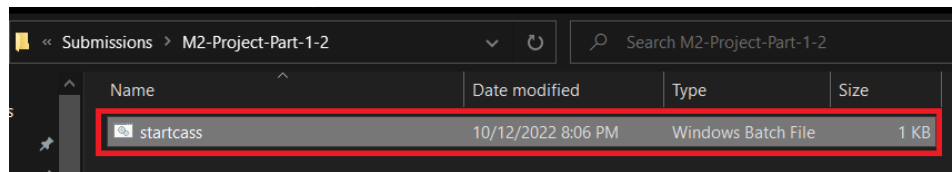## Module 2 Project (Part 1)

### I.    Initializing Cassandra
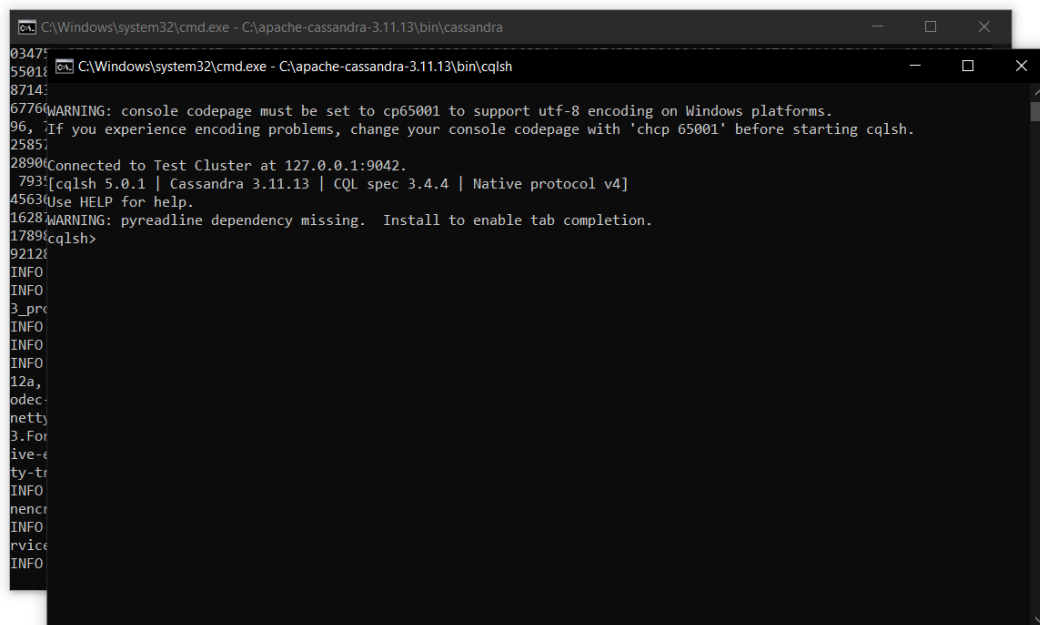
Due to the previous "Module 1 Project" activity, we already have Cassandra set up on our local computer. To interact with the Cassandra database, we first must run the Cassandra executable before running the command-line **cqlsh** to interact with Cassandra using CQL. In this project, we automated the process by creating a .bat file that contains the following:

```
start C:\apache-cassandra-3.11.13\bin\cassandra

echo Wait until cassandra is done initializing before continuing

pause

start C:\apache-cassandra-3.11.13\bin\cqlsh
```

We created this file using a notepad and saved it as **startcass.bat**. Next, we executed this file.



After execution, we had two terminals running in the background: one for the Cassandra database and the other for the CQL shell.

In the terminal where the CQL shell was initialized, we ran the following command, "`describe keyspaces;`", to access the keyspace we made previously in the "Module 1 Project" activity. After verifying that our keyspace is still available, we used the command "`use group23_project;`" to use the keyspace. Next, we ran the command "`describe tables;`" to check that the table containing our CCTV counts still exists.



Our table exists however we need to verify if the data inside is complete. By running this command, "`select * from group23_project_table;`", we can see that our data is ready for use in Talend.

## II.     Opening Talend for Big Data

Talend was also set up on our computer due to the recent "Module 2 Assignment 2 and 3" activities. So, we just need to go to the directory folder where our Talend for Big Data is located to open it.



