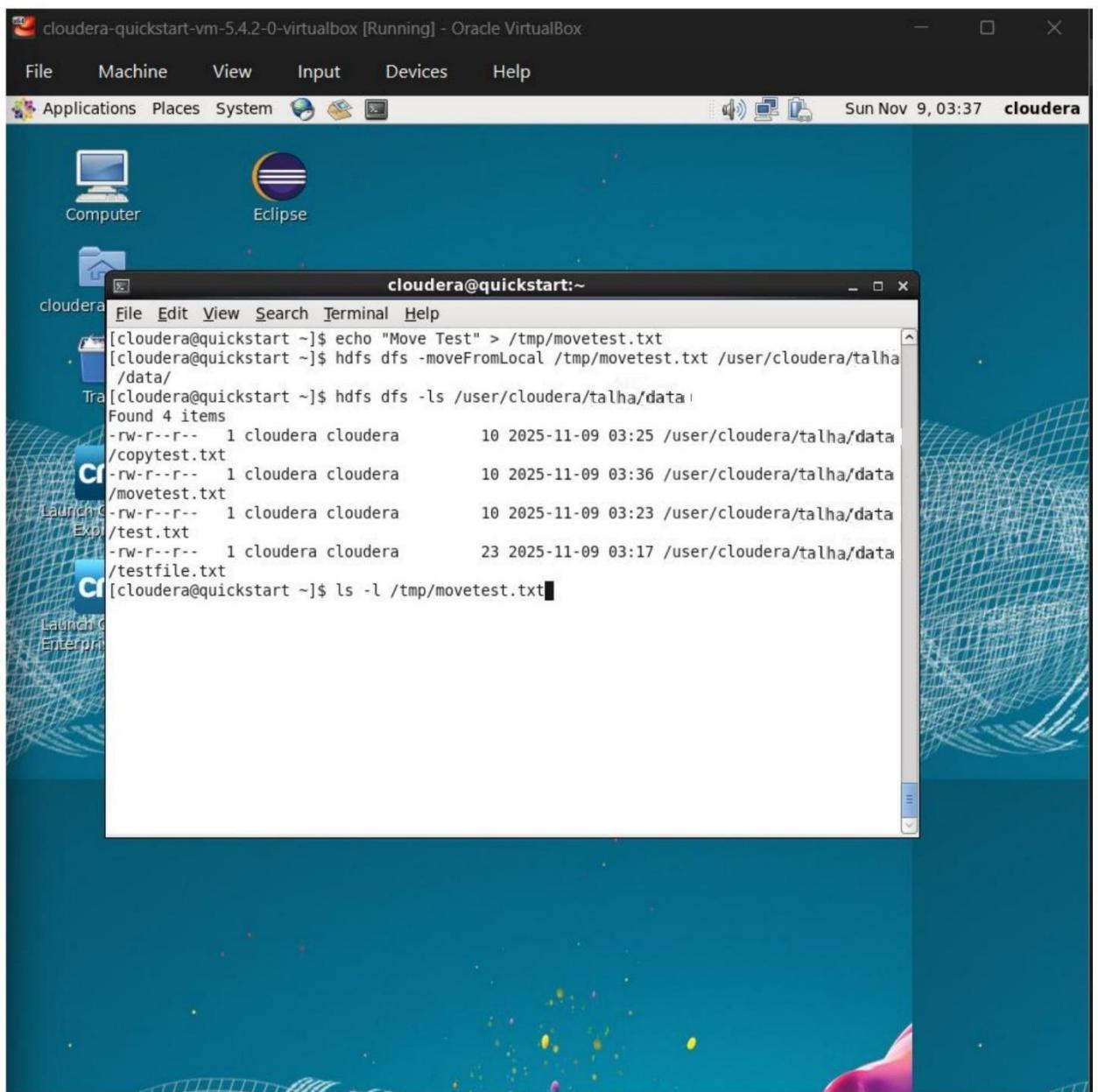
```
hdfs dfs -copyToLocal /user/cloudera/talha/data/copytest.txt  
/tmp/copytest_retrieved.txt
```



8. moveFromLocal

Moves a local file to HDFS (deletes local copy).

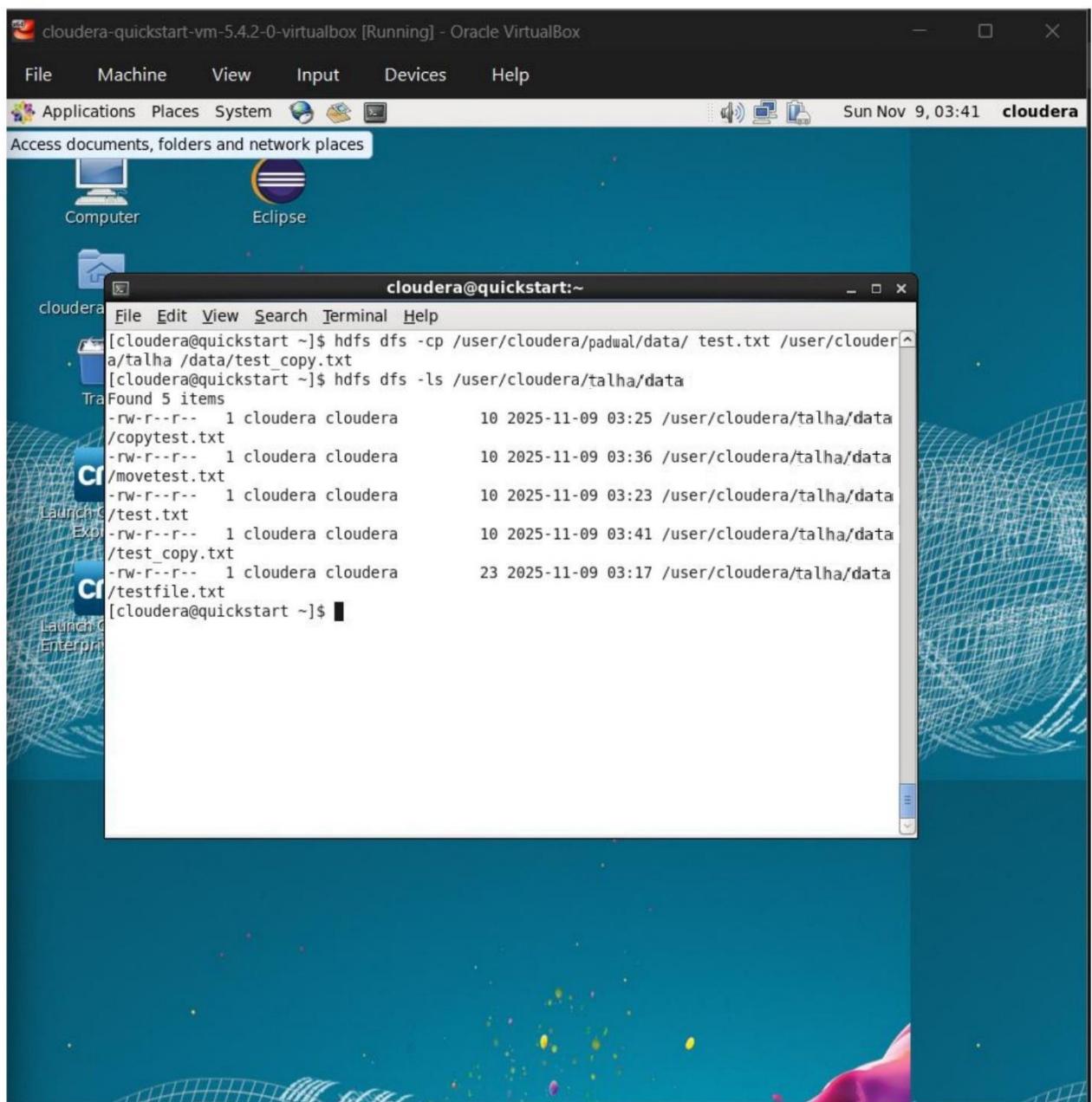
```
hdfs dfs -moveFromLocal /tmp/movetest.txt /user/cloudera/talha/data/
```



9. cp

Copies a file within HDFS directories.

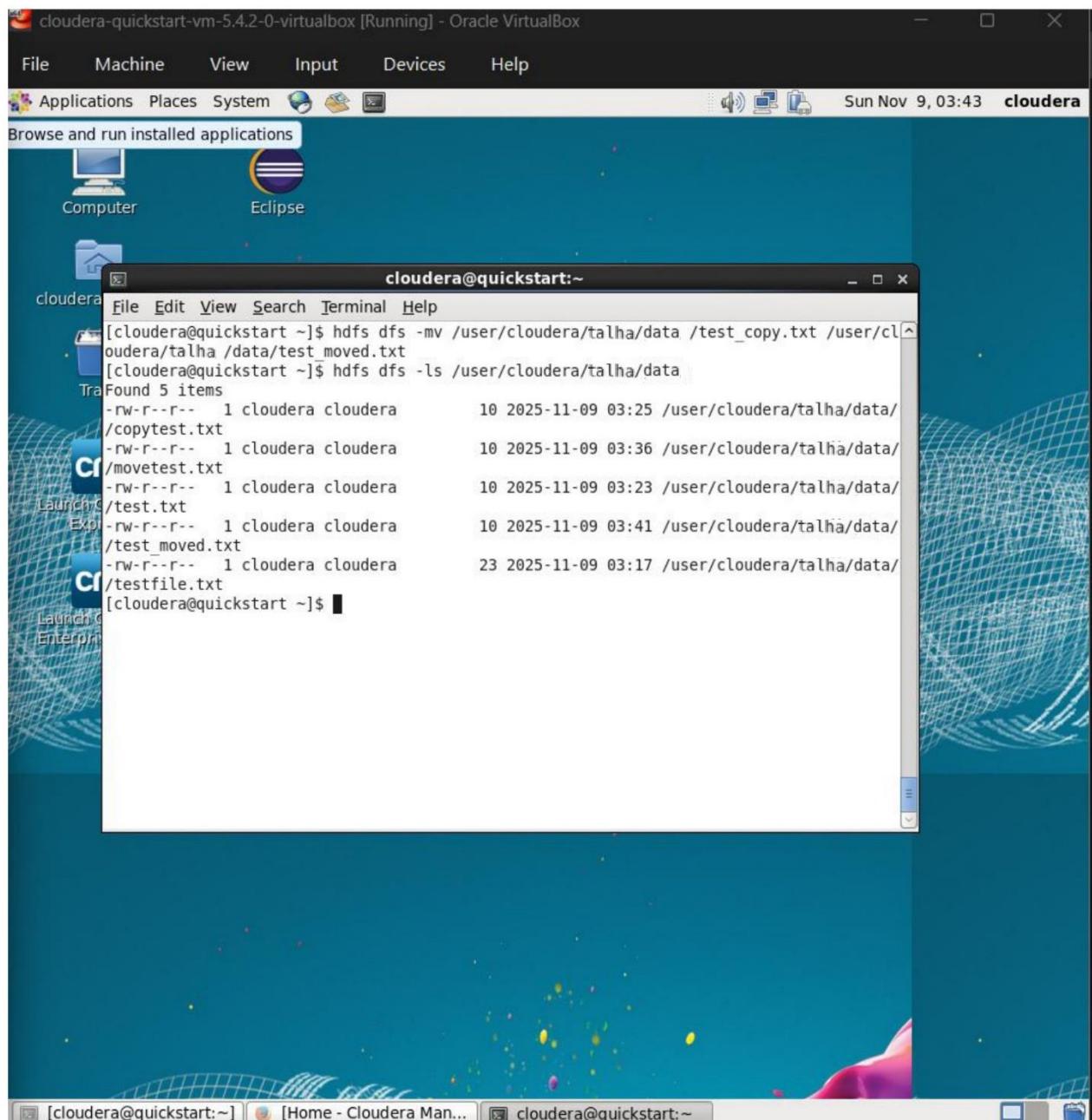
```
hdfs dfs -cp /user/cloudera/talha/data/test.txt
          /user/cloudera/talha/data/test_copy.txt
```



10. mv

Moves or renames a file within HDFS.

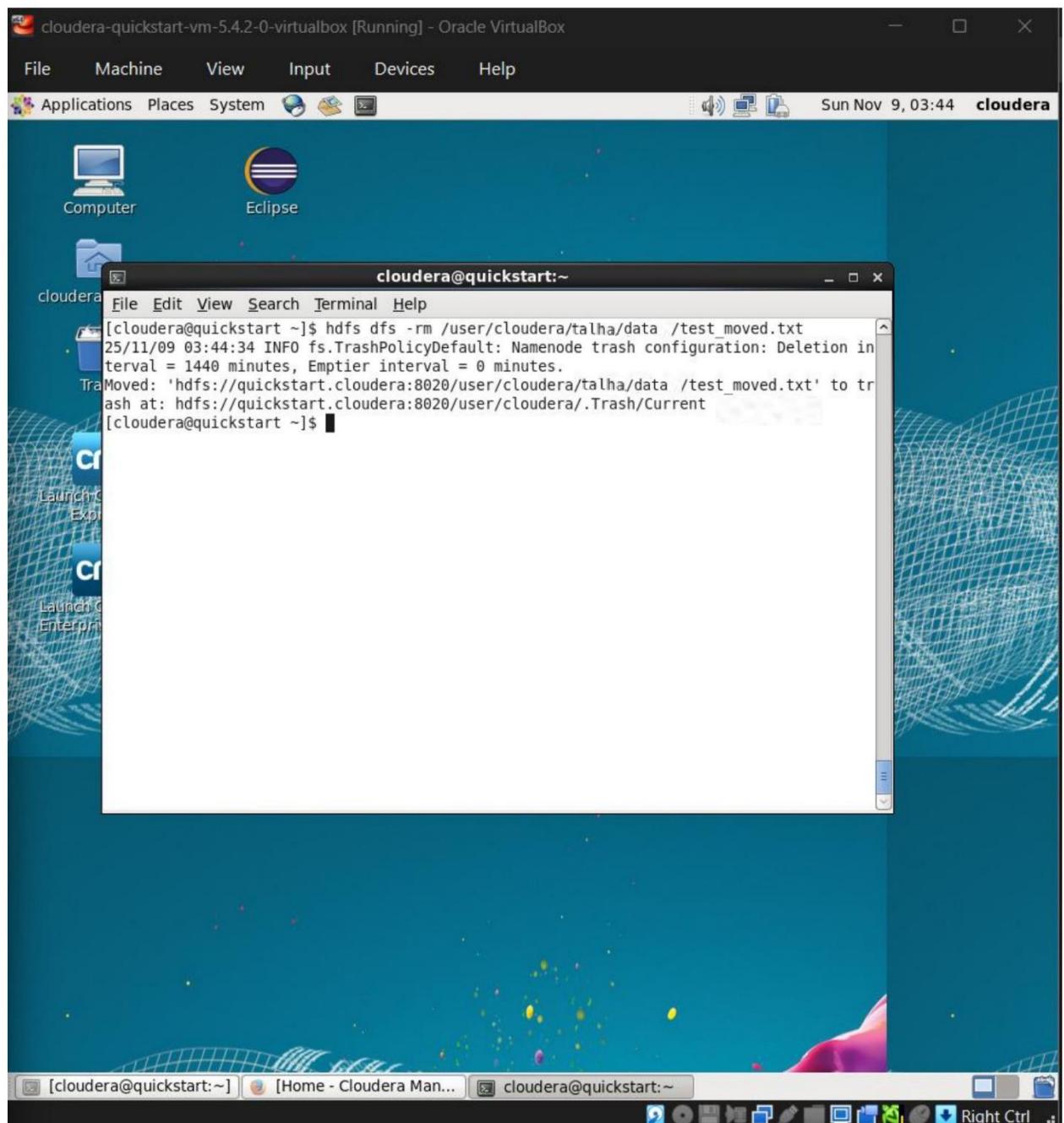
```
hdfs dfs -mv /user/cloudera/talha/data/test_copy.txt
          /user/cloudera/talha/data/test_moved.txt
```



11. rm

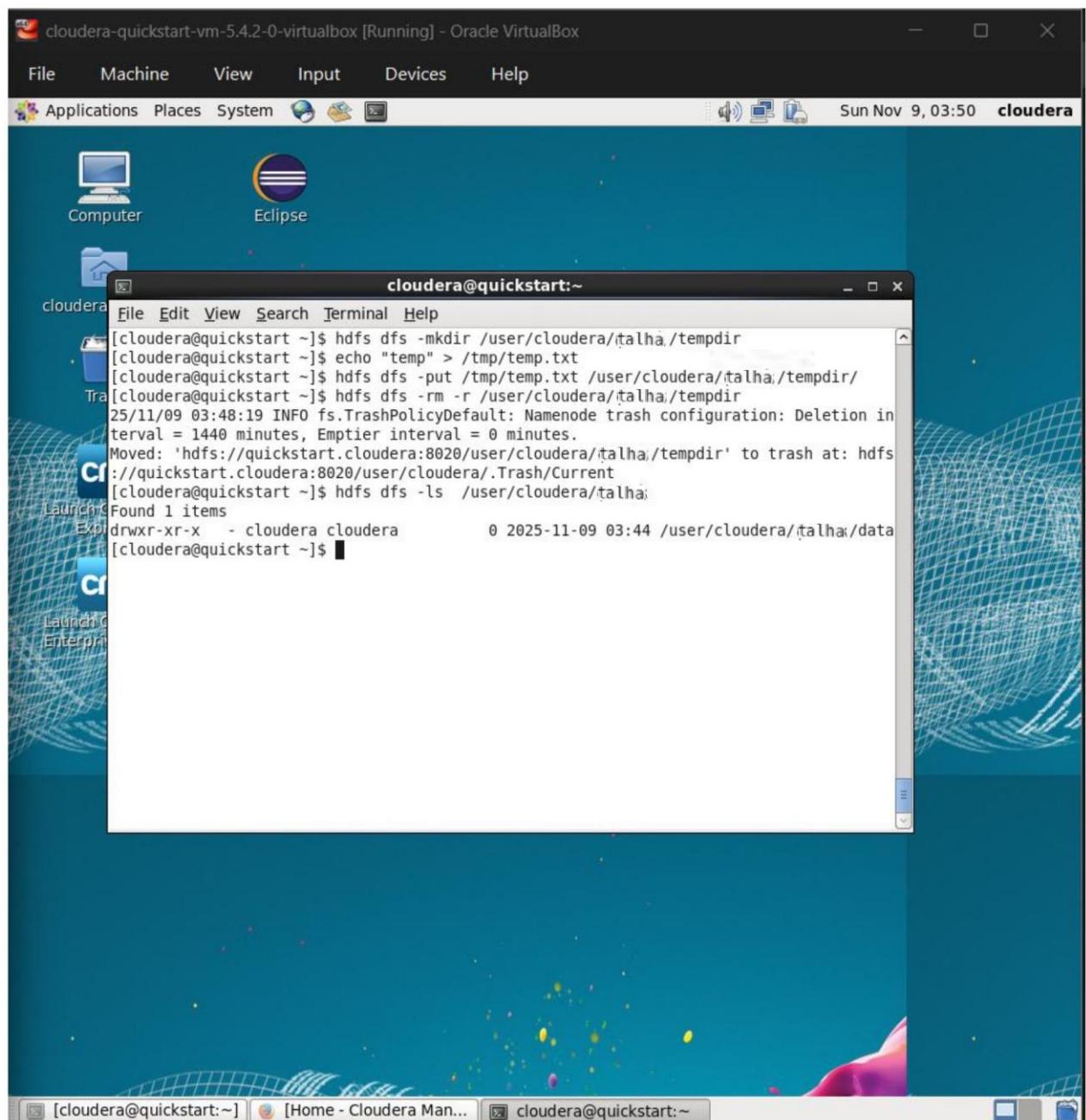
Deletes a specific file from HDFS.

```
hdfs dfs -rm /user/cloudera/talha/data/test_moved.txt
```



12. rm -r

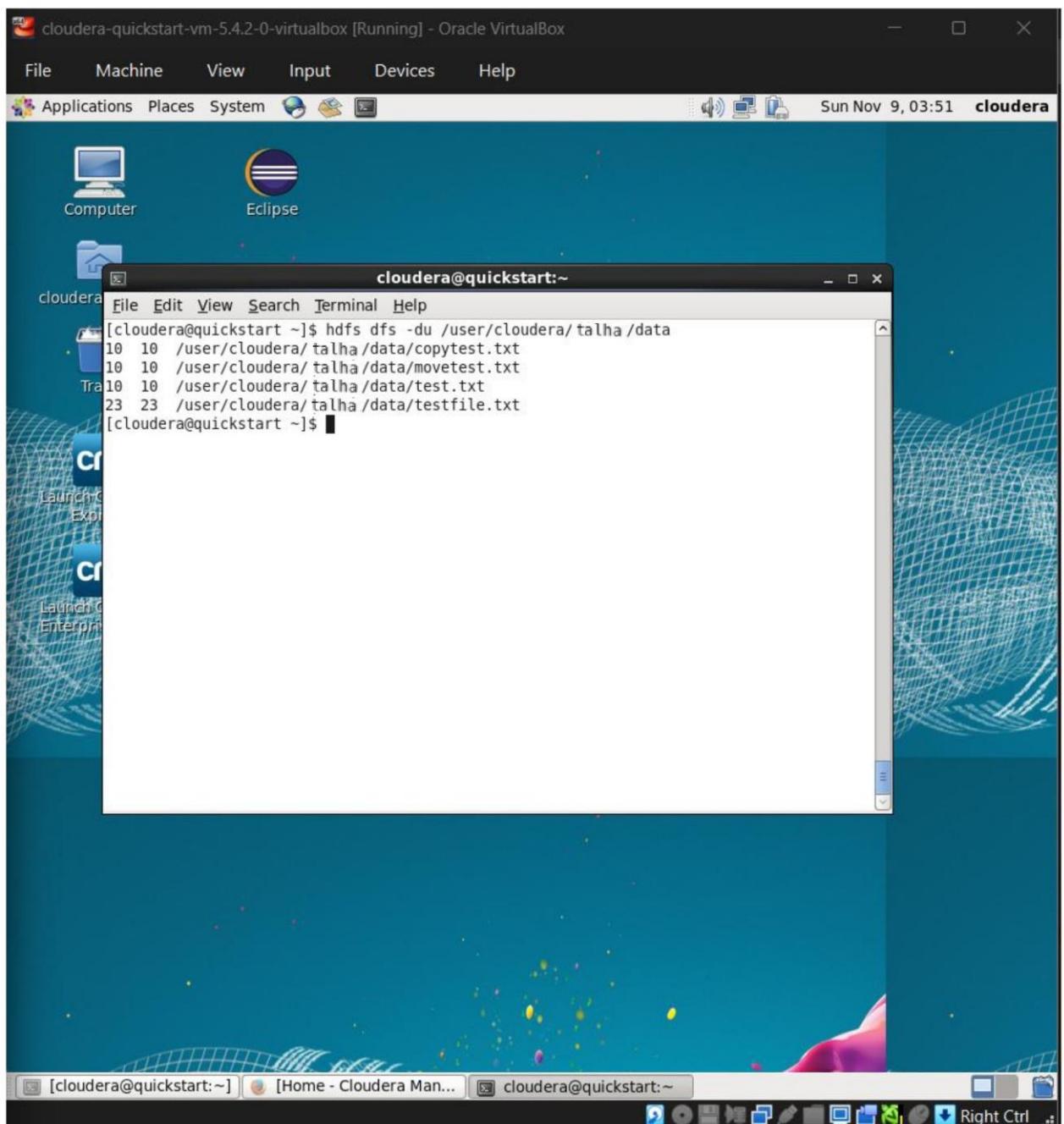
Deletes a directory recursively from HDFS.



13. du

Displays disk usage for each file in a directory.

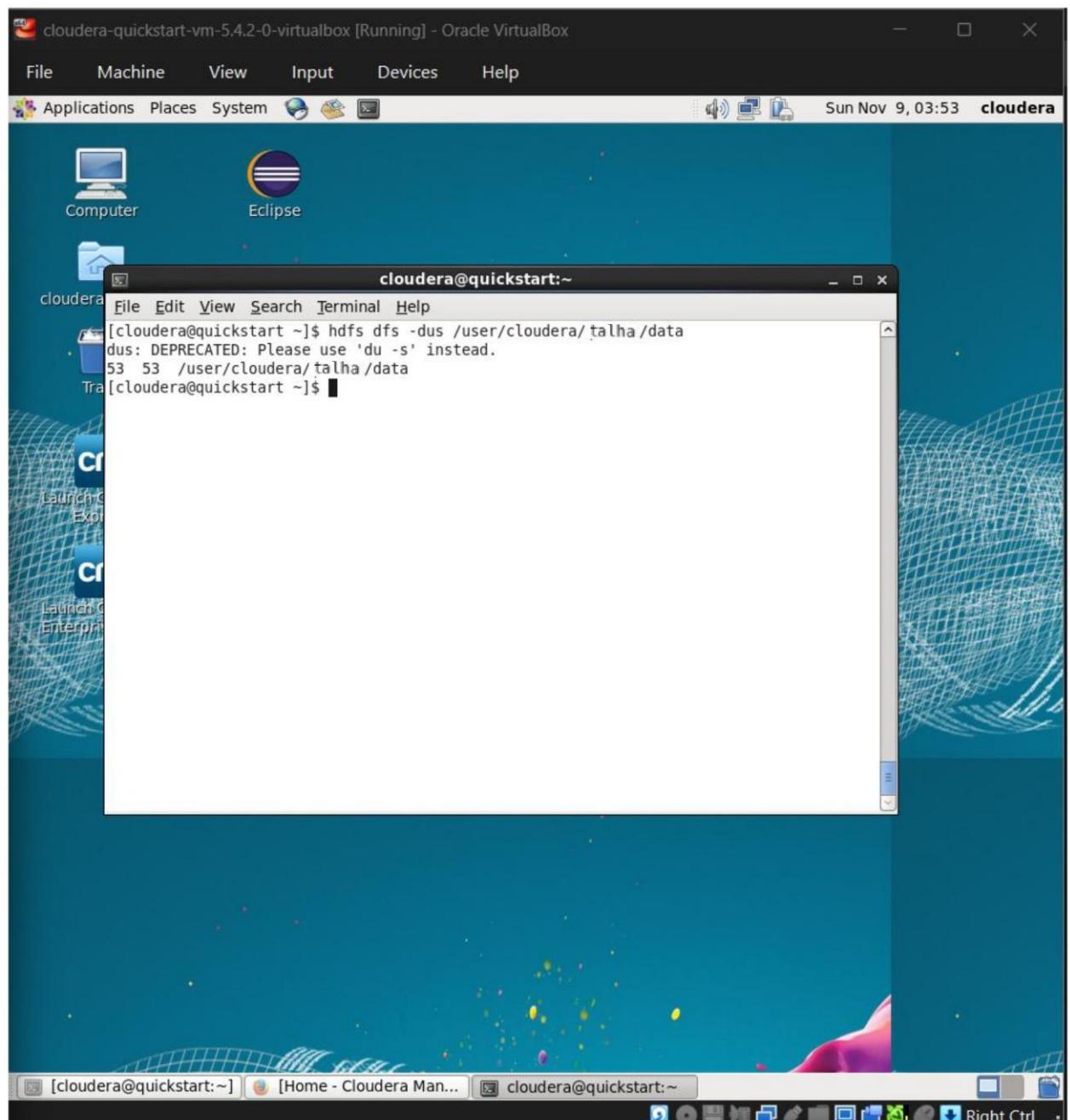
```
hdfs dfs -du /user/cloudera/talha/data
```



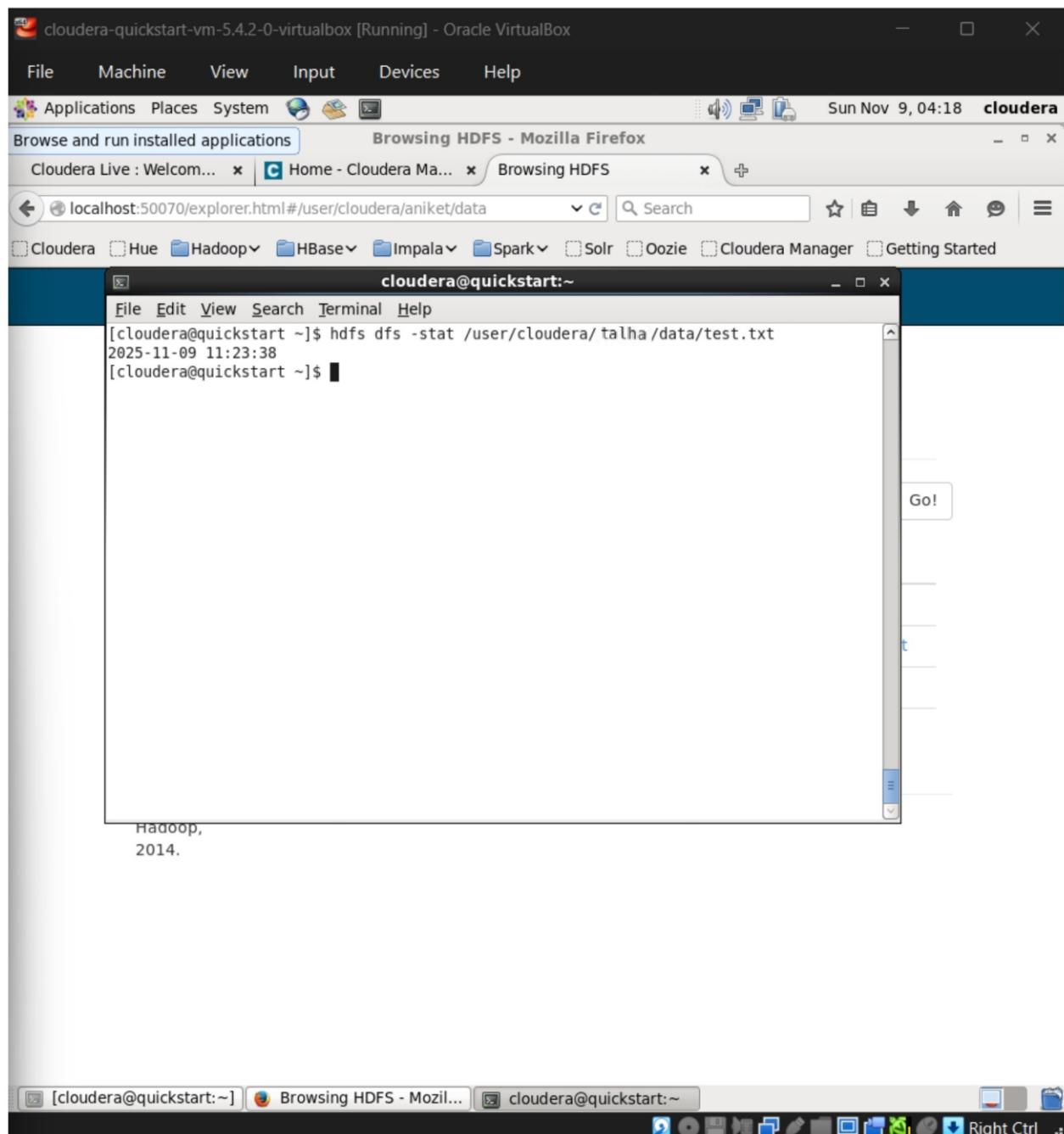
14. **dus**

Displays total disk usage summary for a directory.

```
hdfs dfs -dus /user/cloudera/talha/data
```



15. Displays detailed **status information** about a file or directory in HDFS



The screenshot shows a Firefox browser window titled "Browsing HDFS - Mozilla Firefox". The address bar displays "localhost:50070/explorer.html#". The main content area is titled "Browse Directory" and contains a table listing directory entries:

Permission	Owner	Group	Size	Replication	Block Size	Name
drwxr-xr-x	hbase	supergroup	0 B	0	0 B	hbase
drwxr-xr-x	solr	solr	0 B	0	0 B	solr
drwxrwxrwx	hdfs	supergroup	0 B	0	0 B	tmp
drwxr-xr-x	hdfs	supergroup	0 B	0	0 B	user
drwxr-xr-x	hdfs	supergroup	0 B	0	0 B	var

Below the table, a note reads: "Hadoop, 2014."

Browsing HDFS - Mozilla Firefox

Cloudera Live : Welcom... | C Home - Cloudera Ma... | Browsing HDFS

localhost:50070/explorer.html#/user/cloudera/aniket/data

Search

Cloudera | Hue | Hadoop | HBase | Impala | Spark | Solr | Oozie | Cloudera Manager | Getting Started

Hadoop | Overview | Datanodes | Snapshot | Startup Progress | Utilities

Browse Directory

/user/cloudera/aniket/data

Go!

Permission	Owner	Group	Size	Replication	Block Size	Name
-rw-r--r--	cloudera	cloudera	10 B	1	128 MB	copytest.txt
-rw-r--r--	cloudera	cloudera	10 B	1	128 MB	movetest.txt
-rw-r--r--	cloudera	cloudera	10 B	1	128 MB	test.txt
-rw-r--r--	cloudera	cloudera	23 B	1	128 MB	testfile.txt

Hadoop,
2014.

Practical 2: Map Reduce

1. Write a program in Map Reduce for Word Count operation.

Steps:

Step 1: Open virtual box and then start cloudera quickstart

Step 2: Open eclipse present on the cloudera desktop

Step 3: Create java project

File->New-> Java Project
Give project name:
WordCount Click Next

Step 4: Add Hadoop libraries to project

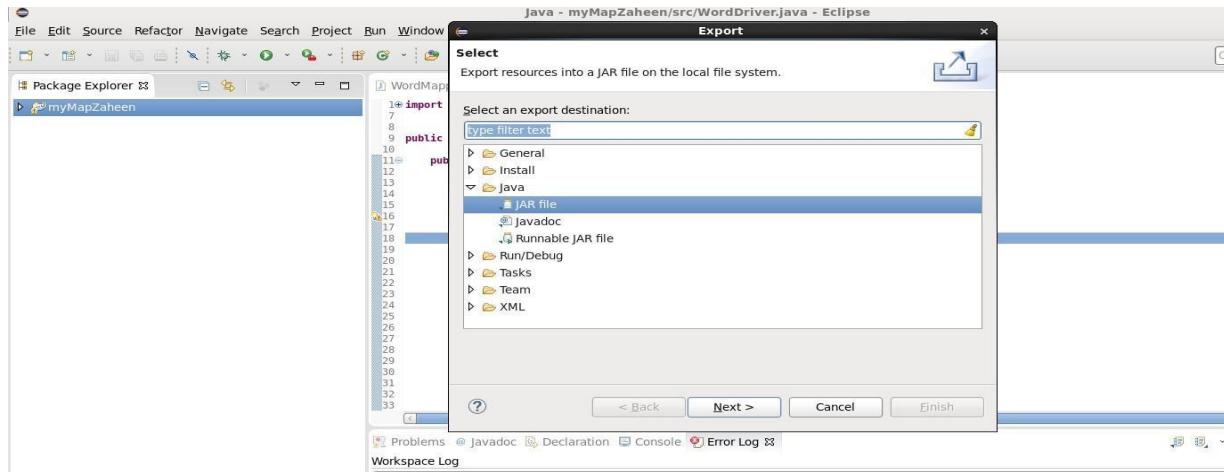
Select Libraries tab->click on Add External Jars File System->user->lib->Hadoop
Select all library(jar) files-
>ok Again click on Add External Jars
File System->user->lib->Hadoop->client
Select all library(jar) files from client->ok->finish

Step 5: Write java code for word count

Right click on src folder of project
WordCount New->class
Write class name
WordMapper Click Finish
Write the code for WordMapper
Same way create classes SumReducer and WordCountDriver

Step 6: Export the project as jar

Right click on project WordCount and select Export>> Java>>JAR file>>Next



Select the export destination-Click browse-give file name wordcount.jar

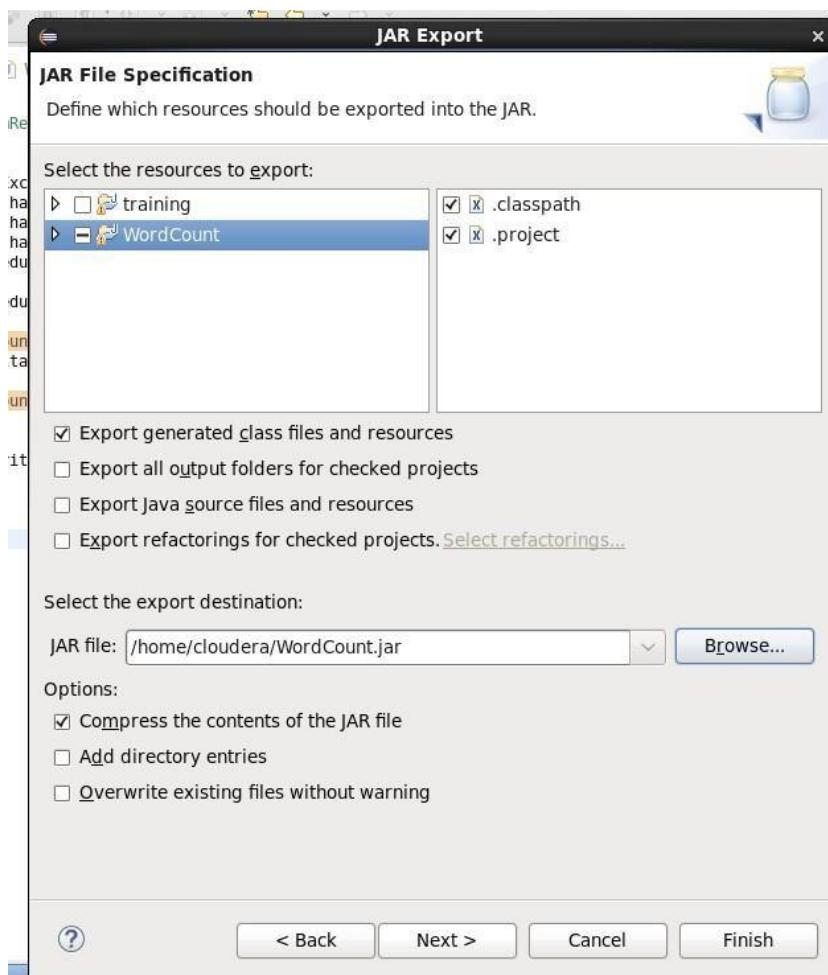
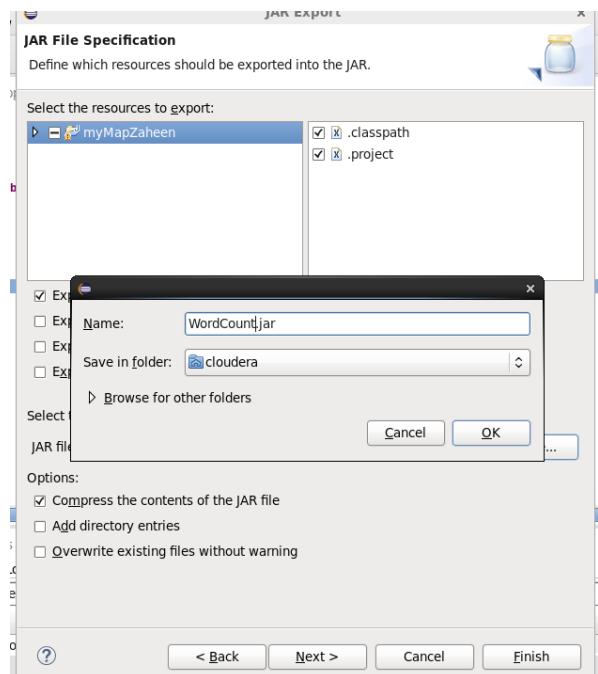
Click ok>>Finish>>ok

Verify the jar file

Verify the jar file through command line open terminal give command
ls

Step 7: Move the jar file to the Hadoop file

system **hdfs dfs -put wordcount.jar**
/user/cloudera hdfs dfs -ls



Step 8: Create the input file for the MapReduce

program Command: **cat > myInputFile.txt**

Welcome to the NMITD

MCA I am persuing MCA

in NMITD

Enter data in input file and press enter and ctrl z

Command: **cat myInputFile.txt**

Step 9: Move the input file to the Hadoop file system

hdfs dfs -put myInputFile.txt

/user/cloudera hdfs dfs -ls

hdfs dfs -ls / (here / indicates root directory of Hadoop file system(hdfs)) step 10: Run mapreduce program on Hadoop

syntax: hadoop jar jarfilename.jar classname inputfilename.txt outputfilename

command: **hadoop jar wordcount.jar WordCountDriver myInputFile.txt myOutput**

Start Yarn if not started and showing connection refused

```
[cloudera@quickstart ~]$ hadoop jar Wordcount.jar newjava.WordCountDriver myInputFile.txt myOutput
25/08/06 22:00:15 INFO client.RMProxy: Connecting to ResourceManager at quickstart.cloudera/127.0.0.1:8032
25/08/06 22:00:15 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
25/08/06 22:00:16 INFO input.FileInputFormat: Total input paths to process : 1
25/08/06 22:00:16 WARN hdfs.DFSClient: Caught exception
java.lang.InterruptedException
        at java.lang.Object.wait(Native Method)
        at java.lang.Thread.join(Thread.java:1281)
        at java.lang.Thread.join(Thread.java:1355)
        at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.closeResponder(DFSOutputStream.java:967)
        at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.endBlock(DFSOutputStream.java:705)
        at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.run(DFSOutputStream.java:894)
25/08/06 22:00:16 INFO mapreduce.JobSubmitter: number of splits:1
25/08/06 22:00:16 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1754542791586_0001
25/08/06 22:00:17 INFO impl.YarnClientImpl: Submitted application application_1754542791586_0001
25/08/06 22:00:17 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1754542791586_0001/
25/08/06 22:00:17 INFO mapreduce.Job: Running job: job_1754542791586_0001
25/08/06 22:00:27 INFO mapreduce.Job: Job job_1754542791586_0001 running in uber mode : false
25/08/06 22:00:27 INFO mapreduce.Job: map 0% reduce 0%
25/08/06 22:00:42 INFO mapreduce.Job: map 100% reduce 0%
25/08/06 22:00:56 INFO mapreduce.Job: map 100% reduce 100%
25/08/06 22:00:57 INFO mapreduce.Job: Job job_1754542791586_0001 completed successfully
25/08/06 22:00:57 INFO mapreduce.Job: Counters: 49
      File System Counters
          FILE: Number of bytes read=144
          FILE: Number of bytes written=295053
          FILE: Number of read operations=0
          FILE: Number of large read operations=0
          FILE: Number of write operations=0
          HDFS: Number of bytes read=232
          HDFS: Number of bytes written=74
          HDFS: Number of read operations=6
          HDFS: Number of large read operations=0
```

```
File System Counters
  FILE: Number of bytes read=144
  FILE: Number of bytes written=295053
  FILE: Number of read operations=0
  FILE: Number of large read operations=0
  FILE: Number of write operations=0
  HDFS: Number of bytes read=232
  HDFS: Number of bytes written=74
  HDFS: Number of read operations=6
  HDFS: Number of large read operations=0
  HDFS: Number of write operations=2

Job Counters
  Launched map tasks=1
  Launched reduce tasks=1
  Data-local map tasks=1
  Total time spent by all maps in occupied slots (ms)=6788608
  Total time spent by all reduces in occupied slots (ms)=5730304
  Total time spent by all map tasks (ms)=13259
  Total time spent by all reduce tasks (ms)=11192
  Total vcore-milliseconds taken by all map tasks=13259
  Total vcore-milliseconds taken by all reduce tasks=11192
  Total megabyte-milliseconds taken by all map tasks=6788608
  Total megabyte-milliseconds taken by all reduce tasks=5730304

Map-Reduce Framework
  Map input records=4
  Map output records=21
  Map output bytes=186
  Map output materialized bytes=140
  Input split bytes=126
  Combine input records=0
  Combine output records=0
  Reduce input groups=10
  Reduce shuffle bytes=140
  Reduce input records=21
  Reduce output records=10
  Spilled Records=42
  Shuffled Maps =1
  Failed Shuffles=0
  Merged Map outputs=1
  GC time elapsed (ms)=795
  CPU time spent (ms)=1520
  Physical memory (bytes) snapshot=240349184
  Virtual memory (bytes) snapshot=1408757760
  Total committed heap usage (bytes)=101449728

Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
```

```

Shuffle Errors
  BAD_ID=0
  CONNECTION=0
  IO_ERROR=0
  WRONG_LENGTH=0
  WRONG_MAP=0
  WRONG_REDUCE=0
File Input Format Counters
  Bytes Read=106
File Output Format Counters
  Bytes Written=74
[cloudera@quickstart ~]$ hdfs dfs -ls /user/cloudera/myOutput
Found 2 items
-rw-r--r--  1 cloudera cloudera      0 2025-08-06 22:00 /user/cloudera/myOutput/_SUCCESS
-rw-r--r--  1 cloudera cloudera    74 2025-08-06 22:00 /user/cloudera/myOutput/part-r-00000
[cloudera@quickstart ~]$ hdfs dfs -cat /user/cloudera/myOutput/part-r-00000
I          2
MCA        4
VMID       3
VMITWelcome   1
Welcome 1
am         2
in         2
persuing    2
the         2
to          2
[cloudera@quickstart ~]$ █

```

Step 11: view output directory**hdfs dfs –ls****hdfs dfs -ls /user/cloudera/myOutput****step 12:** view the output file**hdfs dfs -cat /user/cloudera/myOutput/part-r-00000****Code:****WordMapper.java**

```

import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;
public class WordMapper extends Mapper<LongWritable, Text, Text, IntWritable> { @Override
    public void map(LongWritable key, Text value, Context context) throws
IOException, InterruptedException {
        String line = value.toString();
        for (String word : line.split("\\W+")) {
            if (word.length() > 0) {
                context.write(new Text(word), new IntWritable(1));
            }
        }
    }
}

```

SumReducer.java

```
import java.io.IOException;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Reducer;
public class SumReducer extends Reducer<Text, IntWritable, Text, IntWritable>
{ @Override
    public void reduce(Text key, Iterable<IntWritable> values, Context context)
                      throws
                        IOException, InterruptedException {
        int wordCount = 0;
        for (IntWritable value : values) {
            wordCount += value.get();
        }
        context.write(key, new IntWritable(wordCount));
    }
}
```

WordCountDriver.java

```
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.Job;
public class WordCountDirver {
    public static void main(String[] args) throws
        Exception { if (args.length != 2) {
            System.out.printf("Usage: WordCount <input dir> <output
dir>\n"); System.exit(-1);
        }
        Job job = new Job();
        job.setJarByClass(WordCountDriver.class);
        job.setJobName("Word Count");
        FileInputFormat.setInputPaths(job, new
Path(args[0])); FileOutputFormat.setOutputPath(job,
new Path(args[1]));
        job.setMapperClass(WordMapper.class);
        job.setReducerClass(SumReducer.class);
        job.setMapOutputKeyClass(Text.class);
        job.setMapOutputValueClass(IntWritable.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        boolean success = job.waitForCompletion(true);
        System.exit(success ? 0 : 1);
    }
}
```

3

3

2. Write a program in Map Reduce for Union operation.

Code:

```
import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.Text;
import
org.apache.hadoop.mapreduce.Job;
import
org.apache.hadoop.mapreduce.Mapper;
import
org.apache.hadoop.mapreduce.Reducer;
import
org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import
org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import java.io.IOException;
public class Union {
    // Mapper: Just emits each line as key
    public static class UnionMapper extends Mapper<Object, Text, Text,
    Text> { private static final Text empty = new Text("");
    public void map(Object key, Text value, Context context)
        throws IOException, InterruptedException {
        context.write(value, empty);
    }
}
// Reducer: Writes out each unique key once (union)
```

```
public static class UnionReducer extends Reducer<Text, Text, Text> { public void reduce(Text key, Iterable<Text> values, Context context) throws IOException, InterruptedException { context.write(key, null); } }
```

```
3
```

```
3
```

```
public static void main(String[] args) throws Exception { // Check for input arguments if (args.length < 2) { System.err.println("Usage: Union <input path> <output path>"); System.exit(-1); }
```

```
3
```

```
Configuration conf = new Configuration(); Job job = Job.getInstance(conf, "Union Operation"); job.setJarByClass(Union.class); job.setMapperClass(UnionMapper.class); job.setReducerClass(UnionReducer.class); job.setOutputKeyClass(Text.class); job.setOutputValueClass(Text.class); FileInputFormat.addInputPath(job, new Path(args[0])); FileOutputFormat.setOutputPath(job, new Path(args[1])); System.exit(job.waitForCompletion(true) ? 0 : 1); }
```

```
3
```

```
3
```

Steps to execute:

Step 1: Create two sample files in cloudera.

