# Course name

# Trivedi Maharshi Mayankbhai

# Student ID

<u>**General Theoretical Questions**</u>

**Question 1: List all big data-specific formats that you know.**

**Answer 1:** There are many file formats, which are being used to deal big-data. Plain text file, Sequence file, Avro, Parquet and ORC. Below a short description is provided for each big-data type:

*Plain text file*: Data is saved in *txt* or *csv* formats, they are used in the non-Hadoop environment as well. Heavy disk space/storage can be occupied by these plain text file formats. Thus, a robust compression is required.

*Sequence file*: This file type was mainly developed for MapReduce. Each entry of the data is encoded with a key and a value. Hence the data is stored in a binary format and occupies lesser storage space on the disk than the file formats which are text-based. It is also compatible with block-level compression, which is one of the advantages of the Sequence file.

*Avro file*: While dealing with the big-data files, Avro format is a good choice. On top of a normal file-format, it also supports serialization-deserialization and blocks compression. It is a row based and a splittable format. Avro is only a machine-readable file format.

*Parquet* and *ORC file*: ORC stands for Optimized Row Columnar. Unlike the Avro file, it also stores the data in columnar format. Hence, the horizontal and vertical partition of the data is possible. However, these file formats are also machine-readable only. Because these file formats are columnar, they outperform in terms of storage optimization than any other file formats.

**Question 2: Why the compression of data matters for Hadoop?**

**Answer 2:** Necessity of compression in a typical Hadoop eco-system is because of the huge volumes of data, it must deal with. Reduction in- (1) space needed for storing data and (2) data transfer speed to or from the disk, could be two significant advantages while dealing with big-data. These two advantages could be introduced with data compression. For example, MapReduce job for a compressed-large volume results into a low-latency task.

However, the data compression rate must be tuned with a trade-off between compression and speed of computation. Compressed data first gets decompressed before any other operation on the cluster. Thus, it increases CPU utilization along with the compression. The more compression is set for data, the more resources are used to first decompress it.

<u>**YARN application/commands**</u>

**Question 3: What is YARN? What are YARN's two most important functions?**

**Answer 3:** YARN stands for Yet Another Resource Negotiator. YARN sits in between of Hadoop File System (HDFS) and processing layer. It is essentially used for the two most important function, mentioned below:

   (1) Resource Management: YARN allocates resources such as memory, to the application.

(2) Job Scheduling: It supports multiple scheduling methods for submitting the job to be processed in sequence. Some of the state-of-the-art scheduling methods are FIFO, Capacity Scheduler and Fair Scheduler.

## Question 4: List all running applications

**Answer 4:** The Running application can be visualized with both: YARN UI and HDFS terminal. Figure 1 shows a screenshot of all the running application with YARN UI. Figure 2 lists all the running application with HDFS terminal.



Figure 1. Running the application with YARN UI.



Figure 2. List of all the running application with HDFS terminal (Highlighted yellow is the command).

## HDFS commands

### Question 5: Create a folder/directory 'Lab1_results' in your own HDFS directory.

**Answer 5:** Command used to list the directories/files in a home folder, before generating *Lab1_results* directory: *hdfs dfs -ls*. Figure 3 and Figure 4 shows the listing of all the directories before generating Lab1_results directory. All the commands are implemented in HDFS terminal only.
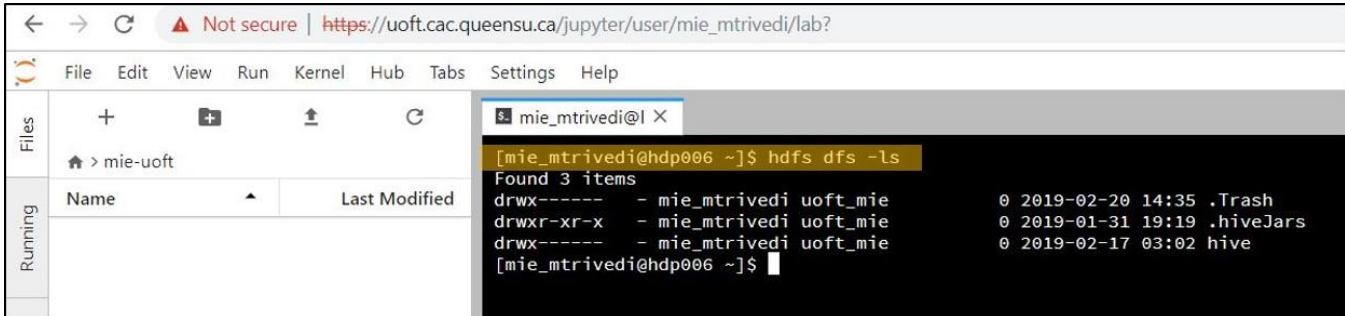