**Note:** If you are still encountering problems such as "Incompatible JVM" then you may need to update the environment variables of your system. Otherwise, proceed here.



However, you may also opt to do the following: Type "cmd" in the address bar of your Talend directory folder.

In the command prompt, we typed "`TOS_BD-win-x86_64 -vm "C:\Program Files\Java\jdk-11.0.16.1\bin"`". Please take note that your Java folder may be different from this example.



Once your Talend for Big Data is open, proceed to creating a new project. In this instance, we named our project "Group-23-M2-Project". Then proceed to select the newly made project and open it in Talend.

## III. Talend Cassandra NoSQL Connection and Job

After creating a new project, the Talend window should now appear. To establish the connection between the Cassandra database and Talend, select *Metadata*. Then, right click *NoSQL Connections > Create connection*.

The new NoSQL Connection window should appear, where the user will be prompted to specify a name, purpose, and description for the new NoSQL connection. Once all required fields are populated, click *Next >*. In this assignment, we named the connection as *NoSQLCass* and left the remaining text field and area blank.

After specifying the name (and purpose and description) of the connection, select the database (DB Type) that you wish to establish a Talend connection with. Since one requirement of this assignment involves Cassandra, select *Cassandra* as the **DB Type** and *Cassandra 3.0.x* as the **DB Version**. Then, specify *localhost* or *127.0.0.1* as the **Server** with *9042* as its **Port**.

From a previous assignment, we created a keyspace called *group23_project*, where we stored the CCTV_counts table (*group23_project_table*). Since the objective of this project is to export that table into a flat file, we specified the **Keyspace** *group23_project* in this connection. Select "Check" to verify all fields.

After clicking "Check" to verify the contents of the fields, the following window will appear to prompt you to download and install the required modules. Click *Download and install all modules available*.



You will be returned to the previous window. The following dialog box should appear to indicate that Connection creation was successful. Click *OK* to close this dialog box. Finally, click *Finish.*

## IV. Exporting CCTV_Counts from Cassandra to an Excel File

To begin exporting the Cassandra table into a flat file, select *Create job* after right-clicking on **Job Designs** under the Repository.



Then, name the new job as *CassExportExcel*.



Next, go to the Repository tab and under *NoSQL Connections*, right click on the newly created NoSql Connection to Cassandra and click on "Retrieve Schema".

In the Schema window prompt, check mark the box beside the name of the table that exists in your Cassandra database. Click on "Next >" to proceed.



There is nothing to change here so click on "Finish".

Drag the NoSQL Connection (*NoSQLCass*) into the job and choose *tCassandraInput*. Click "OK" for it to be placed.



The newly made *tCassandraInput* component should look like this:

Double click the tCassandraInput. Since the **tCassandraInput** subjob requires an external .jar file to be installed (1), select *Download and install all modules available* (2) to do so.



Now, insert the following **Query** on the appropriate text area to show the group23_project_table (CCTV_counts table) and to prepare it for exporting:

```
"select * from group23_project.group23_project_table;"
```

Then, in the Palette tab to your right, search for **tFileOutputExcel**, then select and drag it. This component will perform the flat file export of the Cassandra table. Drag the job to the workspace, as indicated.



Now, drag an arrow from the *NoSQLCass* component that was created earlier through the *tCassandraInput* component to the *tFileOutputExcel* job.

Then, specify the directory where the exported .xls file will be saved. In this assignment, we saved the file in the directory *C:/Program Files (x86)/TOS_BD/workspace/cctv_counts.xls* with a sheet name of *cctv_counts*. Each column from the Cassandra table will be automatically detected by Talend, as indicated in the *Define column auto size* section.

To run the execution, go to the *Run (Job CassExportExcel)* tab and select the *Run* command. The encircled portion of the screenshot should appear, indicating that the exporting is taking place.



Once the execution has finished, go to the specified directory that made earlier and look for a "cctv_counts." The directory should contain the .xls file, as indicated.

