## Question 1

- Import instructions
  - Transferring the files from host OS to HDP access node using the pspc utility in Windows Command Prompt:
    ```
    pscp -P 2222 -pw hadoop Screentime_App_Details.csv
    root@127.0.0.1:/home/cind719/
    pscp -P 2222 -pw hadoop Screentime_App_Ranking.csv
    root@127.0.0.1:/home/cind719/
    pscp -P 2222 -pw hadoop Screentime_Overall_Usage.csv
    root@127.0.0.1:/home/cind719/
    ```

    ```
    Command Prompt

    Microsoft Windows [Version 10.0.19043.2006]
    (c) Microsoft Corporation. All rights reserved.

    C:\Users\annsa>cd c:\data

    c:\data>pscp -P 2222 -pw hadoop Screentime_App_Details.csv root@127.0.0.1:/home/cind719/
    Screentime_App_Details.cs | 1 kB |   1.7 kB/s | ETA: 00:00:00 | 100%

    c:\data>pscp -P 2222 -pw hadoop Screentime_App_Ranking.csv root@127.0.0.1:/home/cind719/
    Screentime_App_Ranking.cs | 0 kB |   1.0 kB/s | ETA: 00:00:00 | 100%

    c:\data>pscp -P 2222 -pw hadoop Screentime_Overall_Usage.csv root@127.0.0.1:/home/cind719/
    Screentime_Overall_Usage. | 0 kB |   0.7 kB/s | ETA: 00:00:00 | 100%

    c:\data>
    ```

  - Creating a new directory using the HDP terminal and moving the files to the HDFS:
    ```
    hadoop fs -mkdir /user/assignment1
    hadoop fs -put /home/cind719/Screentime_App_Details.csv
    /user/assignment1
    hadoop fs -put /home/cind719/Screentime_App_Ranking.csv
    /user/assignment1
    hadoop fs -put /home/cind719/Screentime_Overall_Usage.csv
    user/assignment1
    ```

    ```
    root@sandbox:~

    [root@sandbox ~]# ll /home/cind719/
    total 61168
    -rw-r--r-- 1 root root      115 2022-10-08 14:43 dayofweek.txt
    -rw-r--r-- 1 root root 57016655 2022-10-07 15:29 full_text.txt
    -rw-r--r-- 1 root root     1713 2022-10-08 17:57 Screentime_App_Details.csv
    -rw-r--r-- 1 root root     1009 2022-10-08 17:58 Screentime_App_Ranking.csv
    -rw-r--r-- 1 root root      678 2022-10-08 17:58 Screentime_Overall_Usage.csv
    -rw-r--r-- 1 root root  5589917 2022-10-07 17:55 shakespeare.txt
    -rw-r--r-- 1 root root      323 2022-10-07 16:22 wc_mapper.py
    -rw-r--r-- 1 root root      686 2022-10-07 16:23 wc_reducer.py
    [root@sandbox ~]# hadoop fs -mkdir /user/assignment1
    [root@sandbox ~]# hadoop fs -ls /user
    Found 12 items
    drwxrwx---   - ambari-qa hdfs          0 2015-04-24 12:49 /user/ambari-qa
    drwxr-xr-x   - root      hdfs          0 2022-10-09 16:13 /user/assignment1
    drwxr-xr-x   - guest     guest         0 2015-04-24 13:32 /user/guest
    drwxr-xr-x   - hcat      hdfs          0 2015-04-24 13:13 /user/hcat
    drwx------   - hive      hdfs          0 2015-04-24 13:06 /user/hive
    drwxr-xr-x   - hue       hue           0 2015-04-24 13:32 /user/hue
    drwxr-xr-x   - root      hdfs          0 2022-10-08 14:51 /user/lab
    drwxrwxr-x   - oozie     hdfs          0 2015-04-24 13:10 /user/oozie
    drwx------   - root      hdfs          0 2022-10-07 18:08 /user/root
    drwxr-xr-x   - solr      hdfs          0 2015-04-24 13:25 /user/solr
    drwxrwxr-x   - spark     hdfs          0 2015-04-24 12:59 /user/spark
    drwxr-xr-x   - yarn      yarn          0 2015-04-24 13:33 /user/yarn
    [root@sandbox ~]# hadoop fs -put /home/cind719/Screentime_App_Details.csv /user/assignment1
    [root@sandbox ~]# hadoop fs -put /home/cind719/Screentime_App_Ranking.csv /user/assignment1
    [root@sandbox ~]# hadoop fs -put /home/cind719/Screentime_Overall_Usage.csv /user/assignment1
    [root@sandbox ~]# hadoop fs -ls /user/assignment1
    Found 3 items
    -rw-r--r--   1 root hdfs     1713 2022-10-09 16:14 /user/assignment1/Screentime_App_Details.csv
    -rw-r--r--   1 root hdfs     1009 2022-10-09 16:14 /user/assignment1/Screentime_App_Ranking.csv
    -rw-r--r--   1 root hdfs      678 2022-10-09 16:15 /user/assignment1/Screentime_Overall_Usage.csv
    [root@sandbox ~]#
    ```

