Assignment 2

# BUAN 6346

Big Data Analytics

Dr. Jerry F. Perez
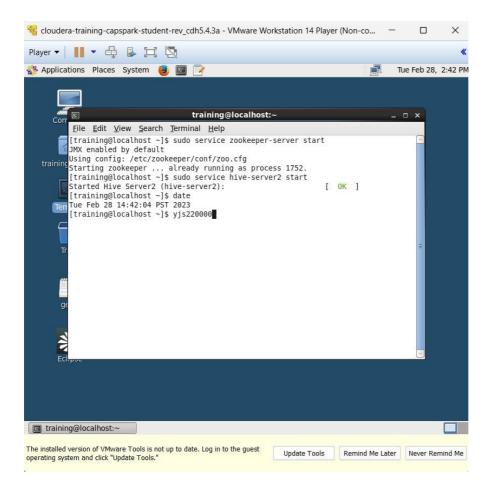
Yug Jyotirmay Singh

YJS220000

# Table of Contents

## Chapter 6

We know that there are a variety of ways to interact with Impala and Hive. In this particular exercise, you will use the Impala or Hive Query Editor in Hue.

1. If you plan to use Hive rather than Impala for this or subsequent exercises, start the Hive server, which is not started by default, by entering the following two commands in a terminal window:

$ sudo service zookeeper-server start

$ sudo service hive-server2 start



2. Now, visit the Hue page in firefox, as described earlier in the "Using HDFS" exercise.

3. Further, open the Impala query editor or Hive query editor, by selecting the editor of your choice from the Query Editors menu.

4. In the query editor page on the right hand side, enter a SQL command to create a table for the webpage data imported in the previous exercises:

CREATE EXTERNAL TABLE webpage

      (page_id SMALLINT,

      name STRING

      assoc_files STRING)

      ROW FORMAT DELIMITED FIELDS TERMINATED BY '\t'

      LOCATION '/loudacre/webpage'



5. Click the EXECUTE button to execute the command.

6. To see the table you just created, refresh the table list on the left hand side

7. Now click on the webpage table to see the column definitions

8. Now click the New Query button, then further enter and execute a test query such as:

SELECT * FROM webpage WHERE name LIKE "ifruit%"



9. Click on the Preview Sample Data icon to view a sampling of the table data.

Data sample for webpage

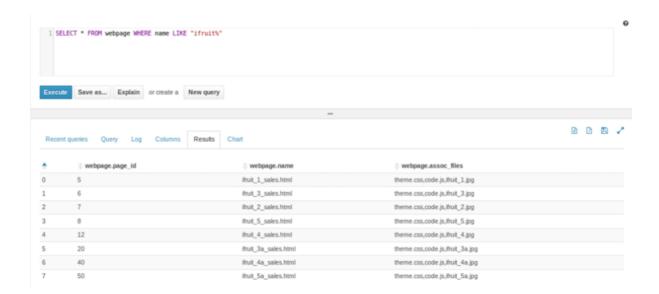| | page_id | name | assoc_files |
|---|---|---|---|
| 0 | 1 | sorrento_f00l_sales.html | theme.css,code.js,sorrento_f00l.jpg |
| 1 | 2 | titanic_2100_sales.html | theme.css,code.js,titanic_2100.jpg |
| 2 | 3 | meetoo_3.0_sales.html | theme.css,code.js,meetoo_3.0.jpg |
| 3 | 4 | meetoo_3.1_sales.html | theme.css,code.js,meetoo_3.1.jpg |
| 4 | 5 | ifruit_1_sales.html | theme.css,code.js,ifruit_1.jpg |
| 5 | 6 | ifruit_3_sales.html | theme.css,code.js,ifruit_3.jpg |
| 6 | 7 | ifruit_2_sales.html | theme.css,code.js,ifruit_2.jpg |
| 7 | 8 | ifruit_5_sales.html | theme.css,code.js,ifruit_5.jpg |
| 8 | 9 | titanic_1000_sales.html | theme.css,code.js,titanic_1000.jpg |
| 9 | 10 | meetoo_1.0_sales.html | theme.css,code.js,meetoo_1.0.jpg |
| 10 | 11 | sorrento_f21l_sales.html | theme.css,code.js,sorrento_f21l.jpg |
| 11 | 12 | ifruit_4_sales.html | theme.css,code.js,ifruit_4.jpg |
| 12 | 13 | sorrento_f23l_sales.html | theme.css,code.js,sorrento_f23l.jpg |

## 6.1 Use Sqoop to Import Directly into Hive and Impala

In this particular section, we will use Sqoop to import data from MySQL into HDFS and also automatically create the corresponding table in the Hive Metastore.