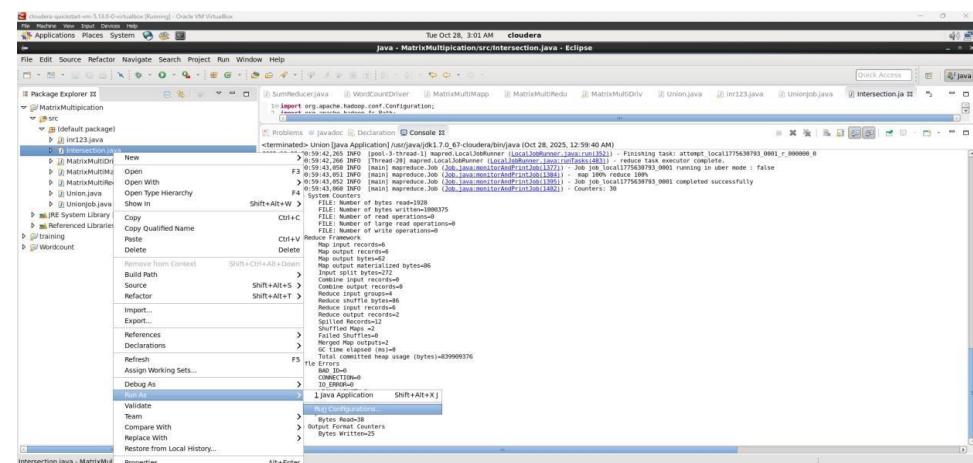
Step 2: Prepare Local Input File. In your Cloudera terminal



Step 3: Run Configuration in Eclipse

Right-click on Intersection.java → Run As → Run

Configurations... In Arguments → Program arguments, enter:



Output:

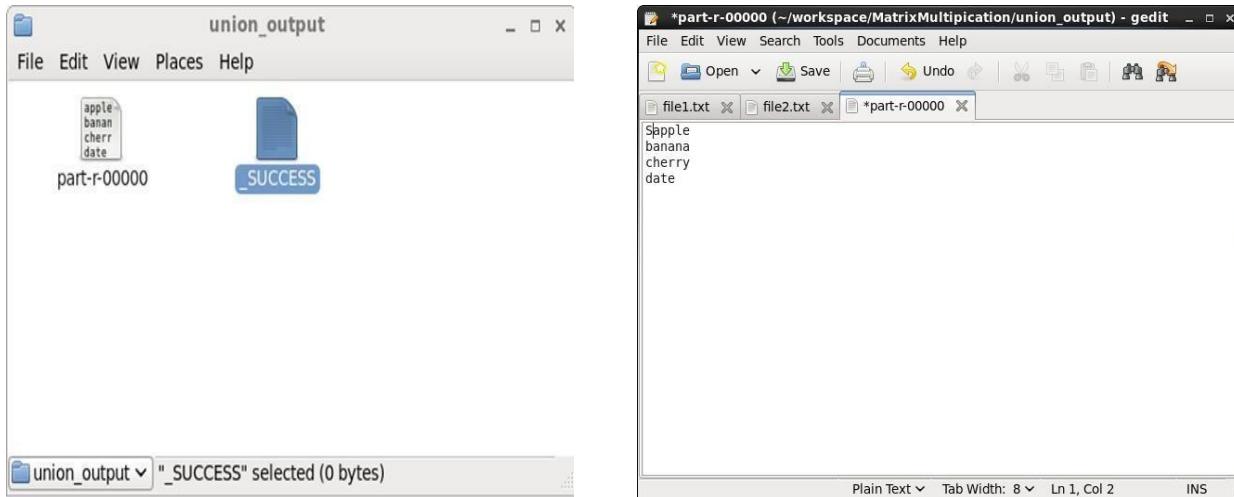
```
<terminated> [Java Application] /usr/java/jdk1.7.0_67-cloudera/bin/java [Oct 28, 2025, 12:44:00 AM]

2025-10-28 00:44:01,614 WARN [main] util.NativeCodeLoader [NativeCodeLoader.java<clinit>(62)] - Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
2025-10-28 00:44:01,964 INFO [main] Configuration.deprecation [Configuration.java<warnDeprecatedIfRepreated>(204)] - session.id is deprecated. Instead, use dfs.metrics.session-id
2025-10-28 00:44:01,965 INFO [main] jvm.JvmMetrics [JvmMetrics.java<init>(76)] - Initializing JVM Metrics with processName=JobTracker, sessionId=
2025-10-28 00:44:02,089 WARN [main] mapreduce.JobResourceUploader [JobResourceUploader.java<uploadFiles>(64)] - Hadoop command-line option parsing not performed. Implement the Tool interface and execute Job with a ToolRunner to access configuration.
2025-10-28 00:44:02,189 INFO [main] InputFileInputFormat [InputFileInputFormat.java<init>(283)] - Total input paths to process: 2
2025-10-28 00:44:02,214 INFO [main] mapreduce.JobSubmitter [JobSubmitter.java<submitJobInternal>(202)] - number of splits:2
2025-10-28 00:44:02,282 INFO [main] mapreduce.JobSubmitter [JobSubmitter.java<printTokens>(291)] - Submitting tokens for job: job_local1773454917_0001
2025-10-28 00:44:02,435 INFO [main] mapreduce.Job [Job.java<submitJobInternal>(133)] - The url to track the job: http://localhost:8088/proxy/application_1773454917_0001/
2025-10-28 00:44:02,436 INFO [main] [Thread-0] mapred.LocalJobRunner [LocalJobRunner.java<createOutputCommitter>(488)] - OutputCommitter set in config null
2025-10-28 00:44:02,436 INFO [main] [Thread-0] mapred.LocalJobRunner [LocalJobRunner.java<createOutputCommitter>(488)] - OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
2025-10-28 00:44:02,439 INFO [main] [Thread-0] output.FileOutputCommitter [FileOutputCommitter.java<init>(124)] - File Output Committer Algorithm version is 1
2025-10-28 00:44:02,439 INFO [main] [Thread-0] output.FileOutputCommitter [FileOutputCommitter.java<init>(139)] - FileOutputCommitter skip cleanup _temporary folders under output directory:false, ignore cl
2025-10-28 00:44:02,440 INFO [main] [Thread-0] mapred.LocalJobRunner [LocalJobRunner.java<createOutputCommitter>(488)] - OutputCommitter is org.apache.hadoop.mapreduce.lib.output.FileOutputCommitter
2025-10-28 00:44:02,466 INFO [LocalJobRunner] [LocalJobRunner.java<run>(251)] - Starting task: attempt_local1773454917_0001_m_000000_0
2025-10-28 00:44:02,484 INFO [LocalJobRunner] [Map Task Executor #0] - Map Task Executor #0
2025-10-28 00:44:02,484 INFO [LocalJobRunner] [Map Task Executor #0] - output.FileOutputCommitter [FileOutputCommitter.java<init>(124)] - File Output Committer Algorithm version is 1
2025-10-28 00:44:02,484 INFO [LocalJobRunner] [Map Task Executor #0] - output.FileOutputCommitter [FileOutputCommitter.java<init>(139)] - FileOutputCommitter skip cleanup _temporary folders under output
2025-10-28 00:44:02,485 INFO [LocalJobRunner] [Map Task Executor #0] - mapred.MapTask [Task.java<initialize>(331)] - Using ResourceCalculatorHint{resourcesFree:[]}
2025-10-28 00:44:02,485 INFO [LocalJobRunner] [Map Task Executor #0] - mapred.MapTask [Task.java<initialize>(331)] - Using ResourceCalculatorHint{resourcesFree:[]}
2025-10-28 00:44:02,765 INFO [LocalJobRunner] [Map Task Executor #0] - mapred.MapTask [MapTask.java<setEquator>(1212)] - (EQUATOR) 0 kv 26214396(104857584)
2025-10-28 00:44:43,441 INFO [main] mapreduce.Job [Job.java<monitorAndPrintJob>(1384)] - map 100% reduce 100%
2025-10-28 00:44:43,442 INFO [main] mapreduce.Job [Job.java<monitorAndPrintJob>(1395)] - Job job_local1773454917_0001 completed successfully
2025-10-28 00:44:03,459 INFO [main] mapreduce.Job [Job.java<monitorAndPrintJob>(1402)] - Counters: 30

File System Counters
    FILE: Number of bytes read=1092
    FILE: Number of bytes written=1000096
    FILE: Number of read operations=0
    FILE: Number of large read operations=0
    FILE: Number of write operations=0

Map-Reduce Framework
    Map input records=6
    Map output records=6
    Map output bytes=44
    Map output materialized bytes=68
    Input split bytes=272
    Combine input records=0
    Combine output records=0
    Reduce input groups=4
    Reduce input bytes=68
    Reduce shuffle bytes=68
    Reduce input records=6
    Reduce output records=4
    Spilled Records=12
    Shuffled Maps =2
    Failed Shuffles=0
    Merged Map partitions=2
    GC time elapsed (ms)=1-A
```

A new output file is created which have all the elements of both the files.



3. Write a program in Map Reduce for Intersection operation.

Code:

```
Import org.apache.hadoop.conf.Configuration;  
  
import org.apache.hadoop.fs.Path;  
  
import org.apache.hadoop.io.IntWritable;  
  
import org.apache.hadoop.io.Text;  
  
import org.apache.hadoop.mapreduce.Job;  
  
import org.apache.hadoop.mapreduce.Mapper;  
  
import org.apache.hadoop.mapreduce.Reducer;  
  
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;  
  
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;  
  
import java.io.IOException;  
  
public class Intersection {  
  
    // Mapper: emit each line as a key,  
    value = 1 public static class  
    IntersectionMapper  
        extends Mapper<Object, Text, Text, IntWritable> {  
  
    private final static IntWritable one = new IntWritable(1);  
  
    public void map(Object key, Text value, Context  
        context) throws IOException, InterruptedException  
    {  
        context.write(value, one);  
    }  
}
```

```
// Reducer: count how many times each key (line) appears

public static class IntersectionReducer
    extends Reducer<Text, IntWritable, Text, Text> {

    public void reduce(Text key, Iterable<IntWritable> values, Context
        context) throws IOException, InterruptedException {

        int count = 0;
        for (IntWritable val : values) {
            count += val.get();
        }

        // Assuming there are 2 files → intersection = lines appearing at least twice
        if (count >= 2) {
            context.write(key, null);
        }
    }

    public static void main(String[] args) throws Exception {
        if (args.length < 2) {
            System.err.println("Usage: Intersection <input path> <output
                path>"); System.exit(-1);
        }

        Configuration conf = new Configuration();

        // Run in local mode for testing (optional)
        conf.set("fs.defaultFS", "file:///");
    }
}
```

```
conf.set("mapreduce.framework.name",
"local");

Job job = Job.getInstance(conf, "Intersection Operation");
job.setJarByClass(Intersection.class);
job.setMapperClass(IntersectionMapper.class);
job.setReducerClass(IntersectionReducer.class)
;
job.setOutputKeyClass(Text.class);
job.setOutputValueClass(IntWritable.class);

FileInputFormat.addInputPath(job, new Path(args[0]));
FileOutputFormat.setOutputPath(job, new Path(args[1]));

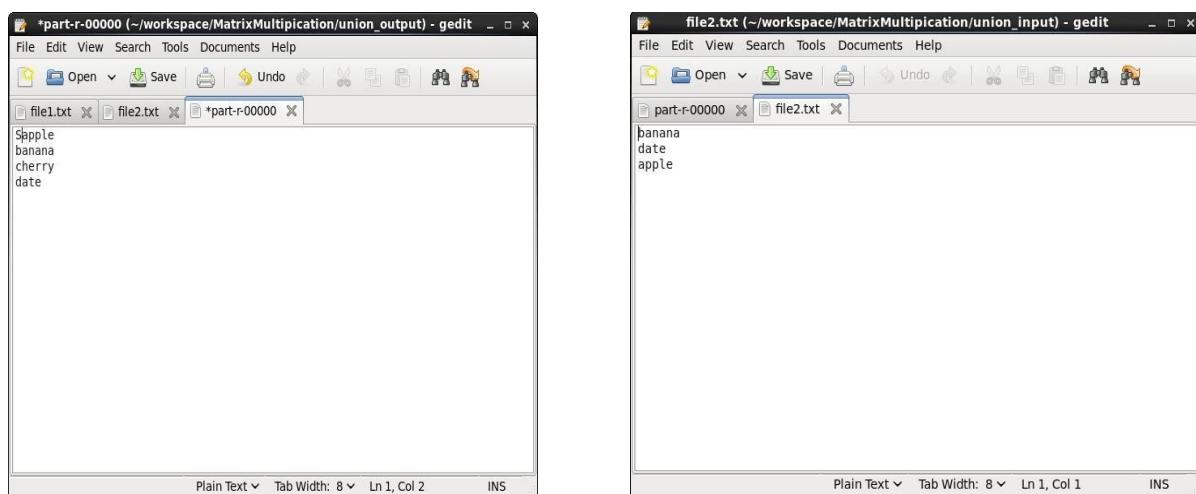
System.exit(job.waitForCompletion(true) ? 0 : 1);
```

3

3

Steps for further implementation:

Step 1: Create two sample files in cloudera.



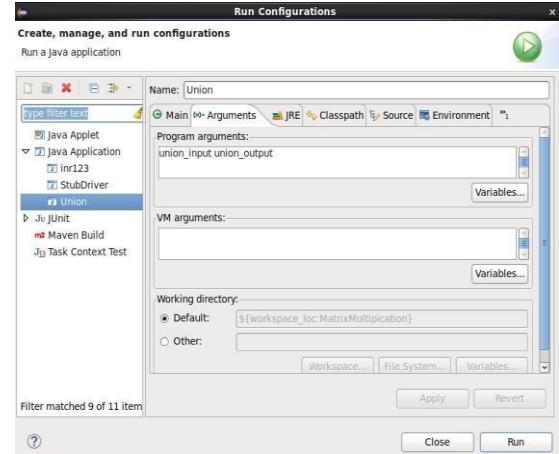
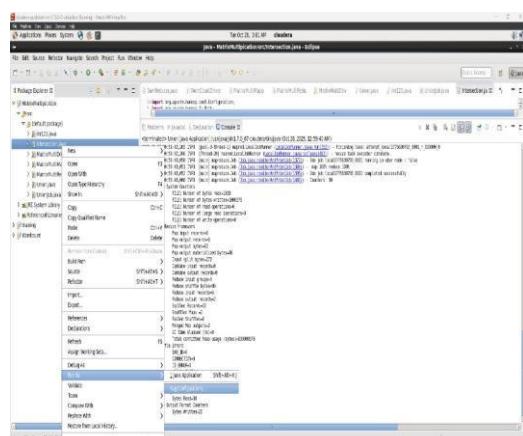
Step 2: Prepare Local Input File. In your Cloudera terminal

```
[cloudera@quickstart ~]$ hdfs dfs -mkdir /user/cloudera/union_input
[cloudera@quickstart ~]$ hdfs dfs -put file1.txt /user/cloudera/union_input/
[cloudera@quickstart ~]$ hdfs dfs -put file2.txt /user/cloudera/union_input/
[cloudera@quickstart ~]$ mkdir -p /home/cloudera/workspace/MatrixMultiplication/u
nion_input
[cloudera@quickstart ~]$
```

Step 3: Run Configuration in Eclipse

Right-click on Intersection.java → Run As → Run

Configurations... In Arguments → Program arguments, enter:



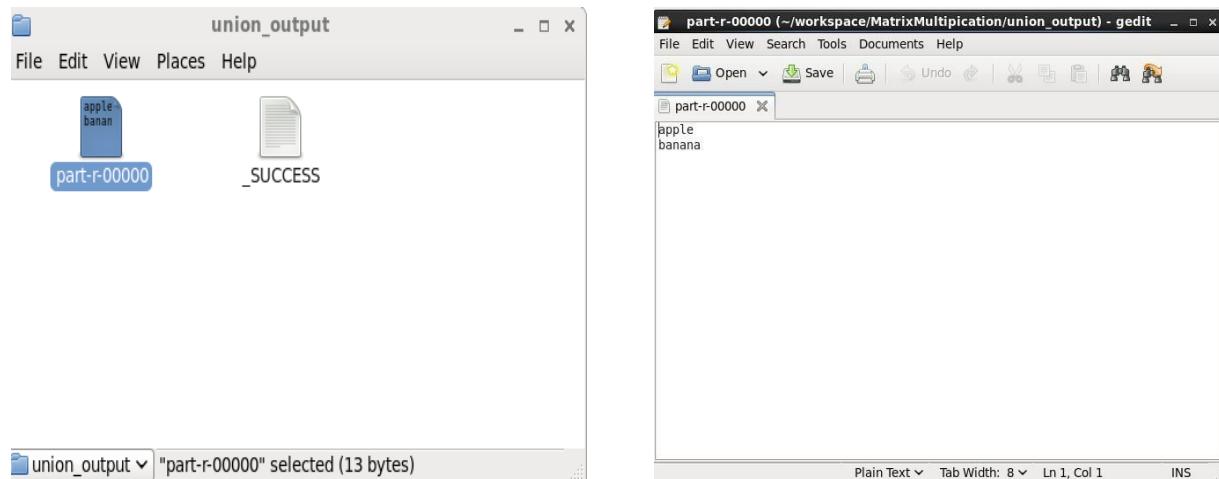
Output:

```

Problems @ Javadoc Declaration Console 
<terminated> Union [Java Application] /usr/java/jdk1.8.0_67-cloudera/bin/java (Oct 28, 2025, 12:44:00 AM)
2025-10-28 00:44:01,614 WARN [main] util.NativeCodeLoader [NativeCodeLoader.java:<clinit>(62)] - Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
2025-10-28 00:44:01,965 INFO [main] Configuration.deprecation [Configuration.java:varNonecfDeprecation(1704)] - session.id is deprecated. Instead, use dfs.metrics.session-id
2025-10-28 00:44:01,965 INFO [main] jvm.JvmMetrics [JvmMetrics.java:init(76)] - Initializing Metrics with preassigned-JvmTracer, session-id=1773454917
2025-10-28 00:44:02,000 INFO [main] JobResourceUploader [JobResourceUploader.java:uploadJobResource(100)] - Hadoop command-line optimization not performed. Implement the Tool interface and execute
2025-10-28 00:44:02,008 WARN [main] mapreduce.JobResourceUploader [JobResourceUploader.java:uploadJobResource(100)] - No job jar file set. User classes may not be found. See Job or Job#setJar(String).
2025-10-28 00:44:02,109 INFO [main] InputFileInputFormat [FileInputFormat.java:listStatus(283)] - Total input paths to process : 2
2025-10-28 00:44:02,156 INFO [main] mapreduce.JobSubmitter [JobSubmission.java:submitJobInternal(202)] - number of splits:2
2025-10-28 00:44:02,282 INFO [main] mapreduce.JobSubmitter [JobSubmission.java:printTokens(291)] - Submitting tokens for job: job_local1773454917_0001
2025-10-28 00:44:02,300 INFO [main] mapreduce.Job [Job.java:monitorAndPrintJob(1356)] - Running job: job_local1773454917_0001
2025-10-28 00:44:02,436 INFO [main] mapreduce.Job [Job.java:monitorAndPrintJob(1356)] - OutputCommitter set in config null
2025-10-28 00:44:02,439 INFO [Thread-20] mapred.LocalJobRunner [LocalJobRunner.java:createOutputCommitter(498)] - OutputCommitter skip cleanup _temporary folders under output directory:false, ignore cl
2025-10-28 00:44:02,439 INFO [Thread-20] mapred.FileOutputCommitter [FileOutputCommitter.java:<init>(124)] - File Output Committer Algorithm version is 1
2025-10-28 00:44:02,439 INFO [Thread-20] mapred.FileOutputCommitter [FileOutputCommitter.java:<init>(124)] - FileOutputCommitter skip cleanup _temporary folders under output
2025-10-28 00:44:02,440 INFO [Thread-20] mapred.LocalJobRunner [LocalJobRunner.java:runTask(1351)] - Using configured calculator classesFree: []
2025-10-28 00:44:02,499 INFO [main] mapreduce.MapTask [MapTask.java:runNewMapper(762)] - Processing split: file:/home/cloudera/workspace/MatrixMultiplication/union_input/fil
2025-10-28 00:44:02,709 INFO [main] MapTask [MapTask.java:setEquator(1212)] - (EQUATOR) 0 kv1 2621496(104857584)
2025-10-28 00:44:03,441 INFO [main] mapreduce.Job [Job.java:monitorAndPrintJob(1384)] - map 100% reduce 100%
2025-10-28 00:44:03,442 INFO [main] mapreduce.Job [Job.java:monitorAndPrintJob(1395)] - Job job_local1773454917_0001 completed successfully
2025-10-28 00:44:03,450 INFO [main] mapreduce.Job [Job.java:monitorAndPrintJob(1402)] - Counters: 39
File System Counters
FILE: Number of bytes read=1892
FILE: Number of bytes written=1000000
FILE: Number of read operations=0
FILE: Number of large read operations=0
FILE: Number of write operations=0
Map-Reduce Framework
Map input records=6
Map output records=6
Map output bytes=40
Map output materialized bytes=68
Input split bytes=272
Combine input records=0
Combine output records=0
Reduce input groups=4
Reduce shuffle bytes=68
Reduce input records=6
Reduce output records=4
Spilled Records=12
Shuffled Maps=2
Failed Shuffles=0
Merged Map outputs=2
File Time balanced (mc1)...A

```

A new output file is created which have all the common elements of the files.



4. Write a program in Map Reduce for Matrix Multiplication.

Steps:

Step 1: Open virtual box and then start cloudera quickstart

Step 2: Open eclipse present on the cloudera desktop

Step 3: Create java project

File->New-> Java Project

Give project name: MatrixMultiplication

Click Next

Step 4: Add Hadoop libraries to project

Select Libraries tab->click on Add External

Jars File System->user->lib->Hadoop

Select all library(jar) files->ok

Again, click on Add External

Jars

File System->user->lib->Hadoop->client

Select all library(jar) files from client->ok->finish

Step 5: Write java code for matrix multiplication using mapreduce Right click on src folder of project

MatrixMultiplication New->class

Write class name MatrixMultiplicationMapper

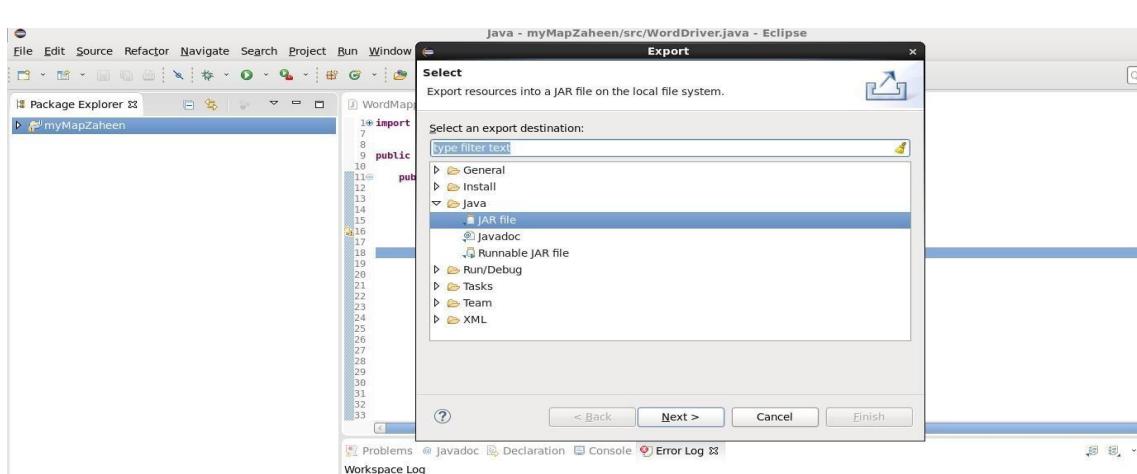
Click Finish

Write the code for MatrixMultiplicationMapper

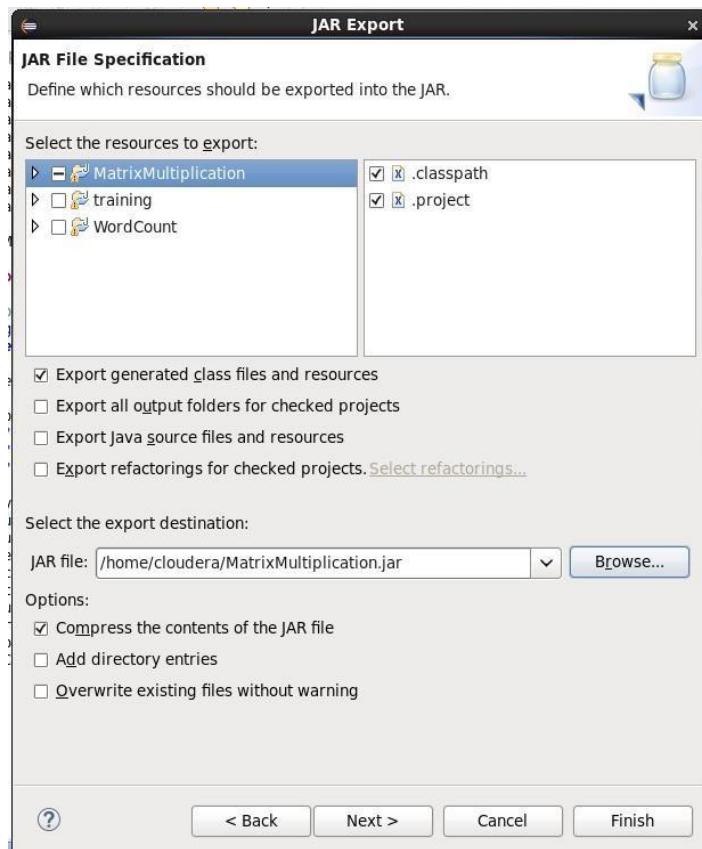
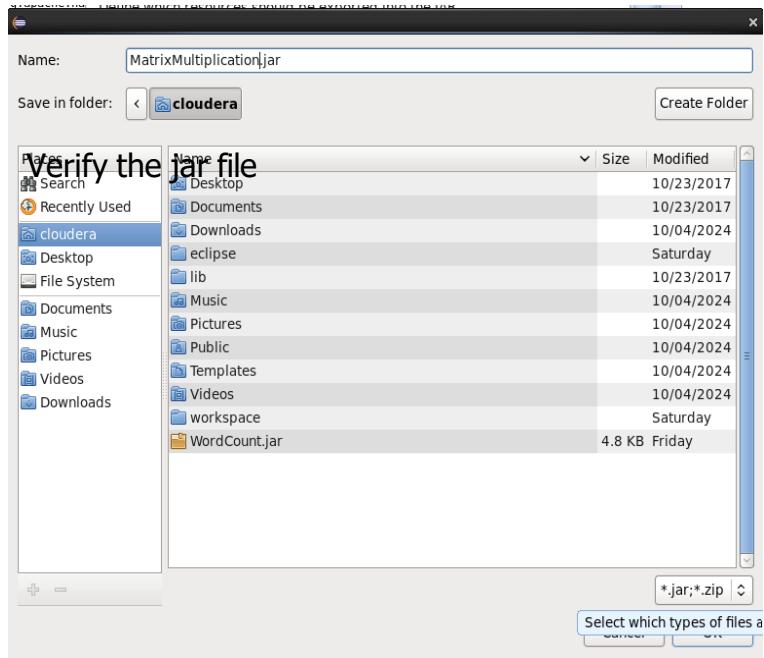
Same way create classes MatrixMultiplicationReducer and

MatrixMultiplicationDriver Step 6: Export the project as jar

Right click on project MatrixMultiplication and select Export>> Java>>JAR file>>Next



Select the export destination-Click browse-give file name matrixmultiplication.jar Click ok>>Finish>>ok



Verify the jar file through command line open terminal give command
ls

step 7: Move the jar file to the Hadoop file system
hdfs dfs -put matrixmultiplication.jar /user/cloudera
hdfs dfs -ls

Step 8: Create the input file for the MapReduce program First input file
Command: **cat > myMMatrix.txt**
Enter data in input file
M,0,0,1
M,0,1,2
M,1,0,3
M,1,1,4
and press enter and ctrl z

Command: **cat myMMatrix.txt**
Second input file
Command: **cat > myNMatrix.txt**
Enter data in input file
N,0,0,3
N,0,1,6
N,1,0,4
N,1,1,2
and press enter and ctrl z
Command: **cat myNMatrix.txt**

Step 9: Create the input directory and move input files into it.

hdfs dfs –mkdir /user/cloudera/matrixInput
hdfs dfs -put myMMatrix.txt
/user/cloudera/matrixInput/ hdfs dfs -put
myNMatrix.txt /user/cloudera/matrixInput/ hdfs dfs -ls

step 10: Run mapreduce program on Hadoop
syntax: hadoop jar jarfilename.jar classname inputfoldername outputfoldername
command: **hadoop jar matrixmultiplication.jar MatrixMultiplicationDriver matrixInput matrixOutput**

step 11: view output directory
hdfs dfs –ls
hdfs dfs -ls /user/cloudera/matrixOutput
step 12: view the output file
hdfs dfs -cat /user/cloudera/matrixOutput/part-r-00000

Code:**MatrixMultiplicationMapper.java**

```
import java.io.IOException;

import org.apache.hadoop.conf.*;
import org.apache.hadoop.io.LongWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Mapper;

public class MatrixMultiplicationMapper extends Mapper<LongWritable, Text,
Text, Text>{

    @Override
        public void map(LongWritable key,
Text value, Context context) throws
IOException,
InterruptedException {
Configuration conf = context.getConfiguration();
int m = Integer.parseInt(conf.get("m"));
int p =
Integer.parseInt(conf.get("p"));
String line = value.toString();
// (M, i, j, Mij);
String[] indicesAndValue = line.split(",");
Text outputKey = new Text();
Text outputValue = new Text();
if (indicesAndValue[0].equals("M")) {
    for (int k = 0; k < p; k++) {
        outputKey.set(indicesAndValue[1] + "," + k); // outputKey.set(i,k);
outputValue.set(indicesAndValue[0] + "," + indicesAndValue[2]
+ "," + indicesAndValue[3]); // outputValue.set(M,j,Mij);
context.write(outputKey, outputValue);
    }
}
else {
    // (N, j, k, Njk);
    for (int i = 0; i < m; i++) {
        outputKey.set(i + "," + indicesAndValue[2]);
outputValue.set("N," + indicesAndValue[1] + ","
+ indicesAndValue[3]);
context.write(outputKey, outputValue);
    }
}
}
```

3

3

3

3

MatrixMultiplicationReducer.java

```
import org.apache.hadoop.io.Text;
import java.io.IOException;
import java.util.HashMap;

public class MatrixMultiplicationReducer extends
org.apache.hadoop.mapreduce.Reducer<Text, Text, Text, Text>{

    @Override
    public void reduce(Text key, Iterable<Text> values, Context context)
        throws IOException, InterruptedException {
        String[] value;
        //key=(i,k),
        //Values = [(M/N,j,V/W),..]
        HashMap<Integer, Float> hashA = new HashMap<Integer,
        Float>();    HashMap<Integer,    Float>    hashB    =    new
        HashMap<Integer, Float>(); for (Text val : values) {
            value = val.toString().split(",");
            if (value[0].equals("M")) {
                hashA.put(Integer.parseInt(value[1]), Float.parseFloat(value[2]));
            } else {
                hashB.put(Integer.parseInt(value[1]), Float.parseFloat(value[2]));
            }
        }
        int n =
        Integer.parseInt(context.getConfiguration().get("n")); float
        result = 0.0f;
        float m_ij;
        float n_jk;
        for (int j = 0; j < n; j++) {
            m_ij = hashA.containsKey(j) ? hashA.get(j) : 0.0f;
            n_jk = hashB.containsKey(j) ? hashB.get(j) : 0.0f;
            result += m_ij * n_jk;
        }
        if (result != 0.0f) {
            context.write(null,
new Text(key.toString() + "," + Float.toString(result)));
        }
    }
}
```

3

MatrixMultiplicationDriver.java

```
import org.apache.hadoop.conf.*;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.*;
import org.apache.hadoop.mapreduce.*;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;

public class MatrixMultiplicationDriver {

    public static void main(String[] args) throws Exception{

        // TODO Auto-generated method stub
        if (args.length != 2) {
            System.err.println("Usage:
                MatrixMultiply <in_dir>
                <out_dir>"); System.exit(2);
        }

        Configuration conf = new Configuration();
        // M is an m-by-n matrix; N is an n-by-p matrix.
        conf.set("m", "1000");
        conf.set("n", "100");
        conf.set("p", "1000");

        Job job = new Job(conf, "MatrixMultiply");
        job.setJarByClass(MatrixMultiplicationDriver.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(Text.class);

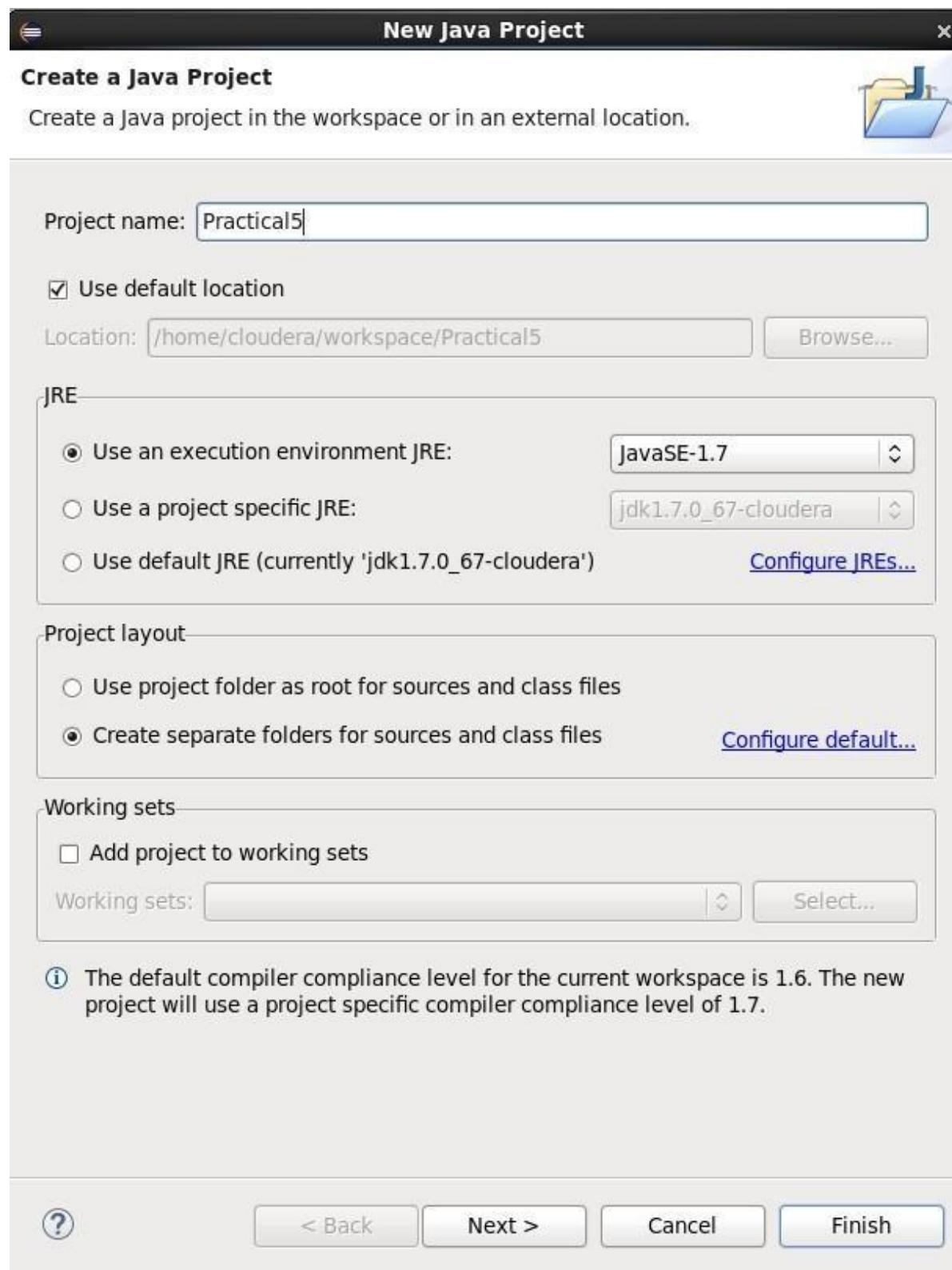
        job.setMapperClass(MatrixMultiplicationMapper.class);
        job.setReducerClass(MatrixMultiplicationReducer.class);

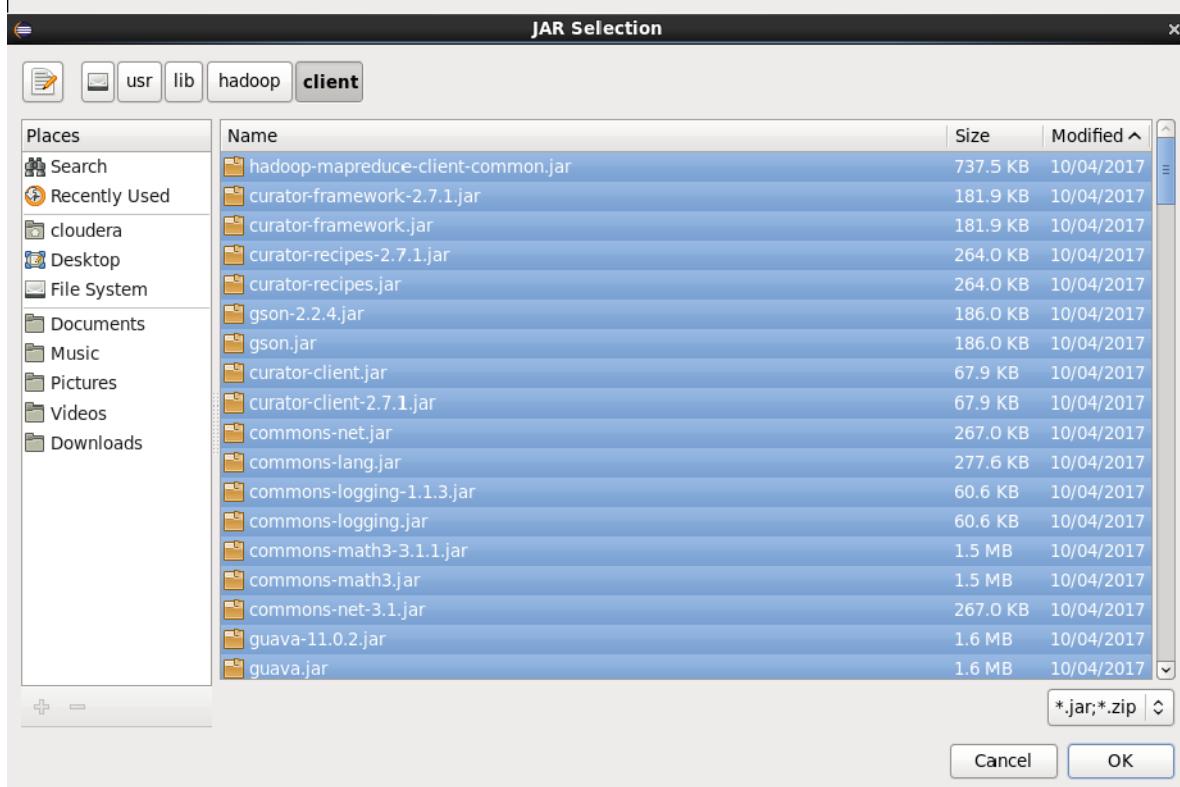
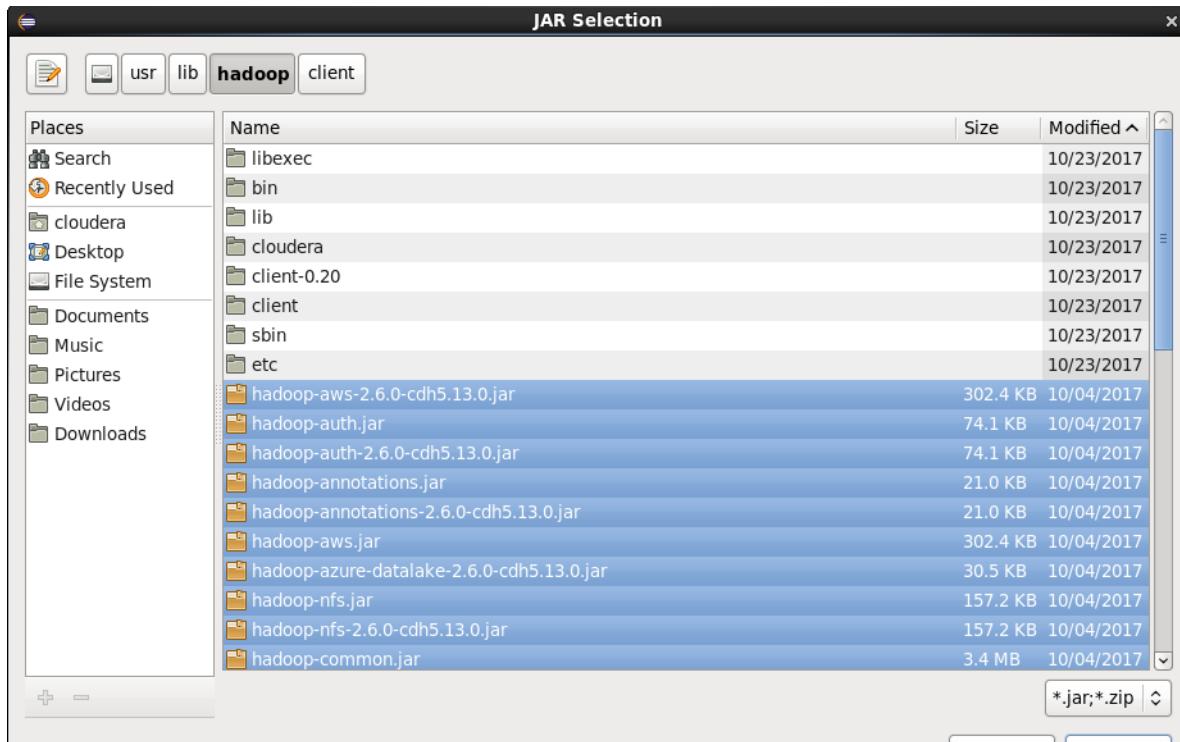
        job.setInputFormatClass(TextInputFormat.class);
        job.setOutputFormatClass(TextOutputFormat.class);

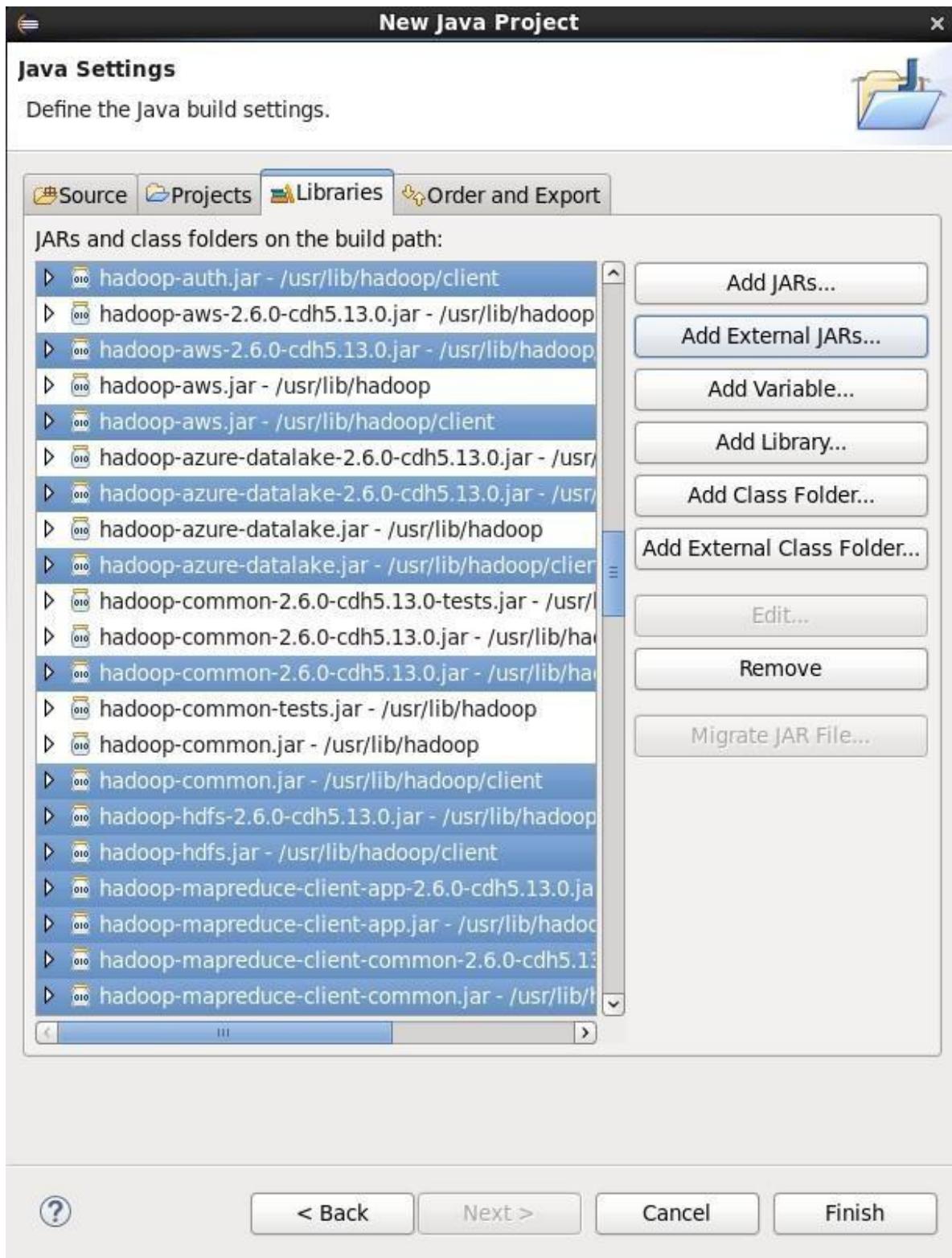
        FileInputFormat.addInputPath(job, new Path(args[0]));
        FileOutputFormat.setOutputPath(job, new Path(args[1]));

        job.waitForCompletion(true);
    }
}
```

Output:







New Java Class

Java Class

 The use of the default package is discouraged.

Source folder:

Package: (default)

Enclosing type:

Name:

Modifiers: public package private protected
 abstract final static

Superclass:

Interfaces:

Which method stubs would you like to create?

public static void main(String[] args)
 Constructors from superclass
 Inherited abstract methods

Do you want to add comments? (Configure templates and default value [here](#))

Generate comments



Eclipse - Java - Practical5/src/MatrixMultiplicationMapper.java

```

1 import java.io.IOException;
2 import org.apache.hadoop.conf.*;
3 import org.apache.hadoop.io.LongWritable;
4 import org.apache.hadoop.io.Text;
5 import org.apache.hadoop.mapreduce.Mapper;
6 public class MatrixMultiplicationMapper extends Mapper<LongWritable, Text, Text, Text> {
7     @Override
8     public void map(LongWritable key, Text value, Context context) throws IOException, InterruptedException{
9         Configuration conf = context.getConfiguration();
10        int m = Integer.parseInt(conf.get("m"));
11        int p = Integer.parseInt(conf.get("p"));
12        String line = value.toString();
13        String[] indicesAndValue = line.split(",");
14        Text outputKey = new Text();
15        Text outputValue = new Text();
16        if(indicesAndValue[0].equals("M")){
17            for(int k = 0; k<p; k++){
18                outputKey.set(indicesAndValue[1]+","+"k");
19                outputValue.set(indicesAndValue[0]+"+"+indicesAndValue[2]+"+"+indicesAndValue[3]);
20                context.write(outputKey, outputValue);
21            }
22        }else{
23            for(int i = 0; i<m; i++){
24                outputKey.set(i+"+"+indicesAndValue[2]);
25                outputValue.set("N,"+indicesAndValue[1]+"+"+indicesAndValue[3]);
26                context.write(outputKey, outputValue);
27            }
28        }
29    }
30 }
31

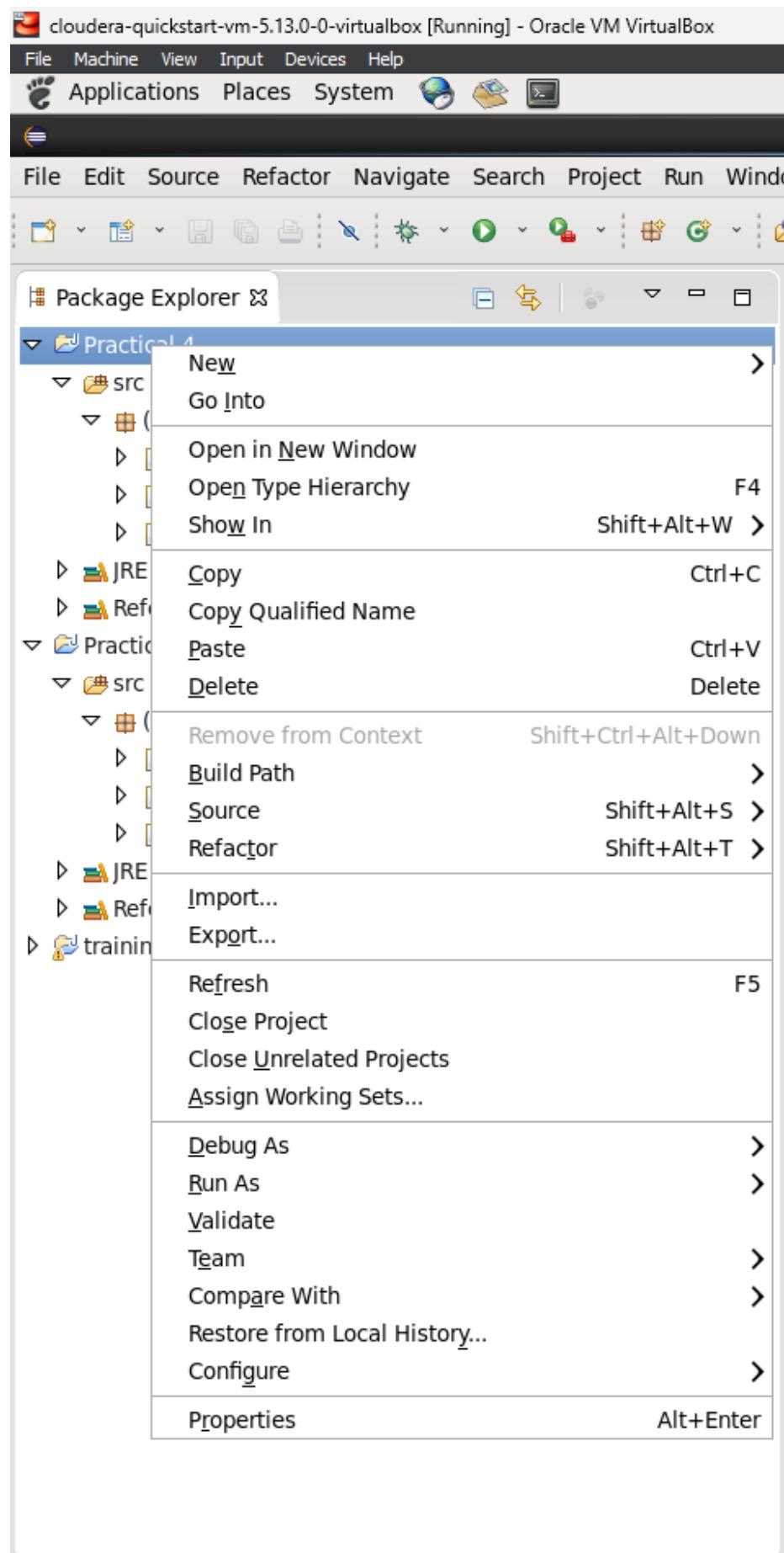
```

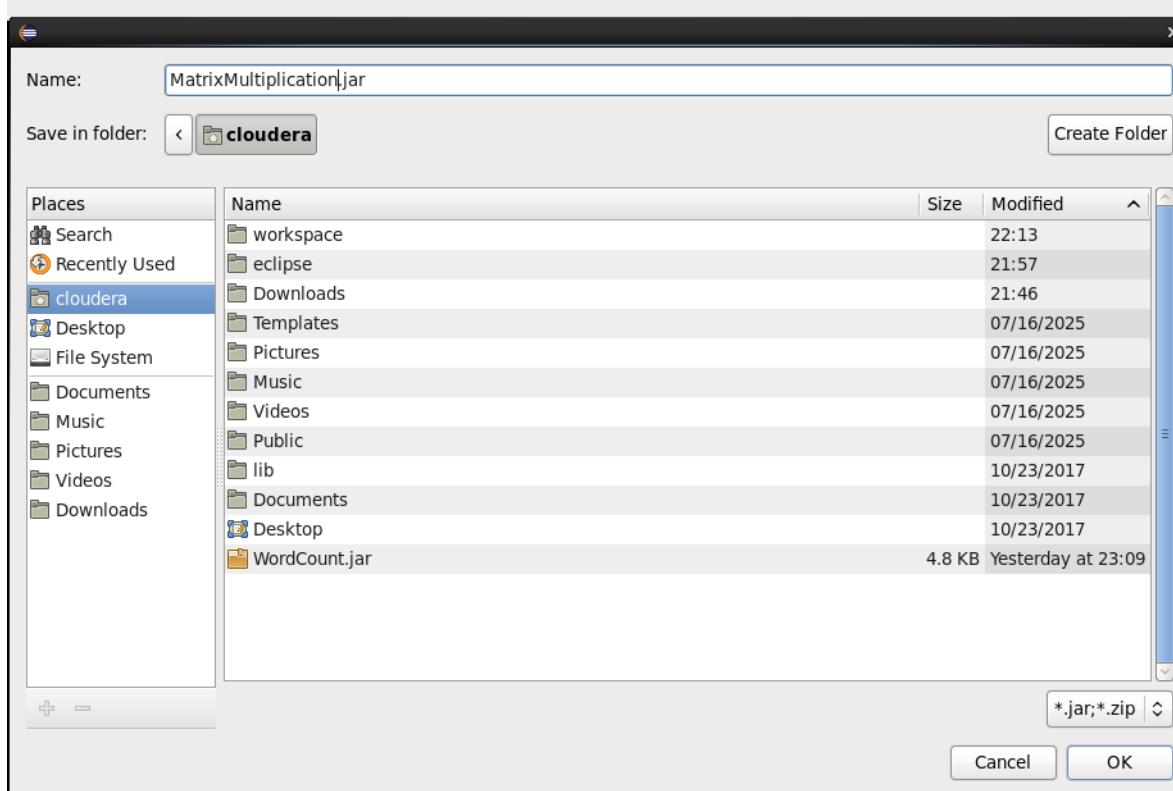
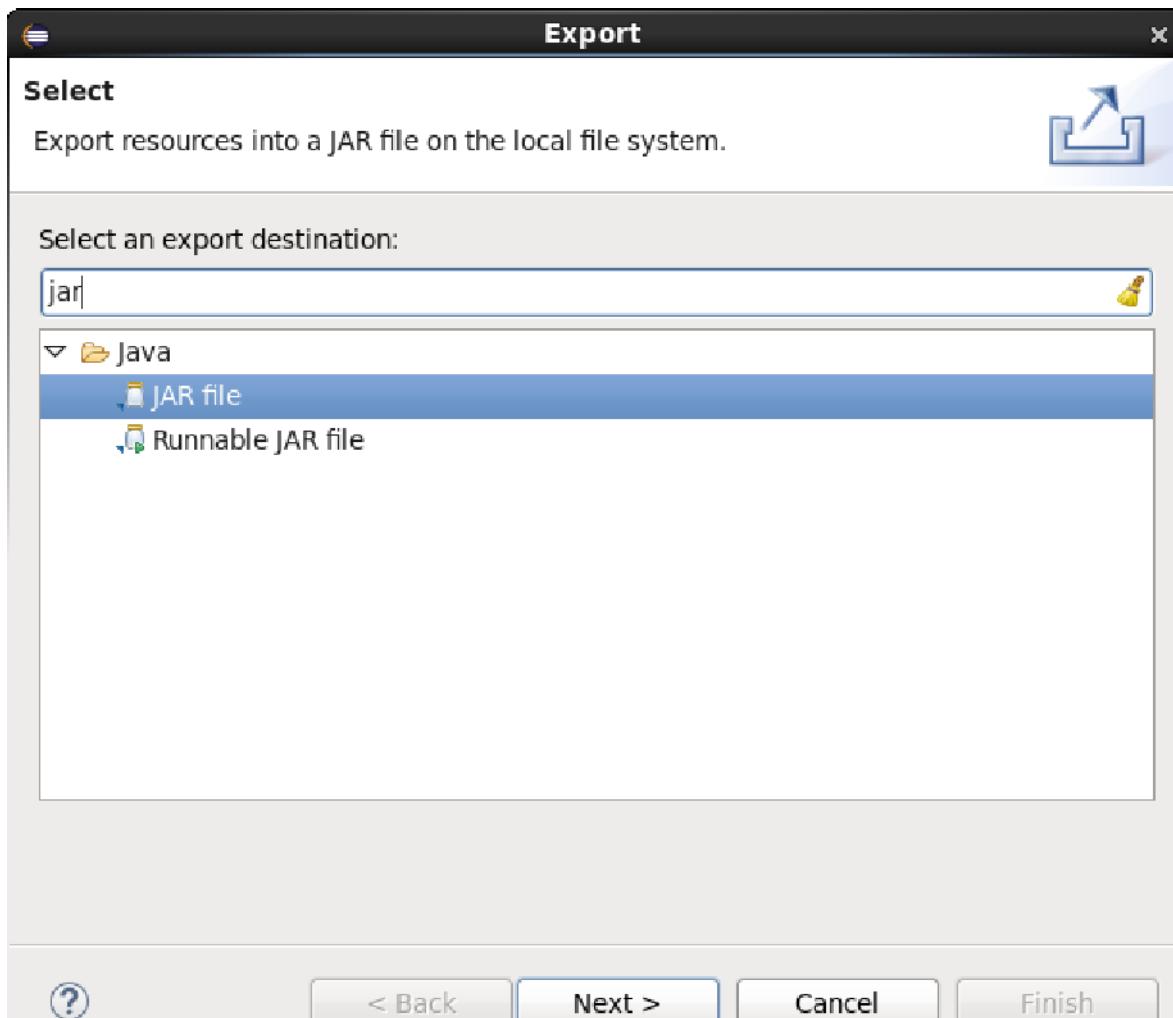
Eclipse - Java - Practical5/src/MatrixMultiplicationReducer.java

