

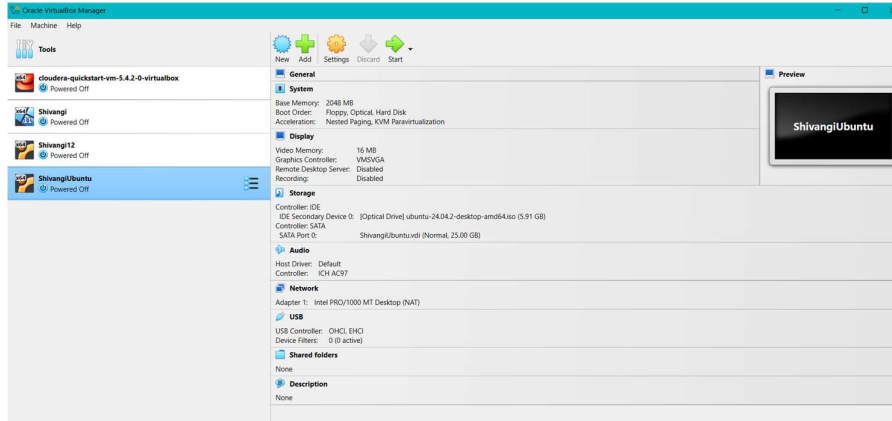
Practical – 2

Steps to install Hadoop

Step 1) From the browser download the Virtual box

link-<https://www.virtualbox.org/wiki/Downloads>

Download the software for the windows hosts



Step 2) Install the HDP Sandbox

visit the cloudera platform & search for the hdp sandbox

Link-<https://www.cloudera.com/>

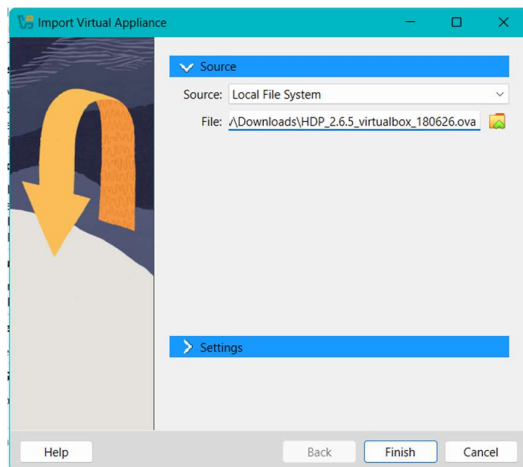
As per new updation since from 2023 the version have been changed so you may get your sandbox through the link provided link below

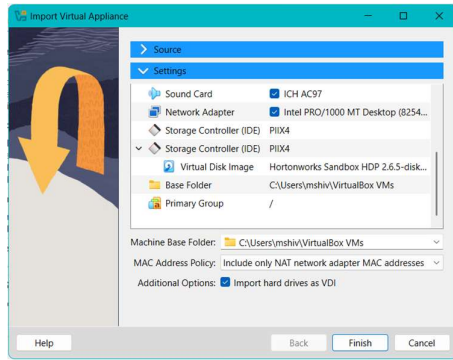
link for hdp sandbox-https://archive.cloudera.com/hwx-sandbox/hdp/hdp-2.6.5/HDP_2.6.5_virtualbox_180626.ova

Step 3) Import the sandbox in virtual box

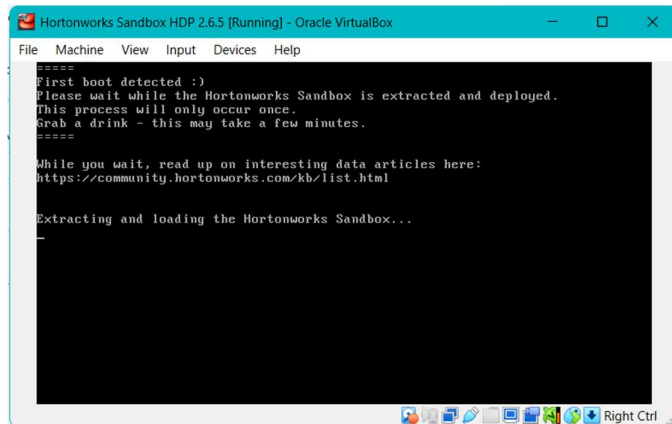
After downloading click on import button on the virtual box and import the file

Hit the start button





Step 4) To visualize what's going on in the Hadoop we may visualize through Ambari
Go & visit on the address provided on the virtual box
Click on the address & launch the dashboard
It requires the username & password
Username – maria_dev
Password- maria_dev



Sign in

Username

Password

[Sign in](#)

Step 5) Small Activity

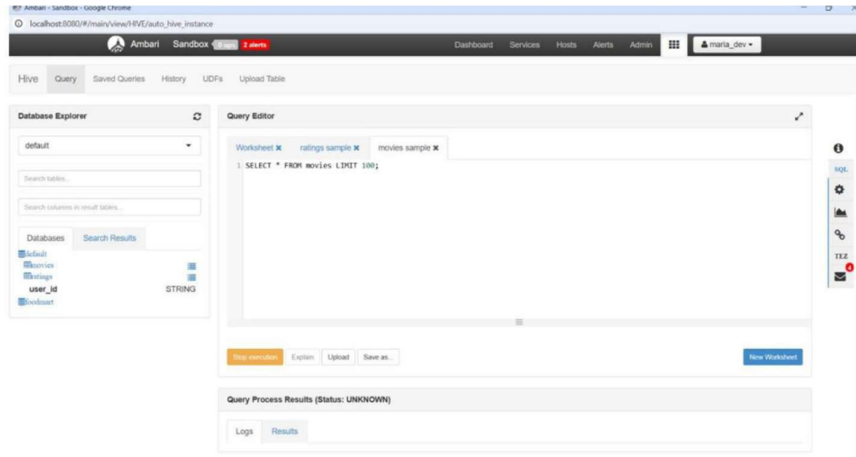
Download the dataset from grouplens

link-<https://grouplens.org/datasets/movielens/>

older dataset is provided on the website

Hit the download button & download ml.100k.zip file

Once you have downloaded extract the data



Step 6) Working on the downloaded dataset

Go into the ambari tool and from the menu go into the hive view

We will import the data from the local file there to import data open the hive view

