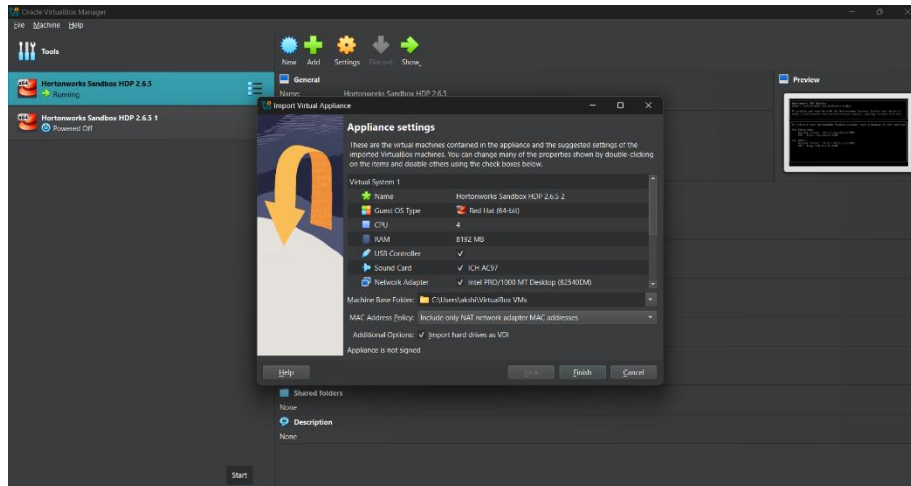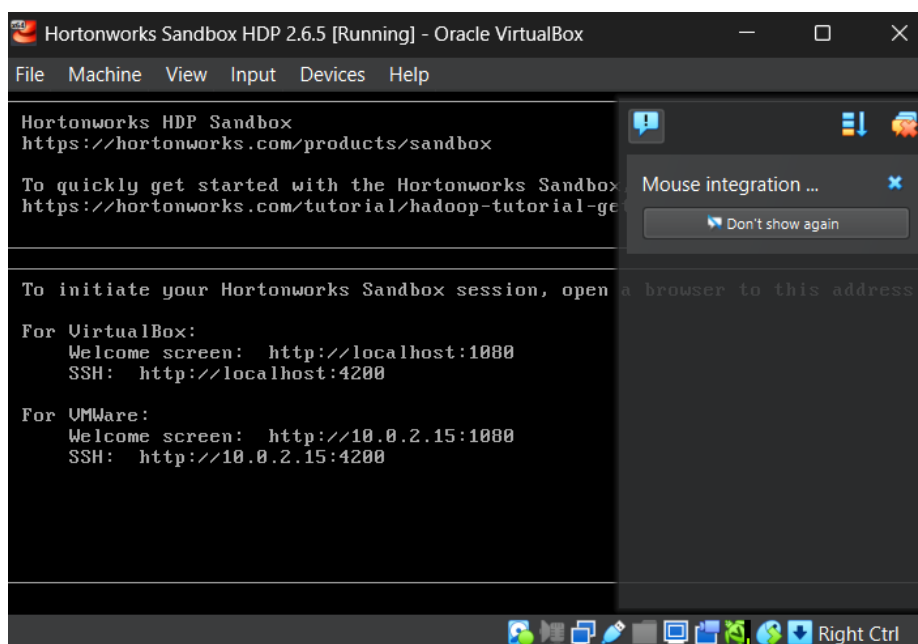# Practical-02

