# CSP 554 Big Data Technologies

# Assignment – #4

# Shiva Sankar Modala(A20517528)

Configure a Hadoop (AWS EMR) environment as you have done for previous assignments.



Download the file TestDataGen.class from the Blackboard. It is one of the attachments to the assignment.

**I have downloaded the required files from the blackboard**

scp the file over to the home directory (/home/hadoop) on your Hadoop VM



Log on to your VM using ssh and execute the file using "java TestDataGen"

This will output a magic number which you should copy down and provide with the results of your assignment.

Magic Number = 65072

When we execute TestDataGen.class file, foodplace<magicnumber>.txt and foodratings<magicnumber>.txt are automatically created.

```
[hadoop@ip-172-31-1-178 ~]$ ls
hql  hql.zip  __MACOSX  TestDataGen.class
[hadoop@ip-172-31-1-178 ~]$ java TestdataGen
Error: Could not find or load main class TestdataGen
[hadoop@ip-172-31-1-178 ~]$ java TestDataGen
Magic Number = 65072
[hadoop@ip-172-31-1-178 ~]$
[hadoop@ip-172-31-1-178 ~]$
[hadoop@ip-172-31-1-178 ~]$ ls
foodplaces65072.txt  foodratings65072.txt  hql  hql.zip  __MACOSX  TestDataGen.class
[hadoop@ip-172-31-1-178 ~]$
```
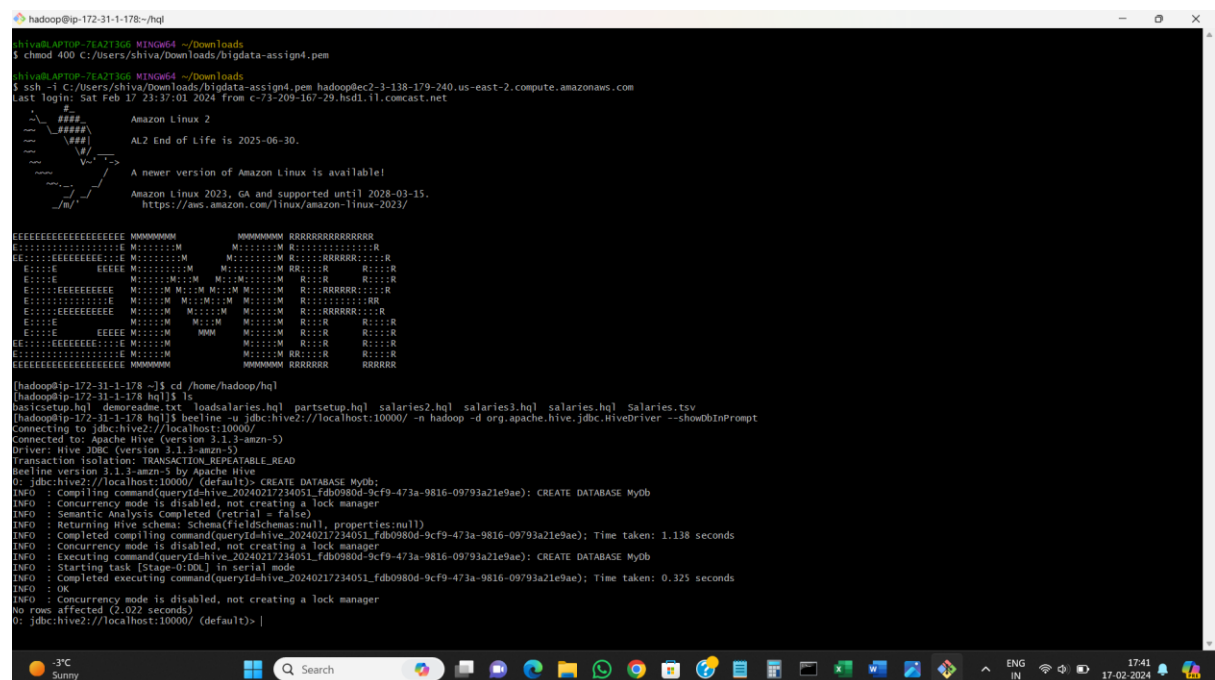
## Starting the Hive beeline

```
[hadoop@ip-172-31-1-178 hql]$ ls
basicsetup.hql  demoreadme.txt  loadsalaries.hql  partsetup.hql  salaries2.hql  salaries3.hql  salaries.hql  Salaries.tsv
[hadoop@ip-172-31-1-178 hql]$ beeline -u jdbc:hive2://localhost:10000/ -n hadoop -d org.apache.hive.jdbc.HiveDriver --showDbInPrompt
Connecting to jdbc:hive2://localhost:10000/
Connected to: Apache Hive (version 3.1.3-amzn-5)
Driver: Hive JDBC (version 3.1.3-amzn-5)
Transaction isolation: TRANSACTION_REPEATABLE_READ
Beeline version 3.1.3-amzn-5 by Apache Hive
0: jdbc:hive2://localhost:10000/ (default)> CREATE DATABASE MyDb;
INFO  : Compiling command(queryId=hive_20240217234051_fdb0980d-9cf9-473a-9816-09793a21e9ae): CREATE DATABASE MyDb
INFO  : Concurrency mode is disabled, not creating a lock manager
INFO  : Semantic Analysis Completed (retrial = false)
INFO  : Returning Hive schema: Schema(fieldSchemas:null, properties:null)
INFO  : Completed compiling command(queryId=hive_20240217234051_fdb0980d-9cf9-473a-9816-09793a21e9ae); Time taken: 1.138 seconds
INFO  : Concurrency mode is disabled, not creating a lock manager
INFO  : Executing command(queryId=hive_20240217234051_fdb0980d-9cf9-473a-9816-09793a21e9ae): CREATE DATABASE MyDb
INFO  : Starting task [Stage-0:DDL] in serial mode
INFO  : Completed executing command(queryId=hive_20240217234051_fdb0980d-9cf9-473a-9816-09793a21e9ae); Time taken: 0.325 seconds
INFO  : OK
INFO  : Concurrency mode is disabled, not creating a lock manager
No rows affected (2.022 seconds)
0: jdbc:hive2://localhost:10000/ (default)>
```