- Database and table creation scripts:

```
CREATE DATABASE screentime;
USE screentime;
set hive.cli.print.current.db=true;

CREATE TABLE detail_app(date string, usage string, notifications
string, time_opened string, app string) ROW FORMAT DELIMITED FIELDS
TERMINATED BY ',' TBLPROPERTIES('skip.header.line.count'='1');

CREATE TABLE detail_ranking(date string, rank1 string, rank2 string,
rank3 string) ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
TBLPROPERTIES('skip.header.line.count'='1');

CREATE TABLE overall_usage(date string, total_usage string,
notifications string, unlocks string) ROW FORMAT DELIMITED FIELDS
TERMINATED BY ',' TBLPROPERTIES('skip.header.line.count'='1');
```
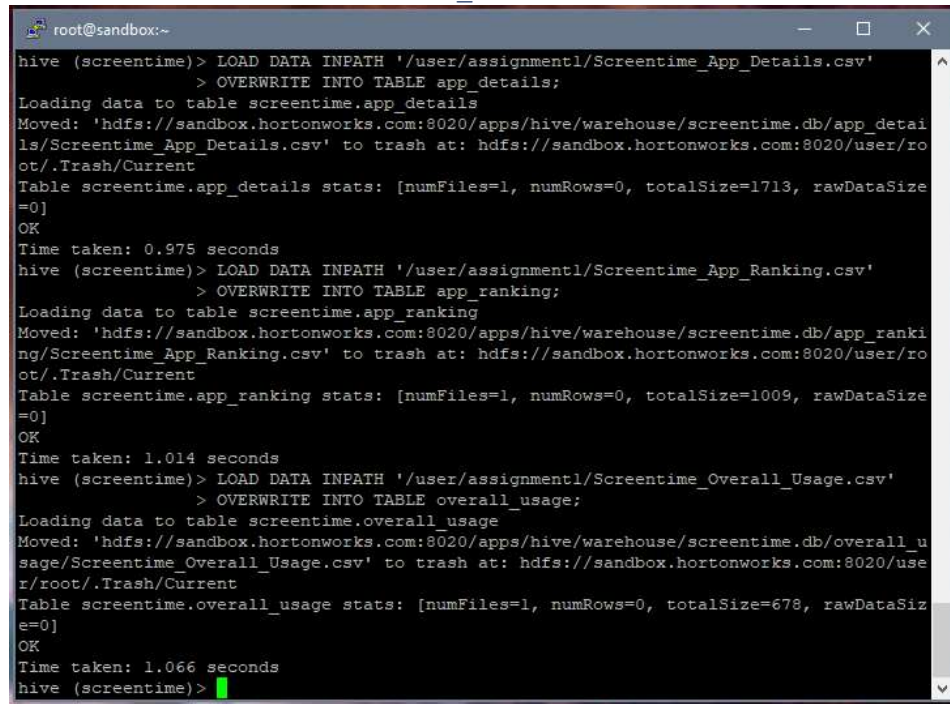
- Importing csv files into tables:
  ```
  LOAD DATA INPATH '/user/assignment1/Screentime_App_Details.csv'
  OVERWRITE INTO TABLE app_details;

  LOAD DATA INPATH '/user/assignment1/Screentime_App_Ranking.csv'
  OVERWRITE INTO TABLE app_ranking;

  LOAD DATA INPATH '/user/assignment1/Screentime_Overall_Usage.csv'
  OVERWRITE INTO TABLE overall_usage;
  ```



- Displaying first 5 rows from each table:
  ```
  SELECT * FROM app_details LIMIT 5;
  SELECT * FROM app_ranking LIMIT 5;
  SELECT * FROM overall_usage LIMIT 5;
  ```

**Question 2**

Find the unique number of applications in Rank 1, Rank 2 and Rank 3 columns. You can write 3 different HIVE commands for each Rank column.

- Unique number of applications in Rank 1: 4
  Unique number of applications in Rank 2: 10
  Unique number of applications in Rank 3: 10

```
SELECT COUNT(DISTINCT rank1) FROM app_ranking;
SELECT COUNT(DISTINCT rank2) FROM app_ranking;
SELECT COUNT(DISTINCT rank3) FROM app_ranking;
```

**Question 3**

Which application has maximum usage?

- Application with the maximum usage: Whatsapp

```
SELECT SUM(ad.usage) as su, ad.app FROM app_details as ad GROUP BY ad.app
ORDER BY su DESC LIMIT 1;
```



**Question 4**

Which application has least no of Notifications?

- Application with the least number of Notification: Instagram

```
SELECT SUM(ad.notifications) as sn, ad.app FROM app_details as as GROUP BY
ad.app ORDER BY sn ASC LIMIT 1;
```

## Question 5

Find the average number of times an Instagram app has been opened up?

- ▪ Average number of times Instagram app is opened up: 32.7 times

```
CREATE VIEW q5 AS SELECT AVG(ad.times_opened), ad.app FROM app_details as ad
GROUP BY ad.app;

SELECT * FROM q5 WHERE apps = 'Instagram';
```

## Question 6

Print total usage for WhatsApp application with dates?

```
SELECT ad.date, ad.usage, ad.app FROM app_details as ad WHERE app =
'Whatsapp';
```