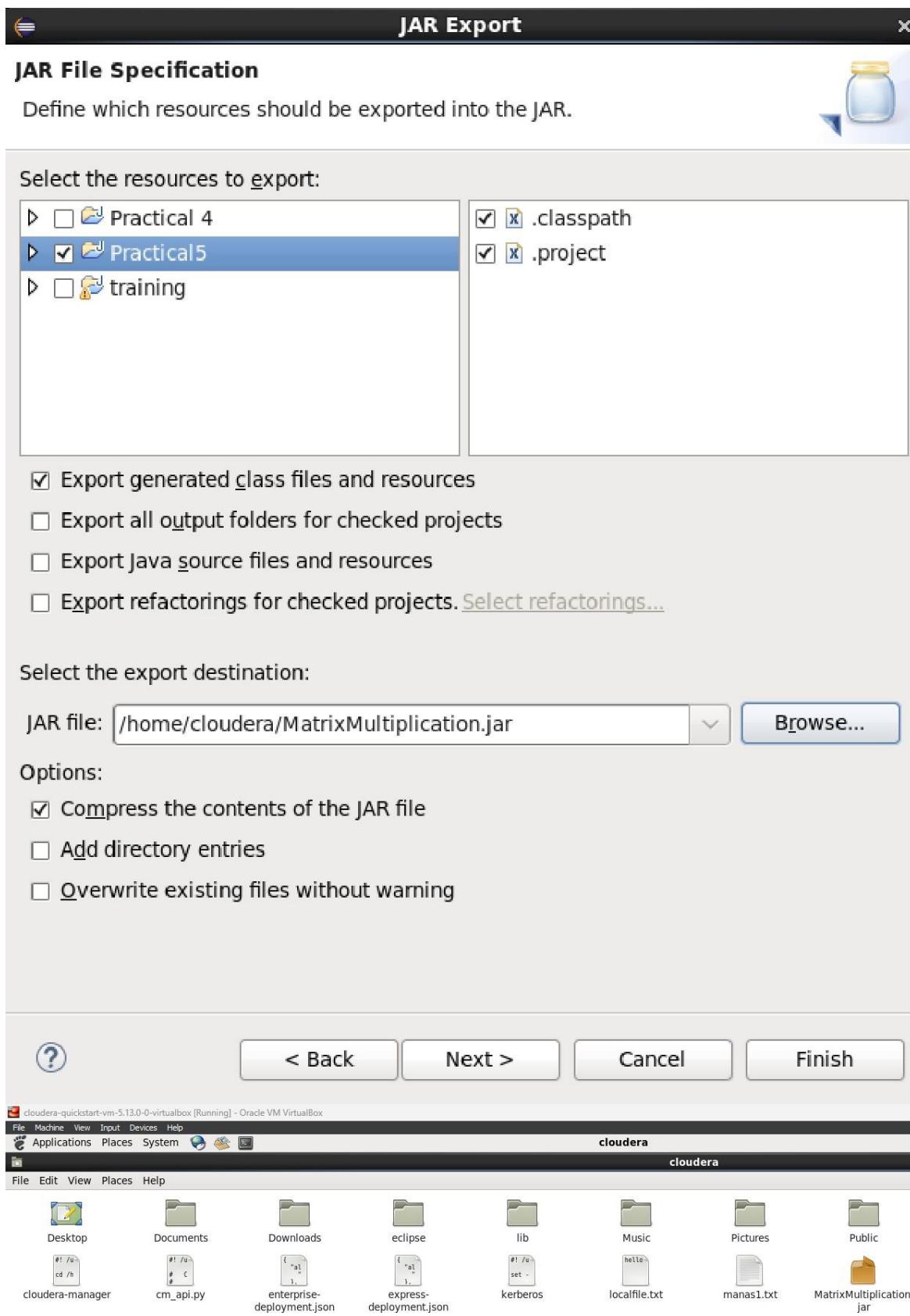
```

1 import java.io.IOException;
2 import java.util.HashMap;
3 import org.apache.hadoop.io.Text;
4 import org.apache.hadoop.mapreduce.Reducer;
5 public class MatrixMultiplicationReducer extends Reducer<Text, Text, Text, Text>{
6     @Override
7     public void reduce(Text key, Iterable<Text> values, Context context) throws IOException, InterruptedException{
8         String[] value;
9         HashMap<Integer, Float> hashA = new HashMap<Integer, Float>();
10        HashMap<Integer, Float> hashB = new HashMap<Integer, Float>();
11        for(Text val : values){
12            value = val.toString().split(",");
13            if(value[0].equals("M")){
14                hashA.put(Integer.parseInt(value[1]), Float.parseFloat(value[2]));
15            }else{
16                hashB.put(Integer.parseInt(value[1]), Float.parseFloat(value[2]));
17            }
18        }
19        int n = Integer.parseInt(context.getConfiguration().get("n"));
20        float result = 0.0f;
21        float m_ij;
22        float n_jk;
23        for(int j = 0; j<n; j++){
24            m_ij = hashA.containsKey(j)?hashA.get(j):0.0f;
25            n_jk = hashB.containsKey(j)?hashB.get(j):0.0f;
26            result += m_ij*n_jk;
27        }
28        if(result != 0.0f){
29            context.write(null, new Text(key.toString()+"+"+Float.toString(result)));
30        }
31    }
32 }
33

```







cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox

File Machine View Input Devices Help
Applications Places System

Java - Practical5/src/MatrixMultiplicationDriver.java

File Edit Source Refactor Navigate Search Project Run Window Help

Package Explorer

- Practical 4
 - (default package)
 - SumReducer.java
 - WordCountDriver.java
 - WordMapper.java
 - JRE System Library [javaSE-1.7]
 - Referenced Libraries
- Practical5
 - src
 - (default package)
 - MatrixMultiplicationDriver.java
 - MatrixMultiplicationMapper.java
 - MatrixMultiplicationReducer.java
 - JRE System Library [javaSE-1.7]
 - Referenced Libraries
- training

```

1 import org.apache.hadoop.fs.Path;
2 import org.apache.hadoop.conf.*;
3 import org.apache.hadoop.io.*;
4 import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
5 import org.apache.hadoop.mapreduce.lib.output.TextOutputFormat;
6 import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
7 import org.apache.hadoop.mapreduce.lib.input.TextInputFormat;
8 import org.apache.hadoop.mapreduce.Job;
9 public class MatrixMultiplicationDriver {
10     public static void main(String[] args) throws Exception{
11         if(args.length !=2){
12             System.err.println("Usage: MatrixMultiply <in_dir> <out_dir>\n");
13             System.exit(2);
14         }
15         Configuration conf = new Configuration();
16         conf.set("m","1000");
17         conf.set("n","100");
18         conf.set("p","1000");
19         @SuppressWarnings("deprecation")
20         Job job = new Job(conf, "MatrixMultiply");
21         job.setJarByClass(MatrixMultiplicationDriver.class);
22         job.setOutputKeyClass(Text.class);
23         job.setOutputValueClass(Text.class);
24         job.setMapperClass(MatrixMultiplicationMapper.class);
25         job.setReducerClass(MatrixMultiplicationReducer.class);
26         job.setInputFormatClass(TextInputFormat.class);
27         job.setOutputFormatClass(TextOutputFormat.class);
28         FileInputFormat.addInputPath(job, new Path(args[0]));
29         FileOutputFormat.setOutputPath(job, new Path(args[1]));
30         job.waitForCompletion(true);
31     }
32 }
33 }
```

cloudera-quickstart-vm-5.13.0-0-virtualbox [Running] - Oracle VM VirtualBox

File Edit View Search Terminal Help

[cloudera@quickstart ~]\$ ls

C24115.txt cm_apl.py Documents eclipse express-deployment.json lib localfile.txt~ MatrixMultiplication.jar myInputfile.4.txt parcels Public Videos workspace

[cloudera@quickstart ~]\$ hdfs dfs -put MatrixMultiplication.jar /user/cloudera

[cloudera@quickstart ~]\$ hdfs dfs -ls

Found 7 items

drwx-----	- cloudera cloudera	0 2025-08-06 22:00 .Trash
drwx-----	- cloudera cloudera	0 2025-08-06 21:44 .staging
-rw-r--r--	1 cloudera cloudera	6555 2025-08-06 23:08 MatrixMultiplication.jar
-rw-r--r--	1 cloudera cloudera	4887 2025-08-05 23:20 WordCount.jar
drwxr-xr-x	- cloudera cloudera	0 2025-07-30 00:31 manas
-rw-r--r--	1 cloudera cloudera	52 2025-08-05 23:25 myInputFile.txt
drwxr-xr-x	- cloudera cloudera	0 2025-08-06 21:43 myoutput

[cloudera@quickstart ~]\$ cat > myMMatrix.txt

M,0,0,1
M,0,1,2
M,1,0,3
M,1,1,4
^Z

[1]+ Stopped cat > myMMatrix.txt

[cloudera@quickstart ~]\$ cat myMMatrix.txt

M,0,0,1
M,0,1,2
M,1,0,3
M,1,1,4

[cloudera@quickstart ~]\$ cat > myNMatrix.txt

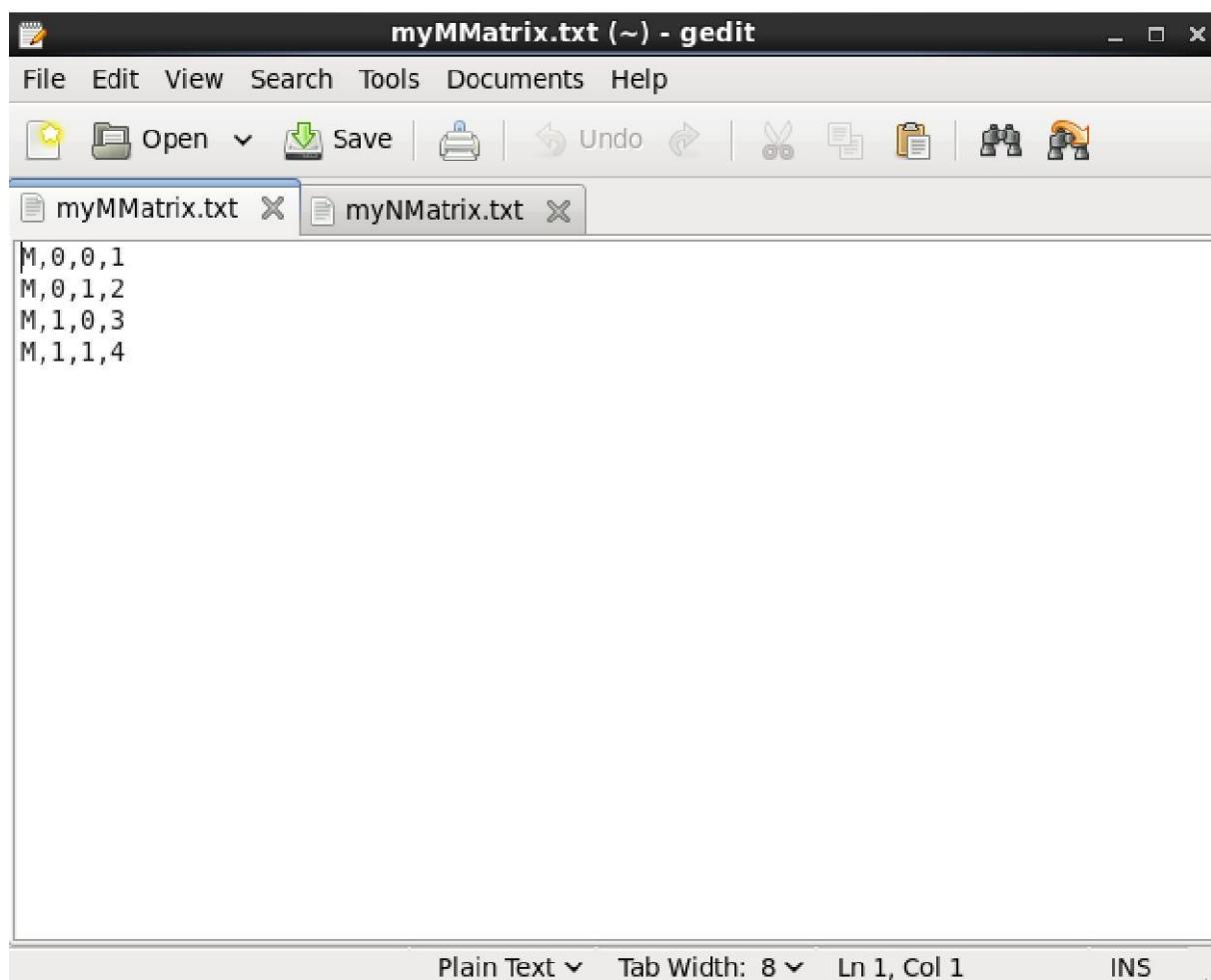
N,0,0,3
N,0,1,6
N,1,0,4
N,1,1,2
^Z

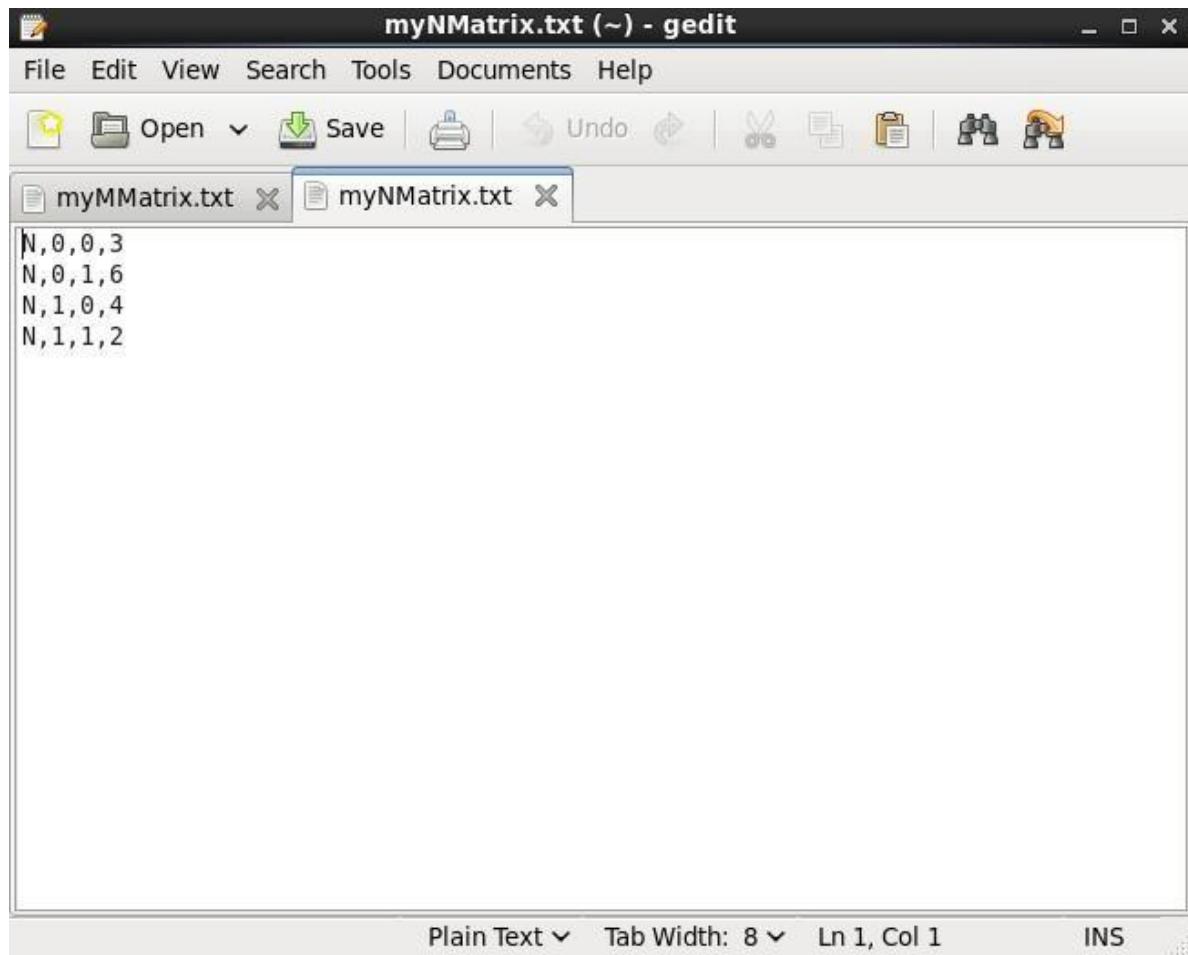
[2]+ Stopped cat > myNMatrix.txt

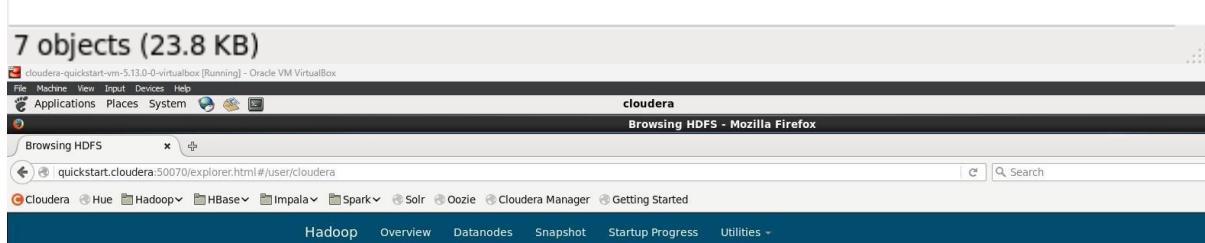
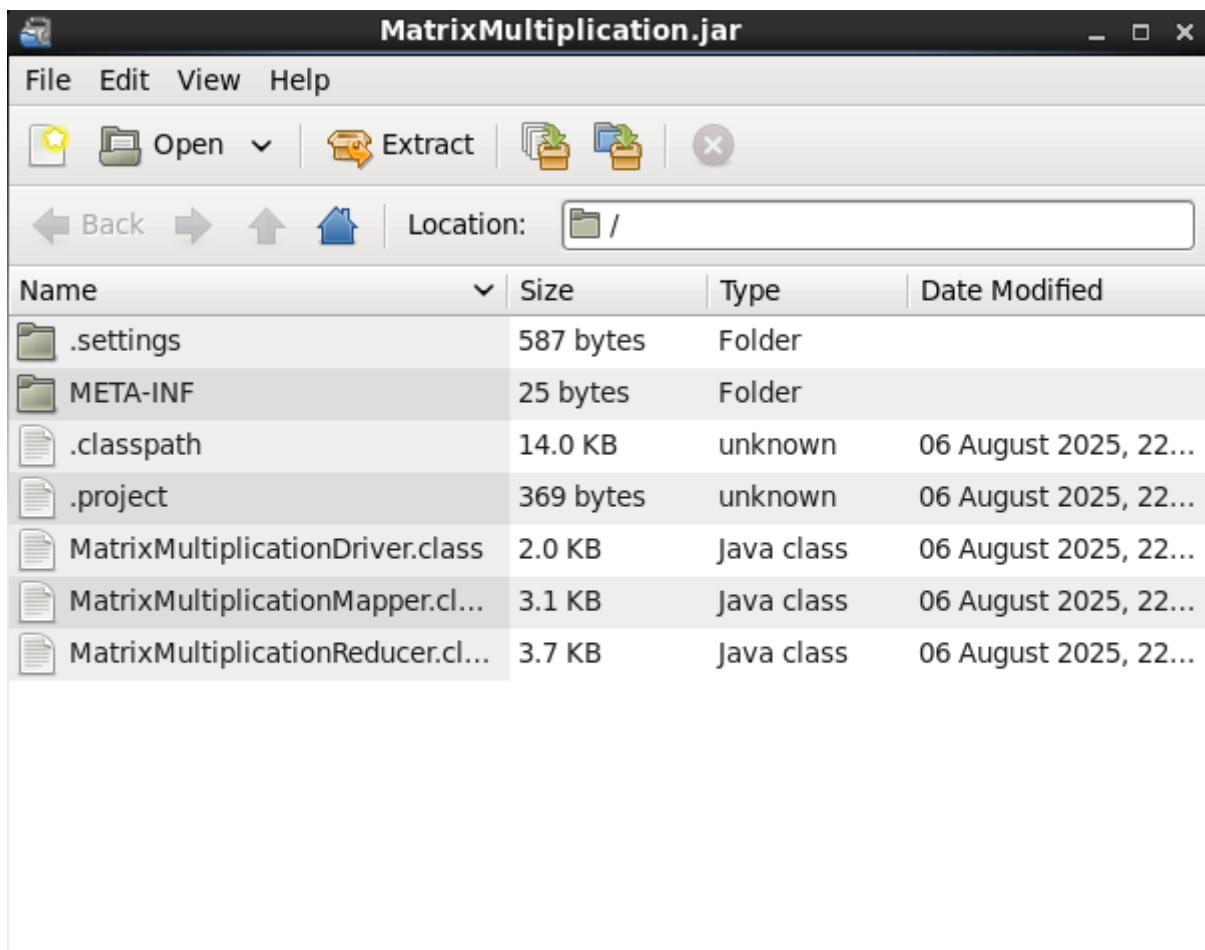
[cloudera@quickstart ~]\$ cat myNMatrix.txt

N,0,0,3
N,0,1,6
N,1,0,4
N,1,1,2

```
[cloudera@quickstart ~]$ hdfs dfs -mkdir /user/cloudera/matrixInput
[cloudera@quickstart ~]$ hdfs dfs -put myMMMatrix.txt /user/cloudera/matrixInput/
[cloudera@quickstart ~]$ hdfs dfs -put myNMatrix.txt /user/cloudera/matrixInput/
[cloudera@quickstart ~]$ hdfs dfs -ls
Found 8 items
drwx-----  cloudera cloudera      0 2025-08-06 22:00 .Trash
drwx-----  cloudera cloudera      0 2025-08-06 21:44 .staging
-rw-r--r--  1 cloudera cloudera  6555 2025-08-06 23:08 MatrixMultiplication.jar
-rw-r--r--  1 cloudera cloudera  4887 2025-08-05 23:20 WordCount.jar
drwxr-xr-x  - cloudera cloudera      0 2025-07-30 00:31 manas
drwxr-xr-x  - cloudera cloudera      0 2025-08-06 23:14 matrixInput
-rw-r--r--  1 cloudera cloudera   52 2025-08-05 23:25 myInputFile.txt
drwxr-xr-x  - cloudera cloudera      0 2025-08-06 21:43 myOutput
[cloudera@quickstart ~]$ hadoop jar MatrixMultiplication.jar MatrixMultiplicationDriver matrixInput matrixOutput
25/08/06 23:15:43 INFO client.RMProxy: Connecting to ResourceManager at quickstart.cloudera/10.6.2.15:8032
25/08/06 23:15:44 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.
25/08/06 23:15:45 INFO input.FileInputFormat: Total input paths to process : 2
25/08/06 23:15:45 INFO mapreduce.JobSubmitter: number of splits:2
25/08/06 23:15:45 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1754541476929_0002
25/08/06 23:15:46 INFO yarn.ClientImpl: Submitted application application_1754541476929_0002
25/08/06 23:15:46 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1754541476929_0002/
25/08/06 23:16:03 INFO mapreduce.Job: Job job_1754541476929_0002 running in uber mode : false
25/08/06 23:16:23 INFO mapreduce.Job: map 0% reduce 0%
25/08/06 23:16:34 INFO mapreduce.Job: map 100% reduce 100%
25/08/06 23:16:35 INFO mapreduce.Job: Job job_1754541476929_0002 completed successfully
25/08/06 23:16:35 INFO mapreduce.Job: Counters
    File Counters
        FILE: Number of bytes read=31169
        FILE: Number of bytes written=505790
        FILE: Number of read operations=0
        FILE: Number of large read operations=0
        FILE: Number of write operations=0
        HDFS: Number of bytes read=336
        HDFS: Number of bytes written=36
        HDFS: Number of read operations=9
        HDFS: Number of large read operations=0
        HDFS: Number of write operations=2
    Job Counters
        Launched map tasks=2
        Launched reduce tasks=1
        Data-local map tasks=2
        Total time spent by all maps in occupied slots (ms)=17917440
        Total time spent by all reduces in occupied slots (ms)=3418112
        Total time spent by all map tasks (ms)=34995
```







Browse Directory

Browse Directory								
/user/cloudera								
Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name	
drwx-----	cloudera	cloudera	0 B	Wed Aug 06 22:00:01 -0700 2025	0	0 B	.Trash	
drwx-----	cloudera	cloudera	0 B	Wed Aug 06 23:16:35 -0700 2025	0	0 B	.staging	
-rw-r--r--	cloudera	cloudera	6.4 KB	Wed Aug 06 23:08:52 -0700 2025	1	128 MB	MatrixMultiplication.jar	
-rw-r--r--	cloudera	cloudera	4.77 KB	Tue Aug 05 23:20:44 -0700 2025	1	128 MB	WordCount.jar	
drwxr-xr-x	cloudera	cloudera	0 B	Wed Jul 30 00:31:25 -0700 2025	0	0 B	manas	
drwxr-xr-x	cloudera	cloudera	0 B	Wed Aug 06 23:14:10 -0700 2025	0	0 B	matrixInput	
drwxr-xr-x	cloudera	cloudera	0 B	Wed Aug 06 23:16:33 -0700 2025	0	0 B	matrixOutput	
-rw-r--r--	cloudera	cloudera	52 B	Tue Aug 05 23:25:23 -0700 2025	1	128 MB	myInputFile.txt	
drwxr-xr-x	cloudera	cloudera	0 B	Wed Aug 06 21:43:57 -0700 2025	0	0 B	myOutput	

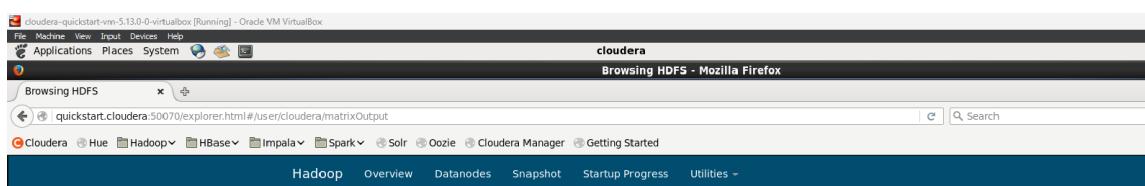
Hadoop, 2017.



Browse Directory

Browse Directory							
/user/cloudera/matrixInput							
Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	cloudera	cloudera	32 B	Wed Aug 06 23:13:51 -0700 2025	1	128 MB	myMMatrix.txt
-rw-r--r--	cloudera	cloudera	32 B	Wed Aug 06 23:14:10 -0700 2025	1	128 MB	myNMatrix

Hadoop, 2017.



Browse Directory

Browse Directory							
/user/cloudera/matrixOutput							
Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rw-r--r--	cloudera	cloudera	0 B	Wed Aug 06 23:16:33 -0700 2025	1	128 MB	_SUCCESS
-rw-r--r--	cloudera	cloudera	36 B	Wed Aug 06 23:16:33 -0700 2025	1	128 MB	part-r-00000

Hadoop, 2017.

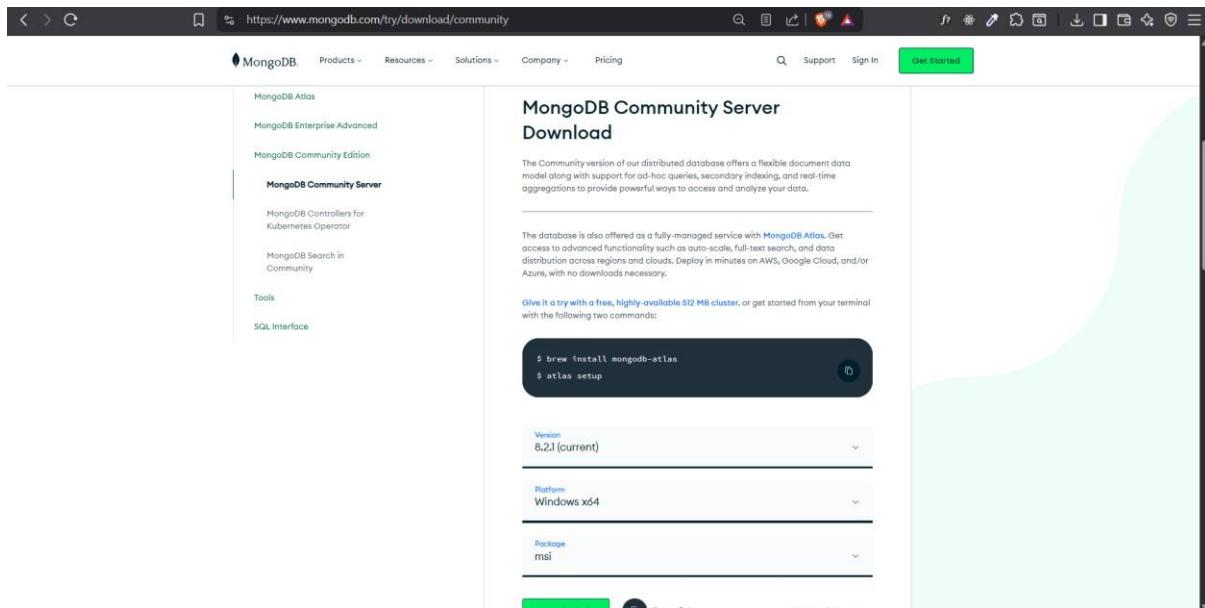
```
Total time spent by all reduce tasks (ms)=6676
Total vcore-milliseconds taken by all map tasks=34995
Total vcore-milliseconds taken by all reduce tasks=6676
Total megabyte-milliseconds taken by all map tasks=17917440
Total megabyte-milliseconds taken by all reduce tasks=3418112
Map-Reduce Framework
    Map input records=8
    Map output records=8000
    Map output bytes=95120
    Map output materialized bytes=31054
    Input split bytes=272
    Combine input records=0
    Combine output records=0
    Reduce input groups=3996
    Reduce shuffle bytes=31054
    Reduce input records=8000
    Reduce output records=4
    Spilled Records=16000
    Shuffled Maps =2
    Failed Shuffles=0
    Merged Map outputs=2
    GC time elapsed (ms)=451
    CPU time spent (ms)=1730
    Physical memory (bytes) snapshot=379375616
    Virtual memory (bytes) snapshot=2108305408
    Total committed heap usage (bytes)=152174592
Shuffle Errors
    BAD_ID=0
    CONNECTION=0
    IO_ERROR=0
    WRONG_LENGTH=0
    WRONG_MAP=0
    WRONG_REDUCE=0
File Input Format Counters
    Bytes Read=64
File Output Format Counters
    Bytes Written=36
[cloudera@quickstart ~]$ hdfs dfs -ls
Found 9 items
drwx-----  - cloudera cloudera      0 2025-08-06 22:00 .Trash
drwx-----  - cloudera cloudera      0 2025-08-06 23:16 .staging
-rw-r--r--  1 cloudera cloudera  6555 2025-08-06 23:08 MatrixMultiplication.jar
-rw-r--r--  1 cloudera cloudera  4887 2025-08-05 23:20 WordCount.jar
drwxr-xr-x  - cloudera cloudera      0 2025-07-30 00:31 manas
drwxr-xr-x  - cloudera cloudera      0 2025-08-06 23:14 matrixInput
drwxr-xr-x  - cloudera cloudera      0 2025-08-06 23:16 matrixOutput
-rw-r--r--  1 cloudera cloudera     52 2025-08-05 23:25 myInputFile.txt
drwxr-xr-x  - cloudera cloudera      0 2025-08-06 21:43 myOutput
[cloudera@quickstart ~]$ hdfs dfs -ls /user/cloudera/matrixOutput
Found 2 items
-rw-r--r--  1 cloudera cloudera      0 2025-08-06 23:16 /user/cloudera/matrixOutput/_SUCCESS
-rw-r--r--  1 cloudera cloudera    36 2025-08-06 23:16 /user/cloudera/matrixOutput/part-r-00000
[cloudera@quickstart ~]$ hdfs dfs -cat /user/cloudera/matrixOutput/part-r-00000
0,0,11.0
0,1,10.0
1,0,25.0
1,1,26.0
[cloudera@quickstart ~]$ █
```

Practical No 3 : MongoDB

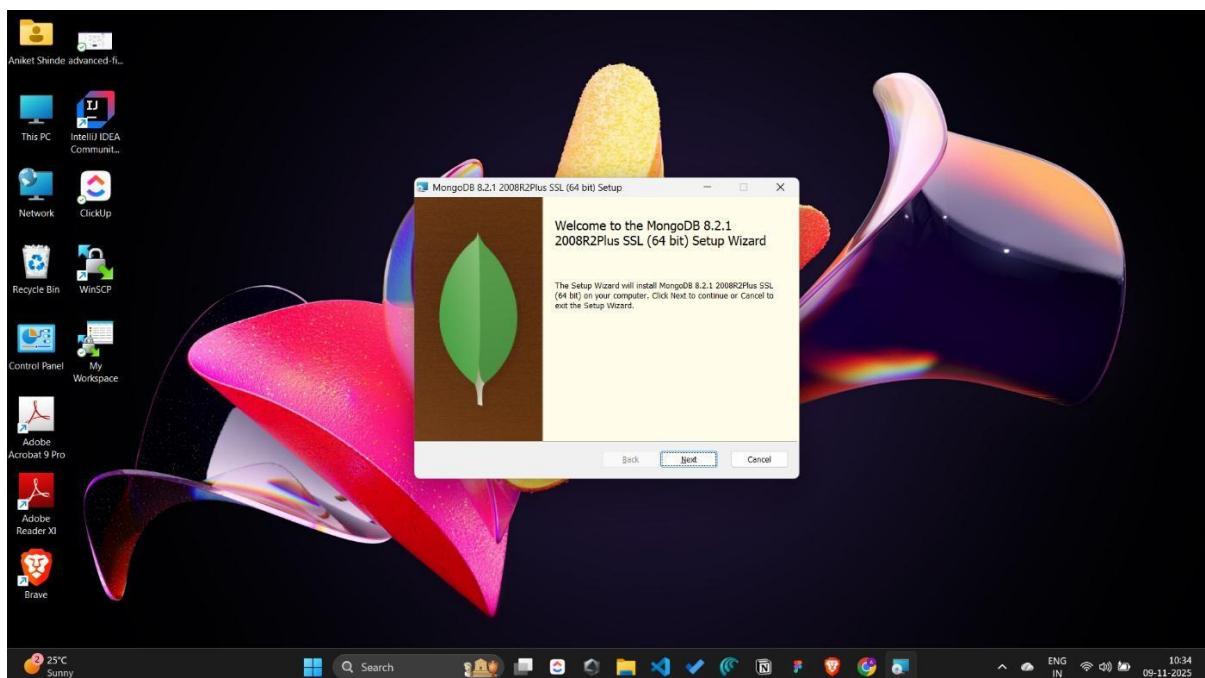
1. Installation

Installing MongoDB on Windows

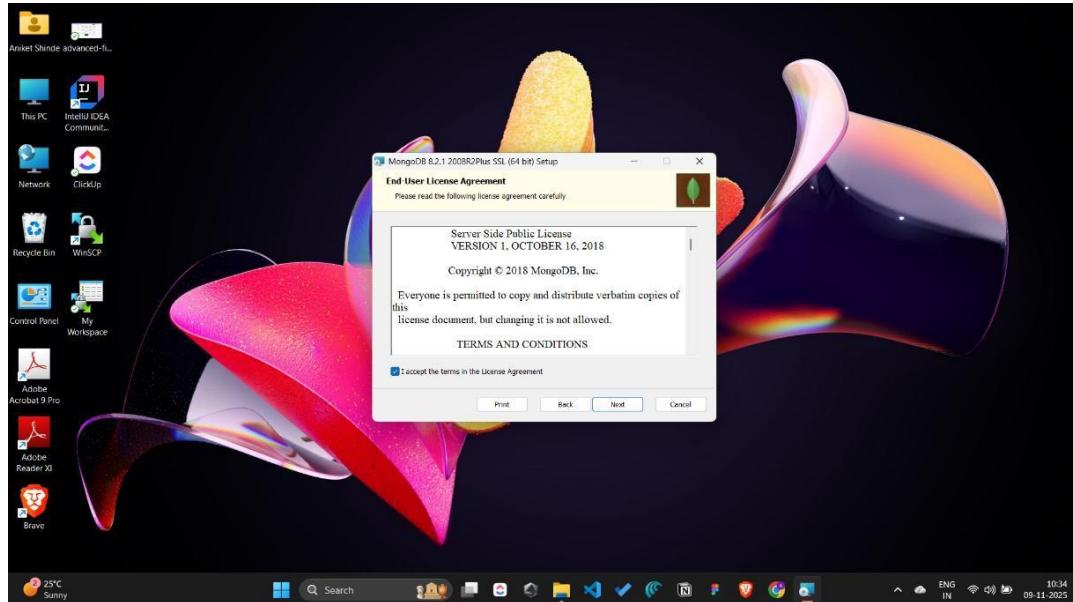
Download MongoDB Community Server. We will install **the 64-bit version** for Windows.



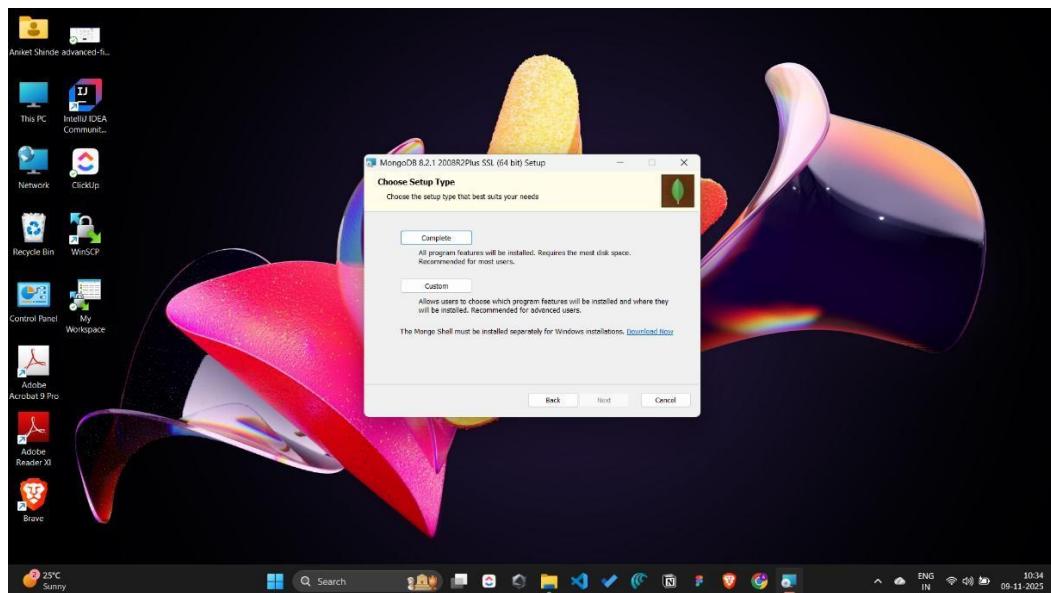
2. Click on Setup. Once download is complete open the MSI file. Click Next in the start-up screen.



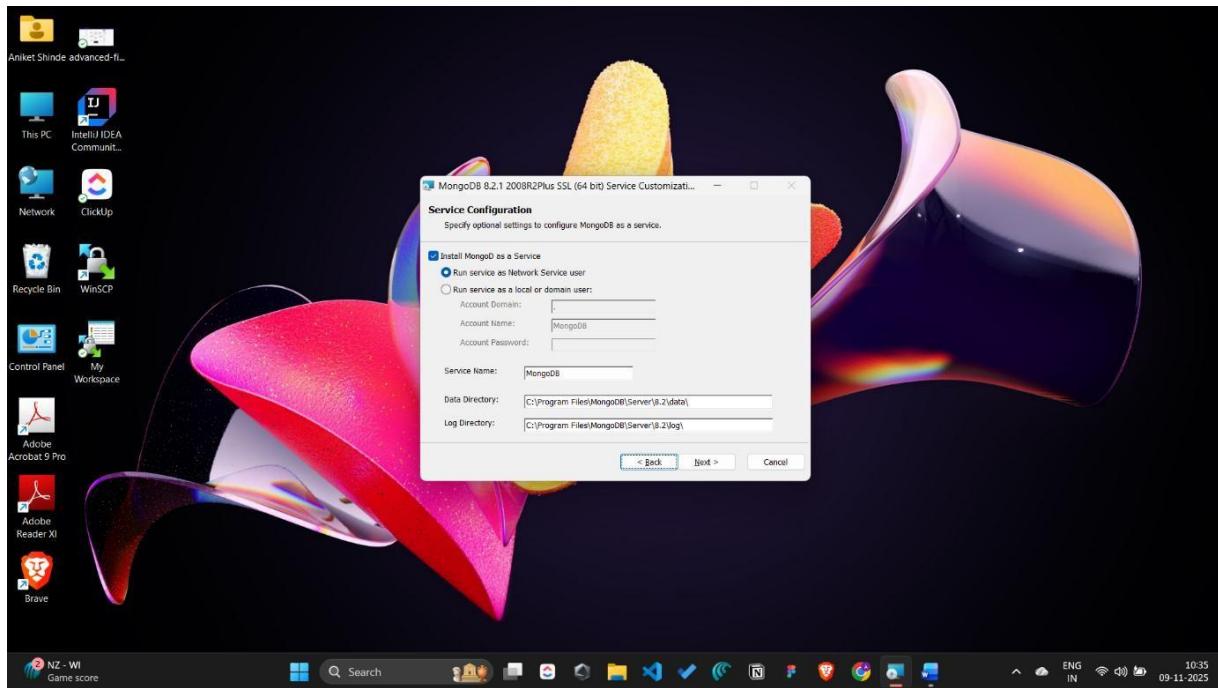
3. Accept the End-User License Agreement & Click Next.



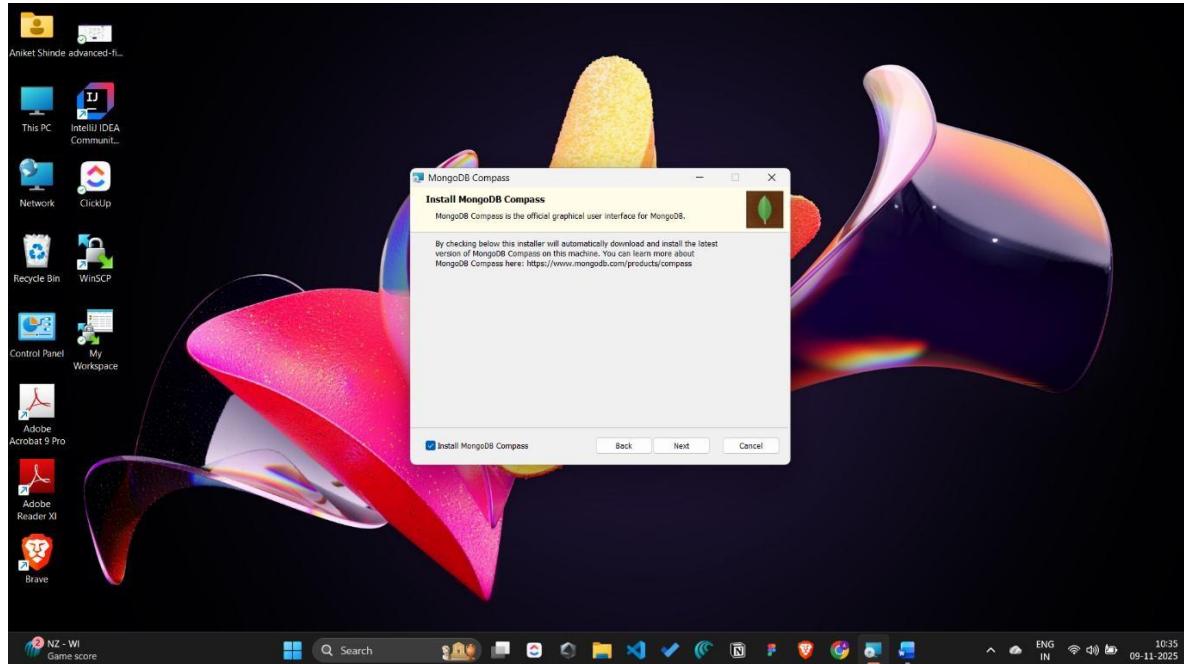
4. Click on the "complete" button to install all of the components. The custom option can be used to install selective components or if you want to change the location of the installation.



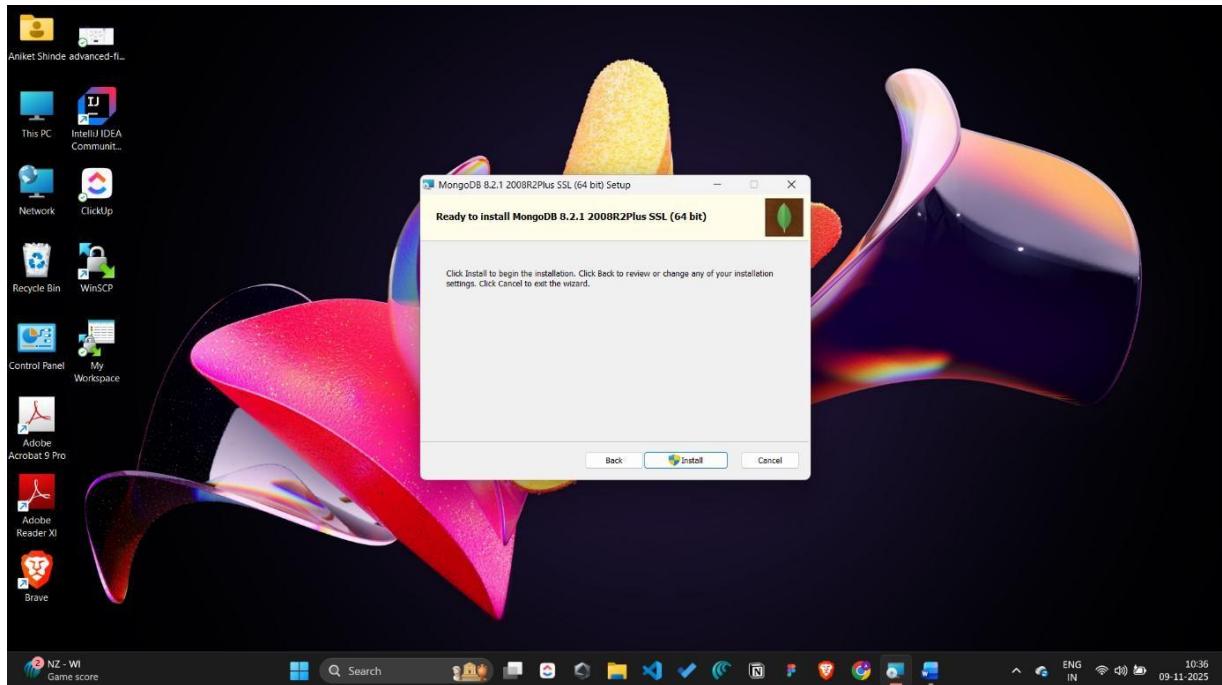
5. Service Configuration: Select “Run service as Network Service user” --> Click Next.



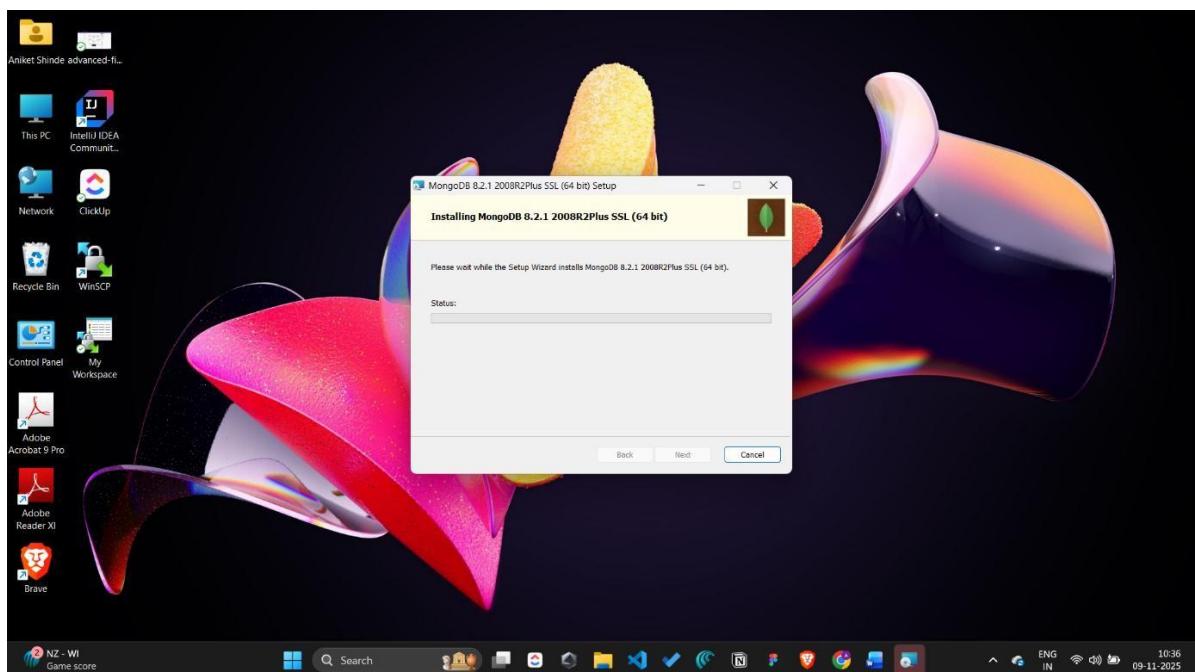
6. Install Mongodb Compass GUI (Optional).