Select the file. The contents of the flat file should be like the screenshot below.

| | A | B | C | D | E | F | G | H | I | J | K | L | M |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | timeuuid_i | bike | bus | car | date_save | jeepney | lgu_code | others | sensor_id | time_save | total | truck | tryke |
| 2 | 166325397 | 5 | 0 | 0 | 09/15/202 | 2 | 1200 | 0 | sensor_09 | 22:58:45 | 9 | 1 | 1 |
| 3 | 166325287 | 4 | 1 | 3 | 09/15/202 | 2 | 1200 | 2 | sensor_06 | 22:40:15 | 15 | 2 | 1 |
| 4 | 166325422 | 2 | 2 | 0 | 09/15/202 | 2 | 1200 | 2 | sensor_02 | 23:03:46 | 9 | 1 | 0 |
| 5 | 166325290 | 3 | 2 | 2 | 09/15/202 | 2 | 1200 | 0 | sensor_07 | 22:41:41 | 12 | 0 | 3 |
| 6 | 166325500 | 3 | 2 | 3 | 09/15/202 | 2 | 1200 | 0 | sensor_10 | 23:18:18 | 14 | 1 | 3 |
| 7 | 166325357 | 4 | 0 | 4 | 09/15/202 | 2 | 1200 | 1 | sensor_03 | 22:53:11 | 13 | 1 | 1 |
| 8 | 166325294 | 1 | 0 | 3 | 09/15/202 | 2 | 1200 | 0 | sensor_05 | 22:42:22 | 8 | 1 | 1 |
| 9 | 166325394 | 4 | 0 | 0 | 09/15/202 | 0 | 1200 | 1 | sensor_04 | 22:59:03 | 8 | 1 | 2 |
| 10 | 166325515 | 5 | 2 | 2 | 09/15/202 | 2 | 1200 | 1 | sensor_02 | 23:19:10 | 16 | 2 | 2 |
| 11 | 166325477 | 2 | 0 | 2 | 09/15/202 | 0 | 1200 | 2 | sensor_08 | 23:12:00 | 10 | 1 | 3 |
| 12 | 166325380 | 0 | 2 | 3 | 09/15/202 | 2 | 1200 | 0 | sensor_06 | 22:57:12 | 10 | 2 | 1 |
| 13 | 166325422 | 1 | 2 | 2 | 09/15/202 | 1 | 1200 | 1 | sensor_02 | 23:03:47 | 9 | 0 | 2 |
| 14 | 166325396 | 4 | 2 | 0 | 09/15/202 | 0 | 1200 | 2 | sensor_01 | 22:59:21 | 11 | 0 | 3 |
| 15 | 166325396 | 3 | 1 | 3 | 09/15/202 | 1 | 1200 | 2 | sensor_01 | 22:59:24 | 13 | 1 | 2 |
| 16 | 166325397 | 0 | 0 | 0 | 09/15/202 | 0 | 1200 | 1 | sensor_07 | 22:58:48 | 5 | 2 | 2 |
| 17 | 166325399 | 1 | 1 | 1 | 09/15/202 | 2 | 1200 | 2 | sensor_06 | 22:59:54 | 10 | 0 | 3 |
| 18 | 166325265 | 5 | 1 | 1 | 09/15/202 | 2 | 1200 | 0 | sensor_09 | 22:36:56 | 11 | 1 | 1 |
| 19 | 166325449 | 1 | 0 | 0 | 09/15/202 | 0 | 1200 | 1 | sensor_03 | 23:08:17 | 5 | 1 | 2 |
| 20 | 166325394 | 4 | 0 | 1 | 09/15/202 | 0 | 1200 | 0 | sensor_06 | 22:59:24 | 5 | 0 | 0 |
| 21 | 166325344 | 4 | 2 | 2 | 09/15/202 | 1 | 1200 | 0 | sensor_06 | 22:51:04 | 12 | 0 | 3 |
| 22 | 166325290 | 3 | 2 | 1 | 09/15/202 | 0 | 1200 | 0 | sensor_09 | 22:43:19 | 10 | 2 | 2 |

cctv_counts

# REFERENCES

*Connect to Cassandra Data and Transfer Data in Talend*. (n.d.). CData Software. Retrieved from
https://www.cdata.com/kb/tech/cassandra-jdbc-talend.rst

*Managing NoSQL metadata*. (n.d.). Talend. Retrieved from https://help.talend.com/r/en-US/7.3/studio-user-guide-big-data/managing-nosql-metadata