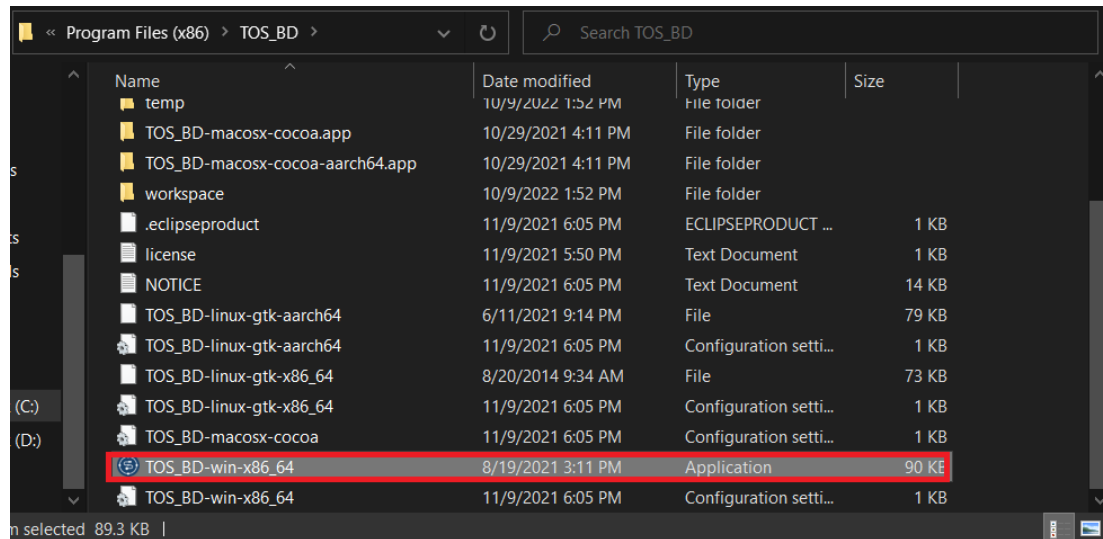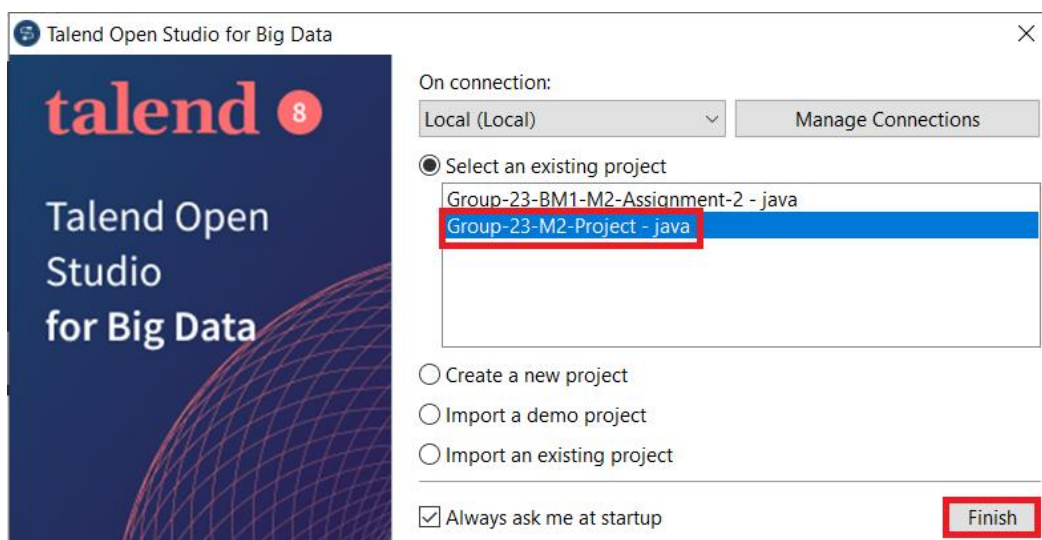## Module 2 Project (Part 2)

### I.    Opening Talend for Big Data

Once again, we will be using Talend for this part of the project. To open it, you simply navigate to the directory folder containing your Talend for Big Data.
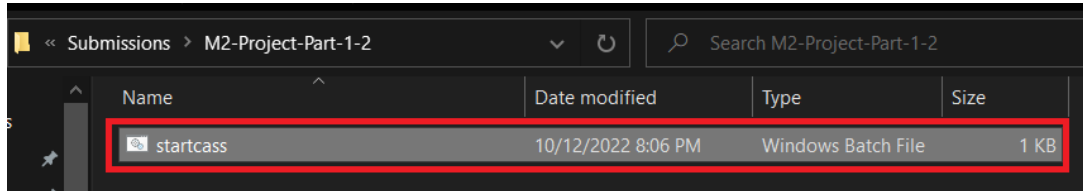


**Note:** If you continue to encounter issues such as "Incompatible JVM," you may need to update your system's environment variables or follow the previous steps from Module 2 Project Part 1.