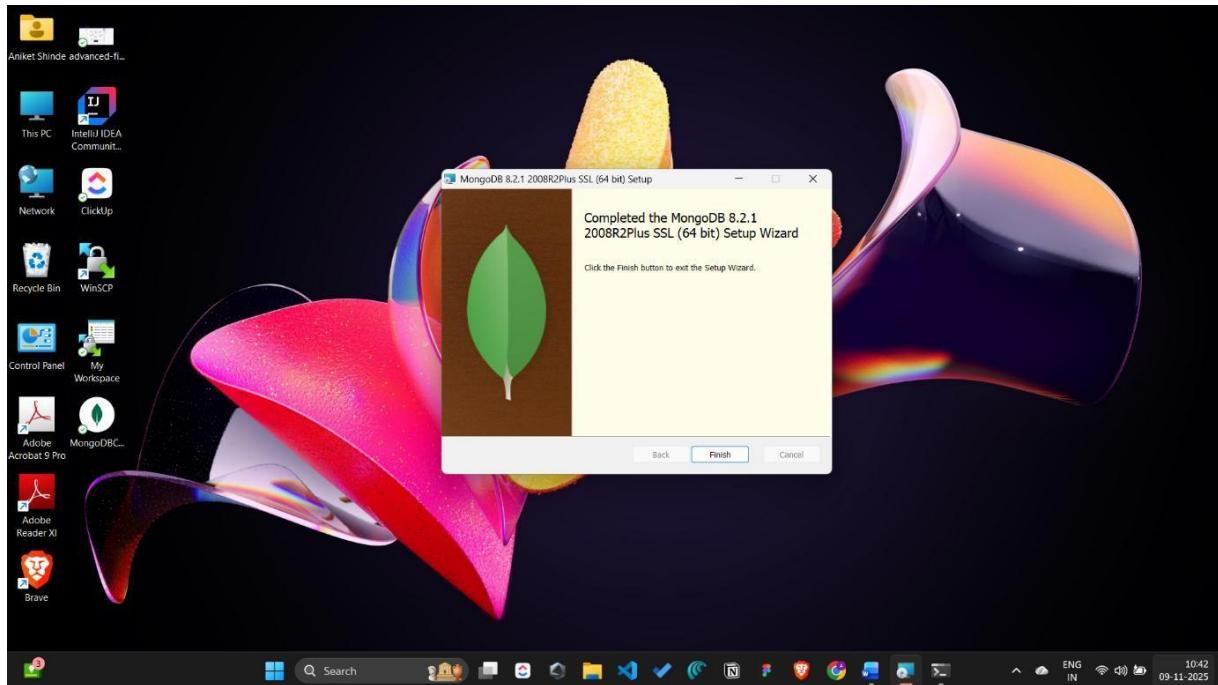
7. Click on the Install button to start the installation.



8. Installation begins. Click Next once completed.



9. Once complete the installation, Click on the Finish button.



2. Sample Database Creation

1. Now create the DB – College and use College

A screenshot of a terminal window titled "mongosh mongodb://127.0.0.1:27017". The command "use College" is entered, followed by "already on db College". The terminal window has a dark background with white text.

3. Query the Sample Database using MongoDB querying commands

1. Create a query & insert a record with insertOne()

A screenshot of a terminal window titled "mongosh mongodb://127.0.0.1:27017". The window contains the following MongoDB shell session:

```
College> db.employee.insertOne({
...   "empcode": "24108",
...   "empfname": "Aniket",
...   "emplname": "Shinde",
...   "job": "SDE",
...   "manager": 2001,
...   "hiredate": "2025-11-09",
...   "salary": 33000,
...   "commission": 0,
...   "depcode": 20
... });
{
  acknowledged: true,
  insertedId: ObjectId('69102e43d5deaae46b63b112')
}
College>
```

2. Now, try to insert more records with - InsertMany()

```
mongosh mongodb://127.0.0.1:27017/College + v

College> db.employee.insertMany([
...   {
...     "empcode": "24109",
...     "empfname": "Sai",
...     "emplname": "Sanas",
...     "job": "SDE",
...     "manager": 2001,
...     "hiredate": "2025-11-09",
...     "salary": 34000,
...     "commission": 0,
...     "depcode": 20
...   },
...   {
...     "empcode": "24110",
...     "empfname": "Vikrant",
...     "emplname": "Paimode",
...     "job": "SDE",
...     "manager": 2001,
...     "hiredate": "2025-11-09",
...     "salary": 35000,
...     "commission": 0,
...     "depcode": 20
...   },
...   {
...     "empcode": "24111",
...     "empfname": "Rutik",
...     "emplname": "Adhav",
...     "job": "SDE",
...     "manager": 2001,
...     "hiredate": "2025-11-09",
...     "salary": 32000,
...     "commission": 0,
...     "depcode": 20
...   },
...   {
...     "empcode": "24112",
...     "empfname": "Alfez",
...     "emplname": "Patel",
...     "job": "SDE",
...     "manager": 2001,
...     "hiredate": "2025-11-09",
...     "salary": 31000,
...     "commission": 0,
...     "depcode": 20
...   }
... ],
... {
...   acknowledged: true,
...   insertedIds: {
...     '0': ObjectId('69102e94d5deaae46b63b113'),
...     '1': ObjectId('69102e94d5deaae46b63b114'),
...     '2': ObjectId('69102e94d5deaae46b63b115'),
...     '3': ObjectId('69102e94d5deaae46b63b116')
...   }
... }
```

3. To view all inserted records use query - find()

```
College> db.employee.find()
[
  {
    _id: ObjectId('69102e43d5deaae46b63b112'),
    empcode: '24108',
    empfname: 'Aniket',
    emplname: 'Shinde',
    job: 'SDE',
    manager: 2001,
    hiredate: '2025-11-09',
    salary: 33000,
    commission: 0,
    depcode: 20
  },
  {
    _id: ObjectId('69102e94d5deaae46b63b113'),
    empcode: '24109',
    empfname: 'Sai',
    emplname: 'Sanas',
    job: 'SDE-I',
    manager: 2001,
    hiredate: '2025-10-15',
    salary: 34500,
    commission: 0,
    depcode: 20
  },
  {
    _id: ObjectId('69102e94d5deaae46b63b114'),
    empcode: '24110',
    empfname: 'Vikrant',
    emplname: 'Paimode',
    job: 'SDE-II',
    manager: 2001,
    hiredate: '2025-09-28',
    salary: 36200,
    commission: 0,
    depcode: 20
  },
  {
    _id: ObjectId('69102e94d5deaae46b63b115'),
    empcode: '24111',
    empfname: 'Rutik',
    emplname: 'Adhav',
    job: 'Backend Engineer',
    manager: 2001,
    hiredate: '2025-08-10',
    salary: 33400,
    commission: 0,
    depcode: 20
  },
  {
    _id: ObjectId('69102e94d5deaae46b63b116'),
    empcode: '24112',
    empfname: 'Alfez',
    emplname: 'Patel',
    job: 'Frontend Engineer',
    manager: 2001,
    hiredate: '2025-07-25',
    salary: 31250,
    commission: 0,
    depcode: 20
  }
]
```

```
College>
```

4. To find one record with empname as Talha – use findOne()

```
        }
    ]
College> db.employee.findOne({empfname: "Aniket"})
{
  _id: ObjectId('69102e43d5deaae46b63b112'),
  empcode: '24108',
  empfname: 'Aniket',
  emplname: 'Shinde',
  job: 'SDE',
  manager: 2001,
  hiredate: '2025-11-09',
  salary: 33000,
  commission: 0,
  depcode: 20
}
College>
```

5. To find employee with job role as Backend Engineer– use find()

```
College> db.employee.find({ job: "Backend Engineer" });
[
  {
    _id: ObjectId('69102e94d5deaae46b63b115'),
    empcode: '24111',
    empfname: 'Rutik',
    emplname: 'Adhav',
    job: 'Backend Engineer',
    manager: 2001,
    hiredate: '2025-08-10',
    salary: 33400,
    commission: 0,
    depcode: 20
  }
]
College> |
```

6. Finding all the records where salary is greater than 33000.

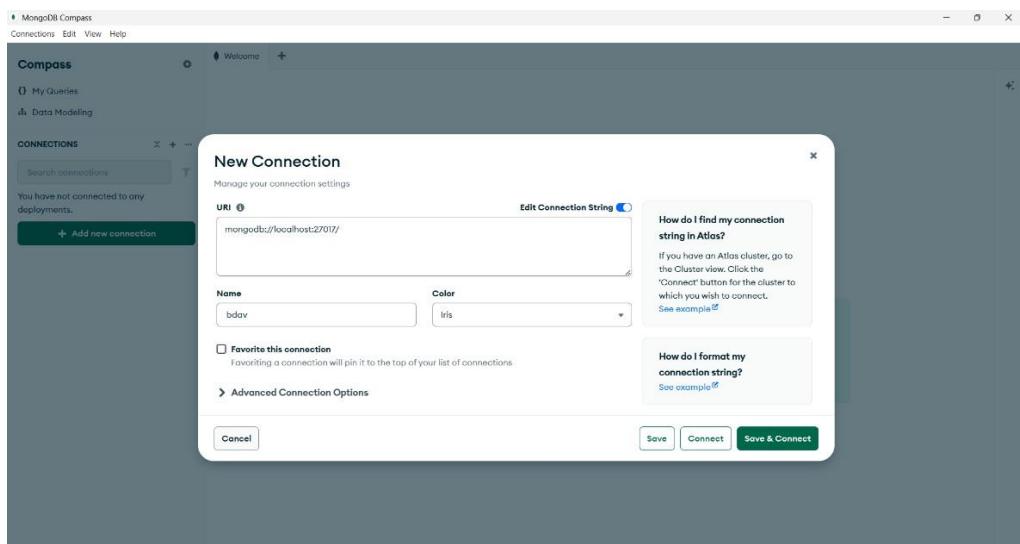
```
College> db.employee.find({ salary: { $gt: 33000 } });
[
  {
    _id: ObjectId('69102e94d5deaae46b63b113'),
    empcode: '24109',
    empfname: 'Sai',
    emplname: 'Sanas',
    job: 'SDE-I',
    manager: 2001,
    hiredate: '2025-10-15',
    salary: 34500,
    commission: 0,
    depcode: 20
  },
  {
    _id: ObjectId('69102e94d5deaae46b63b114'),
    empcode: '24110',
    empfname: 'Vikrant',
    emplname: 'Paimode',
    job: 'SDE-II',
    manager: 2001,
    hiredate: '2025-09-28',
    salary: 36200,
    commission: 0,
    depcode: 20
  },
  {
    _id: ObjectId('69102e94d5deaae46b63b115'),
    empcode: '24111',
    empfname: 'Rutik',
    emplname: 'Adhav',
    job: 'Backend Engineer',
    manager: 2001,
    hiredate: '2025-08-10',
    salary: 33400,
    commission: 0,
    depcode: 20
  }
]
College>
```

7. Finding the records with salary is > 34000 and depcode is 20.

```
College> db.employee.find({salary: { $gt: 34000 }, depcode: 20});
[
  {
    _id: ObjectId('69102e94d5deaae46b63b113'),
    empcode: '24109',
    empfname: 'Sai',
    emplname: 'Sanas',
    job: 'SDE-I',
    manager: 2001,
    hiredate: '2025-10-15',
    salary: 34500,
    commission: 0,
    depcode: 20
  },
  {
    _id: ObjectId('69102e94d5deaae46b63b114'),
    empcode: '24110',
    empfname: 'Vikrant',
    emplname: 'Paimode',
    job: 'SDE-II',
    manager: 2001,
    hiredate: '2025-09-28',
    salary: 36200,
    commission: 0,
    depcode: 20
  }
]
College>
```

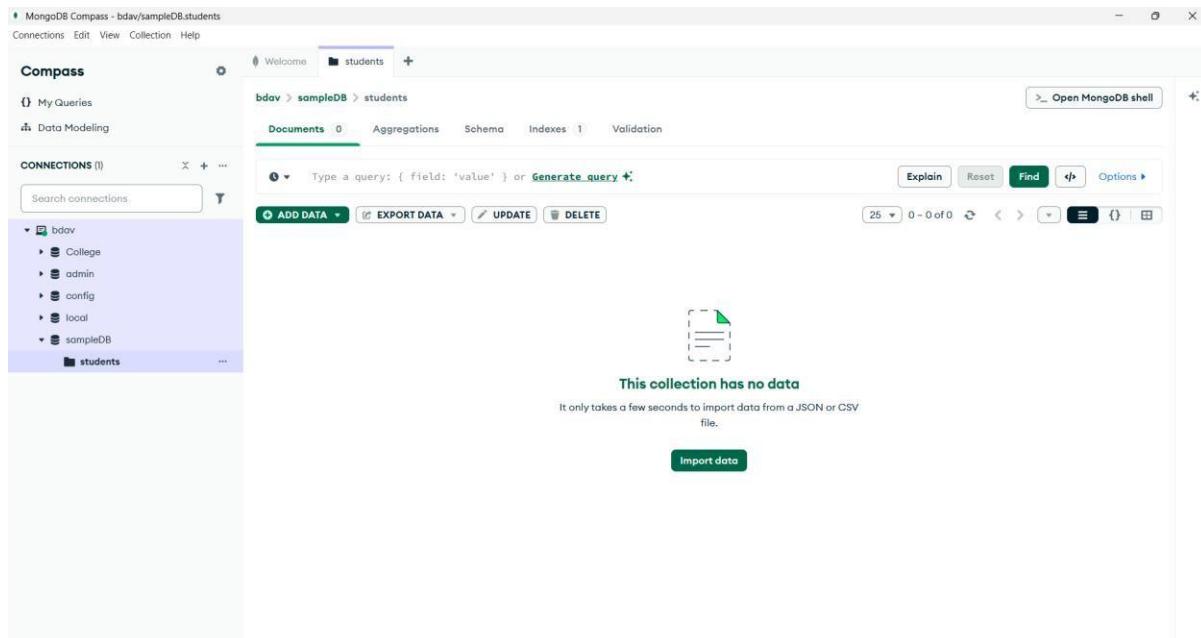
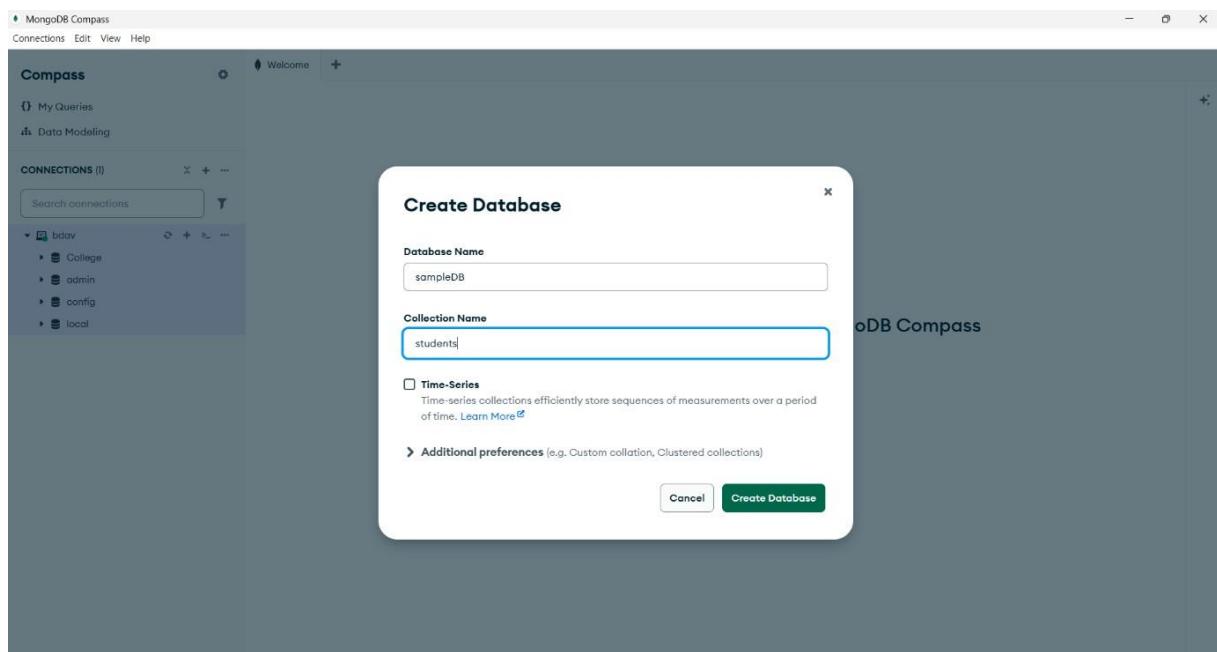
4. Operations on the collection with MongoDB Compass GUI.

For that we have to open MongoDB Compass & create the new connection with following values and make sure to click on Save and connect.



A. Create Collection:

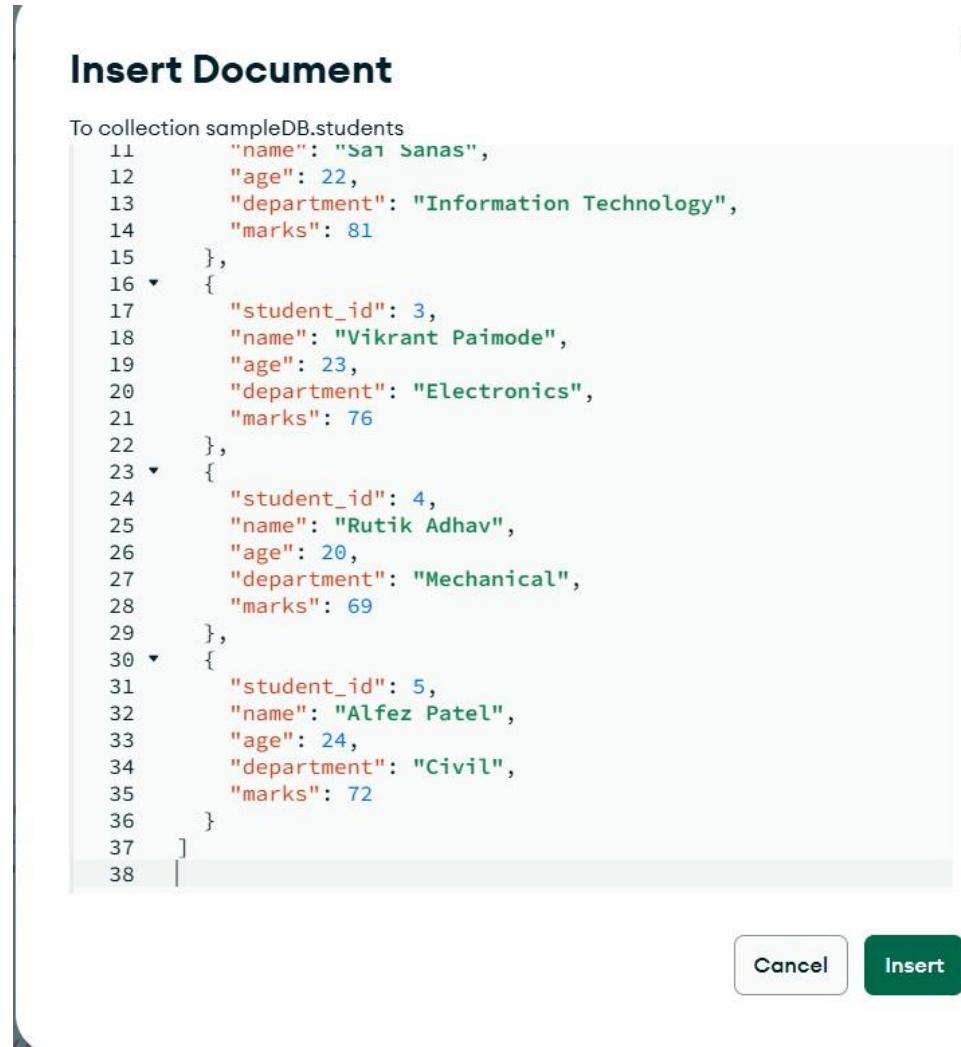
1. Create a database name as Students.



B. Insert Document

1. click on the add data -> and click on the insert document.

1. Insert Multiple values in that document:



_id	student_id	name	age	department	marks
_id: ObjectId('59183484d601ea082cc0c30')	3	Sal Sanas	22	Information Technology	81
_id: ObjectId('59183484d601ea082cc0c31')	3	Vikrant Paimode	23	Electronics	76
_id: ObjectId('59183484d601ea082cc0c32')	4	Rutik Adhav	20	Mechanical	69

C. Query Document

Query the document for finding the name : "Talha Shinde"

The screenshot shows the MongoDB Compass interface. In the top navigation bar, the connection is set to 'bdav' and the database is 'sampleDB'. The collection selected is 'students'. The search bar contains the query '{ "name": "Aniket Shinde" }'. Below the search bar, there are buttons for 'Generate query', 'Explain', 'Reset', 'Find', and 'Options'. The results section displays four documents. The first document is highlighted with a blue border:

```
_id: ObjectId("59163484d601ea682c0c31")
student_id: 1
name : "Aniket Shinde"
age : 21
department : "Computer Science"
marks : 88
```

The other three documents are shown below it:

```
_id: ObjectId("59163484d601ea682c0c31")
student_id: 2
name : "Sai Saneh"
age : 22
department : "Information Technology"
marks : 81
```

```
_id: ObjectId("59163484d601ea682c0c32")
student_id: 3
name : "Vishrant Patmode"
age : 23
department : "Electronics"
marks : 76
```

```
_id: ObjectId("59163484d601ea682c0c33")
student_id: 4
name : "Rutik Adhav"
age : 23
department : "Computer Science"
marks : 88
```

This screenshot is identical to the one above, showing the MongoDB Compass interface with the same search query '{ "name": "Aniket Shinde" }'. The results section displays the same four documents, with the first one highlighted:

```
_id: ObjectId("59163484d601ea682c0c31")
student_id: 1
name : "Aniket Shinde"
age : 21
department : "Computer Science"
marks : 88
```

Query the document and find the Marks greater than 80 and age :20

The screenshot shows the MongoDB Compass interface. The left sidebar displays connections: bdav, College, admin, config, local, startup_log, and sampleDB, with sampleDB selected. The main area shows the 'students' collection under 'sampleDB'. A search bar at the top contains the query: {marks: { \$gt: 80 }, age: 20}. Below the search bar, there are buttons for ADD DATA, EXPORT DATA, UPDATE, and DELETE. The results pane shows one document:

```
_id: ObjectId("59163434d601ea602cd0c31")
student_id: 2
name: "Sai Sana"
age: 22
department: "Information Technology"
marks: 81
```

D. Delete Document

The three screenshots illustrate the step-by-step process of deleting a document from the 'students' collection.

- Screenshot 1:** Shows the 'students' collection with one document. The document details are: _id: ObjectId("60103484d601ea002cd9c32"), student_id: 3, name: "Vikrant Pai mode", age: 23, department: "Electronics", marks: 76. The 'DELETE' button is highlighted.
- Screenshot 2:** Shows the same document selected. A red box highlights the 'Remove document' button in the toolbar above the results pane.
- Screenshot 3:** Shows the document flagged for deletion. A red box highlights the 'DELETE' button in the bottom right corner of the results pane.

The screenshot shows the MongoDB Compass interface. On the left, the 'Connections' sidebar lists 'bdav' (selected), 'College', 'admin', 'config', 'local' (selected), and 'sampleDB'. Under 'local', 'startup_log' is selected. In the center, the 'sampleDB' section shows the 'students' collection with 4 documents. A search bar at the top contains the query '{name : "Vikrant Pai mode"}'. Below the search bar are buttons for 'ADD DATA', 'EXPORT DATA', 'UPDATE', and 'DELETE'. To the right, there are buttons for 'Generate query', 'Explain', 'Reset', 'Find', 'Options', and a dropdown menu. At the bottom, it says '25 1 - 0 of 0' and 'No results' with a note 'Try modifying your query to get results.'

E. Indexing

```
mongosh mongodb://127.0.0.1:27017
College> db.students.createIndex({ student_id: 1 })
student_id_1
College> db.students.getIndexes()
[
  { v: 2, key: { _id: 1 }, name: '_id_' },
  { v: 2, key: { student_id: 1 }, name: 'student_id_1' }
]
College>
```

Practical No 4: Hive

1. Hive Data Types

Primitive Data Types:

- **Numeric Types:**

- TINYINT: 1-byte signed integer.
- SMALLINT: 2-byte signed integer.
- INT / INTEGER: 4-byte signed integer.
- BIGINT: 8-byte signed integer.
- FLOAT: 4-byte single-precision floating-point number.
- DOUBLE: 8-byte double-precision floating-point number.
- DECIMAL: Fixed-point decimal numbers with user-defined precision and scale.

- **Date/Time Types:**

- DATE: Represents a date in 'YYYY-MM-DD' format.
- TIMESTAMP: Represents a point in time with date and time, including optional nanosecond precision.
- INTERVAL: Represents a period of time.

- **String Types:**

- STRING: Variable-length character string.
- VARCHAR(n): Variable-length character string with a maximum length of 'n'.
- CHAR(n): Fixed-length character string with a length of 'n', padded with spaces if shorter.

- Miscellaneous Types:

- BOOLEAN: Stores TRUE or FALSE values.
- BINARY: Stores a sequence of bytes.

- Complex Data Types:

- ARRAY<data_type>:

An ordered collection of elements of the same data_type. Similar to a list in other programming languages.

- MAP<key_data_type, value_data_type>:

An associative array (key-value pairs) where keys are unique and of key_data_type, and values are of value_data_type

- STRUCT<col_name : data_type, ...>:

A record type that groups together fields (columns) of potentially different data types, similar to a struct in C or a record in other systems.

- UNIONTYPE<data_type1, data_type2, ...>:

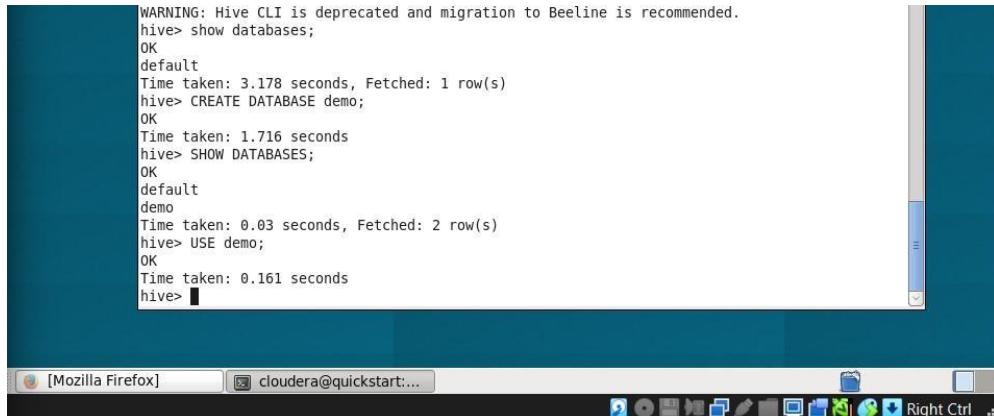
A type that can hold a value of one of several specified data types at any given time.

2. Create Database & Table in Hive

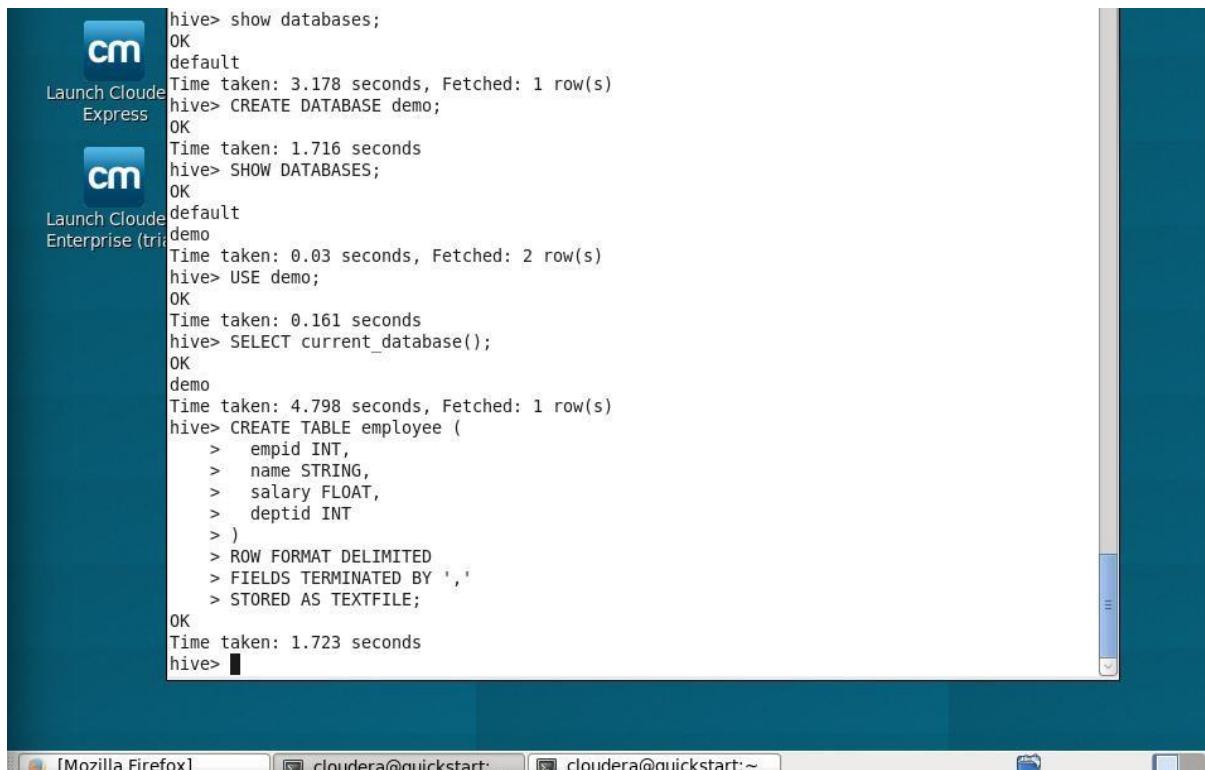
a. show databases - To check default database provided by Hive



```
Logging initialized using configuration in jar:file:/usr/jars/hive-common-1.1.0-cdh5.4.2.jar!/hive-log4j.properties
WARNING: Hive CLI is deprecated and migration to Beeline is recommended.
hive> show databases;
OK
default
Time taken: 3.178 seconds, Fetched: 1 row(s)
hive> █
```

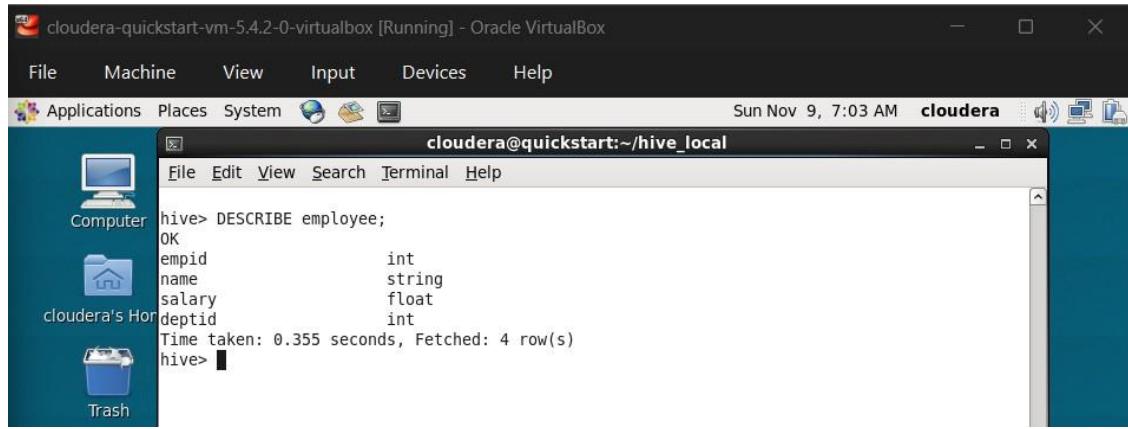
b. create database:- To create a new database

```
WARNING: Hive CLI is deprecated and migration to Beeline is recommended.
hive> show databases;
OK
default
Time taken: 3.178 seconds, Fetched: 1 row(s)
hive> CREATE DATABASE demo;
OK
Time taken: 1.716 seconds
hive> SHOW DATABASES;
OK
default
demo
Time taken: 0.03 seconds, Fetched: 2 row(s)
hive> USE demo;
OK
Time taken: 0.161 seconds
hive> ■
```

c. Create table employee

```
hive> show databases;
OK
default
Time taken: 3.178 seconds, Fetched: 1 row(s)
hive> CREATE DATABASE demo;
OK
Time taken: 1.716 seconds
hive> SHOW DATABASES;
OK
default
demo
Time taken: 0.03 seconds, Fetched: 2 row(s)
hive> USE demo;
OK
Time taken: 0.161 seconds
hive> SELECT current_database();
OK
demo
Time taken: 4.798 seconds, Fetched: 1 row(s)
hive> CREATE TABLE employee (
    >     empid INT,
    >     name STRING,
    >     salary FLOAT,
    >     deptid INT
    > )
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY ','
    > STORED AS TEXTFILE;
OK
Time taken: 1.723 seconds
hive> ■
```

D. Describe the table employee



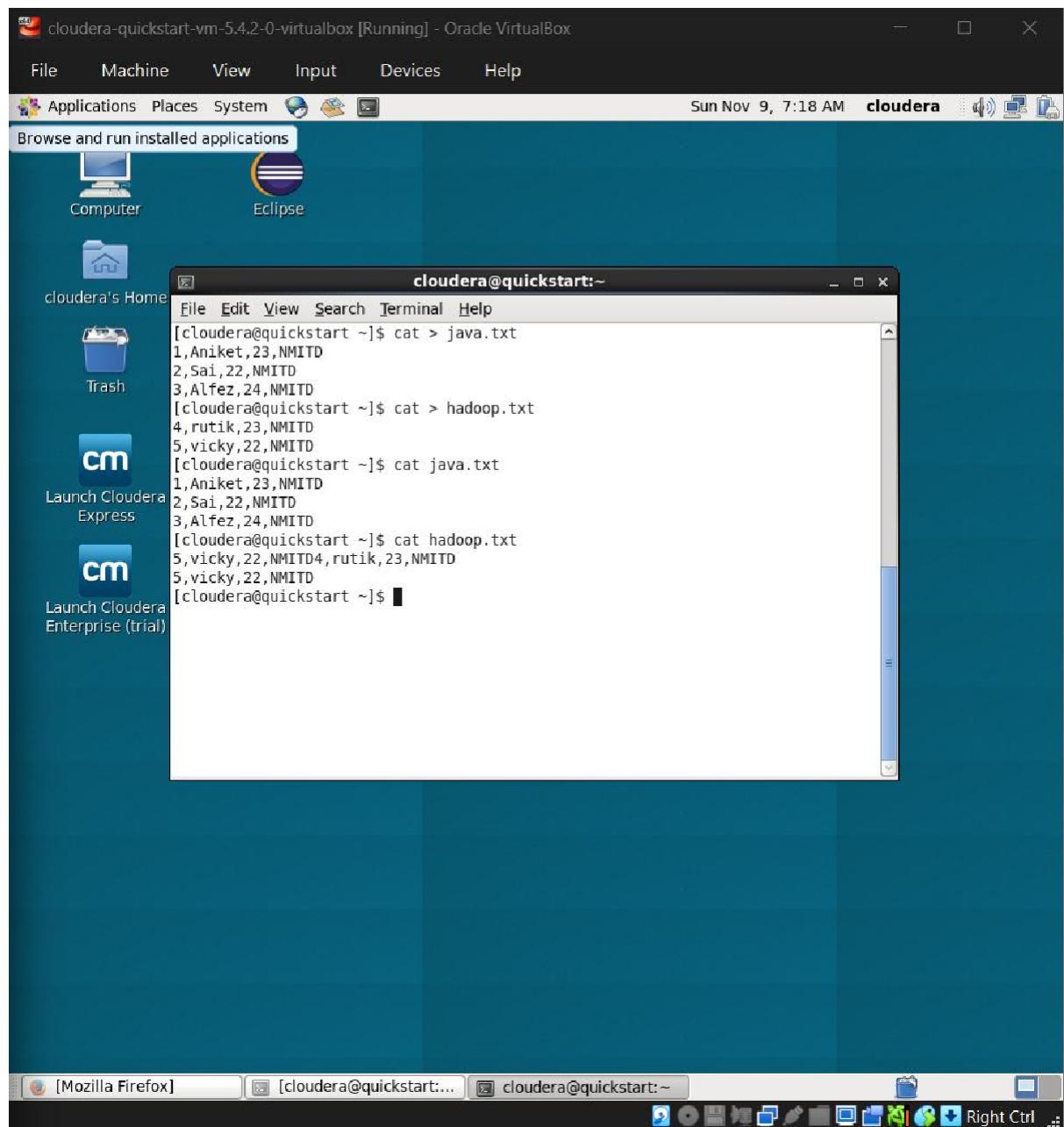
A screenshot of a terminal window titled "cloudera@quickstart:~/hive local". The window is part of an Oracle VirtualBox interface. The terminal shows the following Hive command and its output:

```
hive> DESCRIBE employee;
OK
empid          int
name           string
salary         float
deptid         int
Time taken: 0.355 seconds, Fetched: 4 row(s)
hive> ■
```

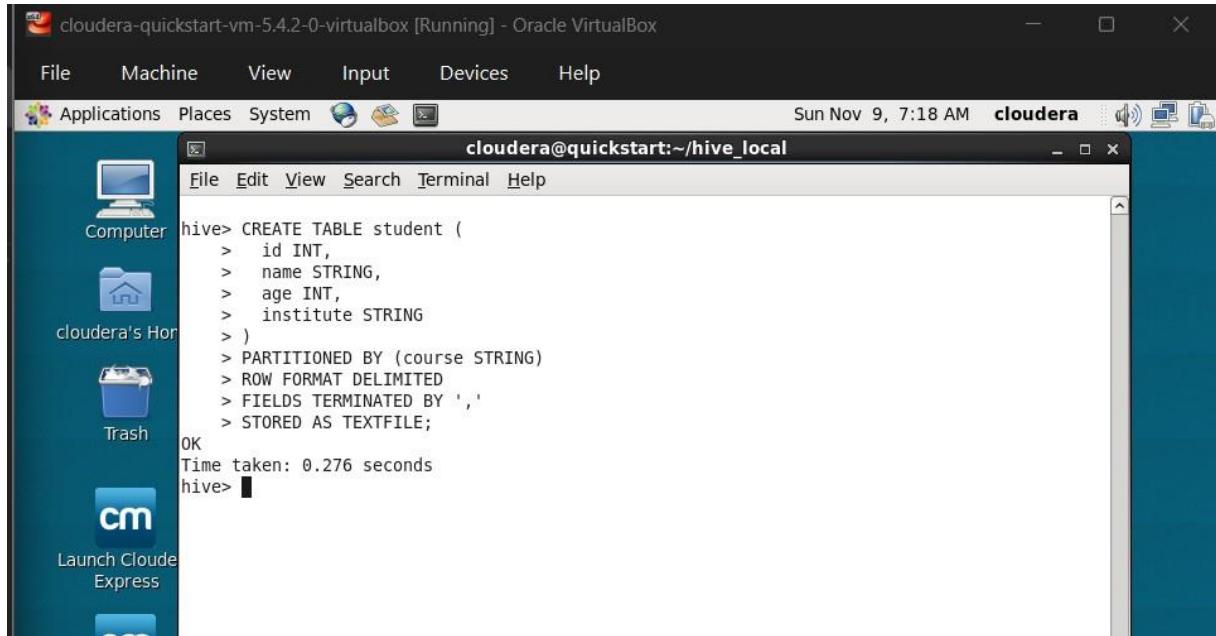
3. Hive Partitioning

A. Static Partition

To learn how to store Hive table data in **separate directory partitions** (for faster queries and organization).



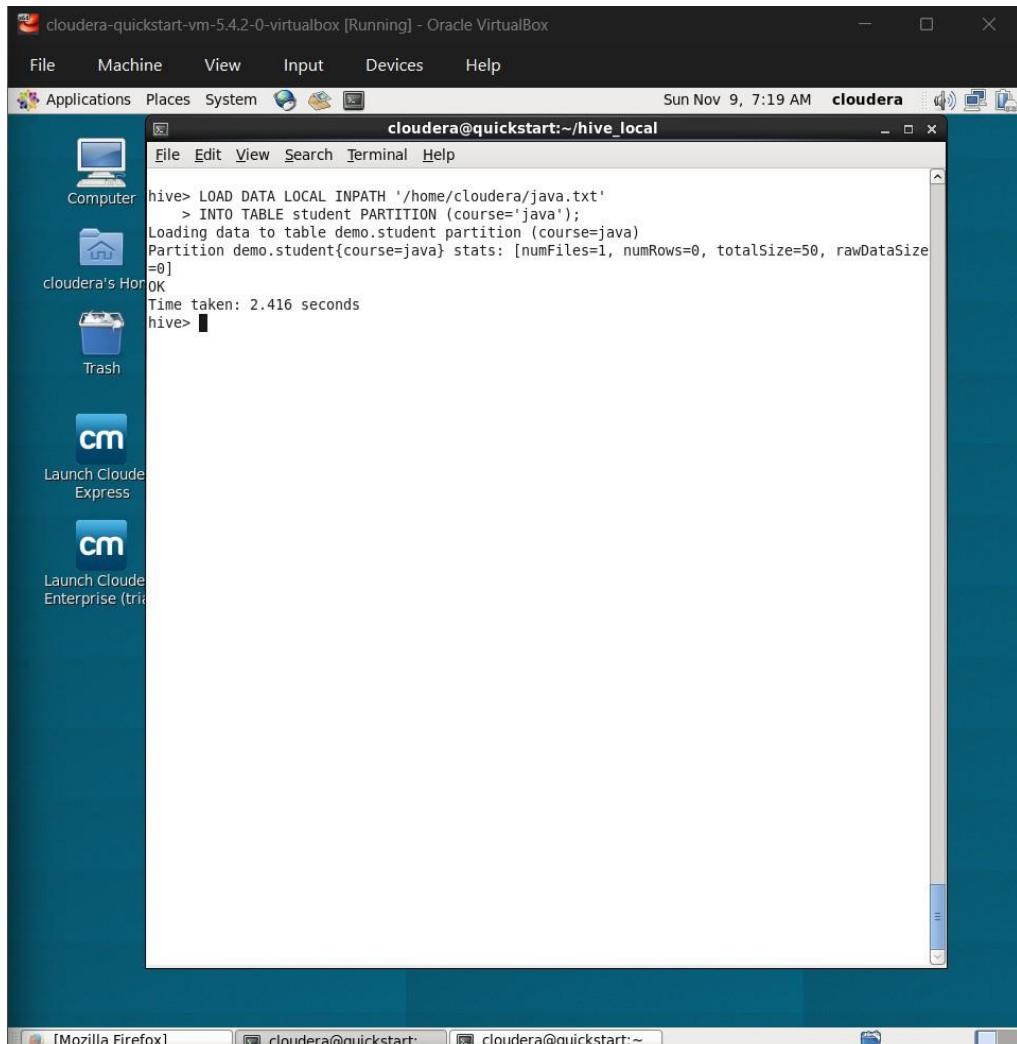
a. For this we are creating table student –



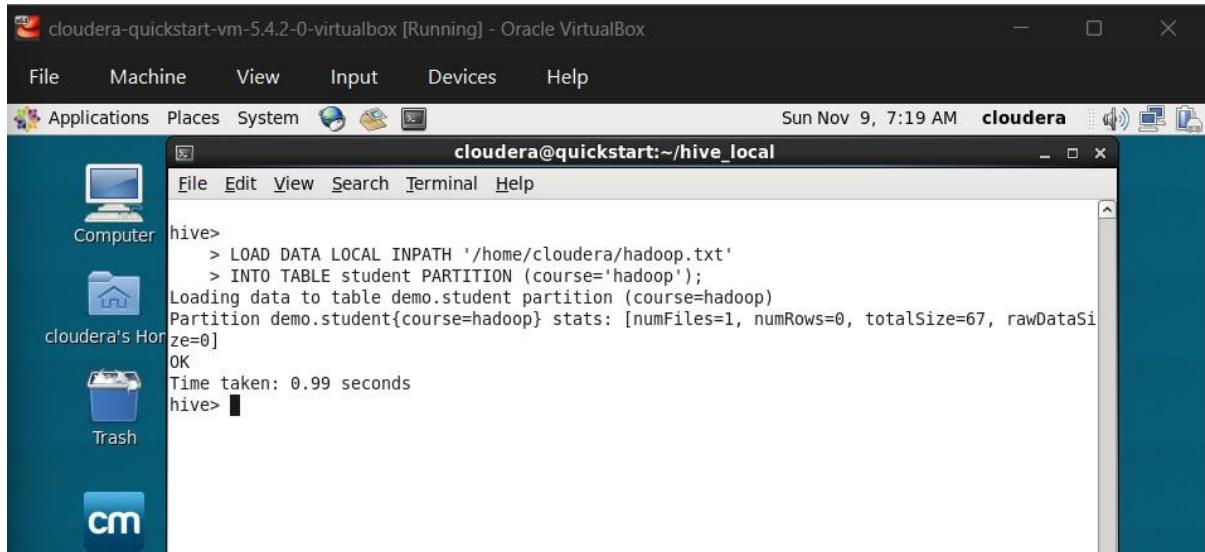
```
hive> CREATE TABLE student (
    >   id INT,
    >   name STRING,
    >   age INT,
    >   institute STRING
    > )
    > PARTITIONED BY (course STRING)
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY ','
    > STORED AS TEXTFILE;
OK
Time taken: 0.276 seconds
hive>
```

b. Load data statically into partitions

Now load each file as a separate partition:



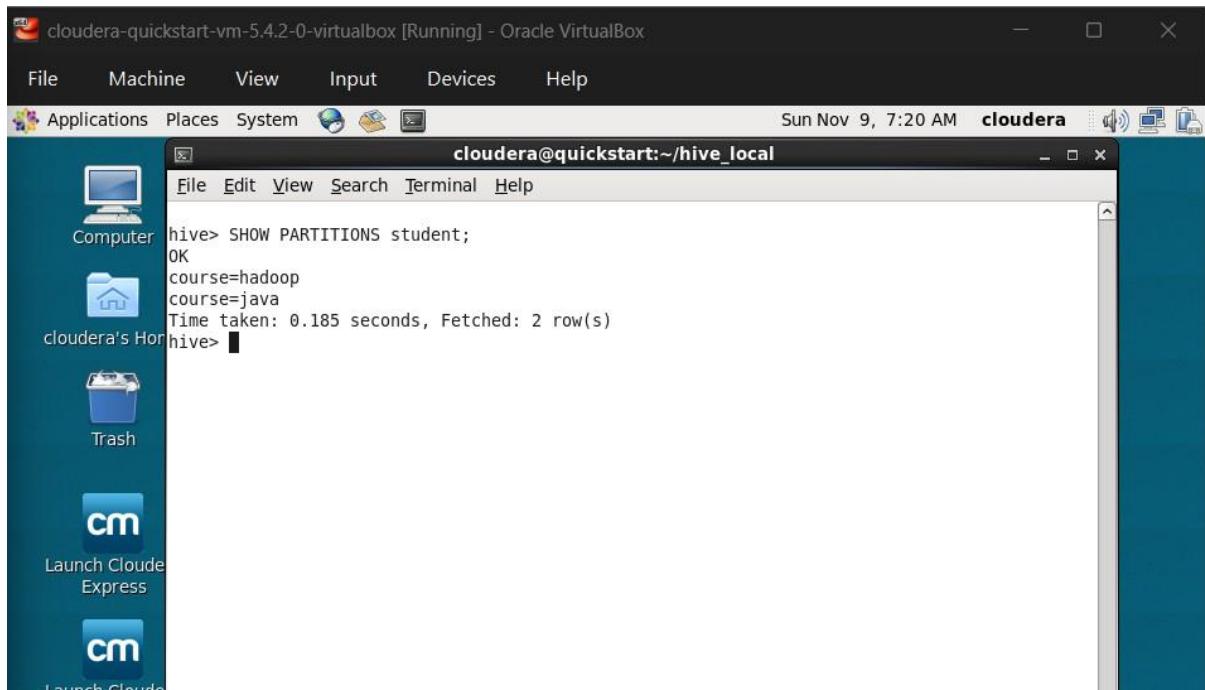
```
hive> LOAD DATA LOCAL INPATH '/home/cloudera/java.txt'
    > INTO TABLE student PARTITION (course='java');
Loading data to table demo.student partition (course=java)
Partition demo.student{course=java} stats: [numFiles=1, numRows=0, totalSize=50, rawDataSize
=0]
OK
Time taken: 2.416 seconds
hive>
```



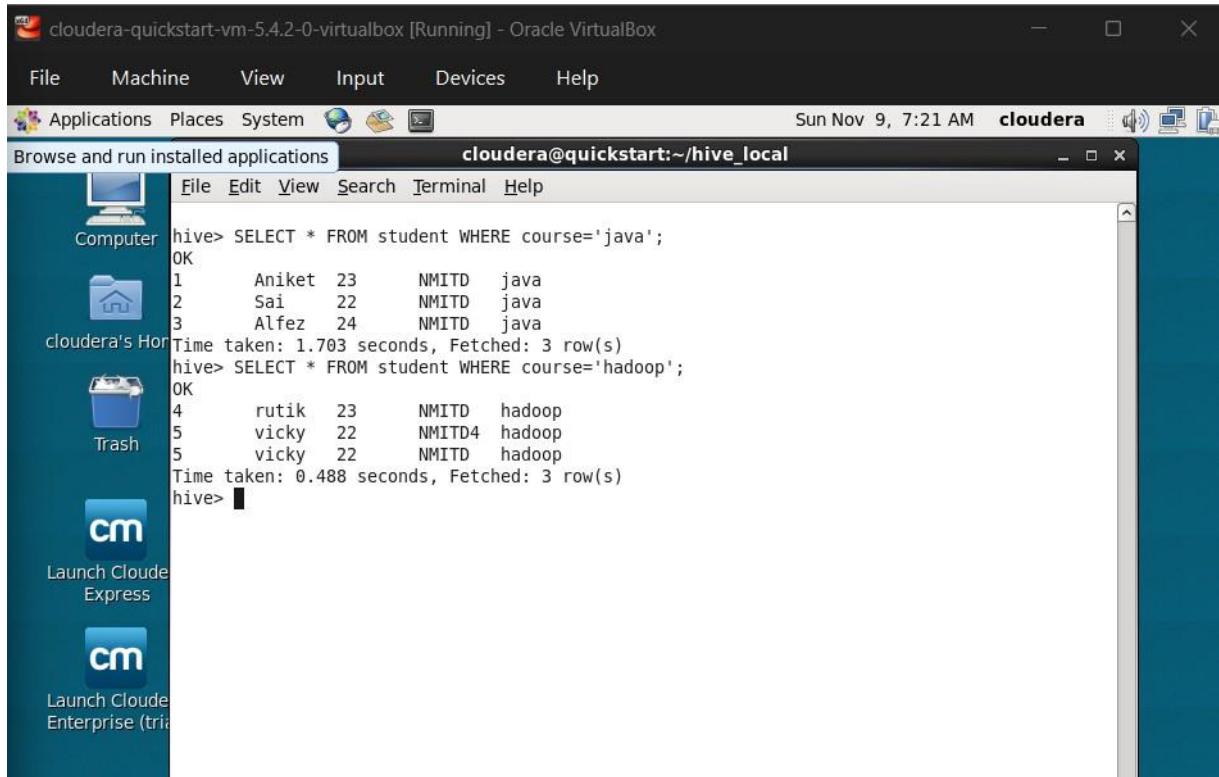
```
hive> > LOAD DATA LOCAL INPATH '/home/cloudera/hadoop.txt'
      > INTO TABLE student PARTITION (course='hadoop');
Loading data to table demo.student partition (course=hadoop)
Partition demo.student{course=hadoop} stats: [numFiles=1, numRows=0, totalSize=67, rawDataSize=0]
OK
Time taken: 0.99 seconds
hive>
```