10. Now in the terminal window, import the device table directly into the Hive Metastore.

$ sqoop import \

--connect jdbc:mysql://localhost/loudacre \

--username training –password training \

--fields-terminated-by '\t' \

--table device \

--hive-import

```
[training@localhost ~]$ sqoop import \
>  --connect jdbc:mysql://localhost/loudacre \
>  --username training --password training \
>  --fields-terminated-by '\t' \
>  --table device \
>  --hive-import
23/02/28 16:15:11 INFO sqoop.Sqoop: Running Sqoop version: 1.4.5-cdh5.4.3
23/02/28 16:15:11 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
23/02/28 16:15:12 INFO manager.MySQLManager: Preparing to use a MySQL streaming resultset.
23/02/28 16:15:12 INFO tool.CodeGenTool: Beginning code generation
23/02/28 16:15:12 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM `device` AS t LIMIT 1
23/02/28 16:15:12 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM `device` AS t LIMIT 1
23/02/28 16:15:12 INFO orm.CompilationManager: HADOOP_MAPRED_HOME is /usr/lib/hadoop-mapreduce
Note: /tmp/sqoop-training/compile/4483863cd7958c15ed3a982a822abdd7/device.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
23/02/28 16:15:14 INFO orm.CompilationManager: Writing jar file: /tmp/sqoop-training/compile/4483863cd7958c15ed3a982a822abdd7/device.jar
23/02/28 16:15:14 WARN manager.MySQLManager: It looks like you are importing from mysql.
23/02/28 16:15:14 WARN manager.MySQLManager: This transfer can be faster! Use the --direct
23/02/28 16:15:14 WARN manager.MySQLManager: option to exercise a MySQL-specific fast path.
23/02/28 16:15:14 INFO manager.MySQLManager: Setting zero DATETIME behavior to convertToNull (mysql)
23/02/28 16:15:14 INFO mapreduce.ImportJobBase: Beginning import of device
23/02/28 16:15:14 INFO Configuration.deprecation: mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
23/02/28 16:15:15 INFO Configuration.deprecation: mapred.jar is deprecated. Instead, use mapreduce.job.jar
23/02/28 16:15:15 INFO Configuration.deprecation: mapred.map.tasks is deprecated. Instead, use mapreduce.job.maps
23/02/28 16:15:15 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
23/02/28 16:15:16 WARN security.UserGroupInformation: PriviledgedActionException as:training (auth:SIMPLE) cause:org.apache.hadoop.mapred.FileAlreadyExistsException: Outpu
t directory hdfs://localhost:8020/user/training/device already exists
```



```
Job Counters
        Launched map tasks=4
        Other local map tasks=4
        Total time spent by all maps in occupied slots (ms)=0
        Total time spent by all reduces in occupied slots (ms)=0
        Total time spent by all map tasks (ms)=13333
        Total vcore-seconds taken by all map tasks=13333
        Total megabyte-seconds taken by all map tasks=3413248
Map-Reduce Framework
        Map input records=50
        Map output records=50
        Input split bytes=464
        Spilled Records=0
        Failed Shuffles=0
        Merged Map outputs=0
        GC time elapsed (ms)=238
        CPU time spent (ms)=2100
        Physical memory (bytes) snapshot=487972864
        Virtual memory (bytes) snapshot=3378053120
        Total committed heap usage (bytes)=191889408
File Input Format Counters
        Bytes Read=0
File Output Format Counters
        Bytes Written=2183
23/02/28 16:42:33 INFO mapreduce.ImportJobBase: Transferred 2.1318 KB in 29.5672 seconds (73.8318 bytes/sec)
23/02/28 16:42:33 INFO mapreduce.ImportJobBase: Retrieved 50 records.
23/02/28 16:42:33 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM `device` AS t LIMIT 1
23/02/28 16:42:33 WARN hive.TableDefWriter: Column release_dt had to be cast to a less precise type in Hive
23/02/28 16:42:33 INFO hive.HiveImport: Loading uploaded data into Hive

Logging initialized using configuration in jar:file:/usr/lib/hive/lib/hive-common-1.1.0-cdh5.4.3.jar!/hive-log4j.properties
OK
Time taken: 2.091 seconds
Loading data to table default.device
chgrp: changing ownership of 'hdfs://localhost:8020/user/hive/warehouse/device/part-m-00000': User does not belong to hive
chgrp: changing ownership of 'hdfs://localhost:8020/user/hive/warehouse/device/part-m-00001': User does not belong to hive
chgrp: changing ownership of 'hdfs://localhost:8020/user/hive/warehouse/device/part-m-00002': User does not belong to hive
chgrp: changing ownership of 'hdfs://localhost:8020/user/hive/warehouse/device/part-m-00003': User does not belong to hive
Table default.device stats: [numFiles=4, totalSize=2183]
OK
Time taken: 0.511 seconds
[training@localhost data-format]$ date
Tue Feb 28 16:42:56 PST 2023
[training@localhost data-format]$ YJS22000
```
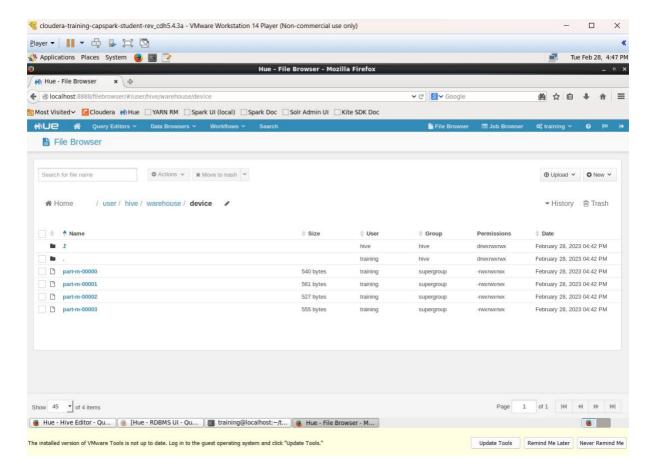
Use –hive-import for either Impala or Hive; this adds metadata to the Metastore, which both tools use.