**Aim:- Install and configure Hadoop using HDP Sandbox in a VirtualBox, and perform Hive queries on a dataset.**
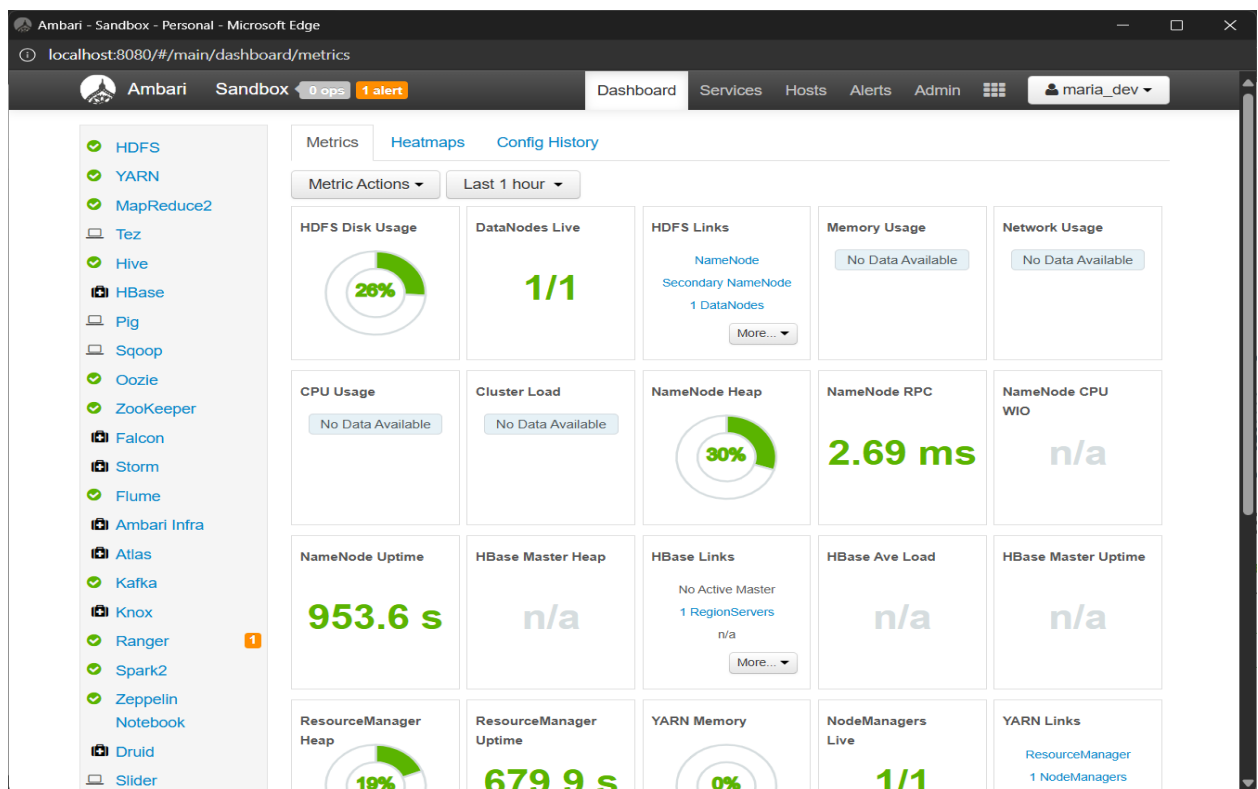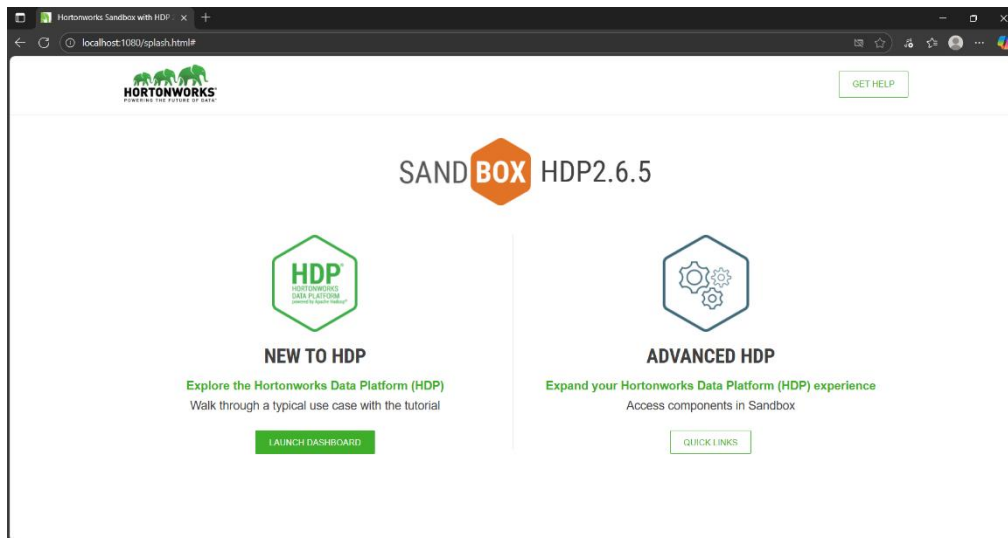
1. Download HDP Sandbox and import into virtualbox.



2. Start the Sandbox VM , it will display the URL to access Ambari.



3. The URL will open Ambari Dashboard , click on launch dashboard it will open a login page, login using:
   - Username : maria_dev
   - Password : maria_dev

4. **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

5. Go into the Ambari tool and from the menu go into the hive view, 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
column name-
movie_id
name
After this hit the upload button

6.  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

Ambari    Sandbox  0 ops   3 alerts          Dashboard    Services    Hosts    Alerts    Admin          maria_dev

Hive    Query    Saved Queries    History    UDFs    Upload Table

**Database Explorer**

default

Search tables...

Databases

default
foodmart

**Query Editor**

Worksheet

```
1 SELECT name
2 FROM movie
3 WHERE movie_id = 50;
```

SQL

TEZ

Execute    Explain    Upload    Save as...                New Worksheet

**Query Process Results (Status: SUCCEEDED)**        Save results...

Logs    Results

---

Execute    Explain    Upload    Save as...                New Worksheet

**Query Process Results (Status: SUCCEEDED)**        Save results...

Logs    Results

Filter columns...                                    previous    next

**name**

Star Wars (1977)

Licensed under the Apache License, Version 2.0.
See third-party tools/resources that Ambari uses and their respective authors