Once your Talend for Big Data is open, you can proceed to selecting or creating a new project. Since we already have a project (Group-23-M2-Project) made for our "Module 2 Project Part 1" we can use that. We selected that project and clicked on Finish to open Talend.

## II.    Initializing Cassandra

In our Module 2 Project Part 1, we created a file named **startcass.bat** that automated the process to opening Cassandra. Since we will be using Cassandra again, we need to execute the file.
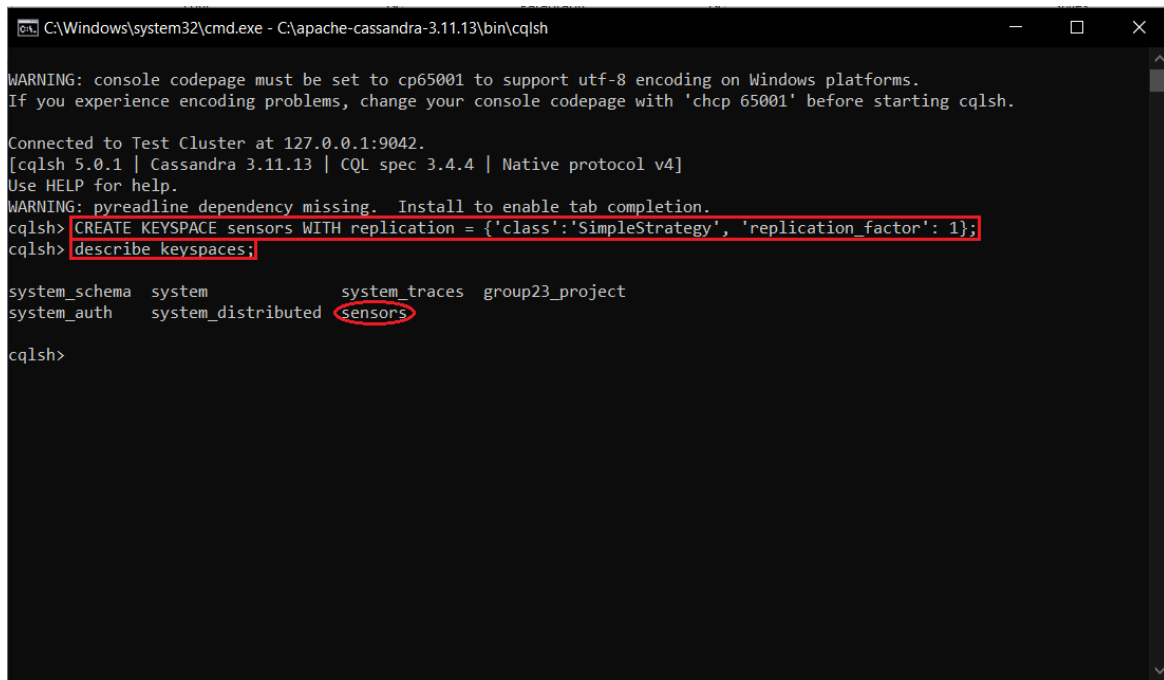


Following execution, two terminals were running in our background: one for the Cassandra database and one for the CQL shell.

In the terminal where the CQL shell was initialized, we ran the following command:

```
CREATE        KEYSPACE        sensors        WITH        replication        =
{'class':'SimpleStrategy', 'replication_factor': 1};
```

This command is used to create a keyspace named "sensors" where we will place the table that will contain data from the dim_sensors.xlxs that was given. To verify if the keyspace was created, we use the command "`describe keyspaces;`".



Cassandra is now ready for the next steps.