Note: There might be a possibility that you may get a warning message that chgrp is unable to change the ownership of the generated files; you can disregard the warning, it does not affect the import.

11. Using Hue or the HDFS command line, review the imported data files. The Hive import copies the data to the default Hive warehouse location:

==/user/hive/warehouse/device==


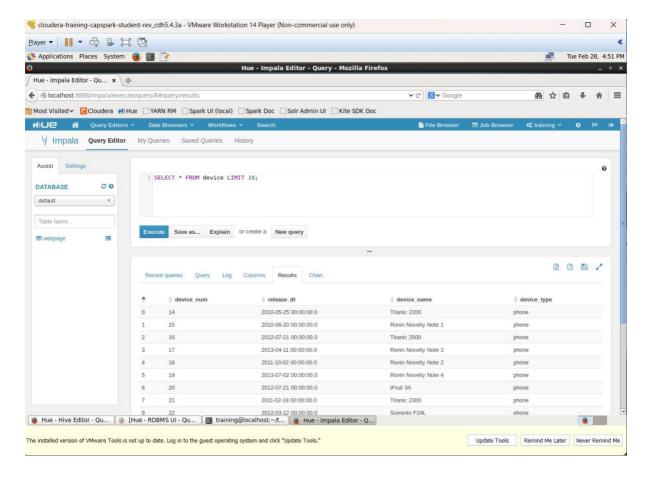
12.        If you are using Impala, make sure you refresh the Impala metadata cache by entering the command in the Hue Impala Query Editor:

==INVALIDATE METADATA==

13.        As in the previous exercise given, view the columns and execute a test query:

==SELECT * FROM device LIMIT 10;==

## Chapter 7

Talking about this exercise, you will you use import data in Avro format and create an Impala/Hive table to access it.