## Files in Hadoop beeline

```
EE:::::EEEEEEEE::::E M:::::M        M:::::M RR::::R     R::::R
E::::::::::::::::::::E M:::::M        M:::::M RR::::R     R::::R
EEEEEEEEEEEEEEEEEEEE MMMMMMM        MMMMMMM RRRRRRRR     RRRRRR

[hadoop@ip-172-31-1-178 ~]$ cd /home/hadoop/hql
[hadoop@ip-172-31-1-178 hql]$ ls
basicsetup.hql  demoreadme.txt  loadsalaries.hql  partsetup.hql  salaries2.hql  salaries3.hql  salaries.hql  Salaries.tsv
[hadoop@ip-172-31-1-178 hql]$
```

## Exercise 1

Create a Hive database called "MyDb".

Now in MyDb create a table with name foodratings having six columns with the name of the first 'name' and the type of the first a string and the names of the remaining columns food1, food2, food3, food4 and id and indicate their types each as an integer. The table should have storage format TEXTFILE and column separator a ",". That is the underlying format should be a CSV file. The table itself and each column should include a comment just to show me you know how to use comments (it does not matter what it says).



Execute a Hive command of 'DESCRIBE FORMATTED MyDb.foodratings;' and capture its output as one of the results of this exercise.

Then in MyDb create a table with name foodplaces having two columns with first called 'id' with the type of the first an integer, and the second column called 'place' with the type of the second a string. This table should also have storage format TEXTFILE and column separator a ",". That is the underlying format should be a CSV file. No comments are needed for this table.

```
18 rows selected (0.551 seconds)
0: jdbc:hive2://localhost:10000/ (default)> CREATE TABLE IF NOT EXISTS mydb.foodratings (
. . . . . . . . . . . . . . . . . . . . . .> id INT,
. . . . . . . . . . . . . . . . . . . . . .> place String
. . . . . . . . . . . . . . . . . . . . . .> )
. . . . . . . . . . . . . . . . . . . . . .> ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
. . . . . . . . . . . . . . . . . . . . . .> STORED AS TEXTFILE;
INFO  : Compiling command(queryId=hive_20240217235011_e241659d-9efc-4bcb-a8bd-bef02c7948a2): CREATE TABLE IF NOT EXISTS mydb.foodratings (
id INT,
place String
)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE
INFO  : Concurrency mode is disabled, not creating a lock manager
INFO  : Semantic Analysis Completed (retrial = false)
INFO  : Returning Hive schema: Schema(fieldSchemas:null, properties:null)
INFO  : Completed compiling command(queryId=hive_20240217235011_e241659d-9efc-4bcb-a8bd-bef02c7948a2); Time taken: 0.022 seconds
INFO  : Concurrency mode is disabled, not creating a lock manager
INFO  : Executing command(queryId=hive_20240217235011_e241659d-9efc-4bcb-a8bd-bef02c7948a2): CREATE TABLE IF NOT EXISTS mydb.foodratings (
id INT,
place String
)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ','
STORED AS TEXTFILE
INFO  : Completed executing command(queryId=hive_20240217235011_e241659d-9efc-4bcb-a8bd-bef02c7948a2); Time taken: 0.002 seconds
INFO  : OK
INFO  : Concurrency mode is disabled, not creating a lock manager
No rows affected (0.037 seconds)
0: jdbc:hive2://localhost:10000/ (default)>
```