C. Verify partition structure

Check which partitions exist:



```
hive> SHOW PARTITIONS student;
OK
course=hadoop
course=java
Time taken: 0.185 seconds, Fetched: 2 row(s)
hive>
```



```

cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> SELECT * FROM student WHERE course='java';
OK
1      Aniket  23      NMITD    java
2      Sai      22      NMITD    java
3      Alfez   24      NMITD    java
Time taken: 1.703 seconds, Fetched: 3 row(s)
hive> SELECT * FROM student WHERE course='hadoop';
OK
4      rutik   23      NMITD    hadoop
5      vicky   22      NMITD4   hadoop
5      vicky   22      NMITD    hadoop
Time taken: 0.488 seconds, Fetched: 3 row(s)
hive>

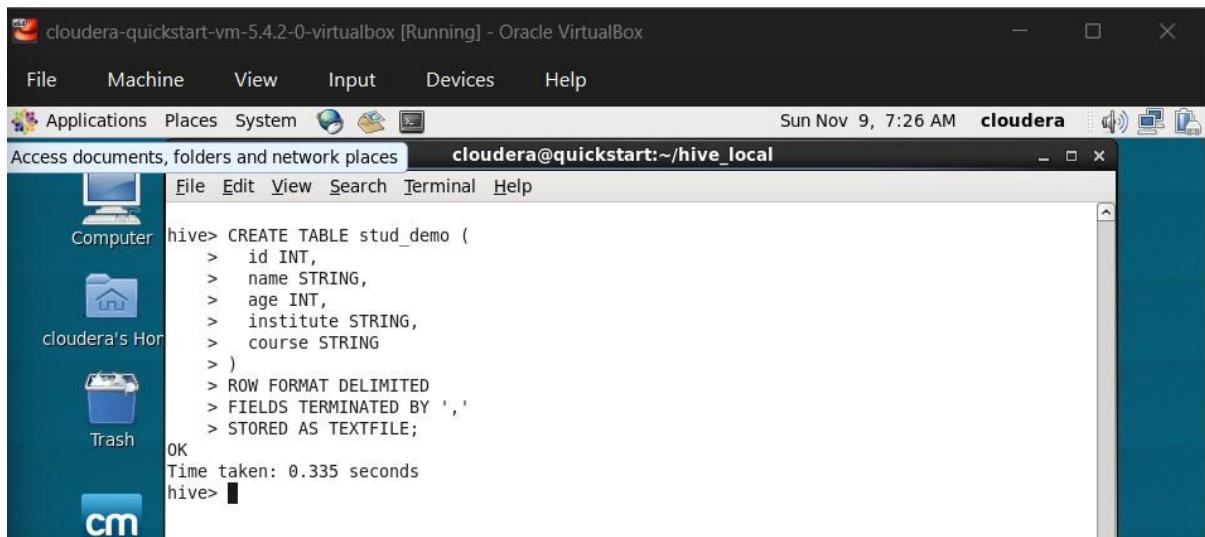
```

B. Dynamic Partitioning

To automatically create Hive table partitions while inserting data — instead of specifying each partition value manually.

a. Create the Staging Table

We'll first create a *non-partitioned* staging table (stud_demo) to hold all data (both Java and Hadoop students).



```

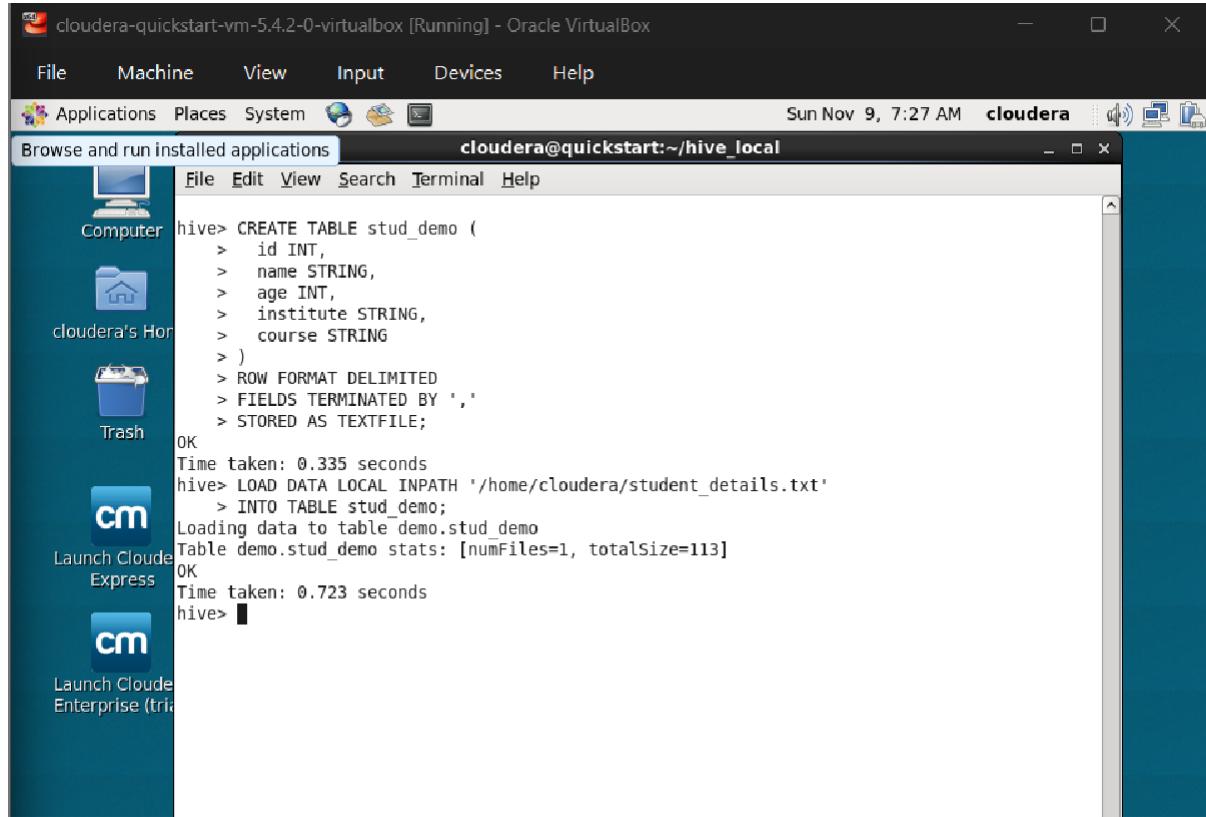
File Edit View Search Terminal Help
Access documents, folders and network places cloudera@quickstart:~/hive_local
hive> CREATE TABLE stud_demo (
  >   id INT,
  >   name STRING,
  >   age INT,
  >   institute STRING,
  >   course STRING
  > )
  > ROW FORMAT DELIMITED
  > FIELDS TERMINATED BY ','
  > STORED AS TEXTFILE;
OK
Time taken: 0.335 seconds
hive>

```

b. Prepare Combined Input File

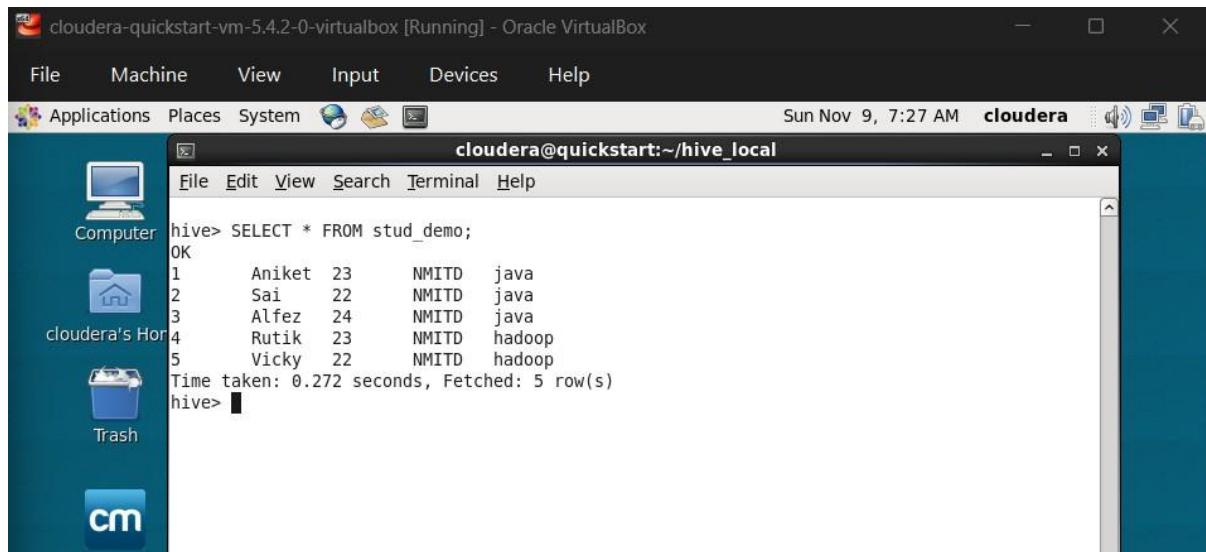
In your **Linux terminal**, create one combined file that includes both course values:

C. Load Data into the Staging Table



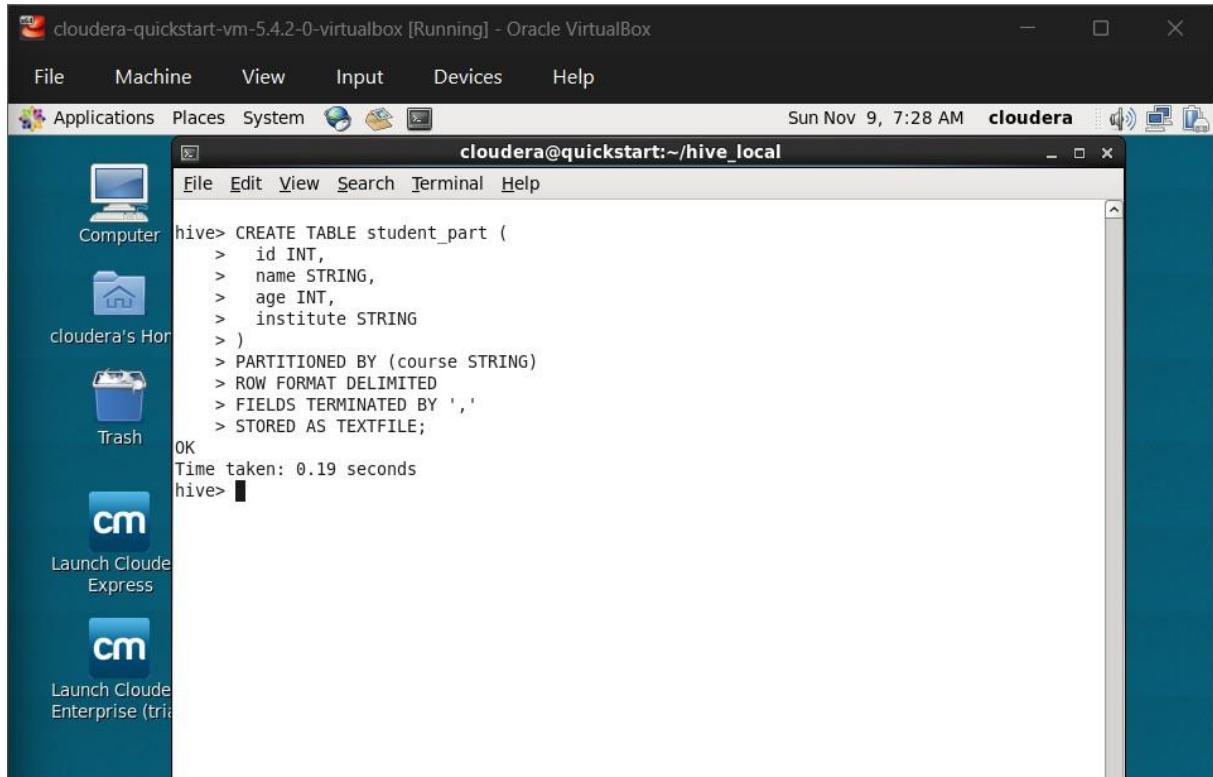
```
hive> CREATE TABLE stud_demo (
    >   id INT,
    >   name STRING,
    >   age INT,
    >   institute STRING,
    >   course STRING
    > )
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY ','
    > STORED AS TEXTFILE;
OK
Time taken: 0.335 seconds
hive> LOAD DATA LOCAL INPATH '/home/cloudera/student_details.txt'
    > INTO TABLE stud_demo;
Loading data to table demo.stud_demo
Table demo.stud_demo stats: [numFiles=1, totalSize=113]
OK
Time taken: 0.723 seconds
hive>
```

D. Show full stud_demo table :



```
hive> SELECT * FROM stud_demo;
OK
1      Aniket  23      NMITD  java
2      Sai     22      NMITD  java
3      Alfez   24      NMITD  java
4      Rutik   23      NMITD  hadoop
5      Vicky   22      NMITD  hadoop
Time taken: 0.272 seconds, Fetched: 5 row(s)
hive>
```

E. Create the Partitioned Target Table



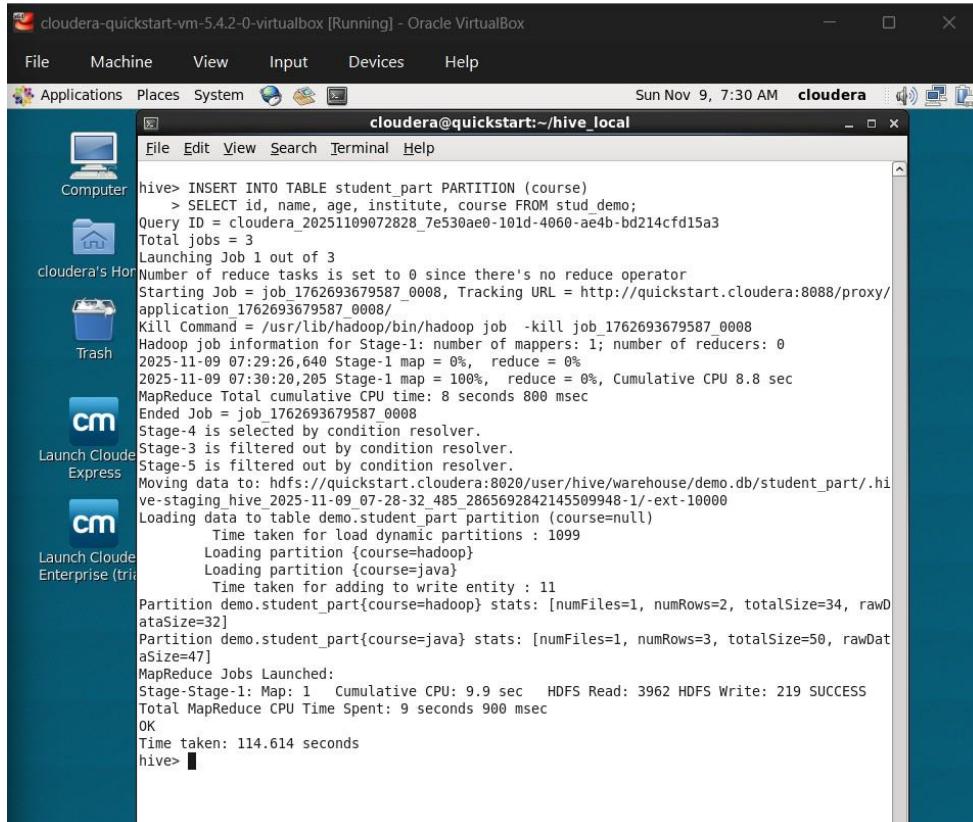
```

cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> CREATE TABLE student_part (
    >   id INT,
    >   name STRING,
    >   age INT,
    >   institute STRING
    > )
    > PARTITIONED BY (course STRING)
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY ','
    > STORED AS TEXTFILE;
OK
Time taken: 0.19 seconds
hive> 

```

F. Insert Data Dynamically

Now execute the dynamic insert statement.



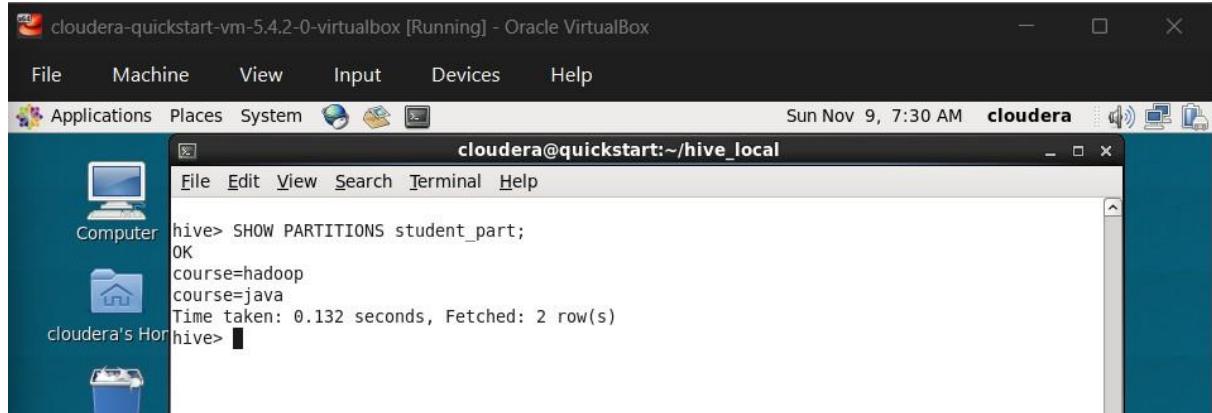
```

cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> INSERT INTO TABLE student_part PARTITION (course)
    > SELECT id, name, age, institute, course FROM stud_demo;
Query ID = cloudera_20251109072828_7e530ae0-101d-4060-ae4b-bd214cf15a3
Total jobs = 3
Launching Job 1 out of 3
Number of reduce tasks is set to 0 since there's no reduce operator
Starting Job = job_1762693679587_0008, Tracking URL = http://quickstart.cloudera:8088/proxy/
application_1762693679587_0008/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1762693679587_0008
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 0
2025-11-09 07:29:26,640 Stage-1 map = 0%, reduce = 0%
2025-11-09 07:30:20,205 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 8.8 sec
MapReduce Total cumulative CPU time: 8 seconds 800 msec
Ended Job = job_1762693679587_0008
Stage-3 is selected by condition resolver.
Stage-3 is filtered out by condition resolver.
Stage-5 is filtered out by condition resolver.
Moving data to: hdfs://quickstart.cloudera:8020/user/hive/warehouse/demo.db/student_part.hive-staging
hive> Moving data to: hdfs://quickstart.cloudera:8020/user/hive/warehouse/demo.db/student_part.hive-staging
hive> Moving data to: hdfs://quickstart.cloudera:8020/user/hive/warehouse/demo.db/student_part.hive-staging
hive> Loading data to table demo.student_part partition (course=null)
    Time taken for load dynamic partitions : 1099
    Loading partition {course=hadoop}
    Loading partition {course=java}
    Time taken for adding to write entity : 11
Partition demo.student_part{course=hadoop} stats: [numFiles=1, numRows=2, totalSize=34, rawDataSize=32]
Partition demo.student_part{course=java} stats: [numFiles=1, numRows=3, totalSize=50, rawDataSize=47]
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Cumulative CPU: 9.9 sec  HDFS Read: 3962 HDFS Write: 219 SUCCESS
Total MapReduce CPU Time Spent: 9 seconds 900 msec
OK
Time taken: 114.614 seconds
hive> 

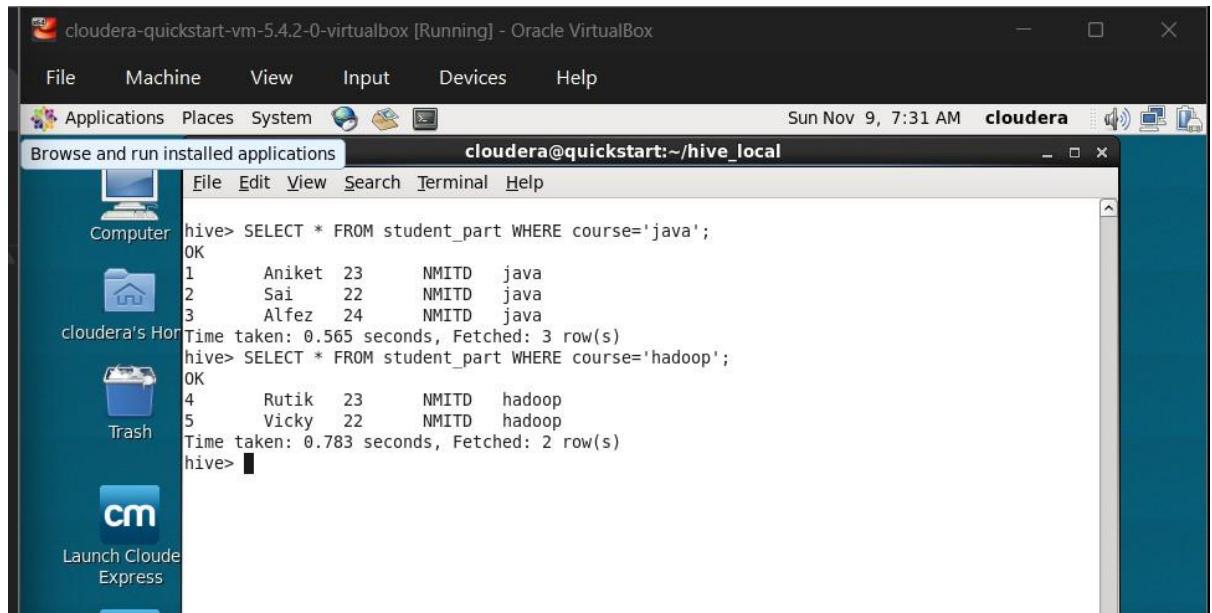
```

G. Verify Dynamic Partitions

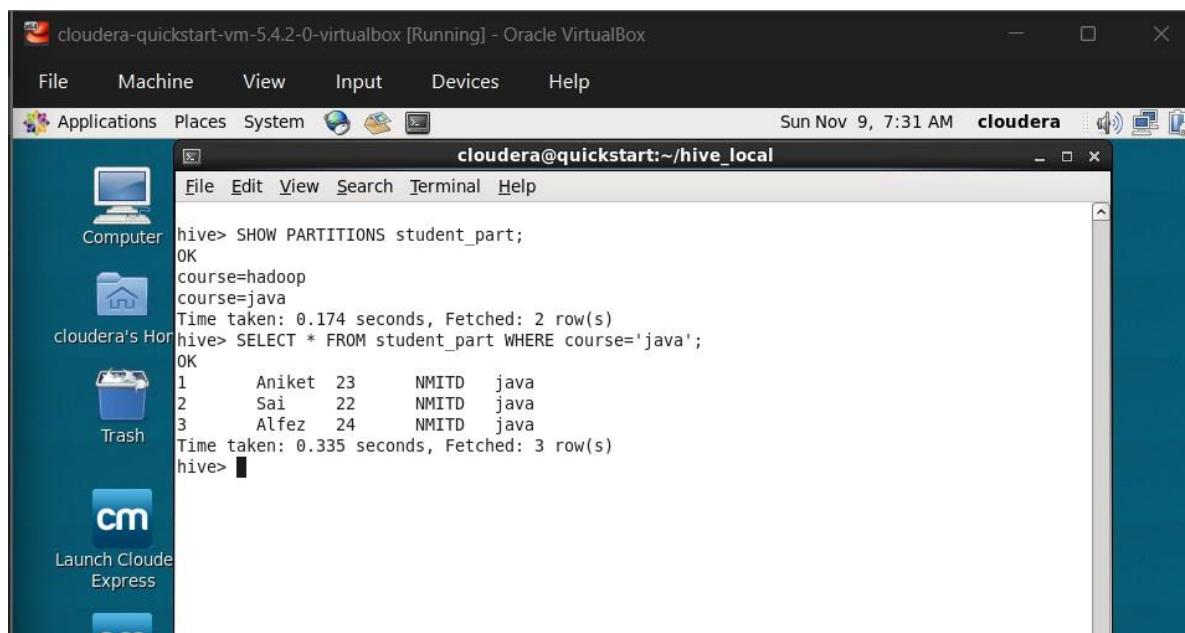
List partitions.



```
cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> SHOW PARTITIONS student_part;
OK
course=hadoop
course=java
Time taken: 0.132 seconds, Fetched: 2 row(s)
hive> ■
```



```
cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> SELECT * FROM student_part WHERE course='java';
OK
1      Aniket  23      NMITD   java
2      Sai     22      NMITD   java
3      Alfez   24      NMITD   java
Time taken: 0.565 seconds, Fetched: 3 row(s)
hive> SELECT * FROM student_part WHERE course='hadoop';
OK
4      Rutik   23      NMITD   hadoop
5      Vicky   22      NMITD   hadoop
Time taken: 0.783 seconds, Fetched: 2 row(s)
hive> ■
```



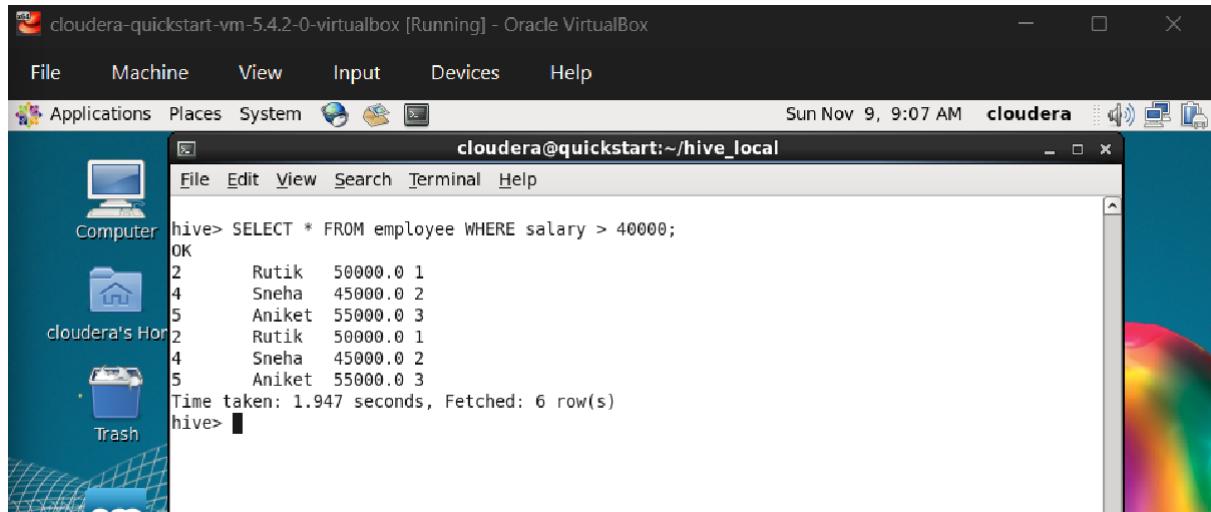
```
cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> SHOW PARTITIONS student_part;
OK
course=hadoop
course=java
Time taken: 0.174 seconds, Fetched: 2 row(s)
hive> SELECT * FROM student_part WHERE course='java';
OK
1      Aniket  23      NMITD   java
2      Sai     22      NMITD   java
3      Alfez   24      NMITD   java
Time taken: 0.335 seconds, Fetched: 3 row(s)
hive> ■
```

4. Hive Built-In Operators

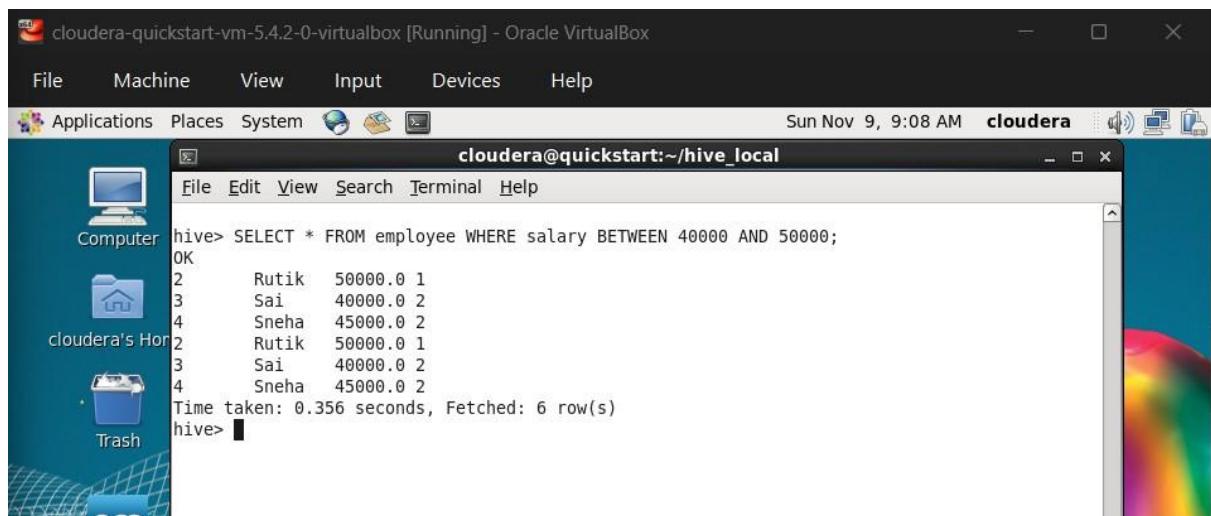
To demonstrate the use of **Hive's built-in operators** for comparison, arithmetic, and logical operations on table data.

A. Comparison Operators

We'll use your existing 'employee' table for this.

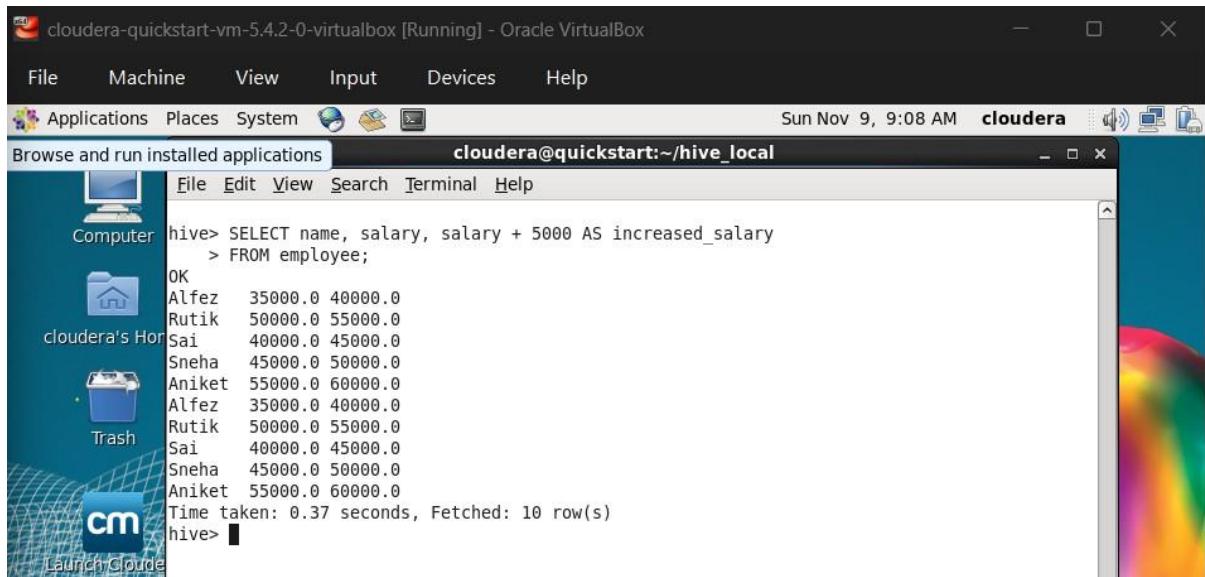


```
File Edit View Search Terminal Help
hive> SELECT * FROM employee WHERE salary > 40000;
OK
2      Rutik  50000.0 1
4      Sneha  45000.0 2
5      Aniket 55000.0 3
2      Rutik  50000.0 1
4      Sneha  45000.0 2
5      Aniket 55000.0 3
Time taken: 1.947 seconds, Fetched: 6 row(s)
hive>
```



```
File Edit View Search Terminal Help
hive> SELECT * FROM employee WHERE salary BETWEEN 40000 AND 50000;
OK
2      Rutik  50000.0 1
3      Sai    40000.0 2
4      Sneha  45000.0 2
2      Rutik  50000.0 1
3      Sai    40000.0 2
4      Sneha  45000.0 2
Time taken: 0.356 seconds, Fetched: 6 row(s)
hive>
```

B. Arithmetic Operators

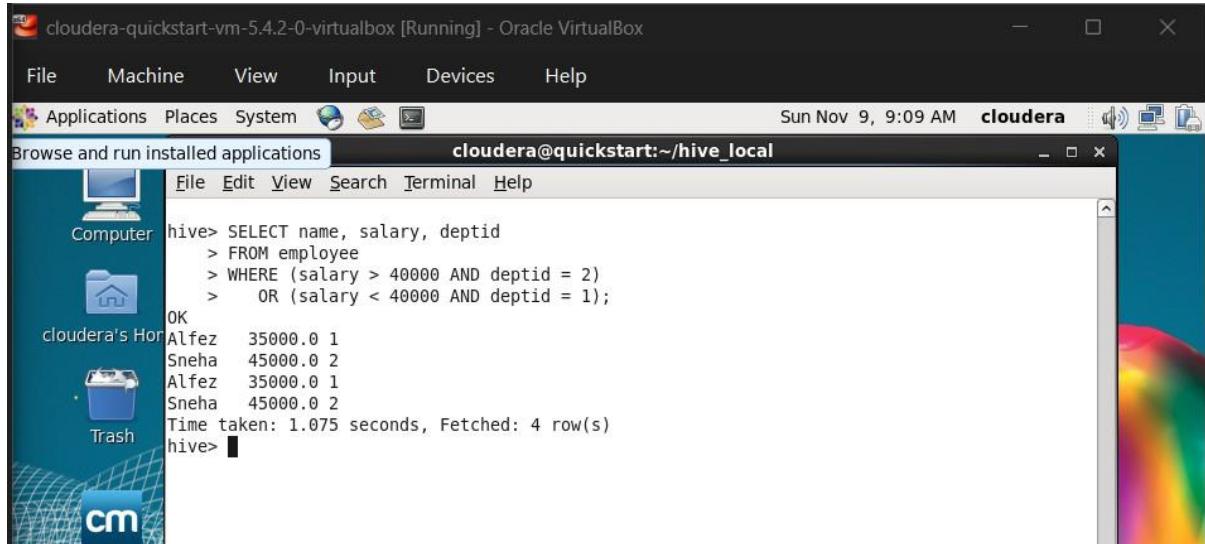


```

cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> SELECT name, salary, salary + 5000 AS increased_salary
      > FROM employee;
OK
Alfez 35000.0 40000.0
Rutik 50000.0 55000.0
Sai 40000.0 45000.0
Sneha 45000.0 50000.0
Aniket 55000.0 60000.0
Alfez 35000.0 40000.0
Rutik 50000.0 55000.0
Sai 40000.0 45000.0
Sneha 45000.0 50000.0
Aniket 55000.0 60000.0
Time taken: 0.37 seconds, Fetched: 10 row(s)
hive>

```

C. Logical Operators



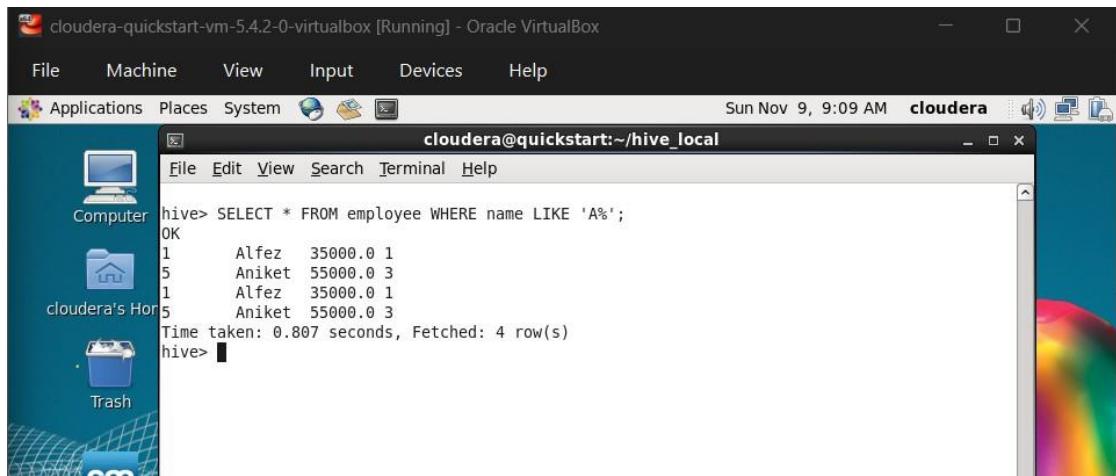
```

cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> SELECT name, salary, deptid
      > FROM employee
      > WHERE (salary > 40000 AND deptid = 2)
      >       OR (salary < 40000 AND deptid = 1);
OK
Alfez 35000.0 1
Sneha 45000.0 2
Alfez 35000.0 1
Sneha 45000.0 2
Time taken: 1.075 seconds, Fetched: 4 row(s)
hive>

```

D. Special Operator Examples

(a) LIKE — pattern matching

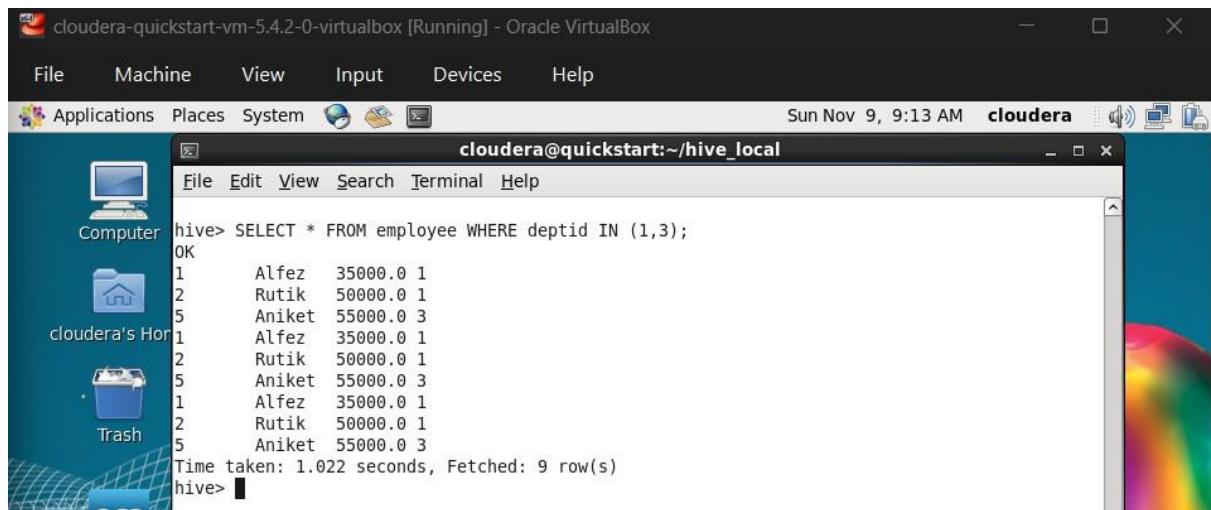


```

cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> SELECT * FROM employee WHERE name LIKE 'A%';
OK
1 Alfez 35000.0 1
5 Aniket 55000.0 3
1 Alfez 35000.0 1
5 Aniket 55000.0 3
Time taken: 0.807 seconds, Fetched: 4 row(s)
hive>

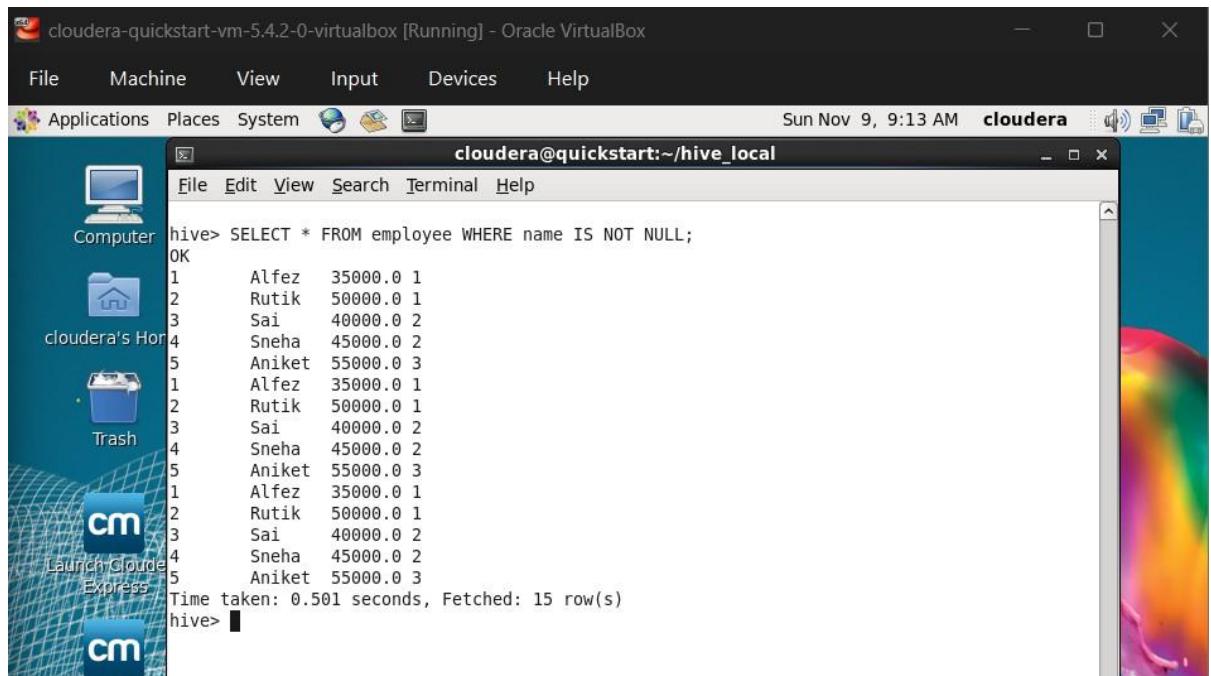
```

(b) IN / NOT IN



```
hive> SELECT * FROM employee WHERE deptid IN (1,3);
OK
1      Alfez  35000.0 1
2      Rutik  50000.0 1
5      Aniket 55000.0 3
1      Alfez  35000.0 1
2      Rutik  50000.0 1
5      Aniket 55000.0 3
1      Alfez  35000.0 1
2      Rutik  50000.0 1
5      Aniket 55000.0 3
Time taken: 1.022 seconds, Fetched: 9 row(s)
hive>
```

(c) IS NULL / IS NOT NULL



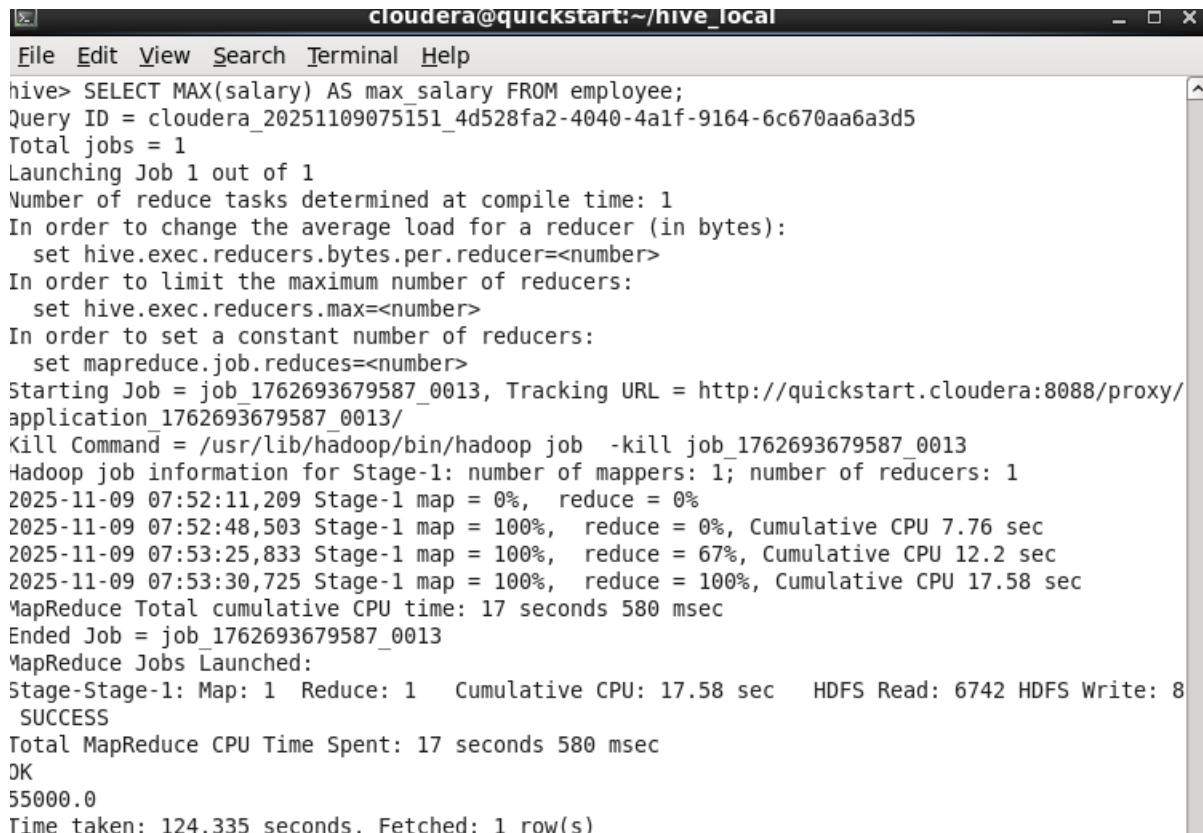
```
hive> SELECT * FROM employee WHERE name IS NOT NULL;
OK
1      Alfez  35000.0 1
2      Rutik  50000.0 1
3      Sai    40000.0 2
4      Sneha  45000.0 2
5      Aniket 55000.0 3
1      Alfez  35000.0 1
2      Rutik  50000.0 1
3      Sai    40000.0 2
4      Sneha  45000.0 2
5      Aniket 55000.0 3
1      Alfez  35000.0 1
2      Rutik  50000.0 1
3      Sai    40000.0 2
4      Sneha  45000.0 2
5      Aniket 55000.0 3
Time taken: 0.501 seconds, Fetched: 15 row(s)
hive>
```

5. Hive Built-In Functions

To demonstrate the use of Hive built-in functions — both aggregate functions (like MAX, MIN, AVG) and string manipulation functions (like UPPER, LOWER, LENGTH, SUBSTR).

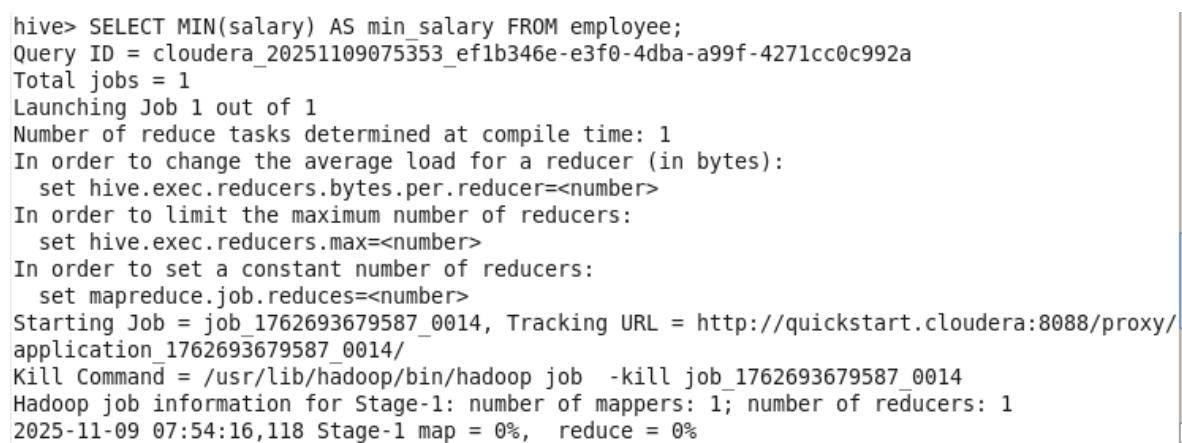
A. Aggregate Functions

```
SELECT MAX(salary) AS max_salary FROM employee;
```



```
cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> SELECT MAX(salary) AS max_salary FROM employee;
Query ID = cloudera_20251109075151_4d528fa2-4040-4a1f-9164-6c670aa6a3d5
Total jobs = 1
Launching Job 1 out of 1
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>
Starting Job = job_1762693679587_0013, Tracking URL = http://quickstart.cloudera:8088/proxy/
application_1762693679587_0013/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1762693679587_0013
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2025-11-09 07:52:11,209 Stage-1 map = 0%, reduce = 0%
2025-11-09 07:52:48,503 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 7.76 sec
2025-11-09 07:53:25,833 Stage-1 map = 100%, reduce = 67%, Cumulative CPU 12.2 sec
2025-11-09 07:53:30,725 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 17.58 sec
MapReduce Total cumulative CPU time: 17 seconds 580 msec
Ended Job = job_1762693679587_0013
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 17.58 sec HDFS Read: 6742 HDFS Write: 8
  SUCCESS
Total MapReduce CPU Time Spent: 17 seconds 580 msec
OK
55000.0
Time taken: 124.335 seconds, Fetched: 1 row(s)
```

B. SELECT MIN(salary) AS min_salary FROM employee;



```
hive> SELECT MIN(salary) AS min_salary FROM employee;
Query ID = cloudera_20251109075353_ef1b346e-e3f0-4dba-a99f-4271cc0c992a
Total jobs = 1
Launching Job 1 out of 1
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>
Starting Job = job_1762693679587_0014, Tracking URL = http://quickstart.cloudera:8088/proxy/
application_1762693679587_0014/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1762693679587_0014
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2025-11-09 07:54:16,118 Stage-1 map = 0%, reduce = 0%
```

C. SELECT AVG(salary) AS avg_salary FROM employee;

```
mapReduce Jobs Launched:  
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 18.62 sec HDFS Read: 6749 HDFS Write: 8  
SUCCESS  
Total MapReduce CPU Time Spent: 18 seconds 620 msec  
OK  
35000.0  
Time taken: 116.092 seconds, Fetched: 1 row(s)  
hive> SELECT AVG(salary) AS avg_salary FROM employee;  
Query ID = cloudera_20251109075555_c907aad7-f1b2-408a-9429-efbd991e2eb4  
Total jobs = 1  
Launching Job 1 out of 1  
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>  
Starting Job = job_1762693679587_0015, Tracking URL = http://quickstart.cloudera:8088/proxy/  
application_1762693679587_0015/  
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1762693679587_0015  
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1  
2025-11-09 07:56:23,417 Stage-1 map = 0%, reduce = 0%  
2025-11-09 07:57:01,849 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 7.94 sec  
2025-11-09 07:57:33,105 Stage-1 map = 100%, reduce = 67%, Cumulative CPU 13.68 sec  
2025-11-09 07:57:37,676 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 18.06 sec  
MapReduce Total cumulative CPU time: 18 seconds 60 msec  
Ended Job = job_1762693679587_0015  
MapReduce Jobs Launched:  
Stage-Stage-1: Map: 1 Reduce: 1 Cumulative CPU: 18.06 sec HDFS Read: 7116 HDFS Write: 8  
SUCCESS  
Total MapReduce CPU Time Spent: 18 seconds 60 msec  
OK  
45000.0  
Time taken: 128.181 seconds, Fetched: 1 row(s)  
hive> S
```

B. String Functions

Hive provides many string manipulation utilities.