1. Change directories to the exercise directory:

   $ cd $DEV1/exercises/data-format/

2. Import the accounts table to an Avro data format.

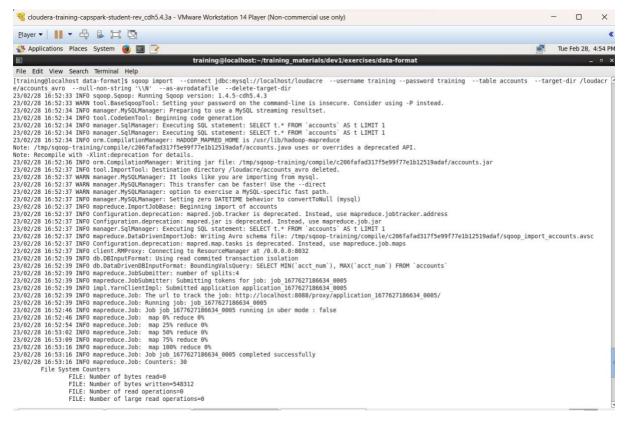$ sqoop import \

--connect jdbc:mysql://localhost/loudacre \
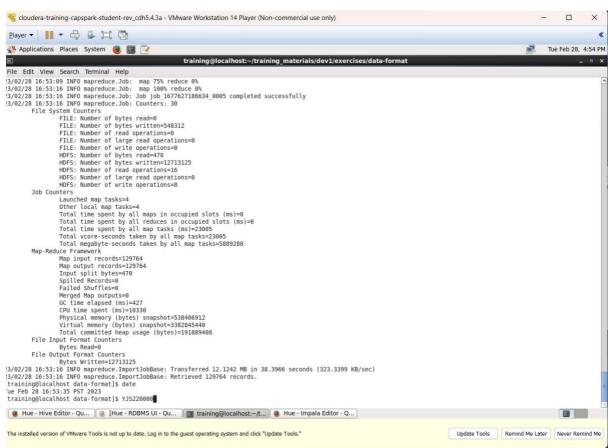
--username training --password training \