Execute a Hive command of 'DESCRIBE FORMATTED MyDb.foodplaces' and capture its output as another of the results of this exercise.

```
18 rows selected (0.093 seconds)
0: jdbc:hive2://localhost:10000/ (default)> DESCRIBE FORMATTED MyDb.foodplaces;
INFO  : Compiling command(queryId=hive_20240217235246_1bffda84-41ed-4ac1-84ee-376865e18b8d): DESCRIBE FORMATTED MyDb.foodplaces
INFO  : Concurrency mode is disabled, not creating a lock manager
INFO  : Semantic Analysis Completed (retrial = false)
INFO  : Returning Hive schema: Schema(fieldSchemas:[FieldSchema(name:col_name, type:string, comment:from deserializer), FieldSchema(name:data_type, type:string, comment:from deserializer), FieldSchema(name:comment, type:string, comment:from deserializer)], properties:null)
INFO  : Completed compiling command(queryId=hive_20240217235246_1bffda84-41ed-4ac1-84ee-376865e18b8d); Time taken: 0.033 seconds
INFO  : Concurrency mode is disabled, not creating a lock manager
INFO  : Executing command(queryId=hive_20240217235246_1bffda84-41ed-4ac1-84ee-376865e18b8d): DESCRIBE FORMATTED MyDb.foodplaces
INFO  : Starting task [Stage-0:DDL] in serial mode
INFO  : Completed executing command(queryId=hive_20240217235246_1bffda84-41ed-4ac1-84ee-376865e18b8d); Time taken: 0.068 seconds
INFO  : OK
INFO  : Concurrency mode is disabled, not creating a lock manager
+-------------------------------+----------------------------------------------------+----------------------------------------------------+
|           col_name            |                     data_type                      |                      comment                       |
+-------------------------------+----------------------------------------------------+----------------------------------------------------+
| # col_name                    | data_type                                          | comment                                            |
| id                            | int                                                |                                                    |
| place                         | string                                             |                                                    |
|                               | NULL                                               | NULL                                               |
| # Detailed Table Information  | NULL                                               | NULL                                               |
| Database:                     | mydb                                               | NULL                                               |
| OwnerType:                    | USER                                               | NULL                                               |
| Owner:                        | hadoop                                             | NULL                                               |
| CreateTime:                   | Sat Feb 17 23:52:13 UTC 2024                       | NULL                                               |
| LastAccessTime:               | UNKNOWN                                            | NULL                                               |
| Retention:                    | 0                                                  | NULL                                               |
| Location:                     | hdfs://ip-172-31-1-178.us-east-2.compute.internal:8020/user/hive/warehouse/mydb.db/foodplaces | NULL |
| Table Type:                   | MANAGED_TABLE                                      | NULL                                               |
| Table Parameters:             | NULL                                               | NULL                                               |
|                               | COLUMN_STATS_ACCURATE                              | {\"BASIC_STATS\":\"true\",\"COLUMN_STATS\":{\"id\":\"true\",\"place\":\"true\"}} |
|                               | bucketing_version                                  | 2                                                  |
|                               | numFiles                                           | 0                                                  |
|                               | numRows                                            | 0                                                  |
|                               | rawDataSize                                        | 0                                                  |
|                               | totalSize                                          | 0                                                  |
|                               | transient_lastDdlTime                              | 1708213933                                         |
|                               | NULL                                               | NULL                                               |
| # Storage Information         | NULL                                               | NULL                                               |
| SerDe Library:                | org.apache.hadoop.hive.serde2.lazy.LazySimpleSerDe | NULL                                               |
| InputFormat:                  | org.apache.hadoop.mapred.TextInputFormat           | NULL                                               |
| OutputFormat:                 | org.apache.hadoop.hive.ql.io.HiveIgnoreKeyTextOutputFormat | NULL                                       |
| Compressed:                   | No                                                 | NULL                                               |
| Num Buckets:                  | -1                                                 | NULL                                               |
| Bucket Columns:               | []                                                 | NULL                                               |
| Sort Columns:                 | []                                                 | NULL                                               |
| Storage Desc Params:          | NULL                                               | NULL                                               |
|                               | field.delim                                        | ,                                                  |
|                               | serialization.format                               | ,                                                  |
+-------------------------------+----------------------------------------------------+----------------------------------------------------+
33 rows selected (0.162 seconds)
0: jdbc:hive2://localhost:10000/ (default)>
```