After hive view click on the upload table option

Select the csv file type & set the file delimiter type to the 9 (i.e horizontal tab)

Choose the file from your local system (i.e u.data file)

Rename the table name-ratings

column name-

user_id

movie_id

rating

rating_time

Hit the upload button

Same for the movie name table

Select the file type as above and set the file delimiter to 124

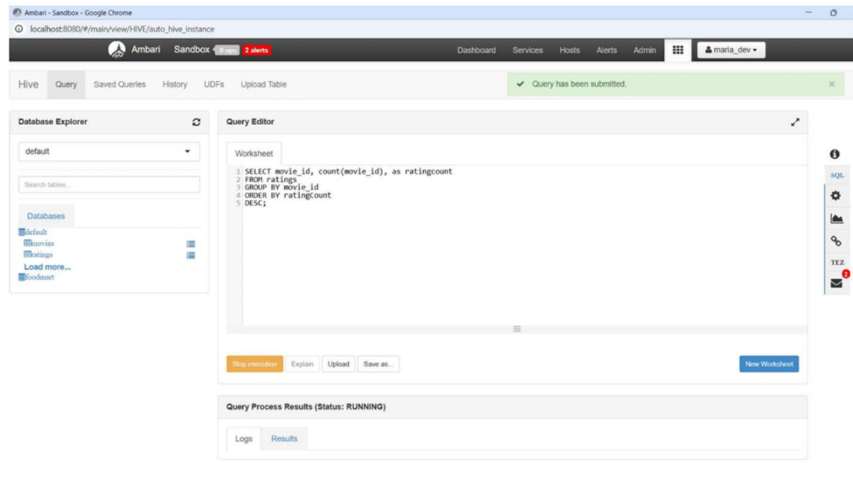
Rename the table name to movie_name

column name-

Movie_id

name

After this hit the upload button



Step 7) Write the SQL Query to perform operations

SQL Query1-

```
SELECT movie_id, count(movie_id), as ratingcount
```

```
FROM ratings
```

```
GROUP BY movie_id
```

```
ORDER BY ratingCount
```

```
DESC
```

After writing the query execute it and see the results

SQL Query2-

```
SELECT name
```

```
FROM movie_names
```

```
WHERE movie_id=50
```

After writing the query execute it and see the results

Ambari - Sandbox - Google Chrome

localhost:8080/#/main/view/HIVE/auto_hive_instance

Query Process Results (Status: SUCCEEDED)

Save results...

LogsResults

Filter columns...

movies.movie_id	movies.movie_name	movies.column3	movies.column4	movies.column5
1	Toy Story (1995)	01-Jan-1995	==	http://us.imdb.com/M/title-exact?Toy%20Story%20(1995)
2	GoldenEye (1995)	01-Jan-1995	==	http://us.imdb.com/M/title-exact?GoldenEye%20(1995)
3	Four Rooms (1995)	01-Jan-1995	==	http://us.imdb.com/M/title-exact?Four%20Rooms%20(1995)
4	Get Shorty (1995)	01-Jan-1995	==	http://us.imdb.com/M/title-exact?Get%20Shorty%20(1995)
5	Copycat (1995)	01-Jan-1995	==	http://us.imdb.com/M/title-exact?Copycat%20(1995)
6	Shanghai Triad (Yao a yao yao dao wai po gao) (1995)	01-Jan-1995	==	http://us.imdb.com/T/Title?Yao+a+yao+yao+dao+wai+po+gao+(1995)
7	Twelve Monkeys (1995)	01-Jan-1995	==	http://us.imdb.com/M/title-exact?Twelve%20Monkeys%20(1995)
8	Babe (1995)	01-Jan-1995	==	http://us.imdb.com/M/title-exact?Babe%20(1995)
9	Dead Man Walking (1995)	01-Jan-1995	==	http://us.imdb.com/M/title-exact?Dead%20Man%20Walking%20(1995)
10	Richard III (1995)	22-Jan-1996	==	http://us.imdb.com/M/title-exact?Richard%20III%20(1995)
11	Seven (Se7en) (1995)	01-Jan-1995	==	http://us.imdb.com/M/title-exact?Se7en%20(1995)
12	Usual Suspects, The (1995)	14-Aug-1995	==	http://us.imdb.com/M/title-exact?Usual%20Suspects,%20The%20(1995)
13	Mighty Aphrodite (1995)	30-Oct-1995	==	http://us.imdb.com/M/title-exact?Mighty%20Aphrodite%20(1995)

Ambari - Sandbox - Google Chrome

localhost:8080/#/main/view/HIVE/auto_hive_instance

Query Process Result

LogsResults

Filter columns...

movies.movie_id	m
1	T
2	G
3	F
4	G
5	C
6	S a q
7	T (1
8	B
9	D (1
10	R
11	S (1
12	U (1