--table accounts \

--target-dir /loudacre/accounts_avro \

--null-non-string '\\N' \

## --as-avrodatafile



```
[training@localhost data-format]$ sqoop import --connect jdbc:mysql://localhost/loudacre --username training --password training --table accounts --target-dir /loudacr
e/accounts_avro --null-non-string '\\N' --as-avrodatafile --delete-target-dir
23/02/28 16:52:33 INFO sqoop.Sqoop: Running Sqoop version: 1.4.5-cdh5.4.3
23/02/28 16:52:33 WARN tool.BaseSqoopTool: Setting your password on the command-line is insecure. Consider using -P instead.
23/02/28 16:52:34 INFO manager.MySQLManager: Preparing to use a MySQL streaming resultset.
23/02/28 16:52:34 INFO tool.CodeGenTool: Beginning code generation
23/02/28 16:52:34 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM `accounts` AS t LIMIT 1
23/02/28 16:52:34 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM `accounts` AS t LIMIT 1
23/02/28 16:52:34 INFO orm.CompilationManager: HADOOP_MAPRED_HOME is /usr/lib/hadoop-mapreduce
Note: /tmp/sqoop-training/compile/c206fafad317f5e99f77e1b12519adaf/accounts.java uses or overrides a deprecated API.
Note: Recompile with -Xlint:deprecation for details.
23/02/28 16:52:36 INFO orm.CompilationManager: Writing jar file: /tmp/sqoop-training/compile/c206fafad317f5e99f77e1b12519adaf/accounts.jar
23/02/28 16:52:37 INFO tool.ImportTool: Destination directory /loudacre/accounts_avro deleted.
23/02/28 16:52:37 WARN manager.MySQLManager: It looks like you are importing from mysql.
23/02/28 16:52:37 WARN manager.MySQLManager: This transfer can be faster! Use the --direct
23/02/28 16:52:37 WARN manager.MySQLManager: option to exercise a MySQL-specific fast path.
23/02/28 16:52:37 INFO manager.MySQLManager: Setting zero DATETIME behavior to convertToNull (mysql)
23/02/28 16:52:37 INFO mapreduce.ImportJobBase: Beginning import of accounts
23/02/28 16:52:37 INFO Configuration.deprecation: mapred.job.tracker is deprecated. Instead, use mapreduce.jobtracker.address
23/02/28 16:52:37 INFO Configuration.deprecation: mapred.jar is deprecated. Instead, use mapreduce.job.jar
23/02/28 16:52:37 INFO manager.SqlManager: Executing SQL statement: SELECT t.* FROM `accounts` AS t LIMIT 1
23/02/28 16:52:37 INFO mapreduce.DataDrivenImportJob: Writing Avro schema file: /tmp/sqoop-training/compile/c206fafad317f5e99f77e1b12519adaf/sqoop_import_accounts.avsc
23/02/28 16:52:37 INFO Configuration.deprecation: mapred.map.tasks is deprecated. Instead, use mapreduce.job.maps
23/02/28 16:52:37 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032
23/02/28 16:52:39 INFO db.DBInputFormat: Using read commited transaction isolation
23/02/28 16:52:39 INFO db.DataDrivenDBInputFormat: BoundingValsQuery: SELECT MIN(`acct_num`), MAX(`acct_num`) FROM `accounts`
23/02/28 16:52:39 INFO mapreduce.JobSubmitter: number of splits:4
23/02/28 16:52:39 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1677627186634_0005
23/02/28 16:52:39 INFO impl.YarnClientImpl: Submitted application application_1677627186634_0005
23/02/28 16:52:39 INFO mapreduce.Job: The url to track the job: http://localhost:8088/proxy/application_1677627186634_0005/
23/02/28 16:52:39 INFO mapreduce.Job: Running job: job_1677627186634_0005
23/02/28 16:52:46 INFO mapreduce.Job: Job job_1677627186634_0005 running in uber mode : false
23/02/28 16:52:46 INFO mapreduce.Job:  map 0% reduce 0%
23/02/28 16:52:54 INFO mapreduce.Job:  map 25% reduce 0%
23/02/28 16:53:02 INFO mapreduce.Job:  map 50% reduce 0%
23/02/28 16:53:09 INFO mapreduce.Job:  map 75% reduce 0%
23/02/28 16:53:16 INFO mapreduce.Job:  map 100% reduce 0%
23/02/28 16:53:16 INFO mapreduce.Job: Job job_1677627186634_0005 completed successfully
23/02/28 16:53:16 INFO mapreduce.Job: Counters: 30
        File System Counters
                FILE: Number of bytes read=0
                FILE: Number of bytes written=548312
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
```
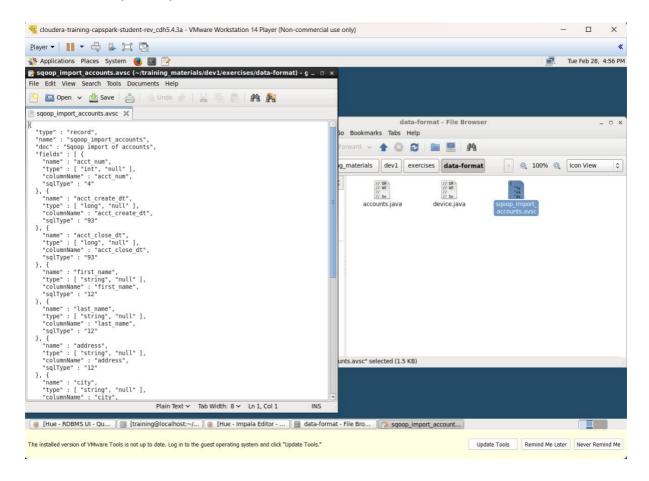


```
23/02/28 16:53:09 INFO mapreduce.Job:  map 75% reduce 0%
23/02/28 16:53:16 INFO mapreduce.Job:  map 100% reduce 0%
23/02/28 16:53:16 INFO mapreduce.Job: Job job_1677627186634_0005 completed successfully
23/02/28 16:53:16 INFO mapreduce.Job: Counters: 30
        File System Counters
                FILE: Number of bytes read=0
                FILE: Number of bytes written=548312
                FILE: Number of read operations=0
                FILE: Number of large read operations=0
                FILE: Number of write operations=0
                HDFS: Number of bytes read=470
                HDFS: Number of bytes written=12713125
                HDFS: Number of read operations=16
                HDFS: Number of large read operations=0
                HDFS: Number of write operations=8
        Job Counters
                Launched map tasks=4
                Other local map tasks=4
                Total time spent by all maps in occupied slots (ms)=0
                Total time spent by all reduces in occupied slots (ms)=0
                Total time spent by all map tasks (ms)=23005
                Total vcore-seconds taken by all map tasks=23005
                Total megabyte-seconds taken by all map tasks=5889280
        Map-Reduce Framework
                Map input records=129764
                Map output records=129764
                Input split bytes=470
                Spilled Records=0
                Failed Shuffles=0
                Merged Map outputs=0
                GC time elapsed (ms)=427
                CPU time spent (ms)=10330
                Physical memory (bytes) snapshot=538406912
                Virtual memory (bytes) snapshot=3382845440
                Total committed heap usage (bytes)=191889408
        File Input Format Counters
                Bytes Read=0
        File Output Format Counters
                Bytes Written=12713125
23/02/28 16:53:16 INFO mapreduce.ImportJobBase: Transferred 12.1242 MB in 38.3966 seconds (323.3399 KB/sec)
23/02/28 16:53:16 INFO mapreduce.ImportJobBase: Retrieved 129764 records.
[training@localhost data-format]$ date
Tue Feb 28 16:53:35 PST 2023
[training@localhost data-format]$ YJS220000
```
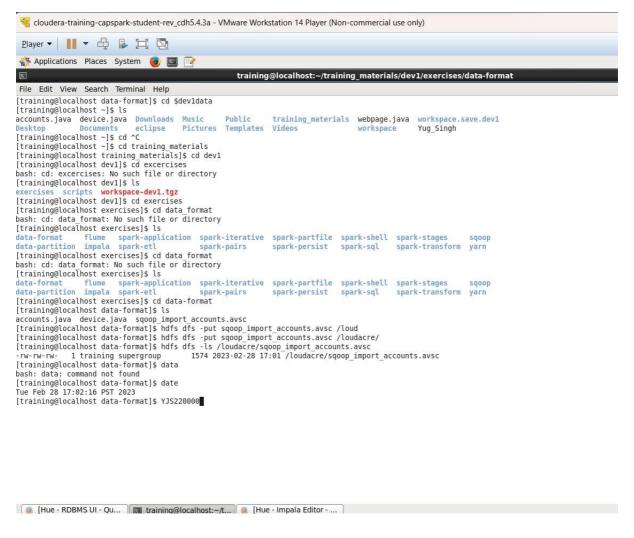
3. Now view the files imported by Sqoop into HDFS. What do you see when you try to view the content of the data files?