Exercise 2

Load the foodratings<magic number>.txt file created using TestDataGen from your local file system into the foodratings table.

```
18 rows selected (0.17 seconds)
0: jdbc:hive2://localhost:10000/ (default)> LOAD DATA INPATH '/home/hadoop/foodratings65072.txt'
. . . . . . . . . . . . . . . . . . . . . .> OVERWRITE INTO TABLE mydb.foodratings;
Error: Error while compiling statement: FAILED: SemanticException Line 1:17 Invalid path ''/home/hadoop/foodratings65072.txt'': No files matching path hdfs://ip-172-31-1-178.us-east-2.compute.internal:8020/home
/hadoop/foodratings65072.txt (state=42000,code=40000)
0: jdbc:hive2://localhost:10000/ (default)> LOAD LOCAL DATA INPATH '/home/hadoop/foodratings65072.txt'
. . . . . . . . . . . . . . . . . . . . . .> OVERWRITE INTO TABLE mydb.foodratings;
Error: Error while compiling statement: FAILED: ParseException line 1:5 extraneous input 'LOCAL' expecting DATA near '<EOF>' (state=42000,code=40000)
0: jdbc:hive2://localhost:10000/ (default)> LOAD DATA LOCAL INPATH '/home/hadoop/foodratings65072.txt'
. . . . . . . . . . . . . . . . . . . . . .> OVERWRITE INTO TABLE mydb.foodratings;
INFO  : Compiling command(queryId=hive_20240218000519_56caa13c-4bf0-46ac-952f-955cf05d32a9): LOAD DATA LOCAL INPATH '/home/hadoop/foodratings65072.txt'
OVERWRITE INTO TABLE mydb.foodratings
INFO  : Concurrency mode is disabled, not creating a lock manager
INFO  : Semantic Analysis Completed (retrial = false)
INFO  : Returning Hive schema: Schema(fieldSchemas:null, properties:null)
INFO  : Completed compiling command(queryId=hive_20240218000519_56caa13c-4bf0-46ac-952f-955cf05d32a9); Time taken: 0.053 seconds
INFO  : Concurrency mode is disabled, not creating a lock manager
INFO  : Executing command(queryId=hive_20240218000519_56caa13c-4bf0-46ac-952f-955cf05d32a9): LOAD DATA LOCAL INPATH '/home/hadoop/foodratings65072.txt'
OVERWRITE INTO TABLE mydb.foodratings
INFO  : Starting task [Stage-0:MOVE] in serial mode
INFO  : Loading data to table mydb.foodratings from file:/home/hadoop/foodratings65072.txt
INFO  : Starting task [Stage-1:STATS] in serial mode
INFO  : Executing stats task
INFO  : Table mydb.foodratings stats: [numFiles=1, numRows=0, totalSize=17472, rawDataSize=0]
INFO  : Completed executing command(queryId=hive_20240218000519_56caa13c-4bf0-46ac-952f-955cf05d32a9); Time taken: 0.513 seconds
INFO  : OK
INFO  : Concurrency mode is disabled, not creating a lock manager
No rows affected (0.584 seconds)
0: jdbc:hive2://localhost:10000/ (default)>
```

Execute a single hive command to output the min, max and average of the values of the food3 column of the foodratings table. This should be one hive command, not three separate ones. A copy of the hive command you wrote, the output of this query and the magic number are the result of this exercise.