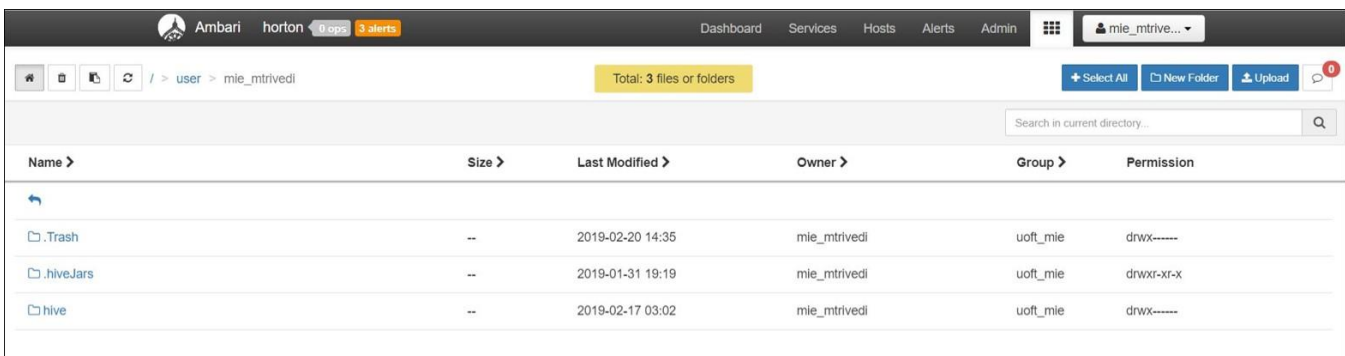Figure 3. Listing of all the directories/files in present HDFS directory with HDFS terminal.



Figure 4. Listing of all the directories/files with HDFS UI.

Directory *Lab1_results* is generated with *hdfs dfs -mkdir Lab1_results* in the present working directory. Figure 5 shows the usage of *mkdir* command to generate *Lab1_results* along with the listing of all the directories after generating *Lab1_results*.
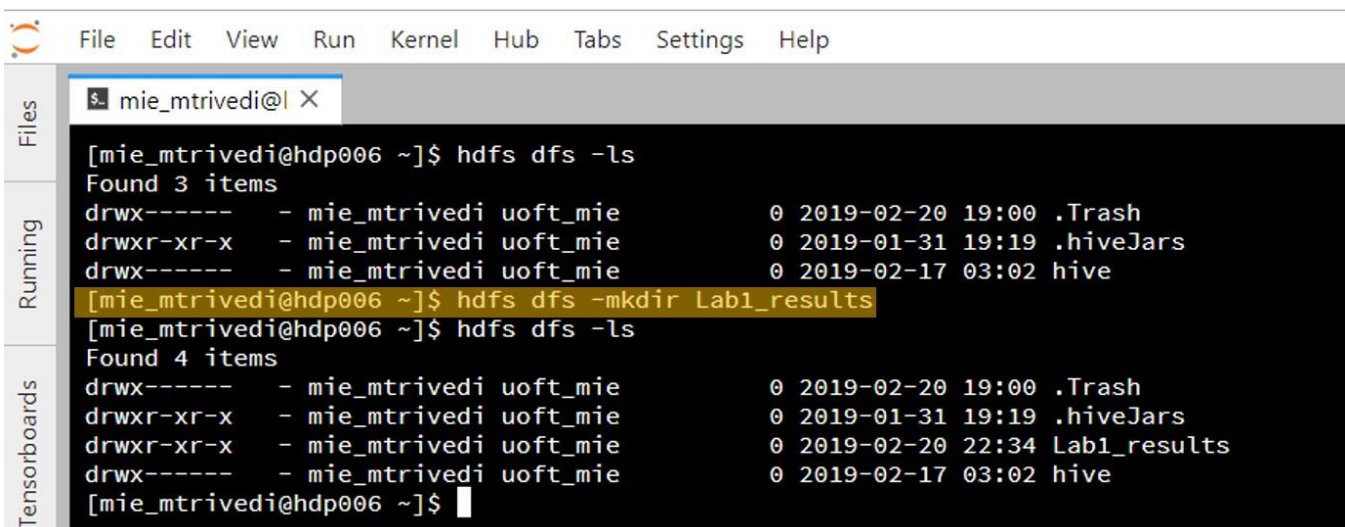