4. Sqoop generated a schema named sqoop_import_accounts.avsc in the current directory. Review this file and then copy it to the /loudacre directory in HDFS.
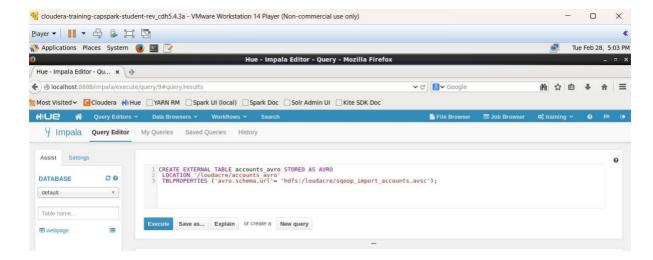
5. In Impala or Hive, create a table using this schema:
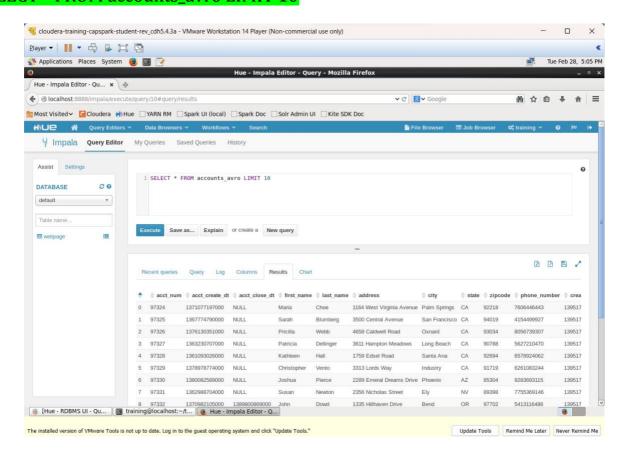
CREATE EXTERNAL TABLE accounts_avro STORED AS AVRO