Let's use your existing employee table to test them:

```

cloudera@quickstart:~/hive_local
hive> SELECT name, UPPER(name) AS upper_name FROM employee;
OK
Alfez    ALFEZ
Rutik    RUTIK
Sai      SAI
Sneha   SNEHA
Aniket  ANIKET
Time taken: 0.325 seconds, Fetched: 5 row(s)
hive> SELECT name, LOWER(name) AS lower_name FROM employee;
OK
Alfez    alfez
Rutik    rutik
Sai      sai
Sneha   sneha
Aniket  aniket
Time taken: 0.299 seconds, Fetched: 5 row(s)
hive> SELECT name, LENGTH(name) AS name_length FROM employee;
OK
Alfez    5
Rutik    5
Sai      3
Sneha   5
Aniket  6
Time taken: 0.252 seconds, Fetched: 5 row(s)
hive> SELECT name, SUBSTR(name, 1, 3) AS first_three_letters FROM employee;
OK
Alfez    Alf
Rutik    Rut
Sai      Sai
Sneha   Sne
Aniket  Ani
Time taken: 0.288 seconds, Fetched: 5 row(s)
hive> ■

```

C. Combining Aggregate & Group

```

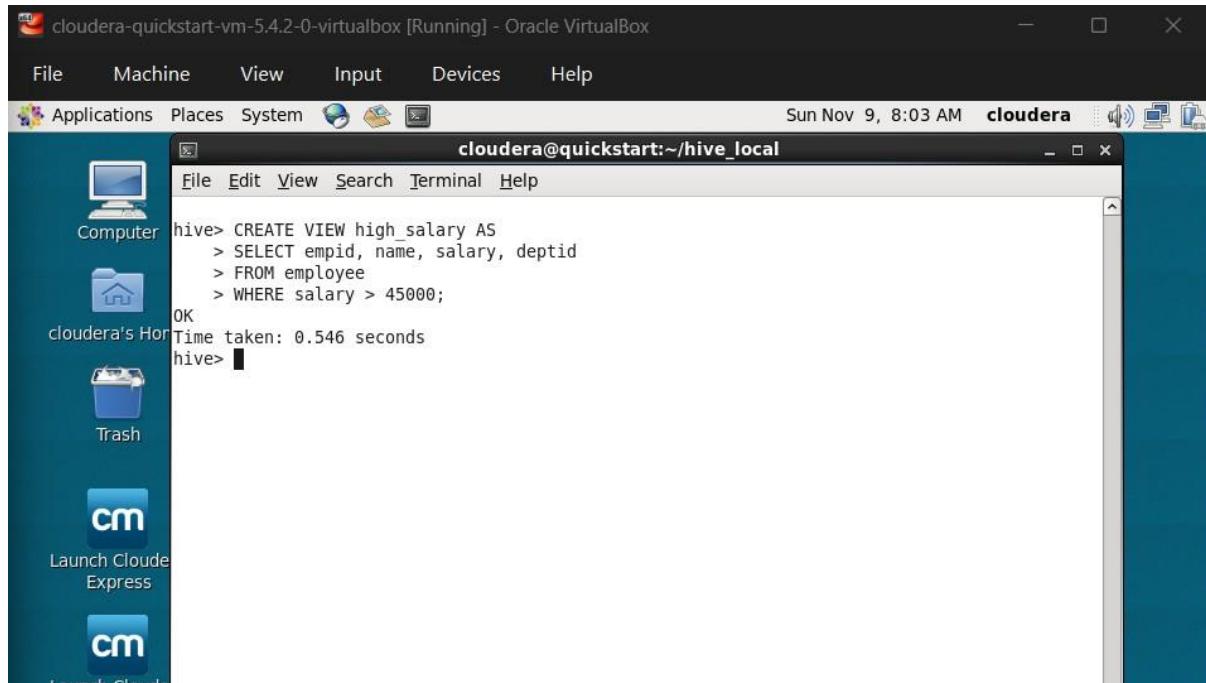
cloudera@quickstart:~/hive_local
hive> SELECT deptid, COUNT(*) AS emp_count, AVG(salary) AS avg_salary
   >   FROM employee
   >   GROUP BY deptid;
Query ID = cloudera_20251109075959_dac80793-eaff-4306-928f-1f5062baf668
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 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>
Starting Job = job_1762693679587_0016, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1762693679587_0016/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1762693679587_0016
Hadoop job information for Stage-1: number of mappers: 1; number of reducers: 1
2025-11-09 08:00:11,423 Stage-1 map = 0%,  reduce = 0%
2025-11-09 08:01:00,491 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 9.42 sec
2025-11-09 08:01:40,834 Stage-1 map = 100%,  reduce = 67%, Cumulative CPU 16.18 sec
2025-11-09 08:01:45,707 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 16.18 sec
MapReduce Total cumulative CPU time: 16 seconds 180 msec
Ended Job = job_1762693679587_0016
MapReduce Jobs Launched:
Stage-Stage-1: Map: 1  Reduce: 1  Cumulative CPU: 21.95 sec  HDFS Read: 7874 HDFS Write: 3
0 SUCCESS
Total MapReduce CPU Time Spent: 21 seconds 950 msec
OK
1      2        42500.0
2      2        42500.0
3      1        55000.0
Time taken: 152.719 seconds, Fetched: 3 row(s)
hive> ■

```

6. Hive Views

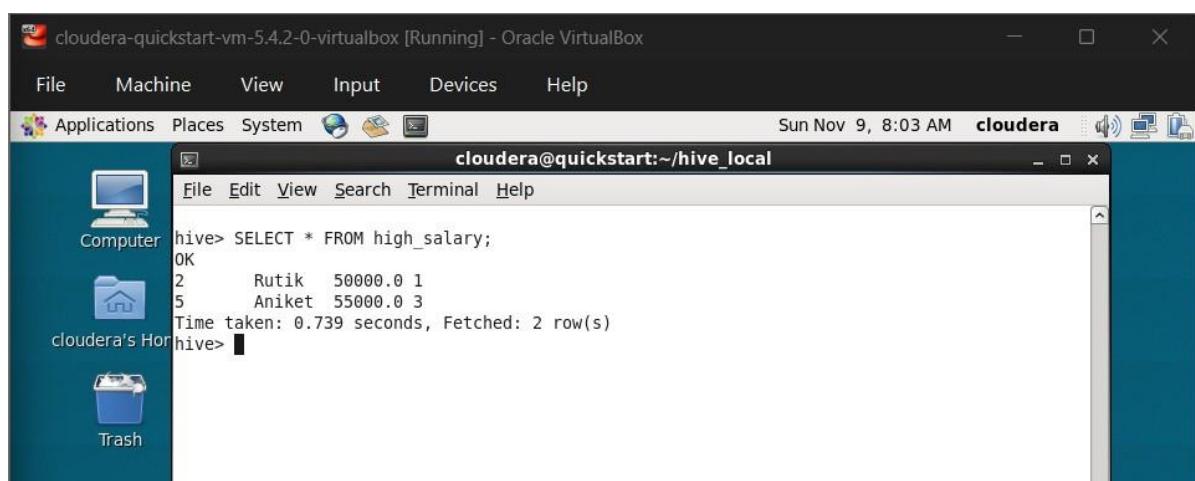
Create and query views (virtual tables) & Understand indexes in Hive (used for query optimization, optional based on VM support).

A. Create a View



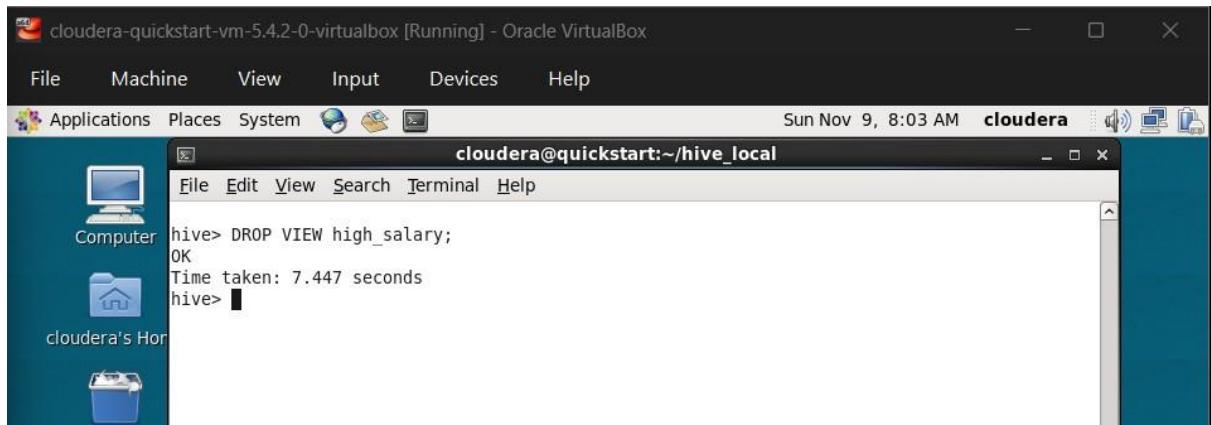
```
hive> CREATE VIEW high_salary AS
>   SELECT empid, name, salary, deptid
>   FROM employee
>   WHERE salary > 45000;
OK
Time taken: 0.546 seconds
hive>
```

B. Query the view



```
hive> SELECT * FROM high_salary;
OK
2      Rutik    50000.0 1
5      Aniket   55000.0 3
Time taken: 0.739 seconds, Fetched: 2 row(s)
hive>
```

C. Drop the View

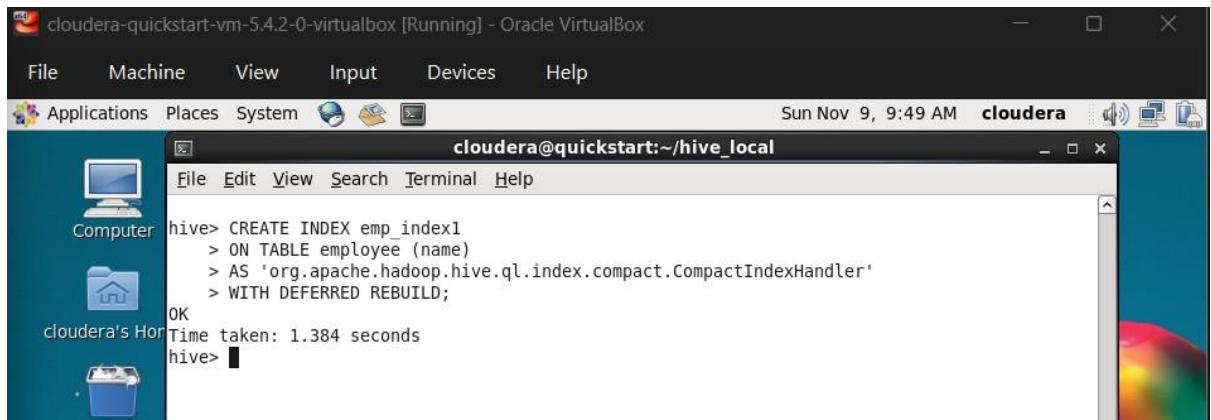


```
File Edit View Search Terminal Help
hive> DROP VIEW high_salary;
OK
Time taken: 7.447 seconds
hive>
```

D. INDEXES

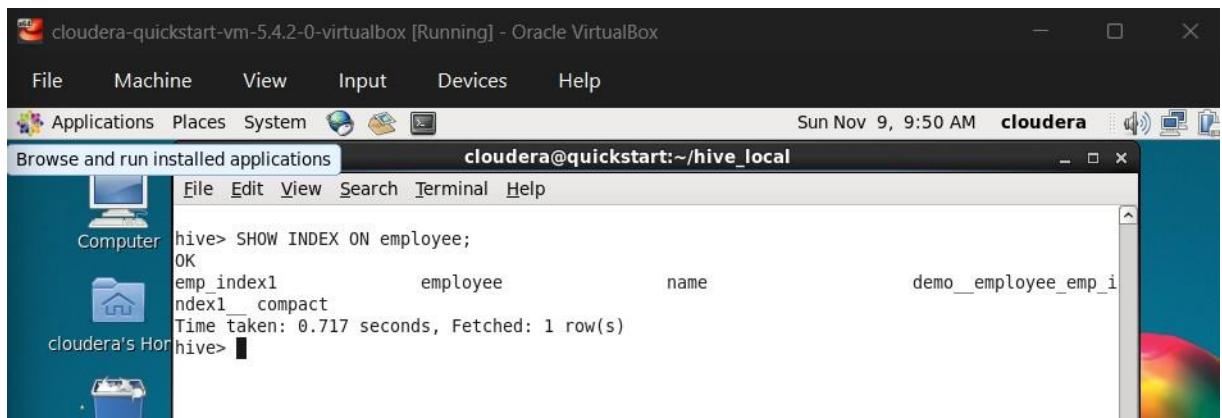
Create an Index - This command creates a metadata index on the name column — intended to speed up queries that filter by name.

The WITH DEFERRED REBUILD part tells Hive to create the index structure later.



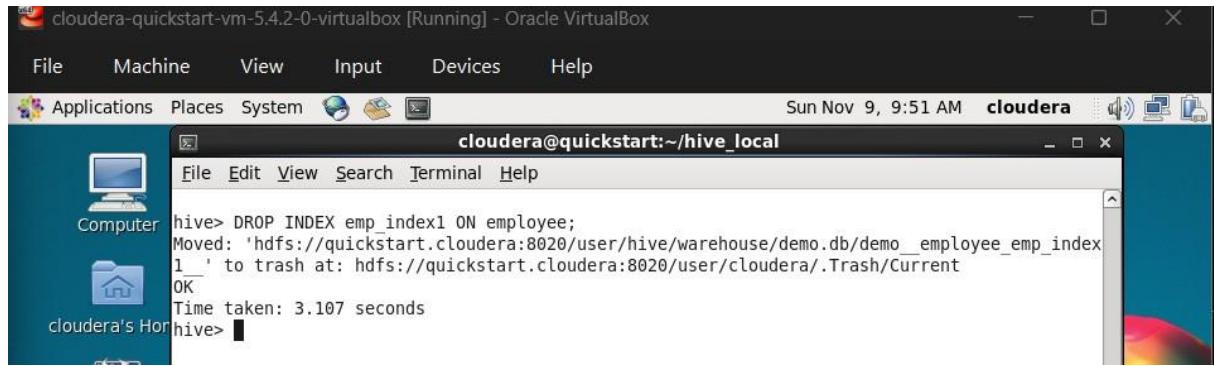
```
File Edit View Search Terminal Help
hive> CREATE INDEX emp_index1
> ON TABLE employee (name)
> AS 'org.apache.hadoop.hive.ql.index.compact.CompactIndexHandler'
> WITH DEFERRED REBUILD;
OK
Time taken: 1.384 seconds
hive>
```

Show Existing Indexes



```
File Edit View Search Terminal Help
Browse and run installed applications
cloudera@quickstart:~/hive_local
hive> SHOW INDEX ON employee;
OK
emp_index1          employee           name
ndex1_compact
Time taken: 0.717 seconds, Fetched: 1 row(s)
hive>
```

Drop the Index

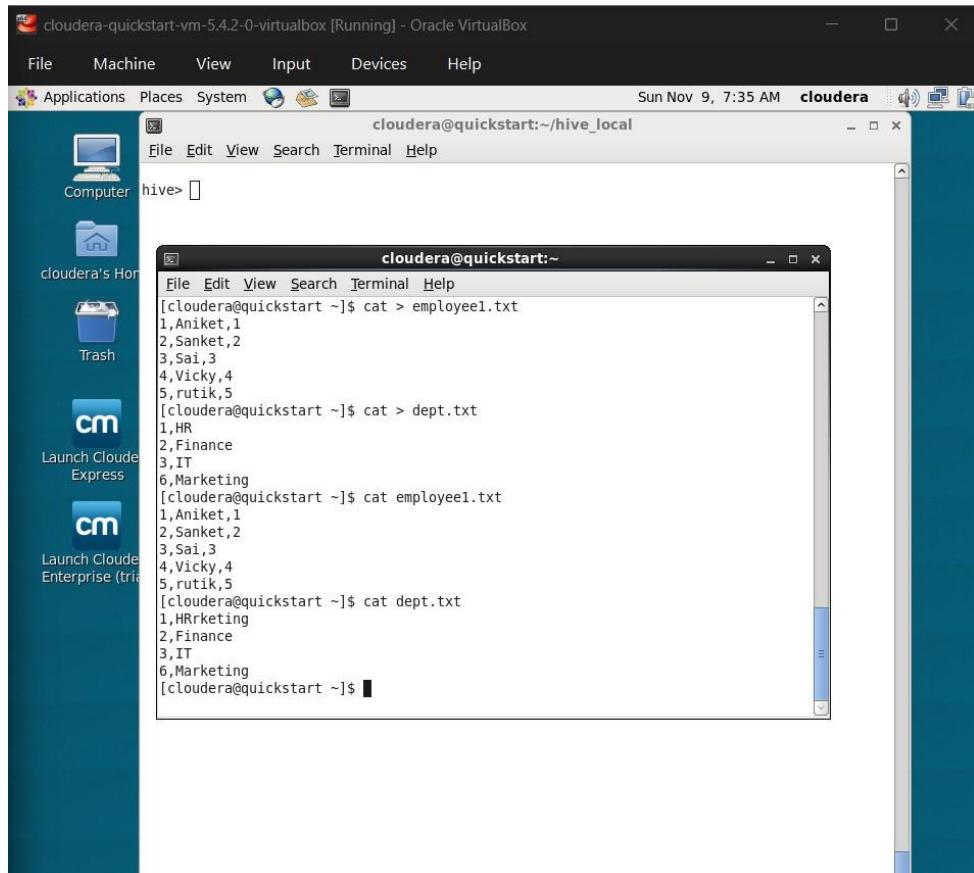


```
cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> DROP INDEX emp_index1 ON employee;
Moved: 'hdfs://quickstart.cloudera:8020/user/hive/warehouse/demo.db/demo_employee_emp_index1' to trash at: hdfs://quickstart.cloudera:8020/user/cloudera/.Trash/Current
OK
Time taken: 3.107 seconds
hive>
```

7. HiveQL: Select Where, Select OrderBy, Select GroupBy, Select Joins

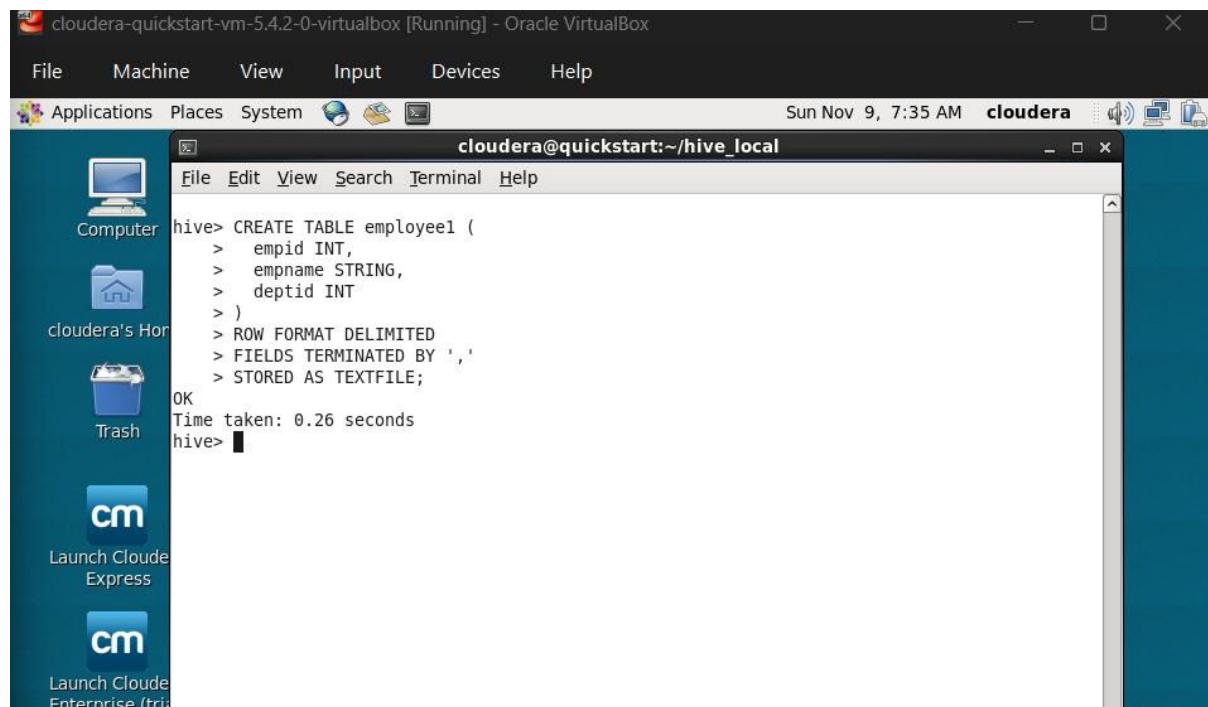
A. Create Sample Input Files

In your **Linux terminal**, create two CSV files:

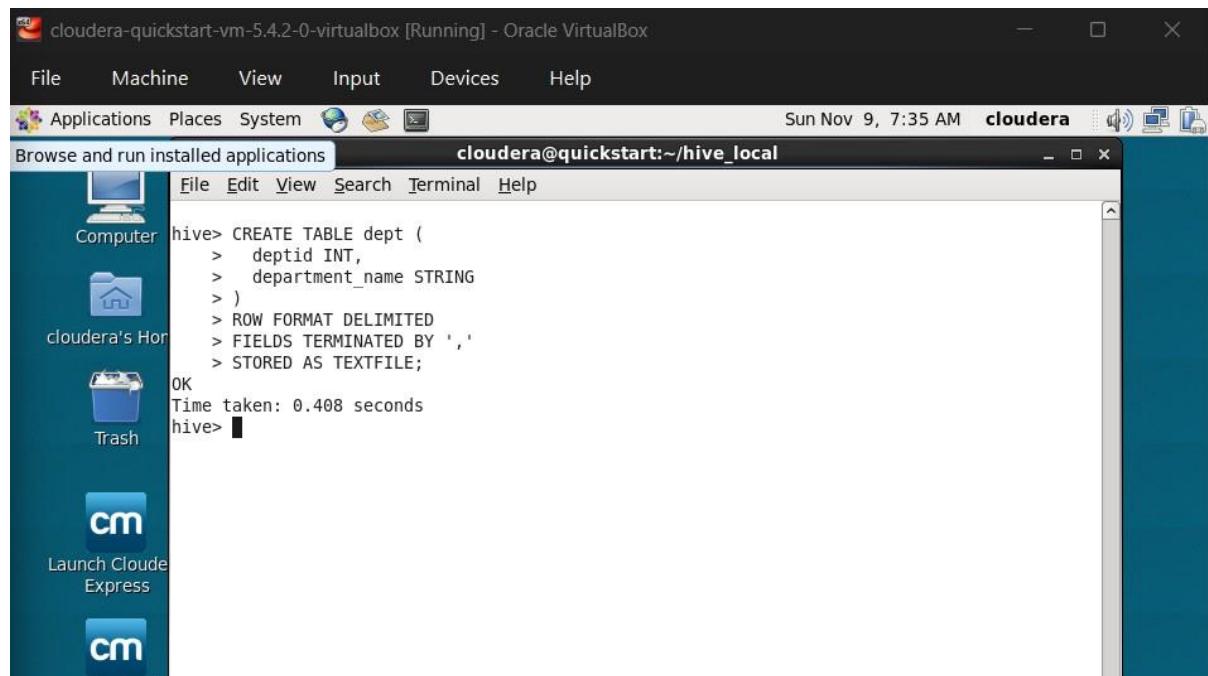


```
cloudera@quickstart:~/hive_local
File Edit View Search Terminal Help
hive> [cloudera@quickstart ~]$ cat > employee1.txt
1,Aniket,1
2,Sanket,2
3,Sai,3
4,Vicky,4
5,rutik,5
[cloudera@quickstart ~]$ cat > dept.txt
1,HR
2,Finance
3,IT
6,Marketing
[cloudera@quickstart ~]$ cat employee1.txt
1,Aniket,1
2,Sanket,2
3,Sai,3
4,Vicky,4
5,rutik,5
[cloudera@quickstart ~]$ cat dept.txt
1,HR
2,Finance
3,IT
6,Marketing
[cloudera@quickstart ~]$
```

B. Create the Tables in Hive

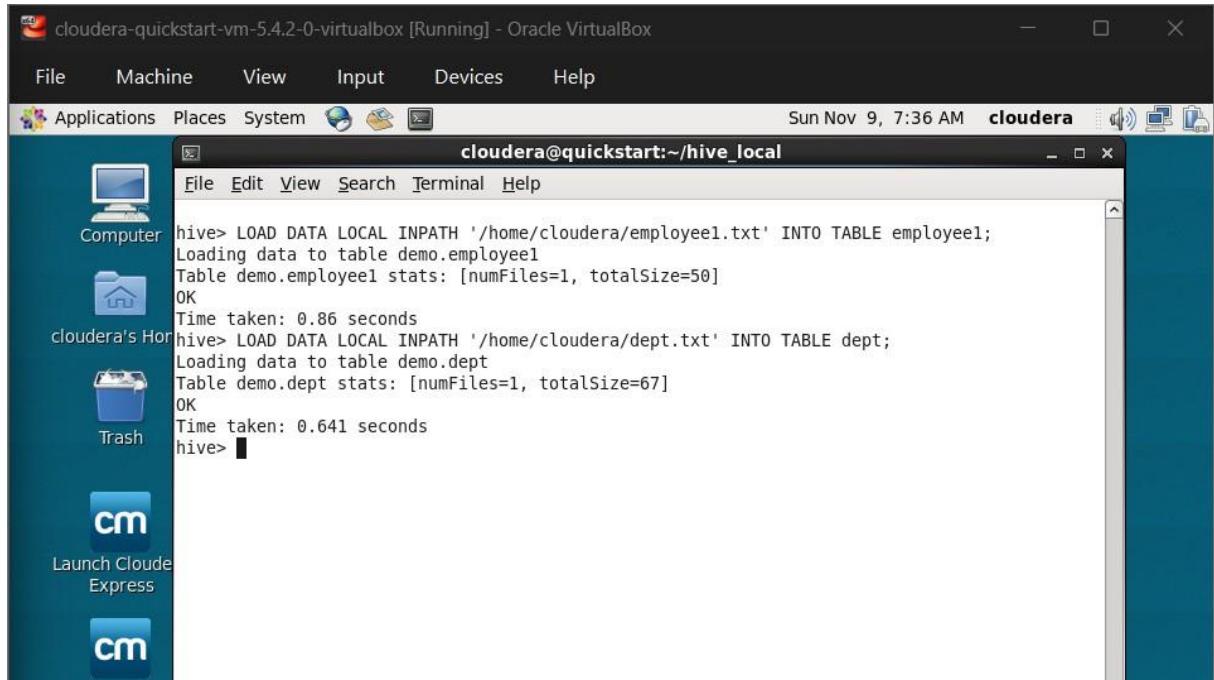


```
hive> CREATE TABLE employee1 (
    >     empid INT,
    >     empname STRING,
    >     deptid INT
    > )
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY ','
    > STORED AS TEXTFILE;
OK
Time taken: 0.26 seconds
hive>
```



```
hive> CREATE TABLE dept (
    >     deptid INT,
    >     department_name STRING
    > )
    > ROW FORMAT DELIMITED
    > FIELDS TERMINATED BY ','
    > STORED AS TEXTFILE;
OK
Time taken: 0.408 seconds
hive>
```

C. Load Data into the Tables



The screenshot shows a Linux desktop environment with a terminal window open. The terminal window title is "cloudera@quickstart:~/hive_local". The terminal content shows two "LOAD DATA LOCAL INPATH" commands being run:

```
hive> LOAD DATA LOCAL INPATH '/home/cloudera/employee1.txt' INTO TABLE employee1;
Loading data to table demo.employee1
Table demo.employee1 stats: [numFiles=1, totalSize=50]
OK
Time taken: 0.86 seconds
hive> LOAD DATA LOCAL INPATH '/home/cloudera/dept.txt' INTO TABLE dept;
Loading data to table demo.dept
Table demo.dept stats: [numFiles=1, totalSize=67]
OK
Time taken: 0.641 seconds
hive>
```

D. Disable Auto Join Optimization

This prevents Hive from converting your join into a map-side broadcast join, which can skip reducers for small tables.

Disabling it ensures you see a MapReduce job plan.

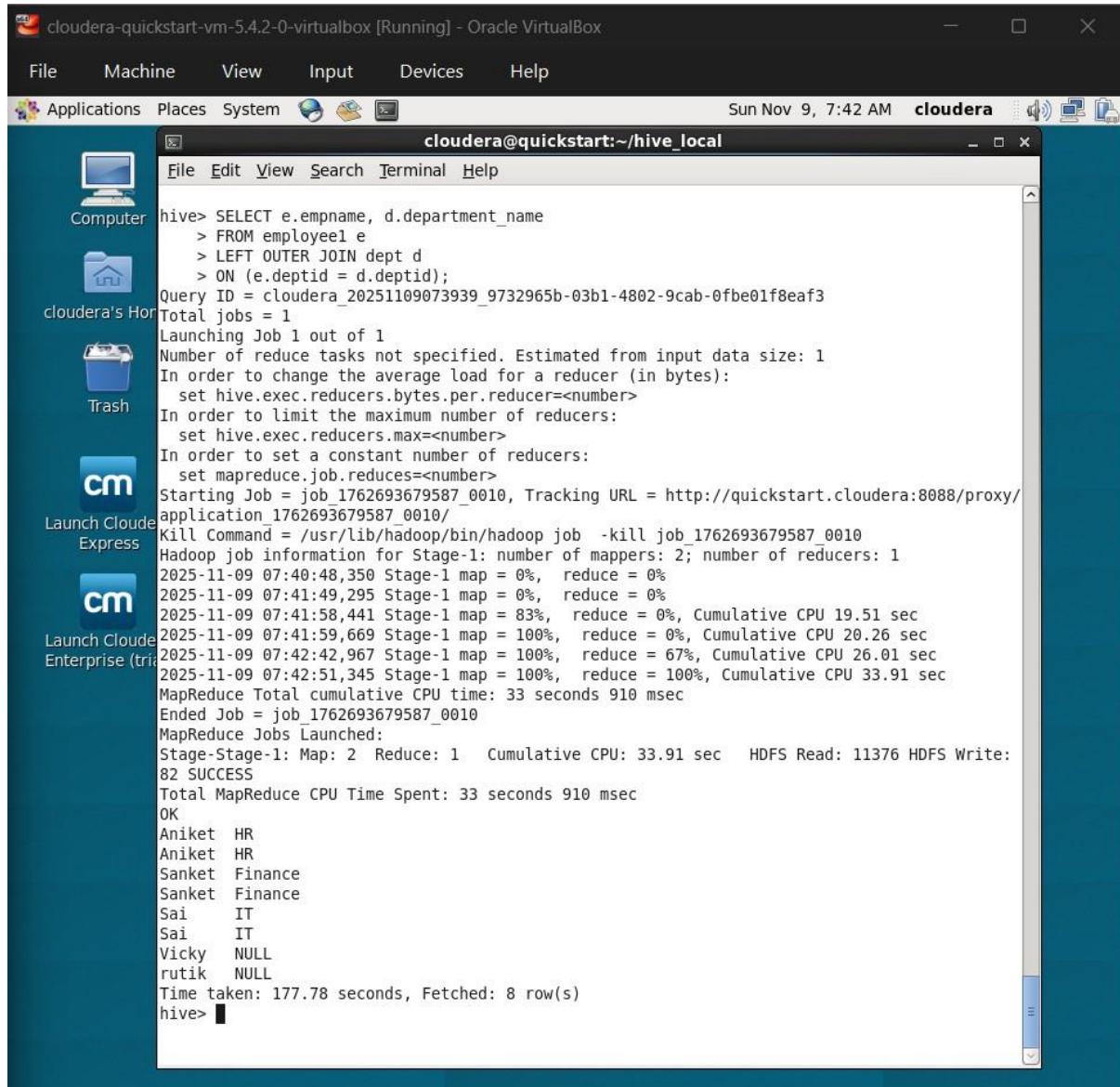
hive.auto.convert.join=false

E. INNER JOIN

The screenshot shows a Linux desktop environment with a terminal window open. The terminal window title is "cloudera@quickstart:~/hive_local". The terminal content displays the following Hive query and its execution results:

```
hive> SELECT e.empname, d.department_name
   > FROM employeee e
   > JOIN dept d
   > ON (e.deptid = d.deptid);
Query ID = cloudera_20251109073636_9cdb24c0-61c7-446e-9f3f-67ab8d39df34
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 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>
Starting Job = job_1762693679587_0009, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1762693679587_0009/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1762693679587_0009
Hadoop job information for Stage-1: number of mappers: 2; number of reducers: 1
2025-11-09 07:37:35,472 Stage-1 map = 0%,  reduce = 0%
2025-11-09 07:38:36,158 Stage-1 map = 0%,  reduce = 0%
2025-11-09 07:38:45,546 Stage-1 map = 50%,  reduce = 0%, Cumulative CPU 19.83 sec
2025-11-09 07:38:48,997 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 29.51 sec
2025-11-09 07:39:28,035 Stage-1 map = 100%,  reduce = 67%, Cumulative CPU 34.87 sec
2025-11-09 07:39:36,208 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 42.89 sec
MapReduce Total cumulative CPU time: 42 seconds 890 msec
Ended Job = job_1762693679587_0009
MapReduce Jobs Launched:
Stage-Stage-1: Map: 2 Reduce: 1 Cumulative CPU: 42.89 sec HDFS Read: 12486 HDFS Write: 64 SUCCESS
Total MapReduce CPU Time Spent: 42 seconds 890 msec
OK
Aniket HR
Aniket HR
Sanket Finance
Sanket Finance
Sai IT
Sai IT
Time taken: 178.419 seconds, Fetched: 6 row(s)
hive>
```

F. LEFT OUTER JOIN



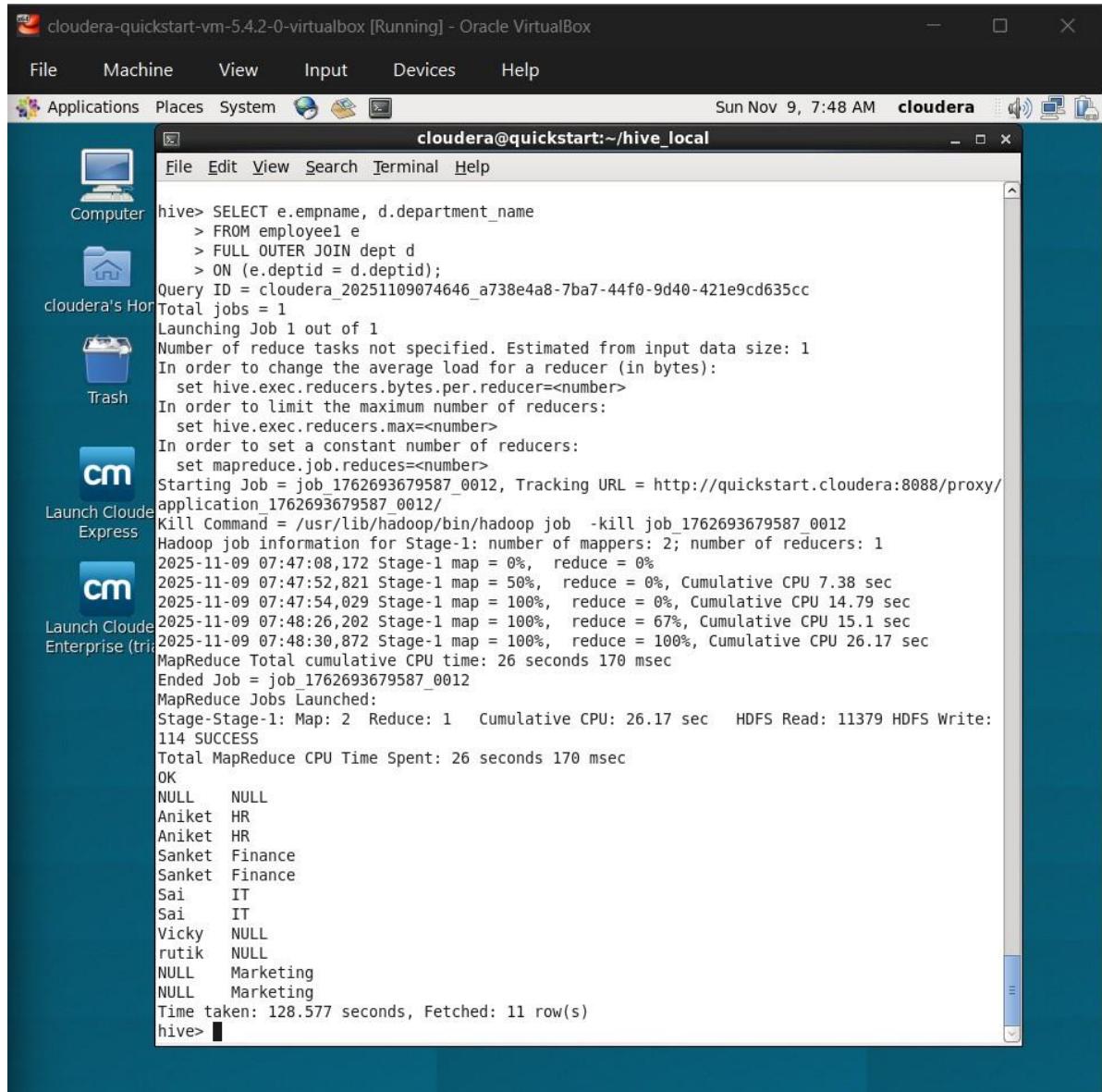
The screenshot shows a Linux desktop environment with a terminal window open. The terminal window title is "cloudera@quickstart:~/hive_local". The terminal content displays a Hive query execution:

```
hive> SELECT e.empname, d.department_name
   > FROM employee1 e
   > LEFT OUTER JOIN dept d
   > ON (e.deptid = d.deptid);
Query ID = cloudera_20251109073939_9732965b-03b1-4802-9cab-0fbe01f8eaf3
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 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>
Starting Job = job_1762693679587_0010, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1762693679587_0010/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1762693679587_0010
Hadoop job information for Stage-1: number of mappers: 2; number of reducers: 1
2025-11-09 07:40:48,350 Stage-1 map = 0%, reduce = 0%
2025-11-09 07:41:49,295 Stage-1 map = 0%, reduce = 0%
2025-11-09 07:41:58,441 Stage-1 map = 83%, reduce = 0%, Cumulative CPU 19.51 sec
2025-11-09 07:41:59,669 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 20.26 sec
2025-11-09 07:42:42,967 Stage-1 map = 100%, reduce = 67%, Cumulative CPU 26.01 sec
2025-11-09 07:42:51,345 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 33.91 sec
MapReduce Total cumulative CPU time: 33 seconds 910 msec
Ended Job = job_1762693679587_0010
MapReduce Jobs Launched:
Stage-Stage-1: Map: 2 Reduce: 1 Cumulative CPU: 33.91 sec HDFS Read: 11376 HDFS Write: 82 SUCCESS
Total MapReduce CPU Time Spent: 33 seconds 910 msec
OK
Aniket  HR
Aniket  HR
Sanket  Finance
Sanket  Finance
Sai     IT
Sai     IT
Vicky   NULL
rutik   NULL
Time taken: 177.78 seconds, Fetched: 8 row(s)
hive>
```

G. RIGHT OUTER JOIN

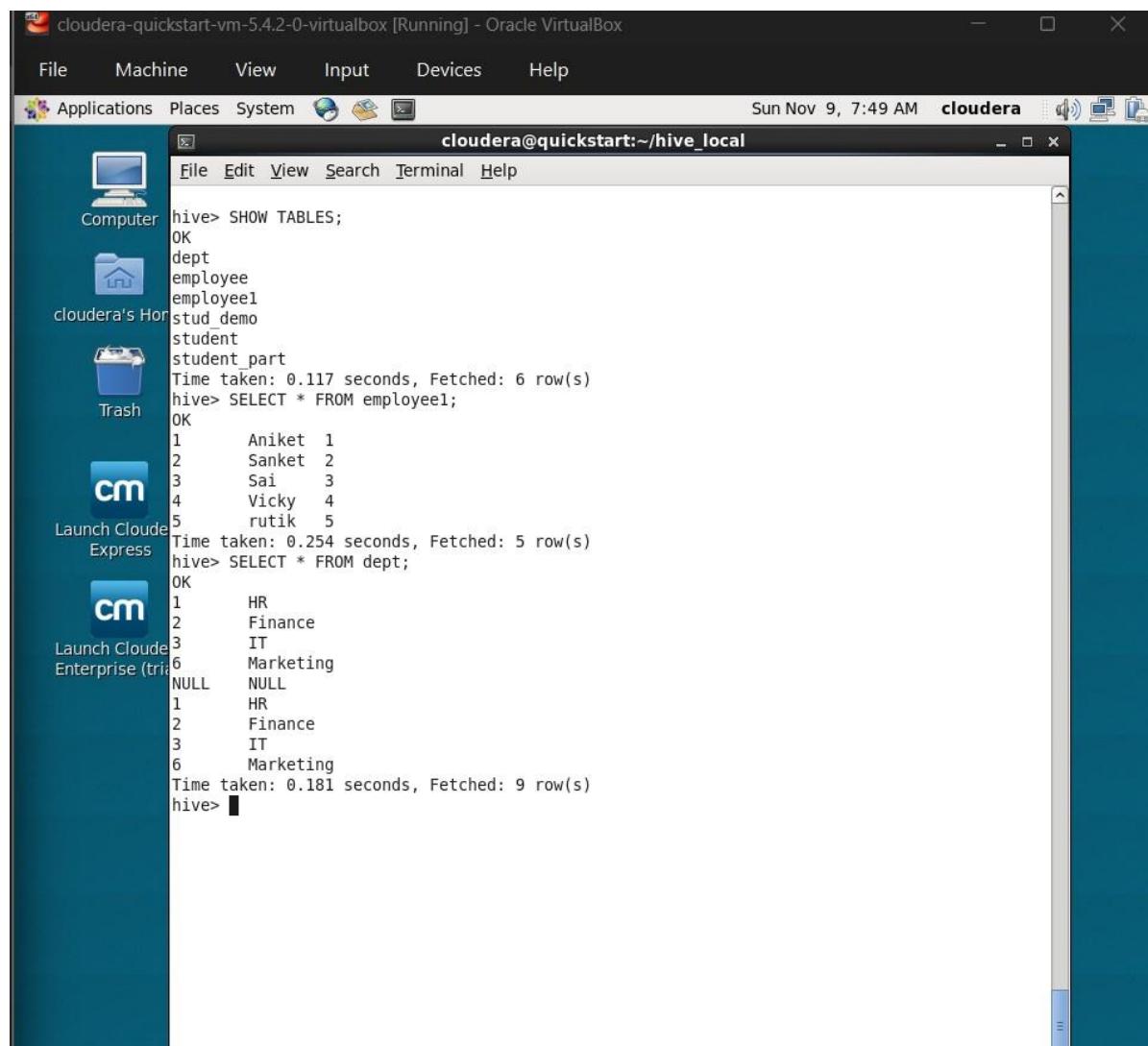
```
hive> SELECT e.empname, d.department_name
   > FROM employee1 e
   > RIGHT OUTER JOIN dept d
   > ON (e.deptid = d.deptid);
Query ID = cloudera_20251109074343_89a48aec-43e4-48c3-8bf6-b0465030b5b3
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 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>
Starting Job = job_1762693679587_0011, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1762693679587_0011/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1762693679587_0011
Hadoop job information for Stage-1: number of mappers: 2; number of reducers: 1
2025-11-09 07:44:13,704 Stage-1 map = 0%, reduce = 0%
2025-11-09 07:45:14,761 Stage-1 map = 0%, reduce = 0%
2025-11-09 07:45:18,603 Stage-1 map = 33%, reduce = 0%, Cumulative CPU 8.55 sec
2025-11-09 07:45:20,394 Stage-1 map = 50%, reduce = 0%, Cumulative CPU 9.25 sec
2025-11-09 07:45:21,534 Stage-1 map = 100%, reduce = 0%, Cumulative CPU 18.56 sec
2025-11-09 07:46:01,364 Stage-1 map = 100%, reduce = 67%, Cumulative CPU 24.12 sec
2025-11-09 07:46:07,792 Stage-1 map = 100%, reduce = 100%, Cumulative CPU 31.25 sec
MapReduce Total cumulative CPU time: 31 seconds 250 msec
Ended Job = job_1762693679587_0011
MapReduce Jobs Launched:
Stage-Stage-1: Map: 2 Reduce: 1 Cumulative CPU: 31.25 sec HDFS Read: 11376 HDFS Write: 96 SUCCESS
Total MapReduce CPU Time Spent: 31 seconds 250 msec
OK
NULL      NULL
Aniket    HR
Aniket    HR
Sanket    Finance
Sanket    Finance
Sai        IT
Sai        IT
NULL      Marketing
NULL      Marketing
Time taken: 181.444 seconds, Fetched: 9 row(s)
hive>
```

H. FULL OUTER JOIN



The screenshot shows a Linux desktop environment with a terminal window open. The terminal window title is "cloudera@quickstart:~/hive_local". The terminal content displays a Hive query execution process:

```
hive> SELECT e.empname, d.department_name
   > FROM employee1 e
   > FULL OUTER JOIN dept d
   > ON (e.deptid = d.deptid);
Query ID = cloudera_20251109074646_a738e4a8-7ba7-44f0-9d40-421e9cd635cc
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks not specified. Estimated from input data size: 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>
Starting Job = job_1762693679587_0012, Tracking URL = http://quickstart.cloudera:8088/proxy/application_1762693679587_0012/
Kill Command = /usr/lib/hadoop/bin/hadoop job -kill job_1762693679587_0012
Hadoop job information for Stage-1: number of mappers: 2; number of reducers: 1
2025-11-09 07:47:08,172 Stage-1 map = 0%,  reduce = 0%
2025-11-09 07:47:52,821 Stage-1 map = 50%,  reduce = 0%, Cumulative CPU 7.38 sec
2025-11-09 07:47:54,029 Stage-1 map = 100%,  reduce = 0%, Cumulative CPU 14.79 sec
2025-11-09 07:48:26,202 Stage-1 map = 100%,  reduce = 67%, Cumulative CPU 15.1 sec
2025-11-09 07:48:30,872 Stage-1 map = 100%,  reduce = 100%, Cumulative CPU 26.17 sec
MapReduce Total cumulative CPU time: 26 seconds 170 msec
Ended Job = job_1762693679587_0012
MapReduce Jobs Launched:
Stage-Stage-1: Map: 2 Reduce: 1 Cumulative CPU: 26.17 sec HDFS Read: 11379 HDFS Write: 114 SUCCESS
Total MapReduce CPU Time Spent: 26 seconds 170 msec
OK
NULL      NULL
Aniket    HR
Aniket    HR
Sanket    Finance
Sanket    Finance
Sai        IT
Sai        IT
Vicky     NULL
rutik    NULL
NULL      Marketing
NULL      Marketing
Time taken: 128.577 seconds, Fetched: 11 row(s)
hive>
```



cloudera@quickstart:~/hive_local

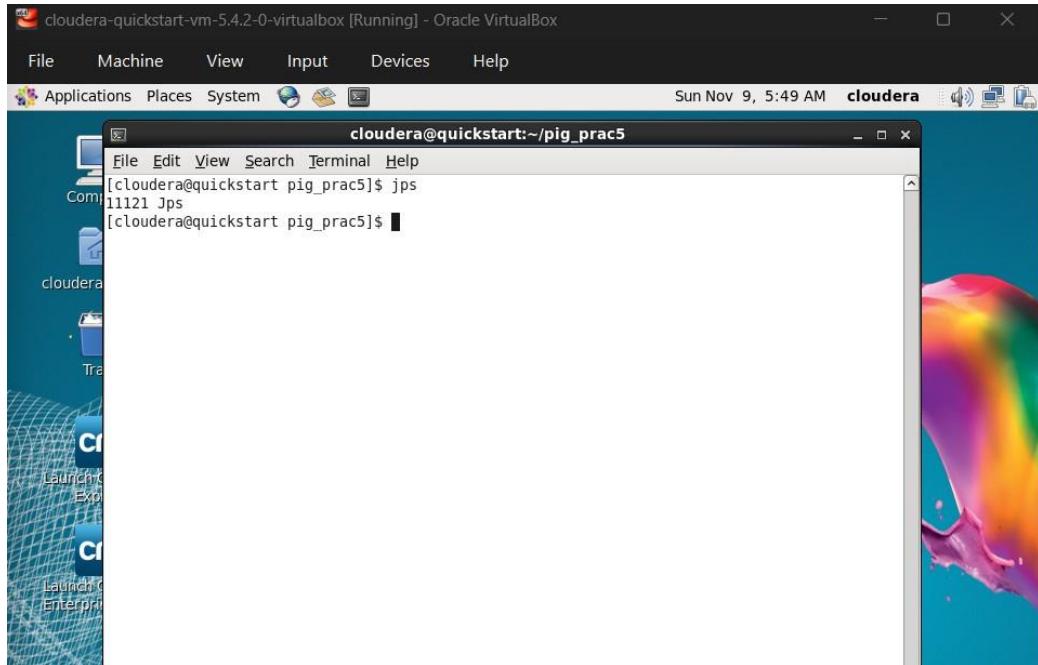
```
hive> SHOW TABLES;
OK
dept
employee
employee1
stud_demo
student
student_part
Time taken: 0.117 seconds, Fetched: 6 row(s)
hive> SELECT * FROM employee1;
OK
1      Aniket  1
2      Sanket  2
3      Sai     3
4      Vicky   4
5      rutik   5
Time taken: 0.254 seconds, Fetched: 5 row(s)
hive> SELECT * FROM dept;
OK
1      HR
2      Finance
3      IT
6      Marketing
NULL  NULL
1      HR
2      Finance
3      IT
6      Marketing
Time taken: 0.181 seconds, Fetched: 9 row(s)
hive>
```

Practical No 5 : Pig

1. Start the Cloudera Quickstart VM

Launch the VirtualBox and start the Cloudera Quickstart virtual machine.

Confirm that Hadoop daemons (NameNode, DataNode, etc.) are running.