Figure 5. Command to generate a directory *Lab1_results*, with *mkdir* command in HDFS.

Figure 6. Checking HDFS UI for Lab1_results after generating it with HDFS command.

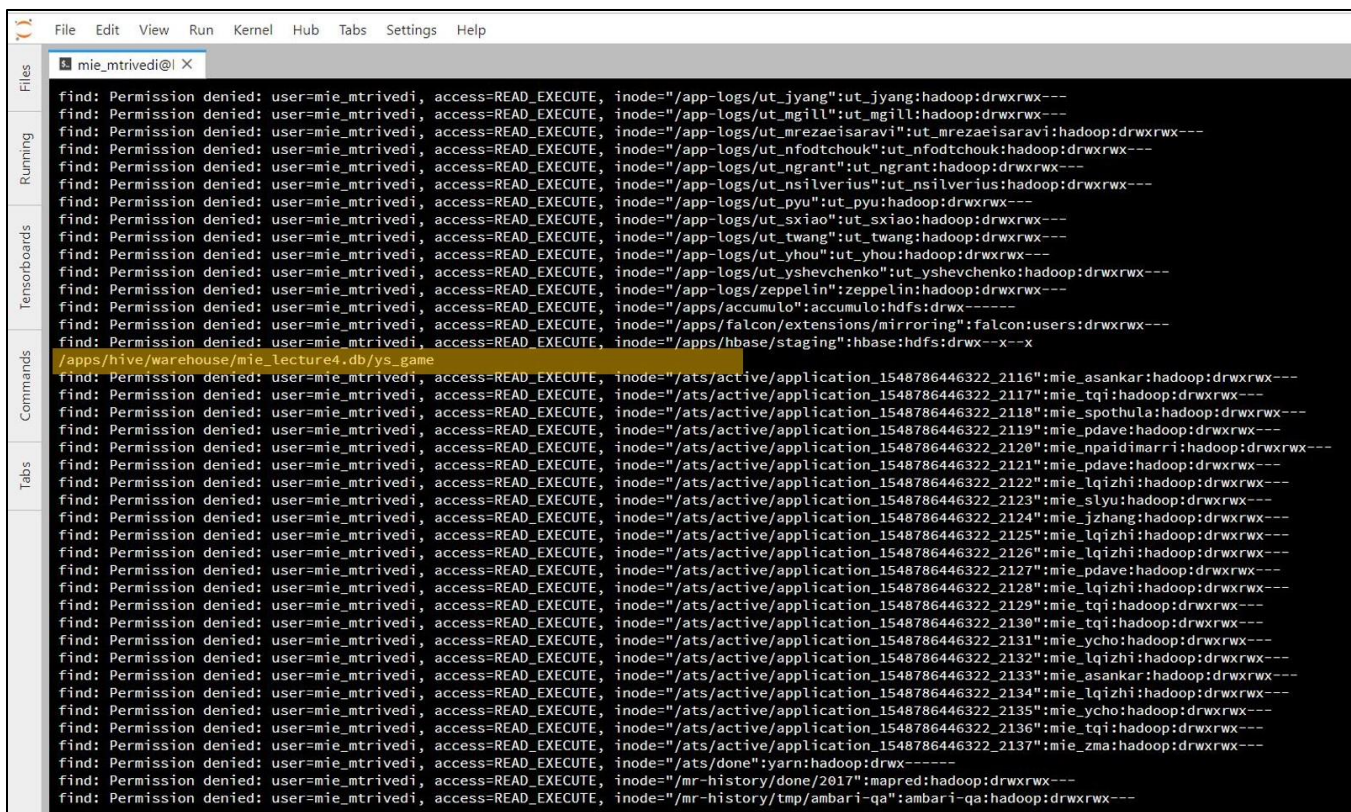## Question 6: Where in HDFS is MIE_Lecture4.ys_game file located? Provide path of the file.