LOCATION '/loudacre/accounts_avro'

TBLPROPERTIES ('avro.schema.url'=
'hdfs:/loudacre/sqoop_import_accounts.avsc');

6. Confirm correct creation of the table by issuing a query such as

SELECT * FROM accounts_avro LIMIT 10



7. Use the DESCRIBE or DESCRIBE FORMATTED command to list the columns and data types of the accounts_avro table created from the Avro schema.

DESCRIBE accounts_avro;

## Chapter 8

Talking about this particular exercise, we will create and load an Impala/Hive table with account data, partitioned by area code.

Talking about the previous exercise, we imported data from the accounts table using Sqoop, into a table called accounts_avro. In this exercise, you will create a new table with some of the account data, partitioned by area code (the first three digits of the phone number).

1. Create a new, empty table in Impala or Hive:

CREATE EXTERNAL TABLE accounts_by_areacode ( acct_num INT,first_name STRING,

last_name STRING, phone_number STRING)

PARTITIONED BY (areacode STRING)

ROW FORMAT DELIMITED

2. In order to populate the new table, we will need to extract the area code from the phone number. We will try executing the following query to demonstrate:

3. Use the SELECT statement above in an INSERT INTO TABLE command to copy the specified columns to the new table, dynamically partitioning by area code.

4. Execute a simple query to confirm that the table was populated correctly, such as

SELECT * FROM accounts_by_areacode LIMIT 10

| | acct_num | first_name | last_name | phone_number | areacode |
|---|---|---|---|---|---|
| 0 | 64920 | Michael | Flynn | 9161476341 | 916 |
| 1 | 64936 | Seth | Williams | 9164802857 | 916 |
| 2 | 64948 | Ruth | Lind | 9161439766 | 916 |
| 3 | 64965 | Sharon | Collier | 9166449005 | 916 |
| 4 | 64967 | Marie | Redding | 9160550780 | 916 |
| 5 | 64993 | Kenneth | Lopez | 9162312933 | 916 |
| 6 | 65010 | Roger | Hall | 9169899934 | 916 |
| 7 | 65013 | Debra | Whittaker | 9168689641 | 916 |
| 8 | 65055 | Philip | Roberts | 9169886979 | 916 |
| 9 | 65081 | Susan | Lozano | 9161441955 | 916 |

5. Using Hue or the hdfs command-line interface, confirm that the directory structure of the accounts_by_areacode table includes partition directories. Review the data in the directories to verify that the partitioning is correct.

## Chapter 9

Talking about this exercise, we will configure Flume to ingest web log data from a local directory to HDFS.

Apache web server logs are generally stored in files on the local machines running the server. In this exercise, we will simulate an Apache server by placing provided web log files into a local spool directory and then using Flume to collect the data. Both the local and HDFS directories must exist before using the spooling directory source.

1. Create a directory in HDFS called /loudacre/weblogs to hold the data files Flume ingests, example:

$ hdfs dfs -mkdir /loudacre/weblogs

2. Create the spool directory into which our weblog simulator will store data files for Flume to ingest. On the local filesystem create */flume/weblogs_spooldir:*

$ sudo mkdir -p /flume/weblogs_spooldir

3. Give all users the permission to write to the */flume/weblogs_spooldir* directory:

$ sudo chmod a+w -R /flume



## 9.1 Configure Flume

In $DEV1/exercises/flume create a Flume configuration file with the characteristics listed below