2: Verify Pig Installation

Check Pig version and environment

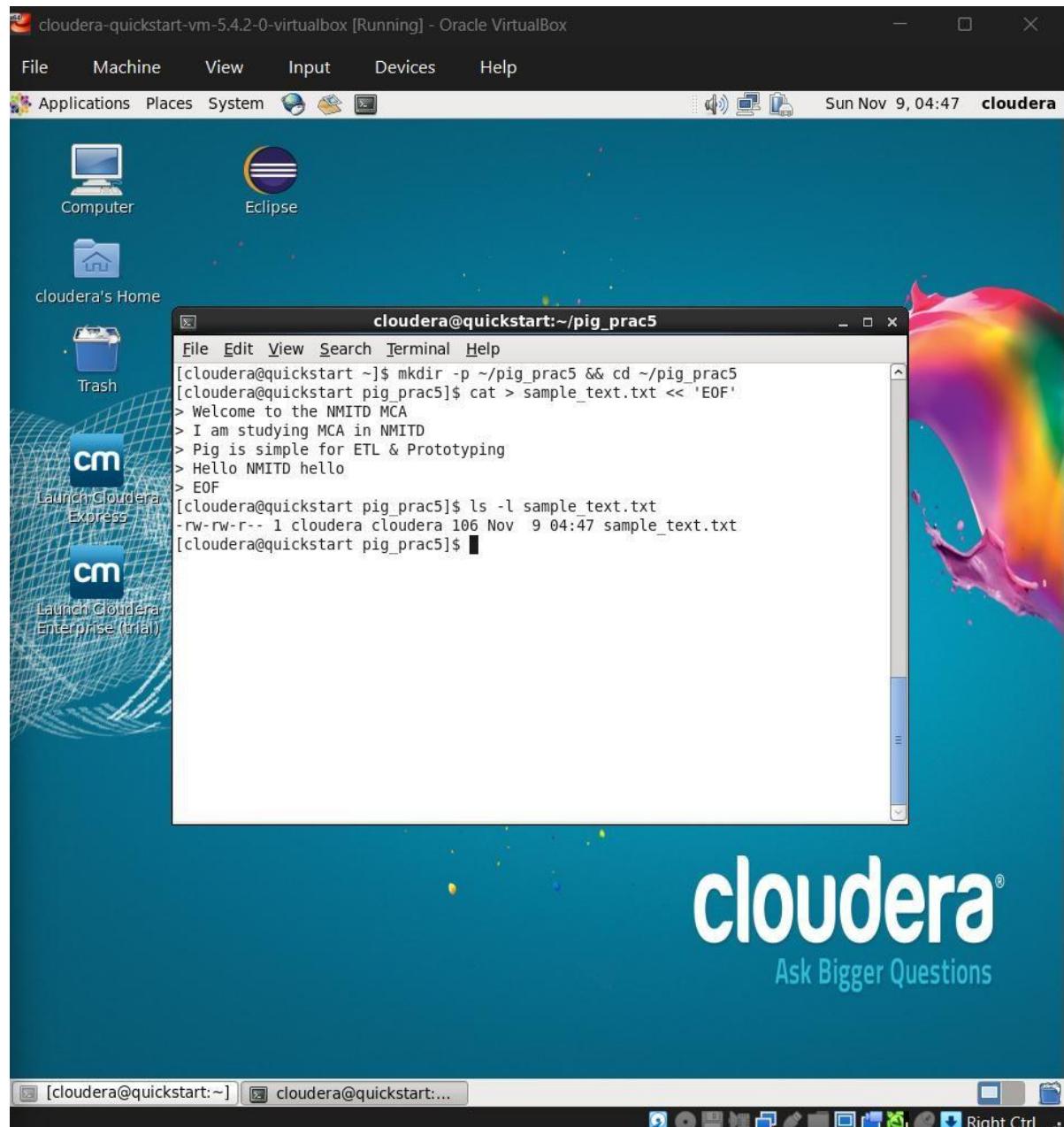


3: Create Working Directory and Dataset

Create a working directory:

Create a sample text file:

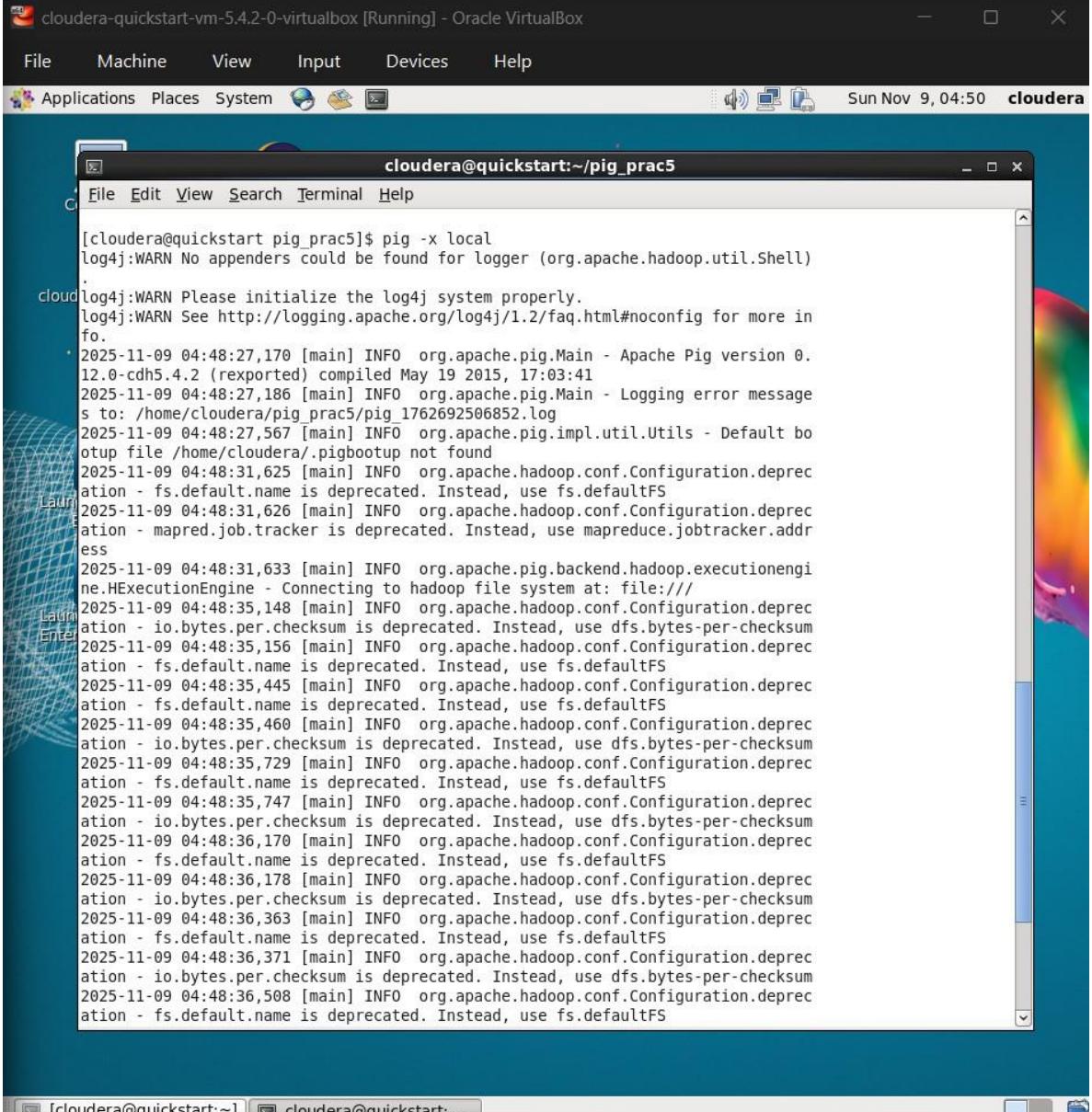
Verify the file



Step 4: Run Pig in Local Mode

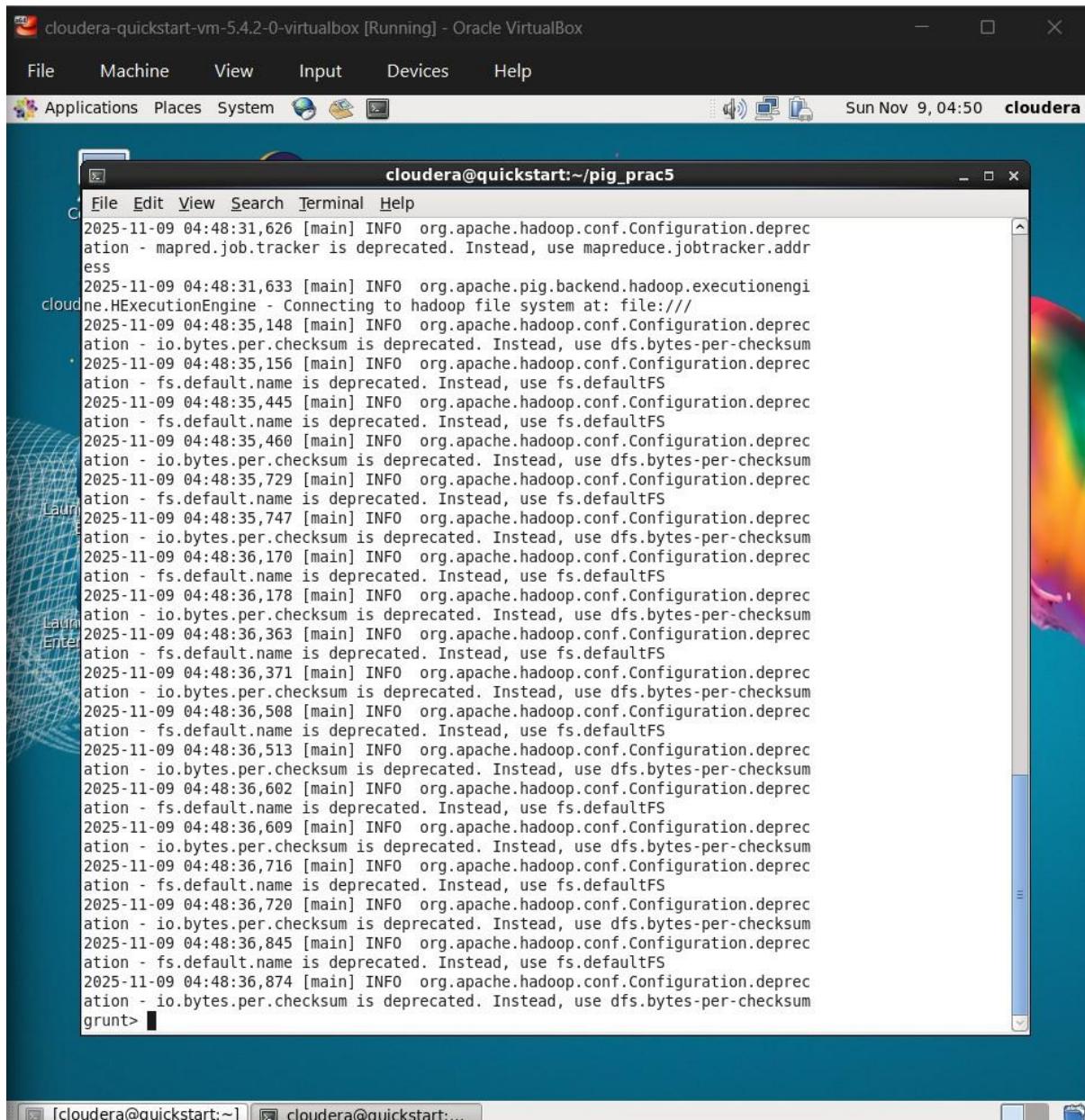
Start Pig shell:

Launch Pig Grunt shell to test basic commands locally.



The screenshot shows a terminal window titled "cloudera@quickstart:~/pig_prac5". The window is running on a Cloudera Quickstart VM (CDH 5.4.2) in Oracle VirtualBox. The terminal displays the output of the command "pig -x local". The log output includes several "WARN" messages from log4j about appenders and configuration. It also shows multiple "INFO" messages from org.apache.pig.Main indicating the version (0.12.0-cdh5.4.2), compilation date (May 19 2015, 17:03:41), and various logging configurations. The terminal window has a standard Linux-style interface with a menu bar (File, Edit, View, Search, Terminal, Help) and a toolbar with icons for Applications, Places, System, and network status. The desktop background is visible behind the window, showing a colorful abstract pattern. The system tray at the bottom right shows the date (Sun Nov 9, 04:50) and the Cloudera logo.

```
[cloudera@quickstart pig_prac5]$ pig -x local
log4j:WARN No appenders could be found for logger (org.apache.hadoop.util.Shell)
.
log4j:WARN Please initialize the log4j system properly.
log4j:WARN See http://logging.apache.org/log4j/1.2/faq.html#noconfig for more in
fo.
2025-11-09 04:48:27,170 [main] INFO  org.apache.pig.Main - Apache Pig version 0.
12.0-cdh5.4.2 (reported) compiled May 19 2015, 17:03:41
2025-11-09 04:48:27,186 [main] INFO  org.apache.pig.Main - Logging error message
s to: /home/cloudera/pig_prac5/pig_1762692506852.log
2025-11-09 04:48:27,567 [main] INFO  org.apache.pig.impl.util.Utils - Default bo
otup file /home/cloudera/.pigbootup not found
2025-11-09 04:48:31,625 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:31,626 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.add
ess
2025-11-09 04:48:31,633 [main] INFO  org.apache.pig.backend.hadoop.executionengi
ne.HExecutionEngine - Connecting to hadoop file system at: file:///-
2025-11-09 04:48:35,148 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:35,156 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:35,445 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:35,460 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:35,729 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:35,747 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:36,170 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:36,178 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:36,363 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:36,371 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:36,508 [main] INFO  org.apache.hadoop.conf.Configuration.deprec
ation - fs.default.name is deprecated. Instead, use fs.defaultFS
```



```
cloudera@quickstart:~/pig_prac5
File Edit View Search Terminal Help
2025-11-09 04:48:31,626 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
2025-11-09 04:48:31,633 [main] INFO org.apache.pig.backend.hadoop.executionengine.HExecutionEngine - Connecting to hadoop file system at: file:///
2025-11-09 04:48:35,148 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:35,156 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:35,445 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:35,460 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:35,729 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:35,747 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:36,170 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:36,178 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:36,363 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:36,371 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:36,508 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:36,513 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:36,602 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:36,609 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:36,716 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:36,720 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:48:36,845 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:48:36,874 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
grunt>
```

Load and view data:

```
ation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
grunt> lines = LOAD 'sample_text.txt' AS (line:chararray);
grunt> DUMP lines;
2025-11-09 04:54:29.083 [main] INFO org.apache.pig.tools.pigstats.ScriptState - Pig features used in
```

```
cloudera@quickstart:~/pig_prac5
File Edit View Search Terminal Help
2025-11-09 04:54:39,909 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Detected Local mode. Stats reported below may be incomplete
2025-11-09 04:54:39,946 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Script Statistics:
cloud
HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.6.0-cdh5.4.2 0.12.0-cdh5.4.2 cloudera 2025-11-09 04:54:32 2025-11-09 04:54:39 UNKNOWN
N

Success!
Job Stats (time in seconds):
Jobid Alias Feature Outputs
job_local2024308346_0001 lines MAP_ONLY file:/tmp/temp1792480495/tmp724326471,
Launch
Input(s):
Successfully read records from: "file:///home/cloudera/pig_prac5/sample_text.txt"
Output(s):
Successfully stored records in: "file:/tmp/temp1792480495/tmp724326471"
Launch
Job DAG:
Enter
job_local2024308346_0001

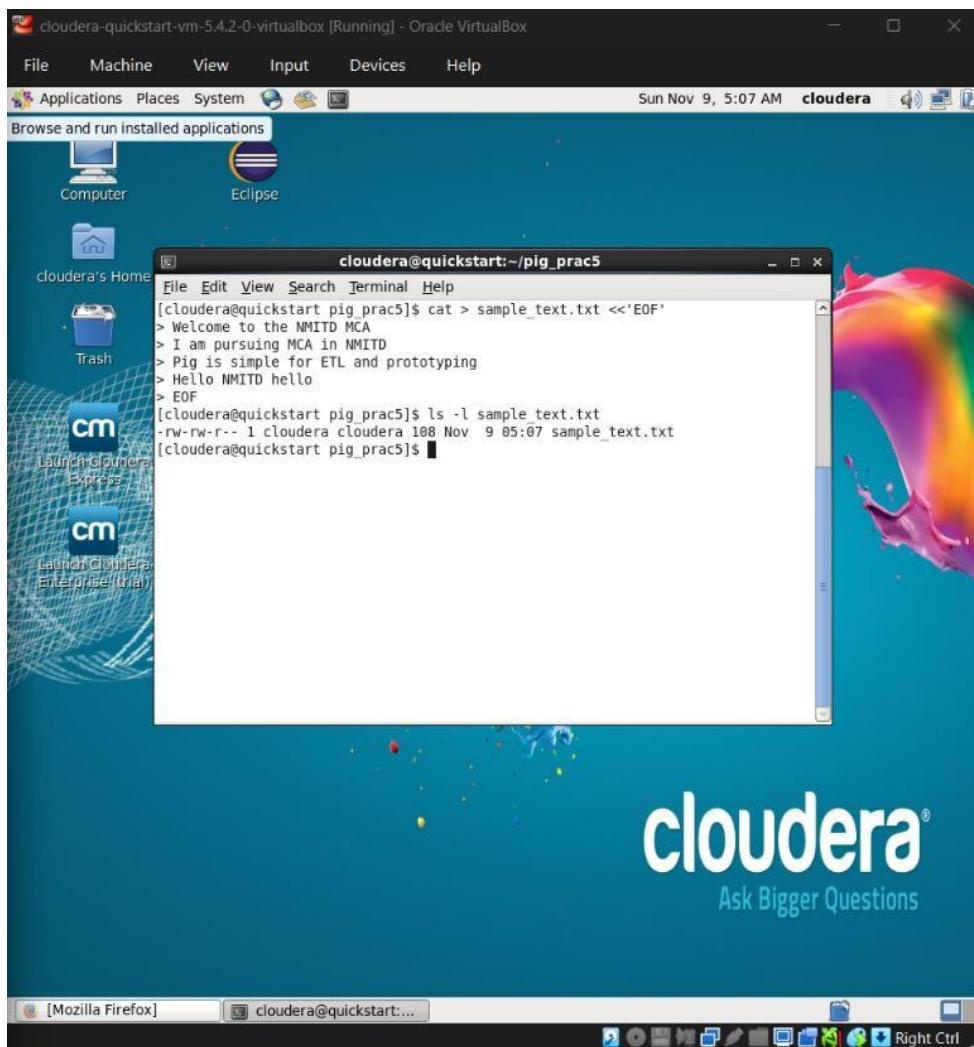
2025-11-09 04:54:39,948 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2025-11-09 04:54:39,970 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 04:54:39,971 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 04:54:39,972 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2025-11-09 04:54:40,023 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2025-11-09 04:54:40,027 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(Welcome to the NMITD MCA)
(I am studying MCA in NMITD)
(Pig is simple for ETL & Prototyping)
>Hello NMITD hello
grunt>
```

5: Write and Execute WordCount Script in Local Mode

Exit Grunt shell using:

quite;

Create script file (Write Pig Latin script for word count in file):



Run that file: `pig -x local wordcount_local.pig`

Output :

The screenshot shows a desktop environment for the Cloudera Quickstart VM. The desktop background is blue with a grid pattern. A terminal window titled "cloudera@quickstart:~/pig_prac5" is open, displaying the output of a Pig Latin script. The terminal window has a dark header bar with the title and a light gray body. The script output includes:

```
Output(s):
Successfully stored records in: "file:/tmp/temp-1690350423/tmp1179811550"

Job DAG:
job_local699795118_0001 ->      job_local189465013_0002,
job_local189465013_0002 ->      job_local2004261248_0003,
job_local2004261248_0003

2025-11-09 05:09:44,617 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - Success!
2025-11-09 05:09:44,712 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - io.bytes.per.checksum is deprecated. Instead, use dfs.bytes-per-checksum
2025-11-09 05:09:44,715 [main] INFO org.apache.hadoop.conf.Configuration.deprecation - fs.default.name is deprecated. Instead, use fs.defaultFS
2025-11-09 05:09:44,716 [main] WARN org.apache.pig.data.SchemaTupleBackend - SchemaTupleBackend has already been initialized
2025-11-09 05:09:44,817 [main] INFO org.apache.hadoop.mapreduce.lib.input.FileInputFormat - Total input paths to process : 1
2025-11-09 05:09:44,818 [main] INFO org.apache.pig.backend.hadoop.executionengine.util.MapRedUtil - Total input paths to process : 1
(nmitd,3)
(hello,2)
(mca,2)
(prototyping,1)
(pursuing,1)
(welcome,1)
(simple,1)
(the,1)
(pig,1)
(for,1)
(etl,1)
(and,1)
(to,1)
(is,1)
(in,1)
(am,1)
(i,1)
grunt>
```

At the bottom of the desktop, there are icons for Mozilla Firefox and the terminal window, along with standard window control buttons.

6. Prepare Data in HDFS (for MapReduce Mode)

Create directory and upload file to HDFS for MapReduce:

The screenshot shows a Linux desktop environment with a terminal window open in the foreground. The terminal window title is "cloudera@quickstart:~/pig_prac5". The terminal session shows the user performing HDFS operations:

```
[cloudera@quickstart pig_prac5]$ hdfs dfs -ls /
Found 5 items
drwxr-xr-x  - hbase supergroup      0 2025-11-09 02:33 /hbase
drwxr-xr-x  - solr   solr          0 2015-06-09 03:38 /solr
drwxrwxrwx  - hdfs supergroup      0 2025-11-09 03:03 /tmp
drwxr-xr-x  - hdfs supergroup      0 2015-06-09 03:38 /user
drwxr-xr-x  - hdfs supergroup      0 2015-06-09 03:36 /var
[cloudera@quickstart pig_prac5]$ hdfs dfs -mkdir -p /user/cloudera/pig_input
[cloudera@quickstart pig_prac5]$ hdfs dfs -ls -d /user/cloudera/pig_input
drwxr-xr-x  - cloudera cloudera    0 2025-11-09 05:21 /user/cloudera/pig_input
[cloudera@quickstart pig_prac5]$ hdfs dfs -put -f sample_text.txt /user/cloudera/pig_input/
[cloudera@quickstart pig_prac5]$ hdfs dfs -ls -h /user/cloudera/pig_input
Found 1 items
-rw-r--r--  1 cloudera cloudera    108 2025-11-09 05:23 /user/cloudera/pig_input/sample_text.txt
[cloudera@quickstart pig_prac5]$
```

Below the terminal, the desktop environment includes a file manager window titled "cloudera@quickstart:~/pig_prac5" showing the same directory structure, and a taskbar at the bottom with icons for Mozilla Firefox and the terminal.

Output :

The screenshot shows a Linux desktop environment with a terminal window open. The terminal window title is "cloudera@quickstart:~/pig_prac5". The terminal session shows the following commands and output:

```
[cloudera@quickstart pig_prac5]$ hdfs dfs -ls /
Found 5 items
drwxr-xr-x  - hbase supergroup      0 2025-11-09 02:33 /hbase
drwxr-xr-x  - solr   solr          0 2015-06-09 03:38 /solr
drwxrwxrwx  - hdfs  supergroup      0 2025-11-09 03:03 /tmp
drwxr-xr-x  - hdfs  supergroup      0 2015-06-09 03:38 /user
drwxr-xr-x  - hdfs  supergroup      0 2015-06-09 03:36 /var
[cloudera@quickstart pig_prac5]$ hdfs dfs -mkdir -p /user/cloudera/pig_input
[cloudera@quickstart pig_prac5]$ hdfs dfs -ls -d /user/cloudera/pig_input
drwxr-xr-x  - cloudera cloudera    0 2025-11-09 05:21 /user/cloudera/pig_input
[cloudera@quickstart pig_prac5]$ hdfs dfs -put -f sample_text.txt /user/cloudera/pig_input/
[cloudera@quickstart pig_prac5]$ hdfs dfs -ls -h /user/cloudera/pig_input
Found 1 items
-rw-r--r--  1 cloudera cloudera    108 2025-11-09 05:23 /user/cloudera/pig_input/sample_text.txt
[cloudera@quickstart pig_prac5]$ hdfs dfs -cat /user/cloudera/pig_input/sample_text.txt
Launching Welcome to the NMITD MCA
I am pursuing MCA in NMITD
Pig is simple for ETL and prototyping
Hello NMITD hello
[cloudera@quickstart pig_prac5]$
```

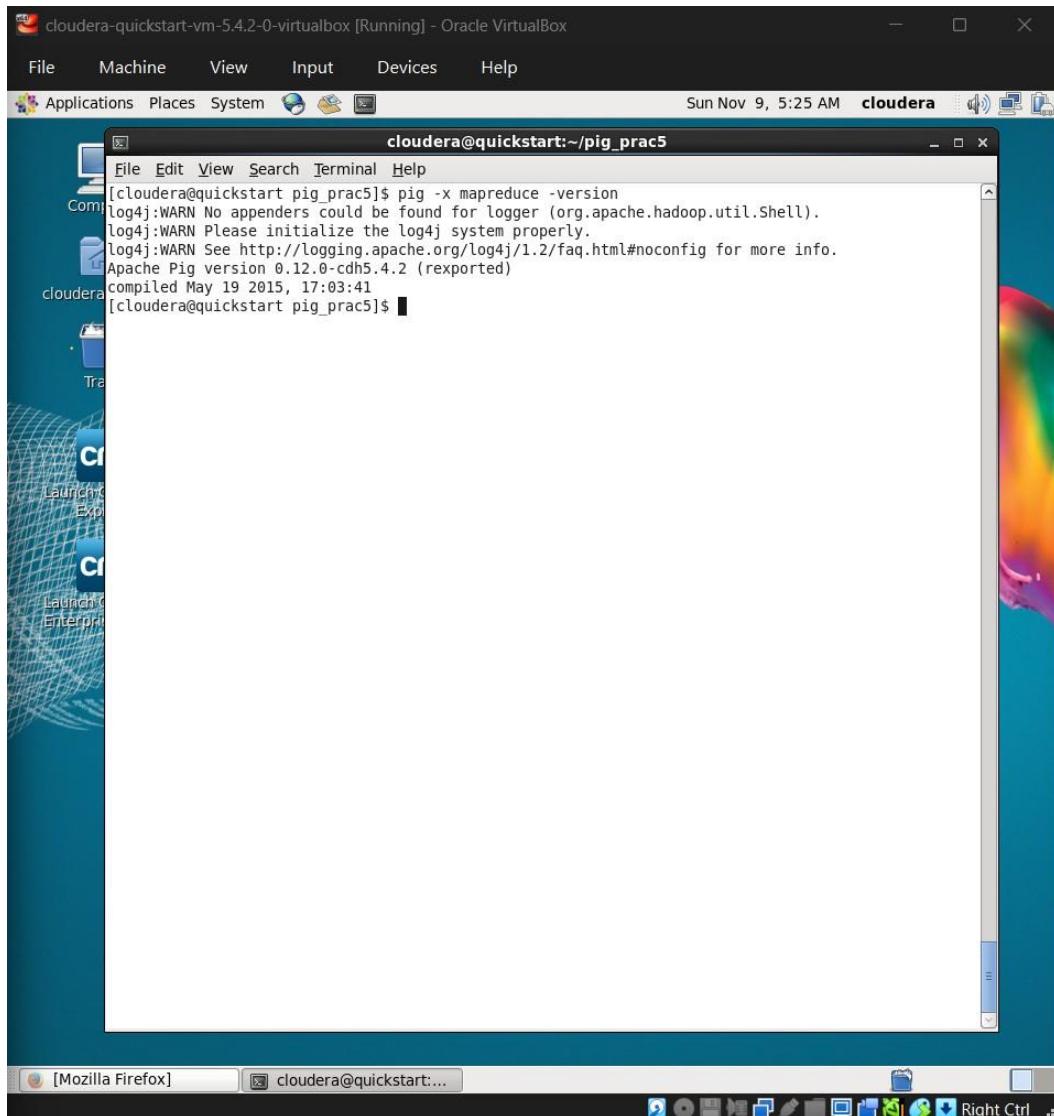
The desktop environment includes a file manager window showing a folder structure with icons for "Compressed (tar)", "cloudera", "Trusted", "Cloudera", and "Launch Enterprise". A taskbar at the bottom shows icons for Mozilla Firefox and the terminal window.

7. Verify Pig MapReduce Mode

Check Pig's MapReduce engine connectivity:

- pig -x mapreduce -version

Output :



8. Create Pig WordCount Script for HDFS

Create the Pig script:

The screenshot shows a terminal window titled "cloudera@quickstart:~/pig_prac5". The terminal content is as follows:

```
[cloudera@quickstart pig_prac5]$ cat > wordcount_mr.pig <<'PIG'
> lines = LOAD '/user/cloudera/pig_input/sample_text.txt' AS (line:chararray);
> words = FOREACH lines GENERATE FLATTEN(TOKENIZE(line)) AS word;
> norm = FOREACH words GENERATE LOWER(word) AS w;
> grp = GROUP norm BY w;
> counts = FOREACH grp GENERATE group AS word, COUNT(norm) AS total;
> ordered = ORDER counts BY total DESC;
> STORE ordered INTO '/user/cloudera/pig_output_wordcount' USING PigStorage('\t');
> PIG
[cloudera@quickstart pig_prac5]$ ls -l wordcount_mr.pig
-rw-rw-r-- 1 cloudera cloudera 398 Nov  9 05:28 wordcount_mr.pig
[cloudera@quickstart pig_prac5]$
```

The terminal window is part of a desktop environment with a blue-themed desktop background. The taskbar at the bottom shows icons for Mozilla Firefox and the terminal window.

9. Run Pig in MapReduce Mode

Execute the script:

`pig wordcount_mr.pig`

Pig launches MapReduce jobs & output path displayed.

Output :

```
cloudera@quickstart:~/pig_pracs
File Edit View Search Terminal Help
at org.apache.hadoop.ipc.Client.getConnection(Client.java:1521)
at org.apache.hadoop.ipc.Client.call(Client.java:1438)
... 38 more
2025-11-09 05:40:05,635 [main] WARN org.apache.pig.tools.pigstats.PigStatsUtil - Failed to
get RunningJob for job job_1762693679587_0003
2025-11-09 05:40:05,645 [main] INFO org.apache.pig.backend.hadoop.executionengine.mapReduceLayer.MapReduceLauncher - 100% complete
2025-11-09 05:40:06,003 [main] INFO org.apache.pig.tools.pigstats.SimplePigStats - Script Statistics:
Tran HadoopVersion PigVersion UserId StartedAt FinishedAt Features
2.6.0-cdh5.4.2 0.12.0-cdh5.4.2 cloudera 2025-11-09 05:29:45 2025-11-09 05:40:05G
ROUP_BY,ORDER_BY
Cl Success!
Launch Job Stats (time in seconds):
Exp JobId Maps Reduces MaxMapTime MinMapTime AvgMapTime MedianMapTime MaxR
duceTime MinReduceTime AvgReduceTime MedianReducetime Alias Feature Outp
uts
job_1762693679587_0001 1 1 0 0 0 0 0 0 0 0 0 0
counts.grp.lines.norm.words GROUP_BY,COMBINER
Launched job_1762693679587_0002 1 1 31 31 31 31 34 34 34 3
EnterFor 4 ordered SAMPLER
job_1762693679587_0003 1 1 0 0 0 0 0 0 0 0 0 0
ordéred ORDER_BY /user/cloudera/pig_output_wordcount,
Input(s):
Successfully read 0 records from: "/user/cloudera/pig_input/sample_text.txt"
Output(s):
Successfully stored 0 records in: "/user/cloudera/pig_output_wordcount"
Counters:
Total records written : 0
Total bytes written : 0
Spillable Memory Manager spill count : 0
Total bags proactively spilled: 0
Total records proactively spilled: 0
Job DAG:
job_1762693679587_0001 -> job_1762693679587_0002,
job_1762693679587_0002 -> job_1762693679587_0003,
```

10. Verify the Output in HDFS

Check output directory:

```
hdfs dfs -ls /user/cloudera/pig_output_wordcount
```

```
hdfs dfs -cat /user/cloudera/pig_output_wordcount/part-r-*
```

Output:

The screenshot shows a Linux desktop environment with a terminal window open. The terminal window title is "cloudera@quickstart:~/pig_prac5". The terminal displays the following command and its output:

```
[cloudera@quickstart pig_prac5]$ hdfs dfs -ls /user/cloudera/pig_output_wordcount
Found 2 items
-rw-r--r-- 1 cloudera cloudera          0 2025-11-09 05:38 /user/cloudera/pig_output_wordcount/_SUCCESS
-rw-r--r-- 1 cloudera cloudera 120 2025-11-09 05:38 /user/cloudera/pig_output_wordcount/part-r-00000
[cloudera@quickstart pig_prac5]$ hdfs dfs -cat /user/cloudera/pig_output_wordcount/part-*
nmitd    3
hello    2
mca     2
Traprototyping   1
pursuing  1
welcome  1
simple   1
the      1
pig      1
Launchfor  1
ExpelL   1
and      1
to       1
Cis      1
in       1
LaunchCam 1
Enterprise 1
[cloudera@quickstart pig_prac5]$
```

The desktop background features a Cloudera logo. The taskbar at the bottom shows icons for Mozilla Firefox and the terminal window.

Practical 6 No - Spark

Apache Spark Commands in Scala	Pair RDD (Key-Value RDD) Operations
Start Spark Shell	Create Pair RDD
Create RDD from Collection	reduceByKey
Read Text File into RDD (from HDFS)	groupByKey — With Output Viewing
Map Transformation	mapValues
Filter Transformation	sortByKey
Reduce Action	join
Collect Action	cogroup
Save RDD to HDFS	aggregateByKey
Cache/Persist RDD	foldByKey
	Read from Local File

1. Start Spark Shell

spark-shell --master local[*] (You can type spark-shell)

Starts the interactive Spark shell using all available CPU cores.

2. Create RDD from Collection

val rdd = sc.parallelize(Seq(1, 2, 3, 4, 5))

Creates an RDD from an in-memory Scala collection.

If there is a error like this

```

scala> val rdds=sc.parallelize(Seq(1,2,3,4,5))
java.lang.IllegalStateException: Cannot call methods on a stopped SparkContext.
This stopped SparkContext was created at:

org.apache.spark.SparkContext.<init>(SparkContext.scala:84)
org.apache.spark.repl.SparkILoop.createSparkContext(SparkILoop.scala:1022)
$iwC$iw$<init>(<console>:15)
$iwC.<init>(<console>:25)
<init>(<console>:27)
.<init>(<console>:31)
.<clinit>(<console>)
.<init>(<console>:7)
.<clinit>(<console>)
$print(<console>)
sun.reflect.NativeMethodAccessorImpl.invoke0(Native Method)
sun.reflect.NativeMethodAccessorImpl.invoke(NativeMethodAccessorImpl.java:57)
sun.reflect.DelegatingMethodAccessorImpl.invoke(DelegatingMethodAccessorImpl.java:43)
java.lang.reflect.Method.invoke(Method.java:606)
org.apache.spark.repl.SparkIMain$ReadEvalPrint.call(SparkIMain.scala:1045)
org.apache.spark.repl.SparkIMain$Request.loadAndRun(SparkIMain.scala:1326)
org.apache.spark.repl.SparkIMain.loadAndRunReq$1(SparkIMain.scala:821)
org.apache.spark.repl.SparkIMain.interpret(SparkIMain.scala:852)
org.apache.spark.repl.SparkIMain.interpret(SparkIMain.scala:800)
org.apache.spark.repl.SparkILoop.reallyInterpret$1(SparkILoop.scala:857)

The currently active SparkContext was created at:
(No active SparkContext.)
```

You need to **create a new `SparkContext`**, since the old one is stopped and can't be used again.

Step-by-step in `spark-shell` or Scala REPL :

1. Stop the existing `SparkContext` explicitly (if not already stopped):

scala

Copy code

```
sc.stop()
```

2. Import necessary Spark classes:

scala

Copy code

```
import org.apache.spark.{SparkConf, SparkContext}
```

3. Create a new `SparkConf` and `SparkContext`:

scala

Copy code

```
val conf = new SparkConf().setAppName("NewApp").setMaster("local[*]")
val sc = new SparkContext(conf)
```

4. Now re-run your command:

scala

Copy code

```
↓
val rdd = sc.parallelize(Seq(1,2,3,4,5))
```

`val rdd = sc.parallelize(Seq(1,2,3,4,5))`

```
scala> val rdd = sc.parallelize(Seq(1,2,3,4,5))
rdd: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[0] at parallelize at <console>:21
```

3. Read Text File into RDD (from HDFS)

`val rdd = sc.textFile("hdfs:///user/cloudera/file.txt")`

Reads a text file from HDFS. Each line becomes an RDD element.

How to Read/View an RDD

Task	Command	Explanation
Print all elements	<code>rdd.collect().foreach(println)</code>	Prints each element. Use only if RDD is small.
Print first 5 lines	<code>rdd.take(5).foreach(println)</code>	Prints first 5 elements. Useful for previews.
Count number of lines	<code>rdd.count()</code>	Shows total number of records.
Print with line numbers	<code>rdd.zipWithIndex().collect().foreach { case (line, i) => println(s"\$i: \$line") }</code>	Shows each line with its index.

```
scala> result.collect()
[Stage 0:>                                (0 + 0) / 2]25/08/23 02:39:39 WARI
; check your cluster UI to ensure that workers are registered and have sufficient resources
res0: Array[Int] = Array(2, 4, 6, 8, 10)

scala> █
```

4. Map Transformation

```
val result = rdd.map(x => x * 2)
```

Applies a function to every element in the RDD.

```
scala> val result=rdd.map(x=>x*2)
result: org.apache.spark.rdd.RDD[Int] = MapPartitionsRDD[1] at map at <console>:23

scala> result.collect()
[Stage 0:>                                         (0 + 0) / 2]25/08/23 02:39:39 WARN cluster.
; check your cluster UI to ensure that workers are registered and have sufficient resources
res0: Array[Int] = Array(2, 4, 6, 8, 10)

scala> ■
```

5. Filter Transformation

```
val result = rdd.filter(x => x > 3)
```

Keeps only the elements that satisfy the condition.

```
scala> val result=rdd.filter(x=>x>3)
result: org.apache.spark.rdd.RDD[Int] = MapPartitionsRDD[2] at filter at <console>:23

scala> result.collect()
res1: Array[Int] = Array(4, 5)

scala> val sum=rdd.reduce((a,b)=>a+b)
sum: Int = 15
```

6. Reduce Action

```
val sum = rdd.reduce((a, b) => a + b)
```

Combines all elements using the specified function.

```
scala> val pairs = sc.parallelize(Seq(("j", 1), ("a", 2), ("j", 3)))
pairs: org.apache.spark.rdd.RDD[(String, Int)] = ParallelCollectionRDD[3] at parallelize at <console>:21

scala> pairs.reduceByKey(_ + _).collect().foreach(println)
(j,4)
(a,2)

scala> ■
```

7. Collect Action

```
rdd.collect()
```

Returns all RDD elements as an array to the driver.

```
scala> val rdd = sc.parallelize(Seq(1,2,3,4,5))
rdd: org.apache.spark.rdd.RDD[Int] = ParallelCollectionRDD[0] at parallelize at <console>:21

scala> val result=rdd.map(x=>x*2)
result: org.apache.spark.rdd.RDD[Int] = MapPartitionsRDD[1] at map at <console>:23

scala> result.collect()
[Stage 0:>                                         (0 + 0) / 2]25/08/23 02:39:39 WARN cluster.
; check your cluster UI to ensure that workers are registered and have sufficient resources
res0: Array[Int] = Array(2, 4, 6, 8, 10)

scala> ■
```

8. Save RDD to HDFS

```
rdd.saveAsTextFile("hdfs:///user/cloudera/output")
```

9. Saves RDD to the given HDFS directory.

The screenshot shows a terminal window with Scala code and a browser window titled "Browse Directory" showing the contents of the HDFS directory "/user/cloudera".

```
scala> rdd.saveAsTextFile("hdfs:///user/cloudera/output.txt")
scala>
```

Browsing HDFS Mozilla Firefox cloudera

Browse Directory

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
drwxr-xr-x	cloudera	cloudera	0 B	Sat Aug 23 02:36:38 -0700 2025	0	0 B	.sparkStaging
drwxr-xr-x	cloudera	cloudera	0 B	Sat Aug 23 02:43:10 -0700 2025	0	0 B	output.txt

Hadoop, 2017.

10. Cache/Persist RDD

```
rdd.cache()
```

Stores the RDD in memory for faster reuse.

```
scala> rdd.cache()
res4: rdd.type = ParallelCollectionRDD[0] at parallelize at <console>:21
scala>
```

Pair RDD (Key-Value RDD) Operations

11. Create Pair RDD

```
val pairs = sc.parallelize(Seq(("a", 1), ("b", 2), ("a", 3)))
scala> val pairs=sc.parallelize(Seq(("a",1),("b",2),("a",3)))
pairs: org.apache.spark.rdd.RDD[(String, Int)] = ParallelCollectionRDD[7] at parallelize at <console>:21
```

```
scala> val pairs=sc.parallelize(Seq(("a",1),("b",2),("c",3)))
pairs: org.apache.spark.rdd.RDD[(String, Int)] = ParallelCollectionRDD[5] at parallelize at <console>:21
```

12. reduceByKey

```
pairs.reduceByKey(_ + _).collect().foreach(println)
```

Adds values with the same key. Efficient for aggregations.

```
scala> pairs.reduceByKey(_+_).collect().foreach(println)
(b,2)
(a,1)
(c,3)

scala> pairs.reduceByKey(_+_).collect().foreach(println)
(b,2)
(a,4)
```

13. groupByKey — With Output Viewing

```
val grouped = pairs.groupByKey()
// View output clearly
grouped.collect().foreach { case (key, values) =>
  println(s"$key -> ${values.mkString(", ")}")
}
```

Explanation:

Groups values by key. mkString makes the Iterable readable.

Output Example:

```
a -> 1, 3
b -> 2
```

```
scala> val grouped = pairs.groupByKey()

grouped.collect().foreach { case (key, values) =>
  println(s"$key -> ${values.mkString(", ")}")
}
j -> 1, 3
a -> 2
grouped: org.apache.spark.rdd.RDD[(String, Iterable[Int])] = ShuffledRDD[5] at groupBy
Key at <console>:23
```

14. mapValues

```
pairs.mapValues(x => x + 1).collect().foreach(println)
```

Applies a function to the value part only.

```
scala> pairs.mapValues(x=>x+1).collect().foreach(println)
(a,2)
(b,3)
(a,4)
```

15. sortByKey

```
pairs.sortByKey().collect().foreach(println)
In) Sorts RDD by key in ascending order.
```

```
scala> pairs.sortByKey().collect().foreach(println)
(a,1)
(a,3)
(b,2)
```

16. Join

```
val rdd1 = sc.parallelize(Seq(("a", 1), ("b", 2)))
val rdd2 = sc.parallelize(Seq(("a", 100), ("b",
200))) val joined = rdd1.join(rdd2)
joined.collect().foreach(println)
```

Performs an inner join on two key-value RDDs.

```
scala> val rdd1=sc.parallelize(Seq(("a",1),("b",2)))
rdd1: org.apache.spark.rdd.RDD[(String, Int)] = ParallelCollectionRDD[14] at parallelize at <console>:21
scala> val rdd2=sc.parallelize(Seq(("a",100),("b",200)))
rdd2: org.apache.spark.rdd.RDD[(String, Int)] = ParallelCollectionRDD[15] at parallelize at <console>:21
scala> val joined=rdd1.join(rdd2)
joined: org.apache.spark.rdd.RDD[(String, (Int, Int))] = MapPartitionsRDD[18] at join at <console>:25
scala> joined.collect().foreach(println)
(b,(2,200))
(a,(1,100))
scala> ■
```

17. Cogroup

```
rdd1.cogroup(rdd2).collect().foreach(println)
Groups values from both RDDs into tuples of lists for each key.
```

```
scala> rdd1.cogroup(rdd2).collect().foreach(println)
(j,(CompactBuffer(3),CompactBuffer(1500)))
(a,(CompactBuffer(5),CompactBuffer(2200)))
scala> ■
```

18. aggregateByKey

```
pairs.aggregateByKey(0)(_ + _, _ + _).collect().foreach(println)
Combines values for each key using custom logic for intra- and inter-partition aggregation.
```

```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
(a,2)  
(j,1)  
(j,3)  
  
scala> val rdd1 = sc.parallelize(Seq(("j", 3), ("a", 5)))  
val rdd2 = sc.parallelize(Seq(("j", 1500), ("a", 2200)))  
val joined = rdd1.join(rdd2)  
joined.collect().foreach(println)  
(j,(3,1500))  
(a,(5,2200))  
rdd1: org.apache.spark.rdd.RDD[(String, Int)] = ParallelCollectionRDD[5] at para  
llelize at <console>:21  
rdd2: org.apache.spark.rdd.RDD[(String, Int)] = ParallelCollectionRDD[6] at para  
llelize at <console>:22  
joined: org.apache.spark.rdd.RDD[(String, (Int, Int))] = MapPartitionsRDD[9] at  
join at <console>:23  
  
scala> rdd1.cogroup(rdd2).collect().foreach(println)  
(j,(CompactBuffer(3),CompactBuffer(1500)))  
(a,(CompactBuffer(5),CompactBuffer(2200)))  
  
scala> pairs.aggregateByKey(0)(_ + _, _ + _).collect().foreach(println)  
(j,4)  
(a,2)  
scala> ■
```

19. **foldByKey**

pairs.foldByKey(0)(_ + _).collect().foreach(println)

Aggregates values with a neutral zero value using a single function.

```
cloudera@quickstart:~  
File Edit View Search Terminal Help  
scala> val rdd1 = sc.parallelize(Seq(("j", 3), ("a", 5)))  
val rdd2 = sc.parallelize(Seq(("j", 1500), ("a", 2200)))  
val joined = rdd1.join(rdd2)  
joined.collect().foreach(println)  
(j,(3,1500))  
(a,(5,2200))  
rdd1: org.apache.spark.rdd.RDD[(String, Int)] = ParallelCollectionRDD[5] at para  
llelize at <console>:21  
rdd2: org.apache.spark.rdd.RDD[(String, Int)] = ParallelCollectionRDD[6] at para  
llelize at <console>:22  
joined: org.apache.spark.rdd.RDD[(String, (Int, Int))] = MapPartitionsRDD[9] at  
join at <console>:23  
  
scala> rdd1.cogroup(rdd2).collect().foreach(println)  
(j,(CompactBuffer(3),CompactBuffer(1500)))  
(a,(CompactBuffer(5),CompactBuffer(2200)))  
  
scala> pairs.aggregateByKey(0)(_ + _, _ + _).collect().foreach(println)  
(j,4)  
(a,2)  
scala> pairs.foldByKey(0)(_ + _).collect().foreach(println)  
(j,4)  
(a,2)  
scala> ■
```

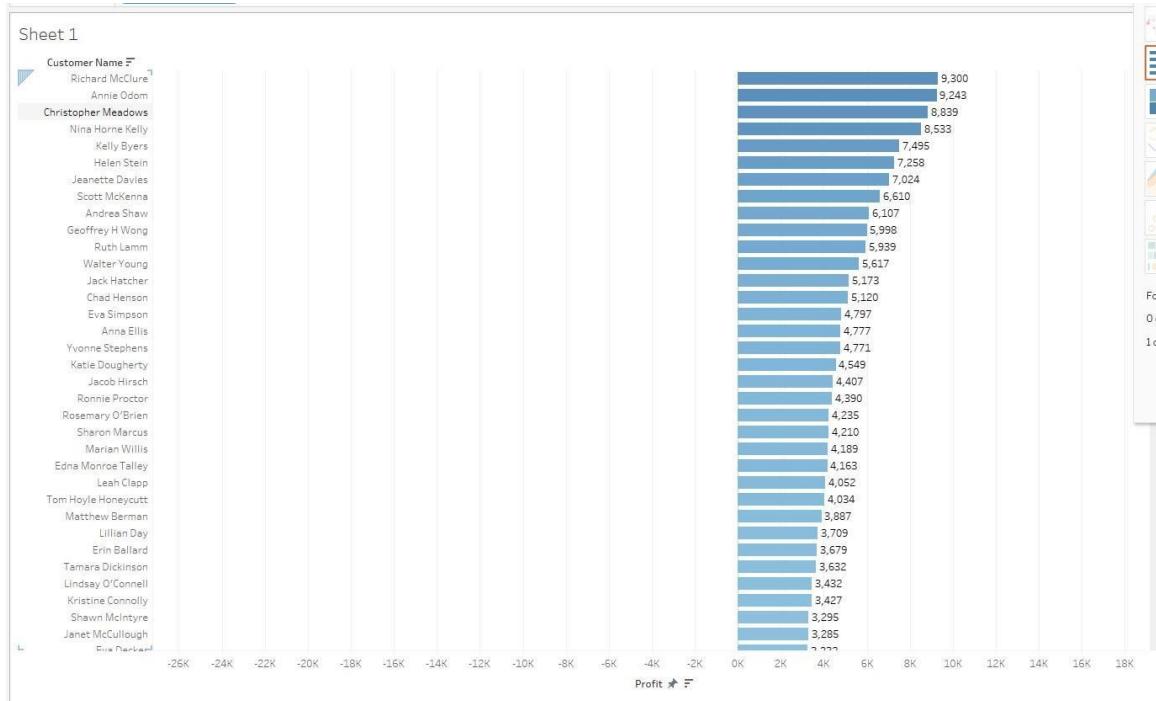
Practical 7 No - Visualization using Tableau:

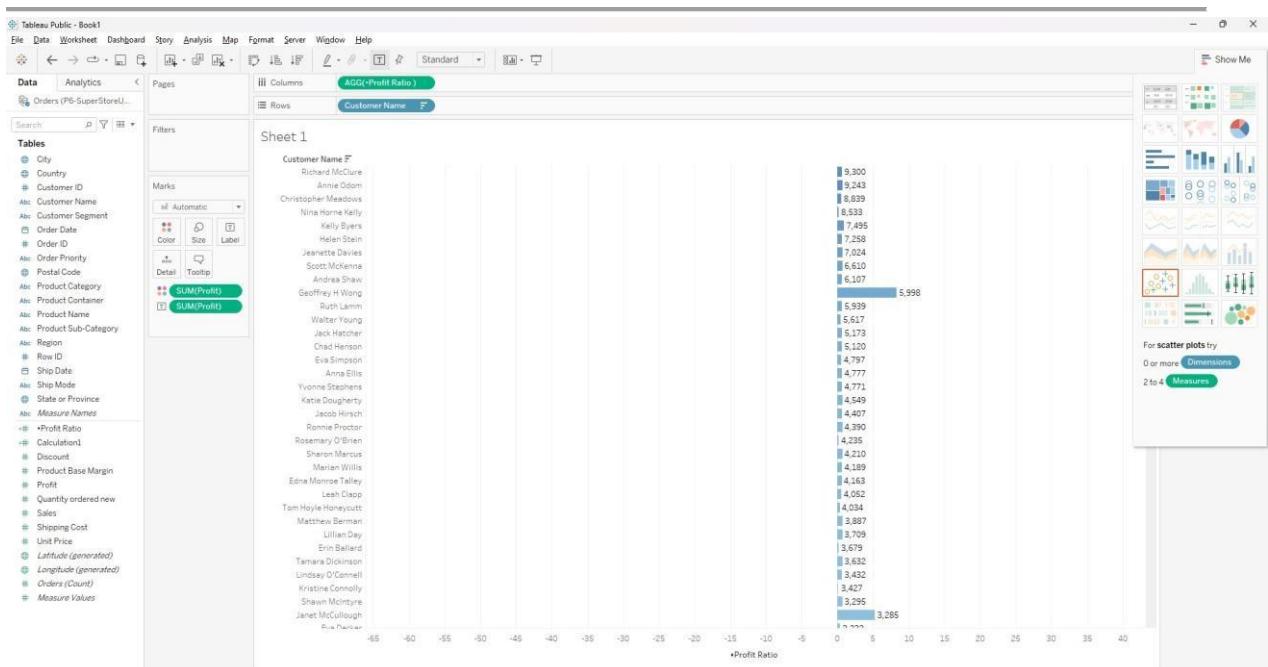
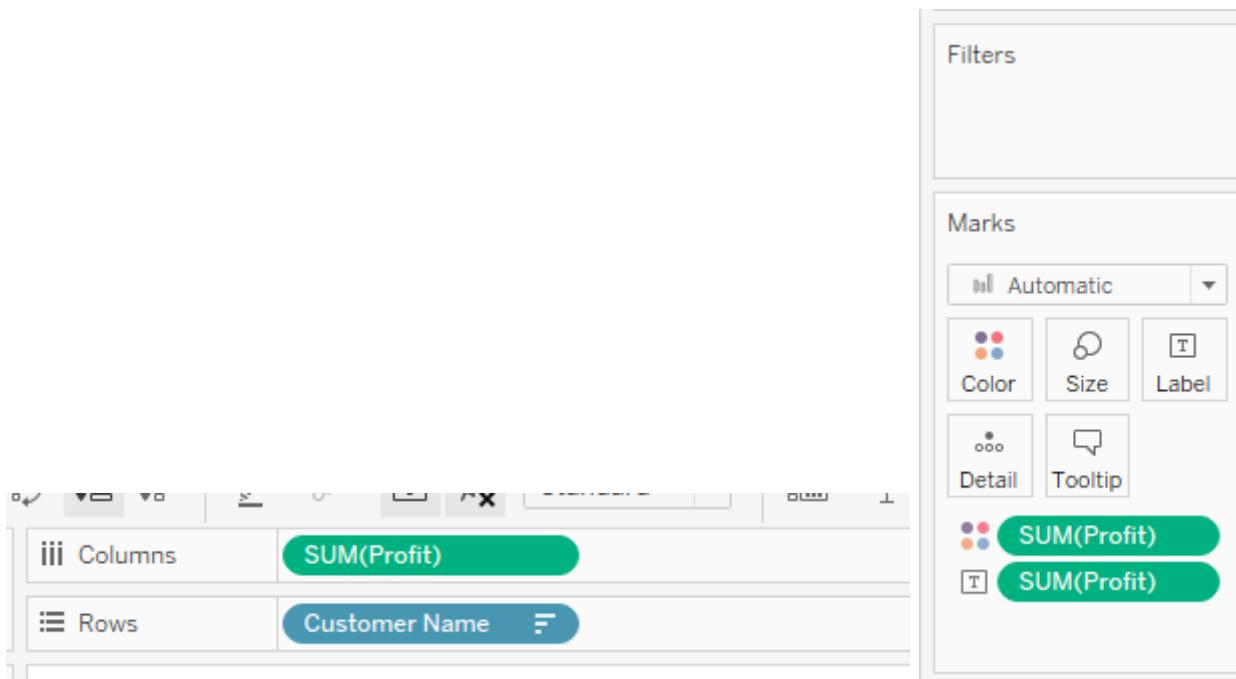
Tableau: Tool Overview, Importing Data, Analyzing with Charts, Creating Dashboards, working with maps

Analysis Operations

Q1. Find the customer with the highest overall profit. What is his/her profit ratio?

- Import superstoreus2015.xlsx dataset.
- Create Calculated Field:
- Profit Ratio = $(\text{SUM}([\text{Profit}]) / \text{SUM}([\text{Sales}]))$
- Drag Customer Name to Rows, Profit to Columns. Sort by Profit.
- Show Profit Ratio in Marks label.