**Answer 6:** Command *hdfs dfs -find / -name ys_game* could be used to find the location of this file with respect to the root directory. Figure 7 shows the usage of *find* command in HDFS terminal.

**Path of the file ys_game:** /apps/hive/warehouse/mie_lecture4.db/ys_game



Figure 7. Command to locate *ys_game* file in HDFS terminal.

The output of *find* command lists the location of all the files with respect to its root directory. Scrolling down the output terminal window, *ys_game* could be located. Figure 8 shows the location of *ys_game* with the highlighted field.

Figure 8. Global path of *ys_game* is highlighted with yellow color.

## Question 7: What format is underlying Hive tables saved in? How can you find the format?

**Answer 7:** Through HDFS, entry to the Hive environment is made.



Figure 9. Entering the Hive with HDFS.

After entering the Hive, all the tables are listed with *show tables*.



Figure 10. Listing of all the tables after entering in Hive environment with HDFS.

Once all the tables are listed, any of them is picked up and file format is checked with *desc formatted Table_name*. Here, table game is chosen for checking the file format and thus the command *desc formatted game* is used.



```
File   Edit   View   Run   Kernel   Hub   Tabs   Settings   Help

  mie_mtrivedi@I ×

hive> desc formatted game;
OK
# col_name                 data_type                  comment

game_id                    int
season                     int
type                       char(2)
date_time                  date
away_team_id               int
home_team_id               int
away_goals                 int
home_goals                 int
outcome                    string
home_rink_side_start       string
venue                      string
venue_link                 string
venue_time_zone_id         string
venue_time_zone_offset     int
venue_time_zone_tz         string

# Detailed Table Information
Database:                  default
Owner:                     mie_mchandraseka
CreateTime:                Thu Jan 31 19:30:30 EST 2019
LastAccessTime:            UNKNOWN
Protect Mode:              None
Retention:                 0
Location:                  hdfs://hdp001.cac.queensu.ca:8020/apps/hive/warehouse/game
Table Type:                MANAGED_TABLE
Table Parameters:
        COLUMN_STATS_ACCURATE      {\"BASIC_STATS\":\"true\"}
        numFiles                   1
        numRows                    7441
        rawDataSize                5513781
        totalSize                  71541
        transient_lastDdlTime      1548981055

# Storage Information
SerDe Library:             org.apache.hadoop.hive.ql.io.orc.OrcSerde
InputFormat:               org.apache.hadoop.hive.ql.io.orc.OrcInputFormat
OutputFormat:              org.apache.hadoop.hive.ql.io.orc.OrcOutputFormat
```

Figure 11. The file format is checked for the *game* table. The file format is bounded with a rectangle which is shown with yellow color.

Above described commands can be used to check the underlying format of all the files. As it can be seen from Figure 11, the file format of the *game* is ORC (Optimized Row Columnar).