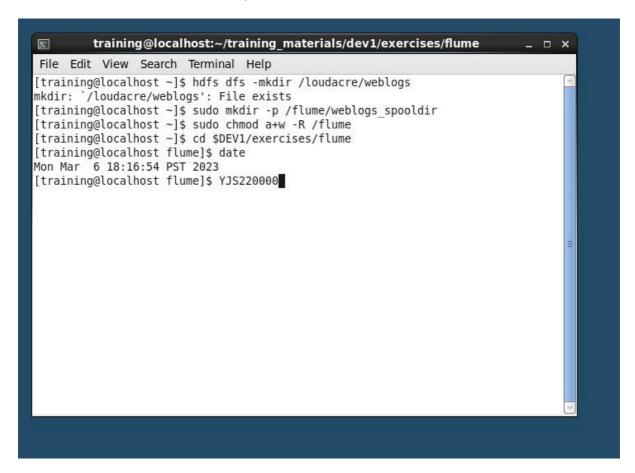• The source is a spooling directory source that pulls from

*/flume/weblogs_spooldir*

• The sink is an HDFS sink that:
- Writes files to the /loudacre/weblogs directory
- Disables time-based file rolling by setting the hdfs.rollInterval property to 0.
- Disables event-based file rolling by setting the hdfs.rollCount property to 0.
- Sets the hdfs.rollSize property as 524288 to enable size-based file rolling at 512KB. - Writes raw text files (instead of SequenceFile format) by setting hdfs.fileType to datastream.

• The channel is a Memory Channel that:
- Can store 10,000 events using the capacity property.
- Has a transaction capacity of 10,000 events using the transactionCapacity property.

4. I changed the directory to the $DEV1/exercises/flume directory as follows:

$ cd $DEV1/exercises/flume

- I created solution directory as it didn't exist.

- I then created the spooldir.conf in the solution directory and we can see that in the /solution directory

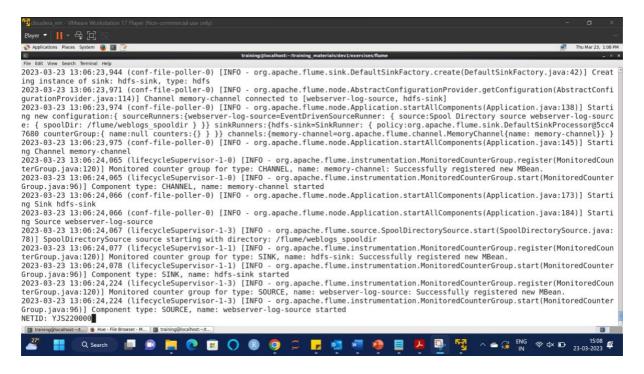5. I started the flume agent that I had created in the last step (since it didn't exist I created it).

```
training@localhost:~/training_materials/dev1/exercises/flume    _ □ ✕
File  Edit  View  Search  Terminal  Help
[training@localhost ~]$ hdfs dfs -mkdir /loudacre/weblogs
mkdir: `/loudacre/weblogs': File exists
[training@localhost ~]$ sudo mkdir -p /flume/weblogs_spooldir
[training@localhost ~]$ sudo chmod a+w -R /flume
[training@localhost ~]$ cd $DEV1/exercises/flume
[training@localhost flume]$ date
Mon Mar  6 18:16:54 PST 2023
[training@localhost flume]$ YJS220000
```

5.1 Starting the flume agent now using the following command:

flume-ng agent --conf /etc/flume-ng/conf \

--conf-file solution/spooldir.conf \

--name agent1 -Dflume.root.logger=INFO,console



Flume agent to start up. I saw the message like: Component type: SOURCE, name: webserver-log-source started

## 9.2 Simulate Apache Web Server Output

7. Open a separate terminal window, and change to the exercise directory. Run the script to place the web log files in the /flume/weblogs_spooldir directory: This script will create a temporary copy of the web log files and move them to the spooldir directory.

I wasn't able to run the following scripts in VMware because of certain error that was popping on my screen several times.

$ cd $DEV1/exercises/flume

$ ./copy-move-weblogs.sh /flume/weblogs_spooldir

8. The other terminal that is running the Flume agent and watch the logging output. The output will give information about the files Flume is putting into HDFS.

9. Once the Flume agent has finished, enter CTRL+C to terminate the process.

10. I then listed the files in HDFS that were added by the flume agent.

$ hdfs dfs -ls /loudacre/weblogs