Exercise 3

Execute a hive command to output the min, max and average of the values of the food1 column grouped by the first column 'name'. This should be one hive command, not three separate ones. The output should look something like:

Mel 10 20 15

Bill 20, 30, 24 …

A copy of the hive command you wrote, the output of this query and the magic number are the result of this exercise.

## Exercise 4

In MyDb create a partitioned table called 'foodratingspart' The partition field should be called 'name' and its type should be a string. The names of the non-partition columns should be food1, food2, food3, food4 and id and their types each an integer. The table should have storage format TEXTFILE and column separator a ",". That is the underlying format should be a CSV file. No comments are needed for this table. Provide the code you have used to create the table as a result of this exercise.



Execute a Hive command of 'DESCRIBE FORMATTED MyDb.foodratingspart;' and capture its output as the result of this exercise.

## Exercise 5

Assume that the number of food critics is relatively small, say less than 10 and the number places to eat is very large, say more than 10,000. In a few short sentences explain why using the (critic) name is a good choice for a partition field while using the place id is not.

Ans: According to the given scenario, if we create a partition table for the name since we have only 10 places, then during the search operation it will only search the query for the particular partition rather than searching whole table with 10000 rows in the database. This dynamic partition will be time efficient and enhances the performances. This method will only work when the number of critics is very less than the number of places. If not, this will be a redundant job.

## Exercise 6

Configure Hive to allow dynamic partition creation as described in the lecture. Now, use a hive command to copy from MyDB.foodratings into MyDB.foodratingspart to create a partitioned table from a non-partitioned one. Hint: The 'name' column from MyDB.foodratings should be mentioned last in this command (whatever it is). Provide a copy of the command you use to load the 'foodratingspart' table as a result of this exercise.

Execute a hive command to output the min, max and average of the values of the food2 column of MyDB.foodratingspart where the food critic 'name' is either Mel or Jill. The query and the output of this query are other results of this exercise, which you must also provide. It should look something like

10 20 15

## Exercise 7

Load the foodplaces<.magic number>.txt  file created using TestDataGen from your local file system into the foodplaces table.



Use a join operation between the two tables (foodratings and foodplaces) to provide the average rating for field food4 for the restaurant 'Soup Bowl' The output of this query is the result of this exercise. It should look something like Soup Bowl 20

Provide the join command and the output as the result of this exercise.

Exercise 8

Read the article "An Introduction to Big Data Formats" found on the blackboard in section "Articles" and provide short (2 to 4 sentence) answers to the following questions:

a) When is the most important consideration when choosing a row format and when a column format for your big data file? Column-based storage is most beneficial for running analytics queries that involve only a subset of columns to be analyzed across massive sets of data. Row-based storage is more appropriate if the majority or all of the columns in each row of data must be accessible for your queries.
Ans: At the highest level, column-based storage is most useful when performing analytics queries that require only a subset of columns examined over very large data sets. If your queries require access to all or most of the columns of each row of data, row-based storage will be better suited to your needs.

b) What is "splittability" for a column file format and why is it important when processing large volumes of data? Big data is huge. In order to handle huge datasets efficiently, it is typically necessary to divide the work into segments that can be assigned to different processors. If the query computation is only concerned with one column at a time, a column-based format will be easier to divide into distinct jobs. A batch of rows is taken and stored in columnar format using row-columnar columnar formats. Once split, these batches form boundaries.
Ans: Big data is HUGE. Processing such datasets efficiently usually requires breaking the job into parts that can be farmed out to separate processors. A column-based format will be more amenable to splitting into separate jobs if the query calculation is concerned with a single column at a time. The columnar formats are row-columnar, which means they take a batch of rows and store that batch in columnar format. These batches then become split boundaries.

c) What can files stored in column format achieve better compression than those stored in row format?
Ans: Compression uses encoding for frequently repeating data to achieve this reduction. Columnar data can achieve better compression rates than row-based data. Storing values by column, with the same type next to each other, allows you to do more efficient compression on them than if

you're storing rows of data. For example, storing all dates together in memory allows for more efficient compression than storing data of various types next to each other—such as string, number, date, string, date.

d) Under what circumstances would it be the best choice to use the "Parquet" column file format?
Ans: Parquet is commonly used with Apache Impala, an analytics database for Hadoop. Parquet is especially adept at analysing wide datasets with many columns. Each Parquet file contains binary data organized by "row group." For each row group, the data values are organized by column. This enables the compression benefits that we described above. Parquet is a good choice for read-heavy workloads.