## Hive functionality, file manipulation in Hive via Ambari

**Question 8: Calculate the average score of the away team, and the home team for every season. Show results by ordering by the season.**

**Hive Query**

-- Seasonwise average score for both, home and away teams

select season,

avg(home_team_id) Average_score_away,

avg(home_goals) Average_score_home

from game

group by season

sort by season;

-- End of the query

**Result**

| season | average_score_away | average_score_home |
|---|---|---|
| 20122013 | 16.588089330024815 | 2.8647642679900074 |
| 20132014 | 16.804232804232804 | 2.88435374149965987 |
| 20142015 | 17.589082638362395 | 2.83775587566633814 |
| 20152016 | 17.649507948523844 | 2.805450416351249 |
| 20162017 | 17.538344722854973 | 2.915717539863326 |
| 20172018 | 19.206642066420663 | 3.1202952029520294 |

**Question 9: For every season, find the highest scoring home team. Include the score in your answer.**

**Hive Query**

-- List all the team ids which scored highest goals in a season

create temporary table if not exists season_go_max as

select season,

max(home_goals) goals_max

from game

group by season;

select game.season,game.home_team_id,game.home_goals

from game

join season_go_max

on game.season=season_go_max.season and game.home_goals=season_go_max.goals_max

-- End of the query

**Result**

| game.season | game.home_team_id | game.home_goals |
|---|---|---|
| 20122013 | 5 | 8 |
| 20122013 | 14 | 8 |
| 20132014 | 24 | 9 |
| 20132014 | 28 | 9 |
| 20122013 | 17 | 8 |
| 20142015 | 52 | 8 |
| 20142015 | 26 | 8 |
| 20142015 | 5 | 8 |
| 20152016 | 2 | 8 |
| 20172018 | 8 | 10 |
| 20162017 | 8 | 10 |
| 20152016 | 24 | 8 |
| 20162017 | 29 | 10 |
| 20172018 | 16 | 10 |

**Question 10: Provide team_id and name for the team that played as home team at TD Garden.**

**Hive Query**

-- Team id and Team name which played at TD Garden location

select team_id,teamname from team_info

where team_info.team_id in (select distinct home_team_id from game

where venue='TD Garden')

-- End of the query

**Result**

| team_id | teamname |
|---------|----------|
| 6 | Bruins |

**Question 11: Create a new table that lists all games that happened at TD Garden or Madison Square Garden. Add a column that summarizes away and home goals.**

**Hive Query**

-- Summarize goals of home and away team which played at TD Garden or Medison Square Garden

create table if not exists mie_mtrivedi_goals as

select *, concat('Away Goals + Home Goals: ',away_goals+home_goals) as goal_summary

from game

where venue='TD Garden' or venue='Madison Square Garden';


select * from mie_mtrivedi_goals

-- End of the query

**Result**

| mie_mtrivedi_goals.venue_time_zone_id | mie_mtrivedi_goals.venue_time_zone_offset | mie_mtrivedi_goals.venue_time_zone_tz | mie_mtrivedi_goals.goal_summary |
|---|---|---|---|
| America/New_York | -4 | EDT | Away Goals + Home Goals: 5 |
| America/New_York | -4 | EDT | Away Goals + Home Goals: 7 |
| America/New_York | -4 | EDT | Away Goals + Home Goals: 3 |
| America/New_York | -4 | EDT | Away Goals + Home Goals: 7 |
| America/New_York | -4 | EDT | Away Goals + Home Goals: 4 |
| America/New_York | -4 | EDT | Away Goals + Home Goals: 3 |
| America/New_York | -4 | EDT | Away Goals + Home Goals: 1 |
| America/New_York | -4 | EDT | Away Goals + Home Goals: 7 |
| America/New_York | -4 | EDT | Away Goals + Home Goals: 7 |
| America/New_York | -4 | EDT | Away Goals + Home Goals: 1 |
| America/New_York | -4 | EDT | Away Goals + Home Goals: 4 |

**Question 12: Use subquery and count unique team names of teams that played as away team at TD Garden center and scored > 6. Provide the code and the answer.**

**Hive Query**

-- Use of subquery and counting unique team names of teams that played as away teams at TD Garden center and scored >6

select count(distinct(team_id)) as Distinct_count_teams from team_info

where team_id in (select away_team_id from game

where away_goals>6 and venue='TD Garden')

-- End of the query

**Result**

| distinct_count_teams |
|---|
| 2 |