Sheet 1

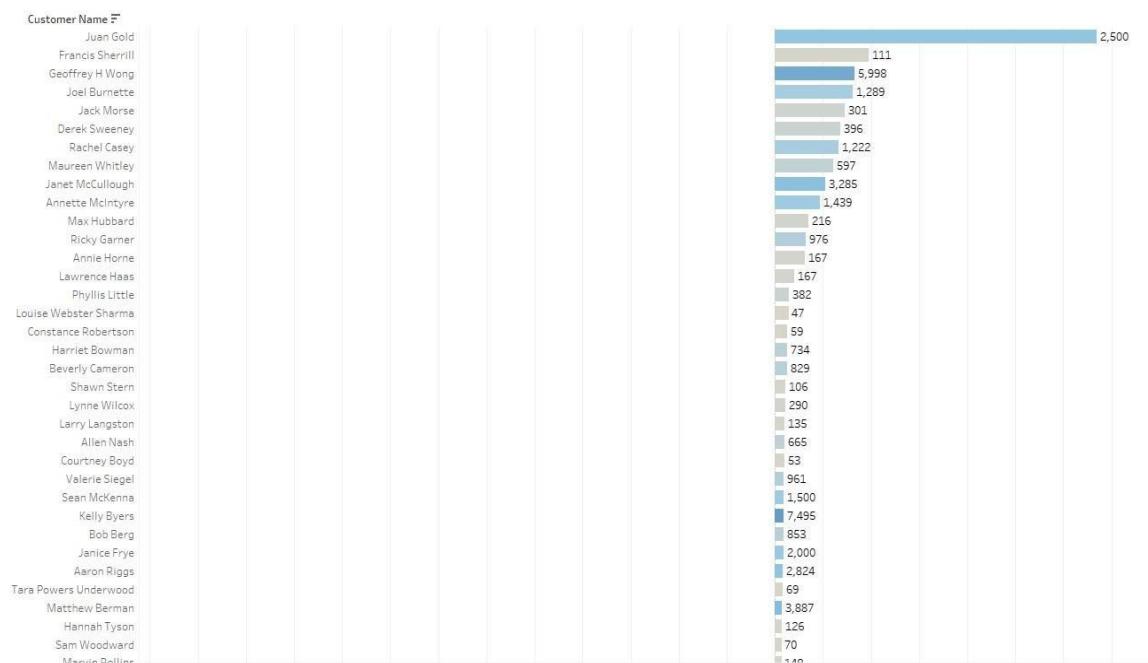


Tableau Public - Book1

File Data Window Help

Connections P6-SuperStoreUS-2015

Sheets Orders

Use Data Interpreter

Orders Returns Users

New Union New Table Extension

Orders

Need more data? Drag tables here to relate them. [Learn more](#)

•Profit Ratio

(SUM([Profit]) / SUM([Sales]))

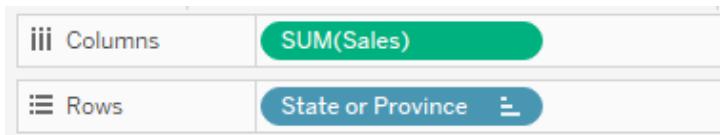
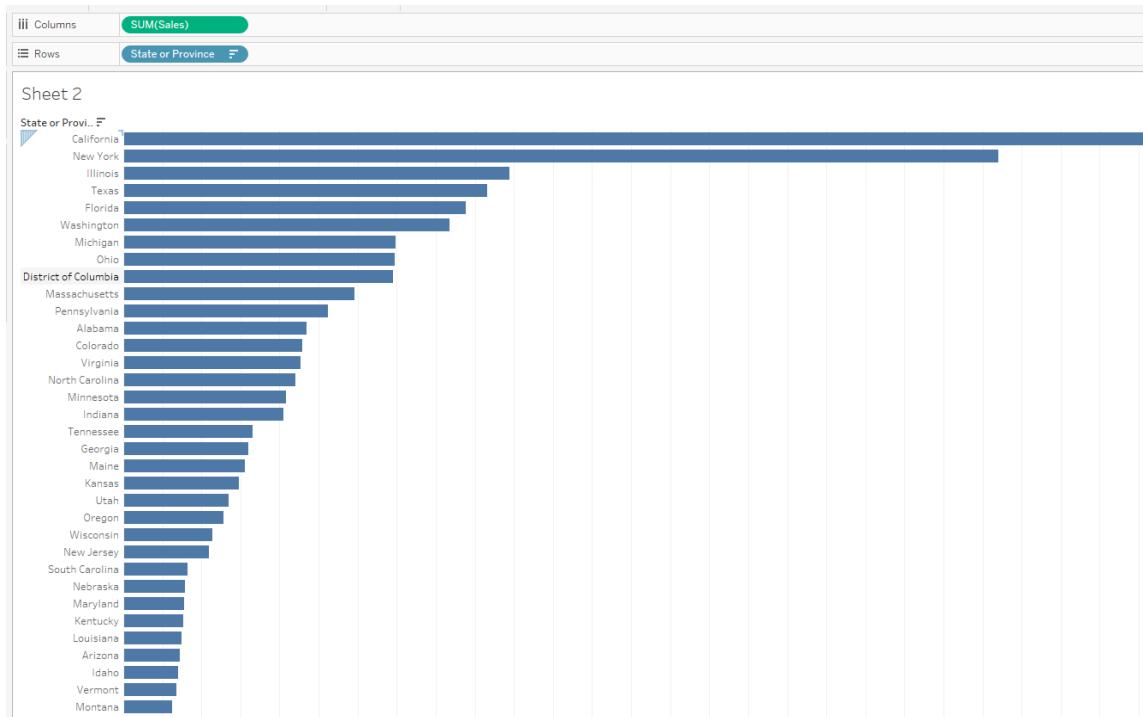
The calculation is valid.

1 Dependency

Apply OK

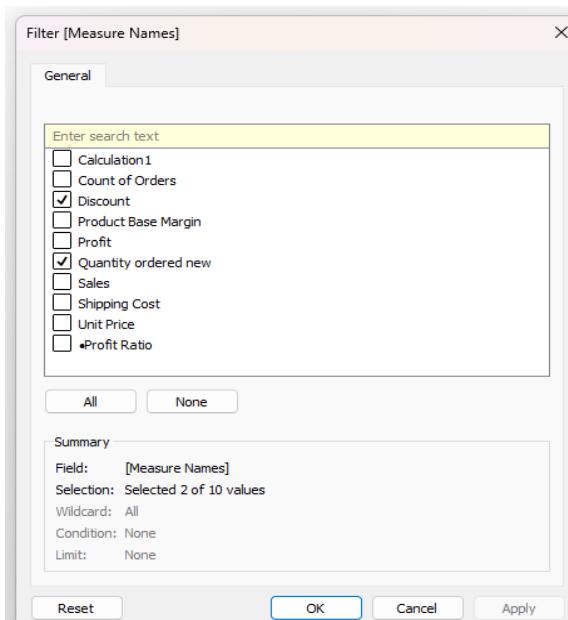
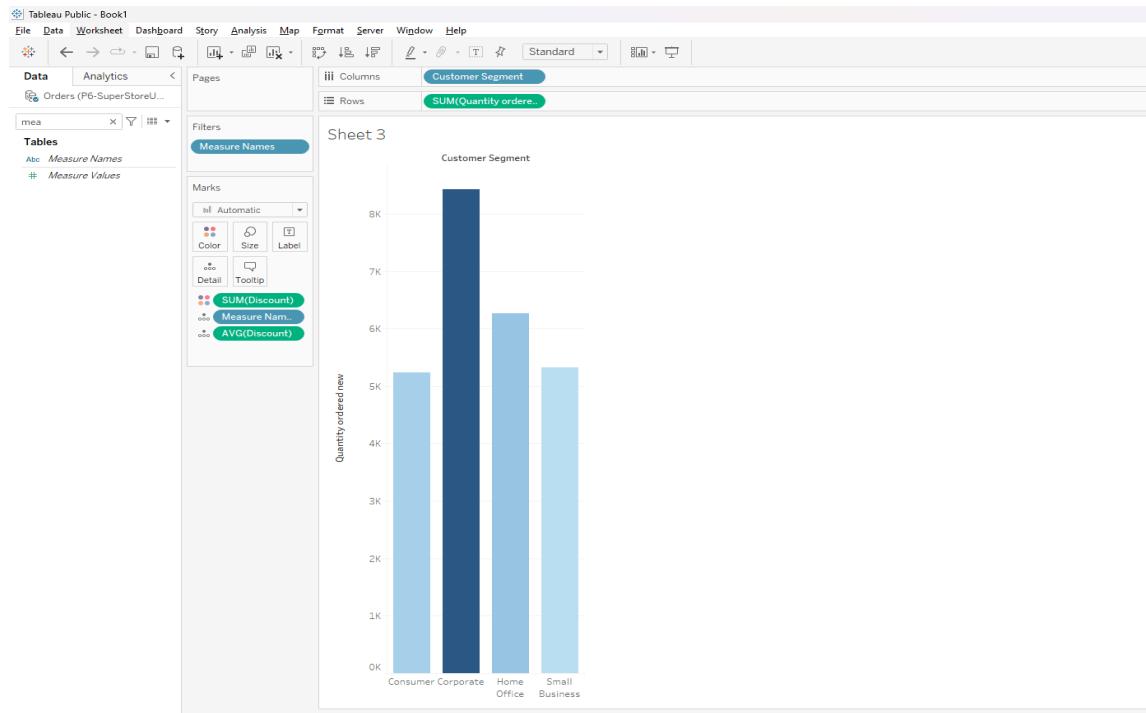
Q2. Which state has the highest Sales (Sum)? What is the total Sales for that state?

- Drag State to Rows, Sales to Columns.
- Sort descending.



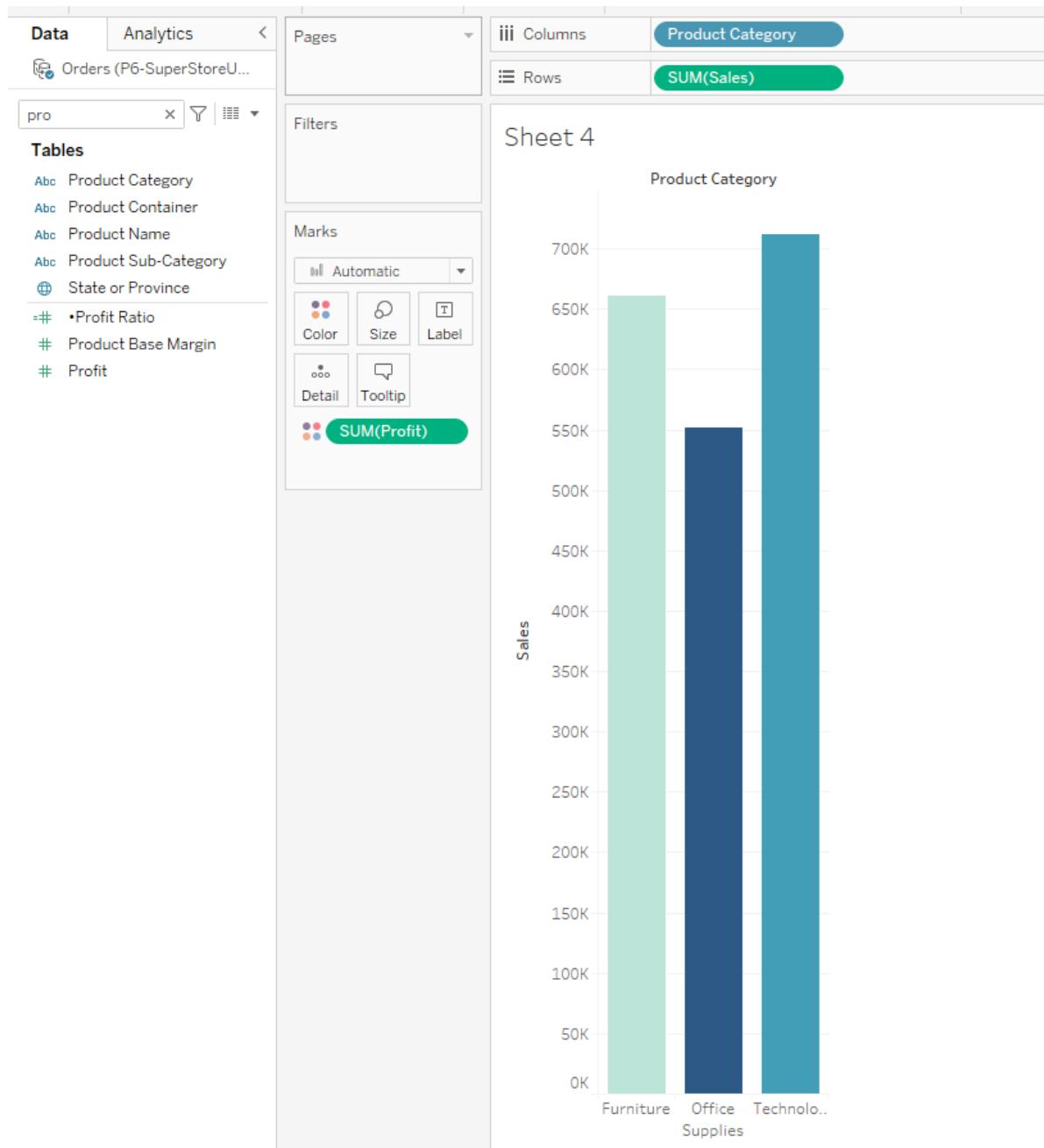
Q3. Which customer segment has both the highest order quantity and average discount rate?

- Drag Segment to Rows.
- Add Order Quantity and Discount to Measures.
- Analyze visually.



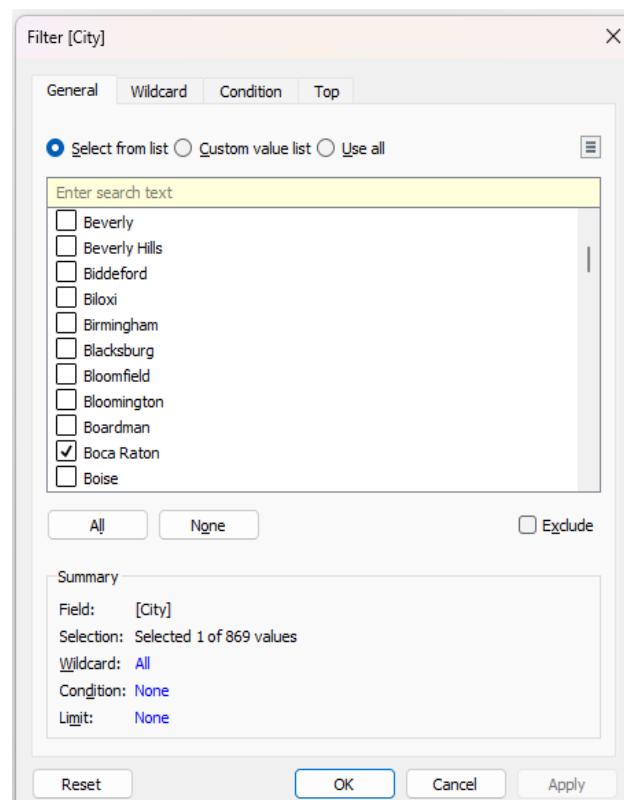
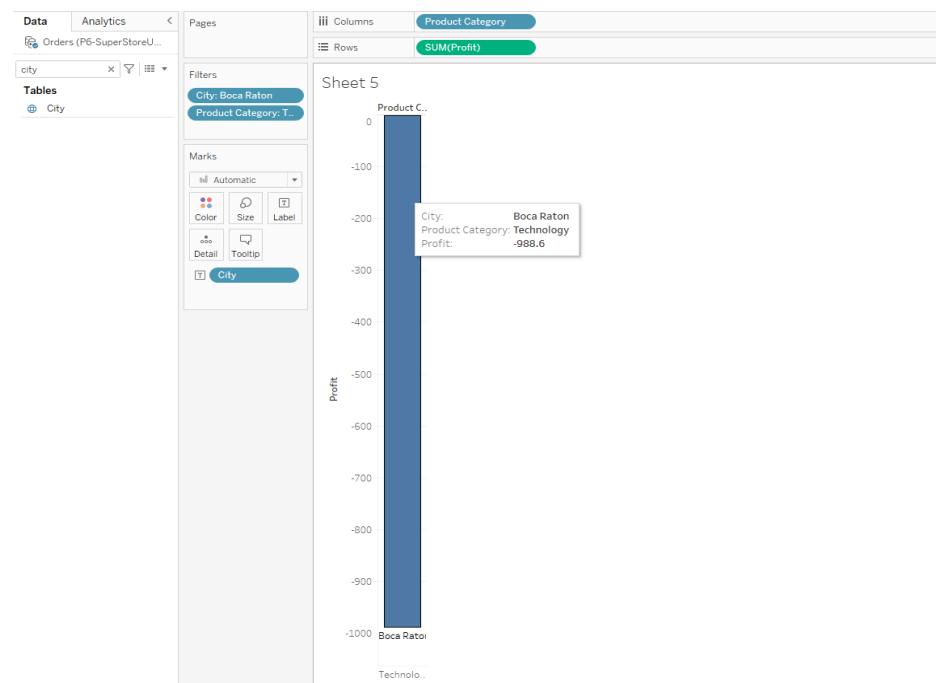
Q4. Which Product Category has the highest total Sales? Which Product Category has the worst Profit?

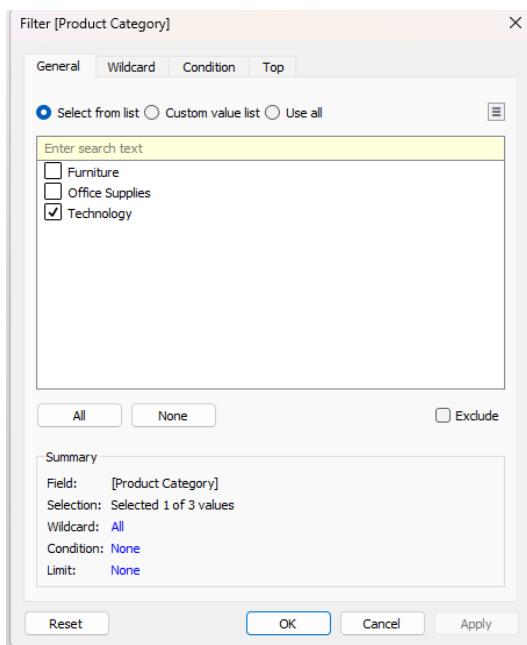
- Create a Bar Chart with Category VS Sales.
- Add Profit as color scale.



Q5. What was the Profit on Technology in Boca Raton (City)?

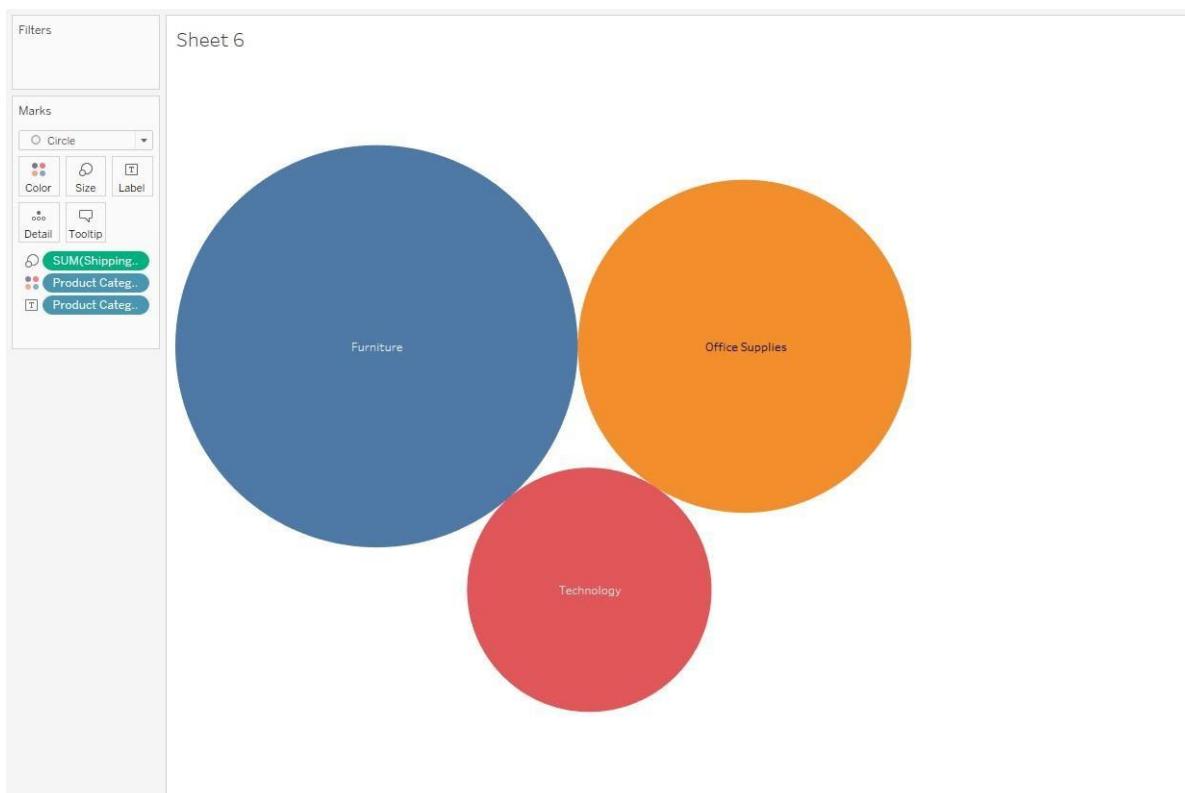
- Add City filter.
- Filter for “Boca Raton”.
- Show Profit.





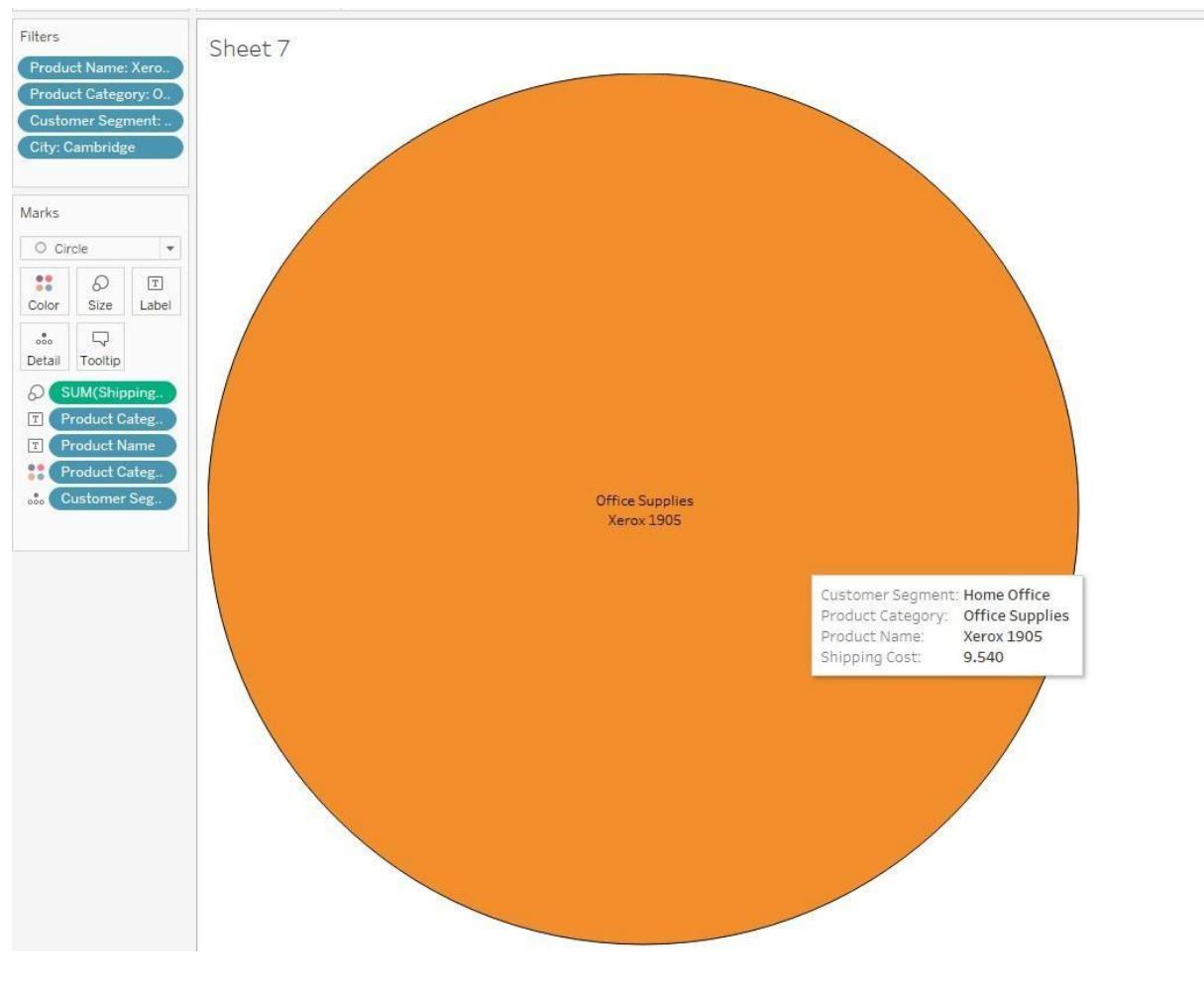
Q6. Which Product Department has the highest Shipping Costs?

- Create **Packed Bubble Chart** with Department.
- Use Shipping Cost as bubble size.



Q7. What was the shipping cost of Office Supplies for Xerox 1905 in the Home Customer Segment in Cambridge?

- Apply filters: Product → Xerox 1905, Category → Office Supplies, Segment → Home, City → Cambridge.



Preparing Maps

Dataset: **Superstore**

1. Create Geographic Hierarchy (Country → State → City → Postal Code)

1. In the **Data Pane**, right-click on Country → select **Hierarchy** → **Create Hierarchy**.
2. Name it (e.g., *Geography*).
3. Drag State, City, and Postal Code into the hierarchy (below Country in order).
 - Now you can drill down from Country → State → City → Postal Code.



Filters

Product Name: Grip ...
AGG(Calculation2)

Marks

Map

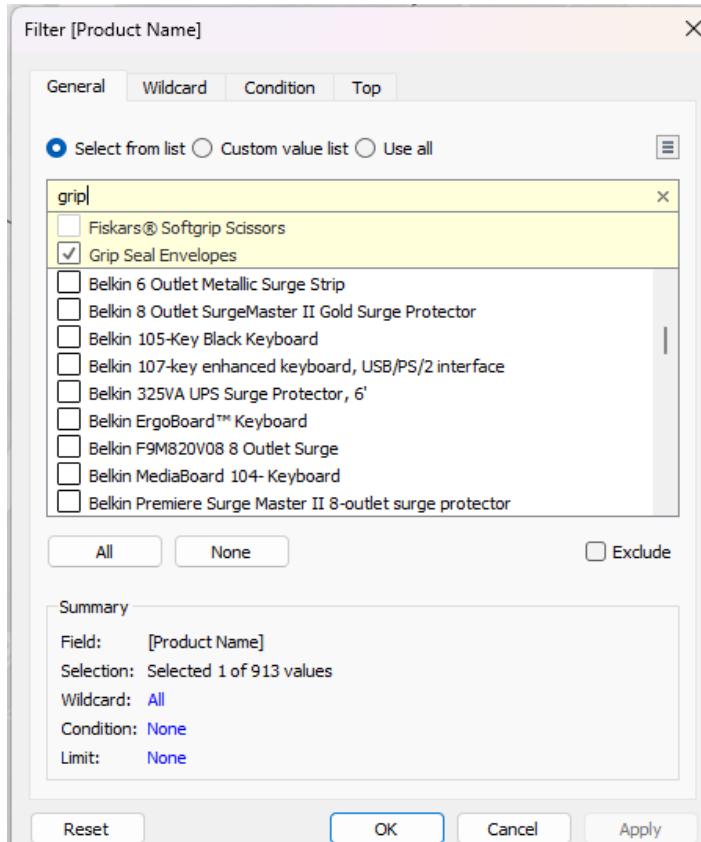
Color Size Label
Detail Tooltip

AGG(Calculation2)
Country State or Pr...
City Postal Code

Calculation2

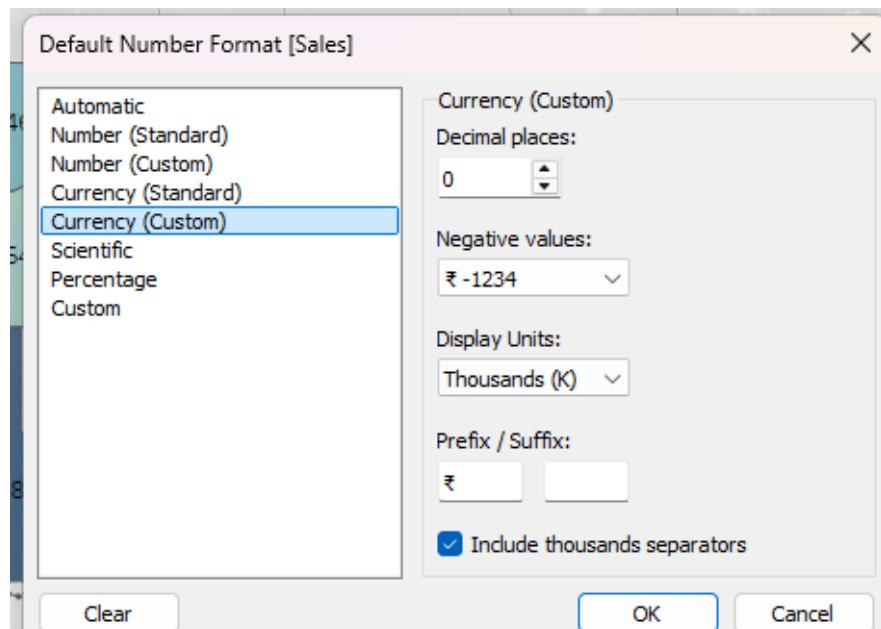
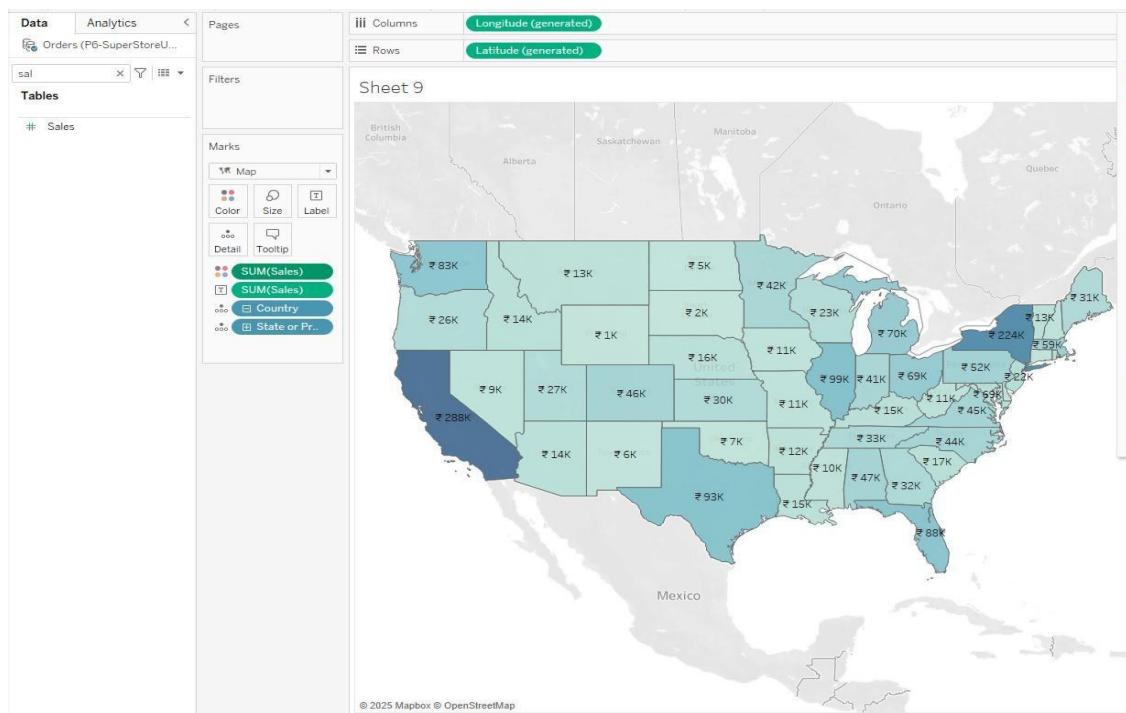
(SUM([Profit])/SUM([Sales]))

The calculation is valid. 1 Dependency ▾ Apply OK



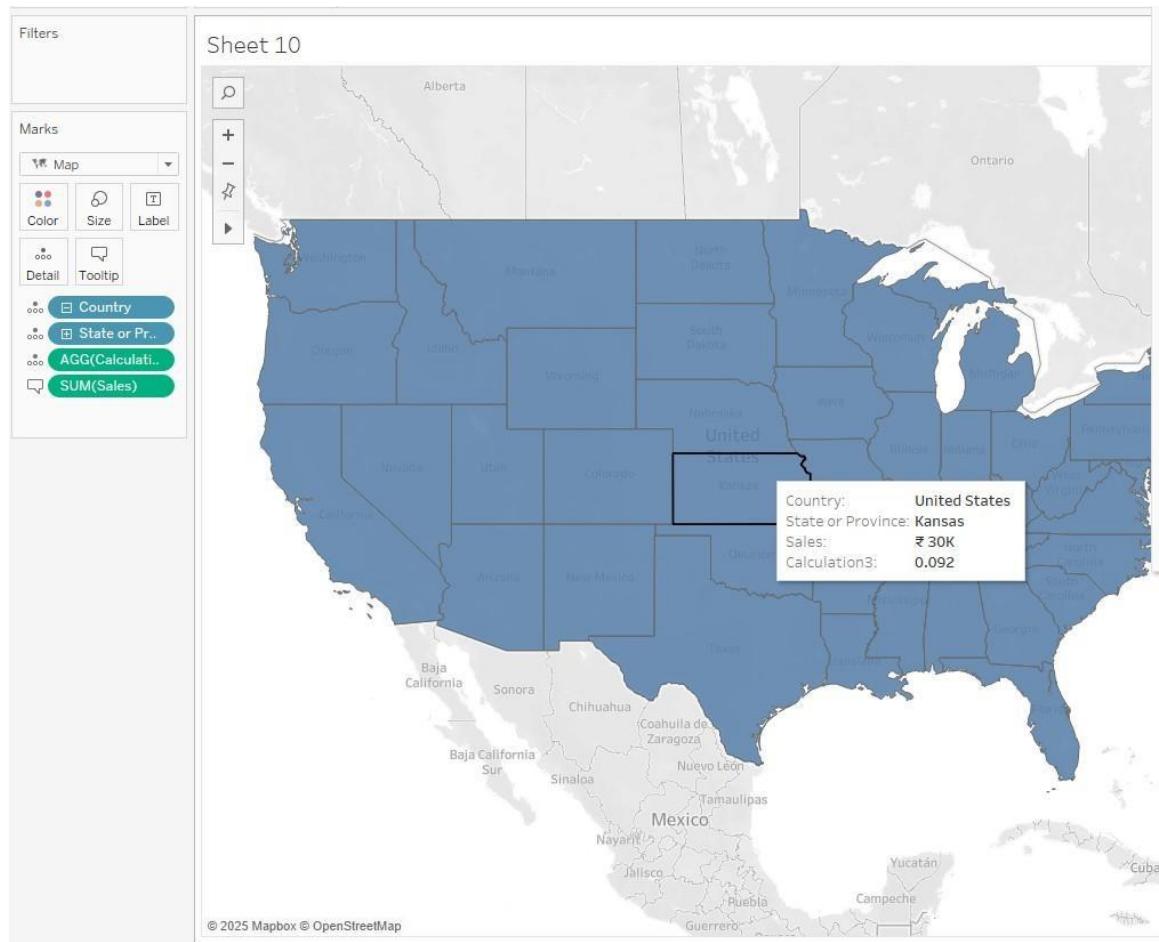
2. Build a Basic Map of Sales by State

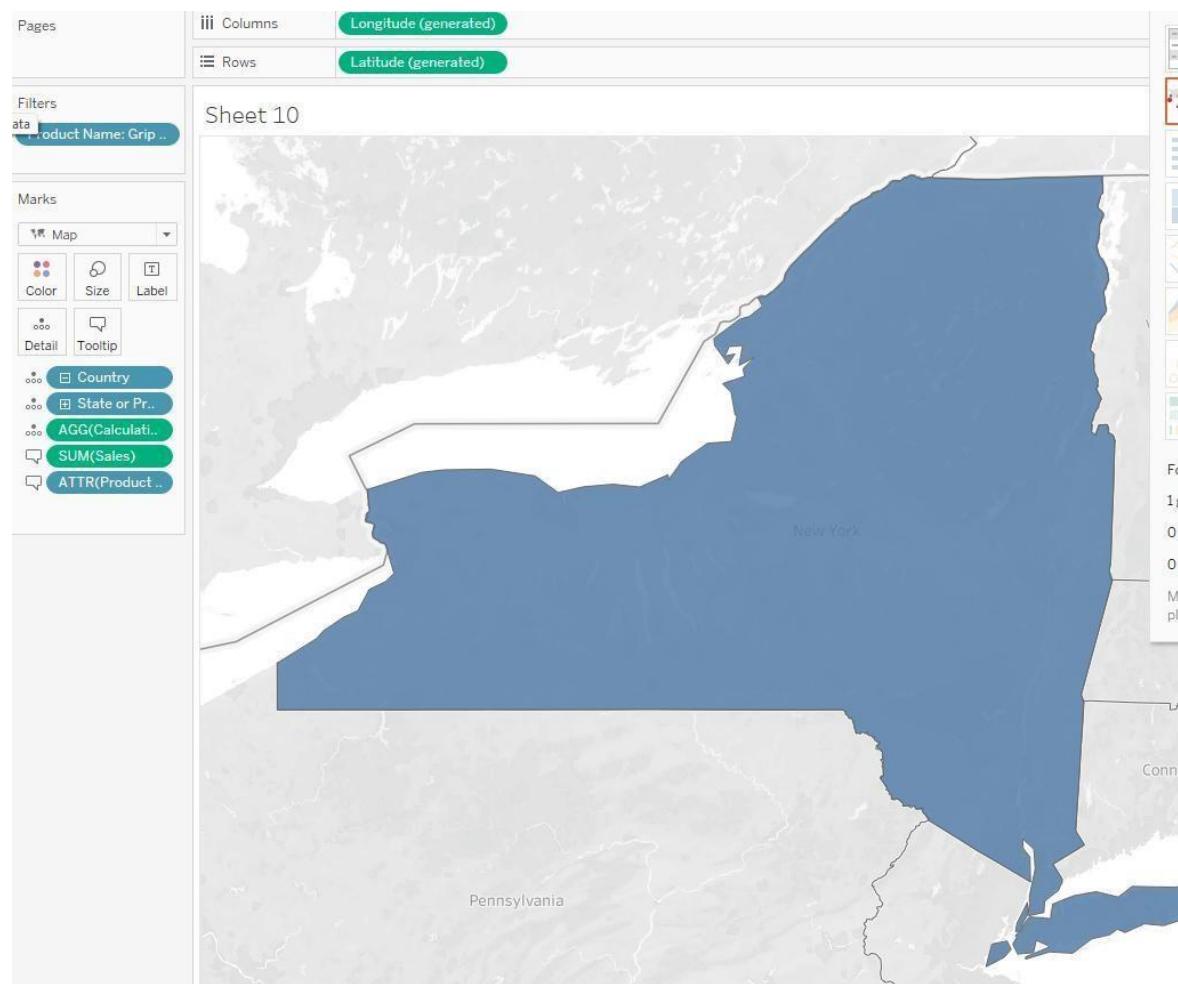
1. Double-click on State → Tableau will create a map.
2. From **Measures**, drag Sales → drop it on **Color** (Marks card).
 - Now each state is colored by its Sales.
3. Drag Sales again → drop it on **Label** (so values appear on the map).
4. Right-click on Sales → **Default Properties** → **Number Format** → **Currency / Custom (No decimals, in K)**.

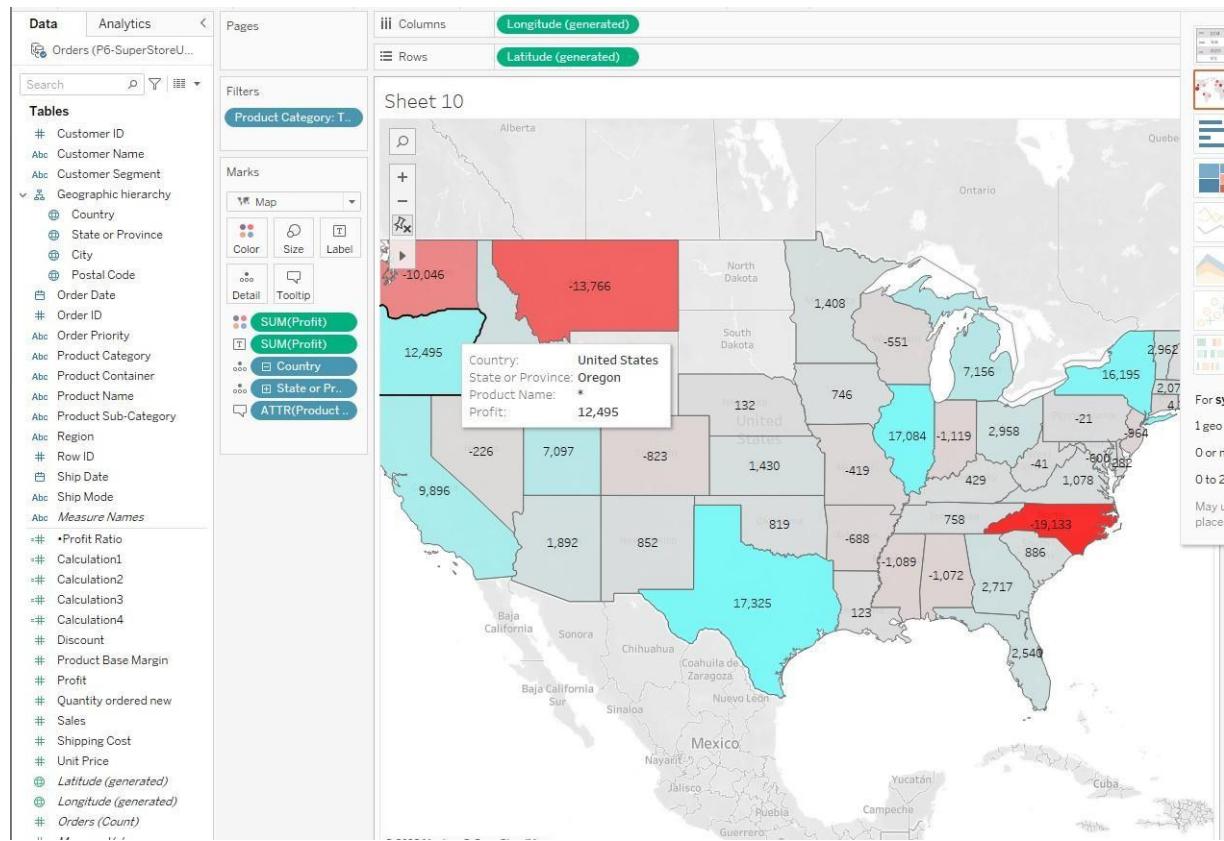


3. Add Profit Ratio as Tooltip

1. Create a calculated field:
2. Profit Ratio = $\text{SUM}([\text{Profit}]) / \text{SUM}([\text{Sales}])$
3. Drag this field onto the **Tooltip** shelf (Marks card).
4. Now when you hover over a state, you will see both Sales and Profit Ratio.
 - Show Profit Ratio of each state as tooltip on map
 - Show Profit ratio for Grip Envelop products
 - Identify Unprofitable States in Technology Surrounded by Profitable Ones (e.g., Nevada)
 1. Drag Category to **Filters** → select only **Technology**.
 2. Drag Profit to **Color** (Marks card).
 - Positive profit → blue/green, Negative profit → red.
 3. Add Profit also to **Label** if you want exact values.
 4. Look at the map → Nevada will appear **red (loss)** but surrounded by **green states (profit)**.







Preparing Reports

Dataset: **Superstore**

- Report: Product category-wise sales.
 - Report: Region-wise product sales.
 - Report: State-wise sales.
 - Example Question: *What is % of total Sales for Home Office segment in 2015?*
 - Example Question: *Find Top 10 Product Names by Sales in each region. Which product is ranked #2 in Central & West in 2015?*

Conclusion: Successfully performed Visualization in Tableau with charts, dashboards, stories, maps, and reports.

Pages iii Columns

Rows State or Province City

Filters

Marks

Automatic

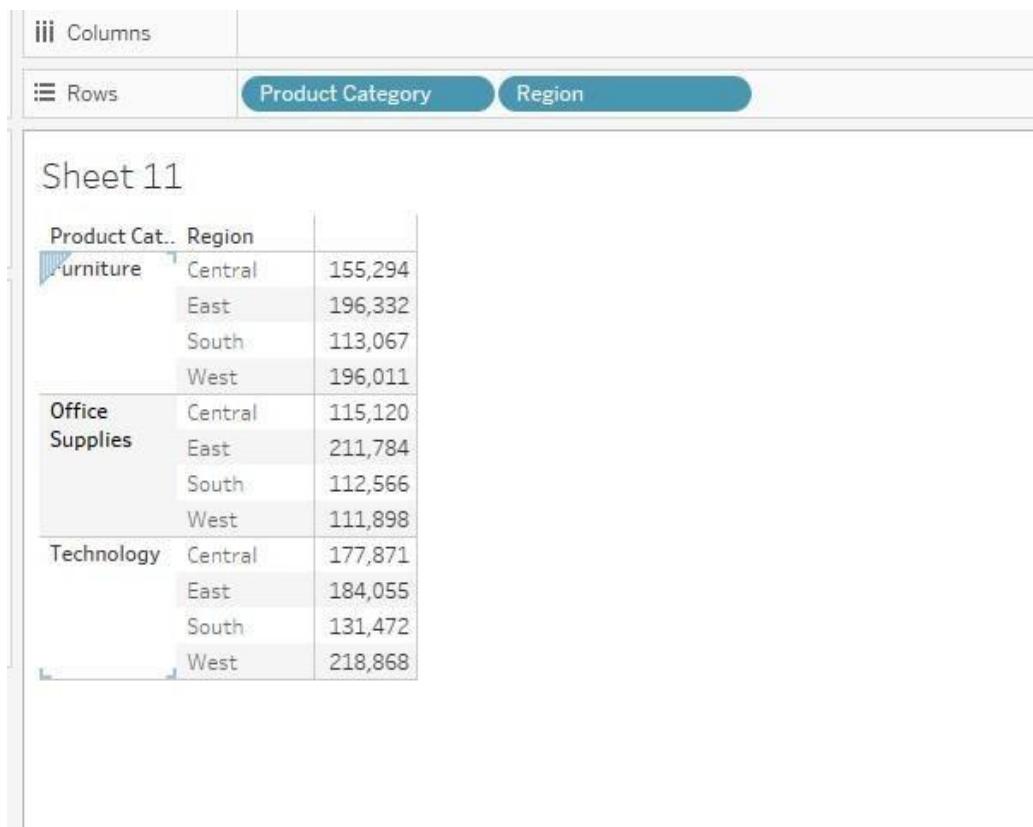
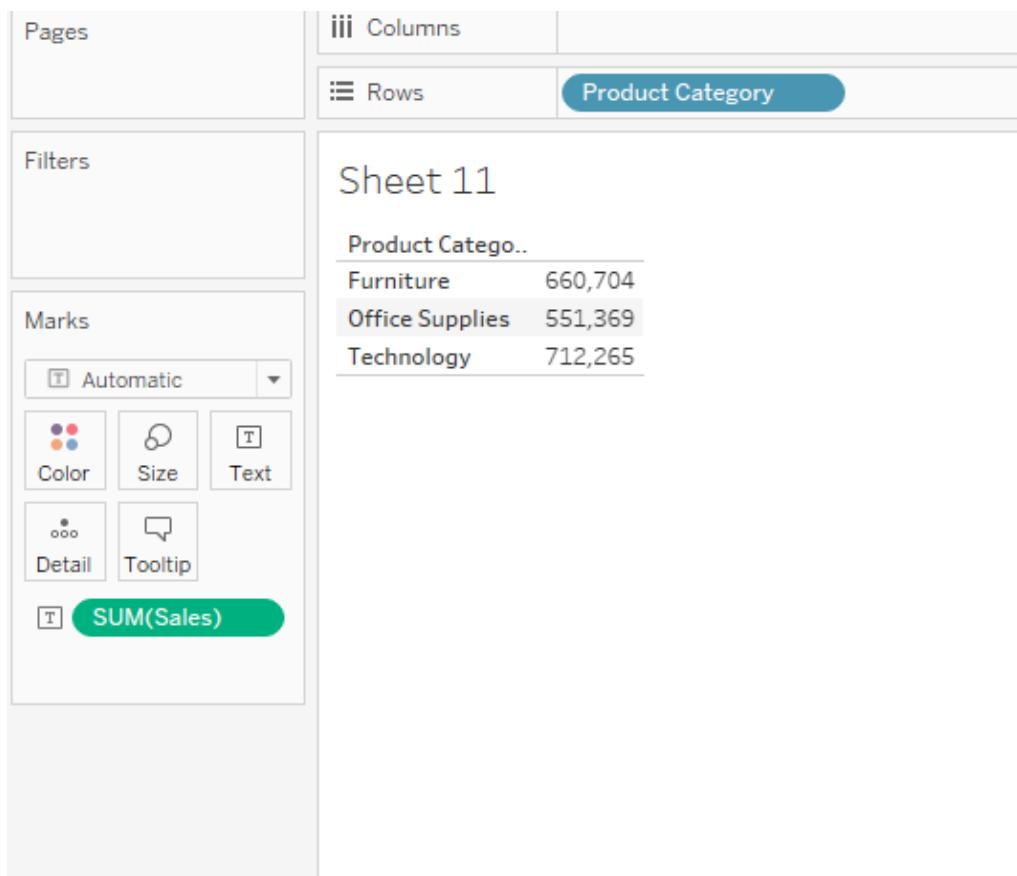
Color Size Text

Detail Tooltip

SUM(Sales)

Sheet 11

State or Province	City	
Alabama	Auburn	10,630
	Birmingham	116
	Decatur	2,874
	Enterprise	443
	Hoover	3,550
	Madison	20,640
	Mobile	2,199
	Northport	128
	Opelika	2,835
	Phenix City	201
Arizona	Tuscaloosa	613
	Vestavia Hills	2,599
	Bullhead City	1,424
	Chandler	2,875
	El Mirage	59
	Flagstaff	229
	Gilbert	1,542
	Glendale	190
	Kingman	1,715
	Mesa	324
Arkansas	Oro Valley	3,762
	Peoria	353
	Prescott	159
	Prescott Valley	1,149
	Scottsdale	55
	Surprise	254
	Tucson	278
	Bryant	658
	Cabot	370
	Hot Springs	130
	Jonesboro	465
	Little Rock	1,095
	North Little Rock	343
	Paragould	168
	Pine Bluff	883
	Rogers	100
	Searcy	2,848



Sheet 11

Product Cat..	Region	State or Province	
Furniture	Central	Illinois	20,376
		Indiana	15,094
		Iowa	2,342
		Kansas	17,074
		Michigan	24,884
		Minnesota	25,998
		Missouri	1,907
		Nebraska	5,835
		North Dakota	68
		Oklahoma	3,401
		South Dakota	131
		Texas	33,026
		Wisconsin	5,160
	East	Connecticut	3,664
		District of Columbia	9,595
		Maine	16,230
		Maryland	8,796
		Massachusetts	16,248
		New Hampshire	4,458
		New Jersey	10,933
		New York	72,852
		Ohio	20,665
		Pennsylvania	19,850
		Vermont	3,116
		West Virginia	9,928
	South	Alabama	4,547
		Arkansas	3,112
		Florida	46,019
		Georgia	12,053
		Kentucky	5,888
		Louisiana	2,860
		Mississippi	2,795
		North Carolina	11,657
		South Carolina	3,106
		Tennessee	7,146
		Virginia	13,883
		Arizona	8,055