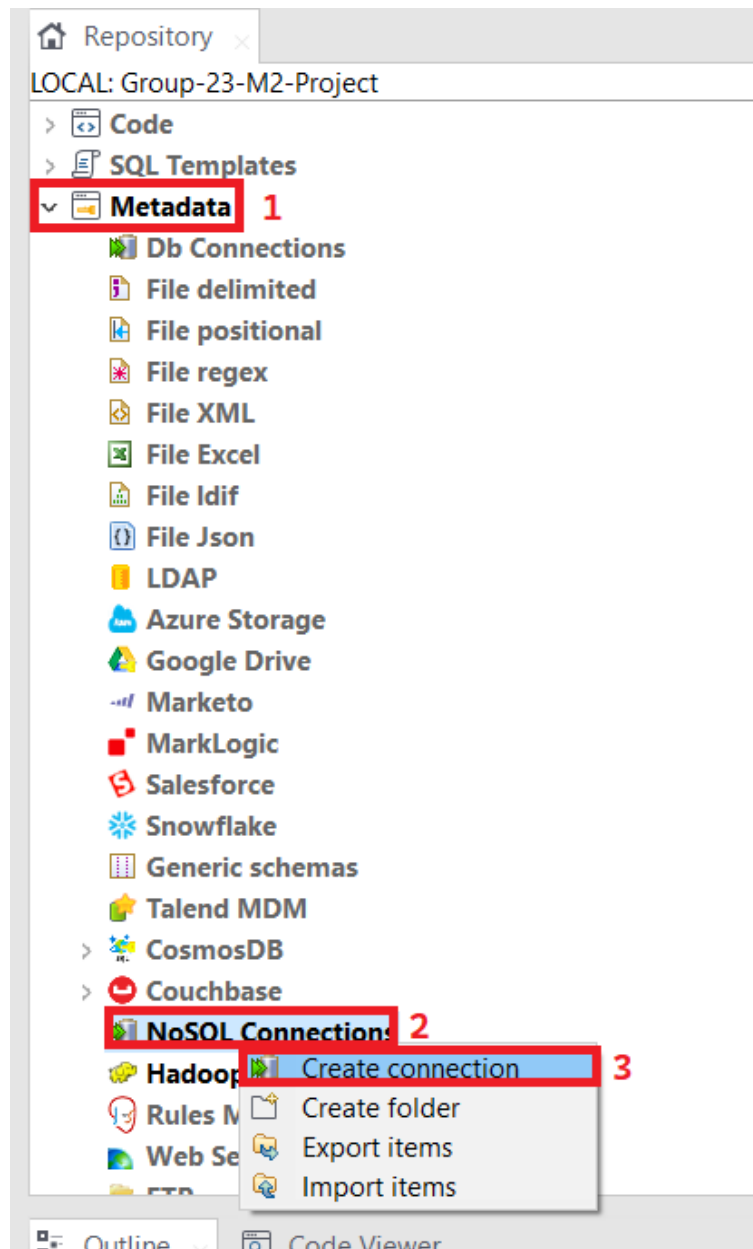## III.     Creating a Job and Establishing a NoSql Connection

Go back to your open Talend Studio for Big Data. Under the Repository tab, right click *Job Designs* and click on *Create Job*.

In our example, we named our job *ImporttoCass* but you may name it however you wish.

Previously we made a NoSql Connection in this Talend project; however, for the sake of documentation, we will reinstate it again. To establish the connection between the Cassandra database and Talend, select *Metadata* > Right click *NoSQL Connections* > *Create connection*.

The new NoSQL Connection window should appear, prompting the user to enter a name, purpose, and description for the new NoSQL connection. After filling out all of the required fields, click "Next >." For this project, we named the connection *NoSQLCass* and left the remaining text field and area blank.

The next step is to select the database (DB Type) that you wish to establish a Talend connection with. Select *Cassandra* as the **DB Type** and *Cassandra 3.0.x* as the **DB Version**. Specify *localhost* or *127.0.0.1* as the **Server** with *9042* as its **Port**. For the **Keyspace**, you may name it however you like. Once everything is finished, click on "Check" to test the connection.

You will then encounter a message that states that some list of modules should be downloaded for you to continue. Just click on "Download and install all modules available".



After downloading the modules needed, your connection should be successful. Click "Finish" once everything is all set.

## IV.    Importing Data to Cassandra using Talend

To import the provided dataset into Cassandra using Talend, select the ">" next to **Metadata**, then right click on **File Excel** > **Create file Excel**.

The *New Excel File* window should appear. Here we specified "dim_sensors" as the **Name**, however, you may choose a different one. Then, click *Next >*.

Then, select the "Browse…" button on the right of the *File* text field, and find the *dim_sensors* file. The file's directory should appear like the one in the screenshot below. Select the *Read excel2007 file format(xlsx)* checkbox to ensure that all Excel versions of the dataset will be read. The File Viewer should automatically detect the columns and its contents, as shown in the bottom half of the window. No changes are required so click *Next >* to proceed.

For the next step, make sure to check mark the box next to *Set heading row as column names*. After, click on "Refresh Preview" to check if the excel output looks right. Click "Finish" once everything is done.

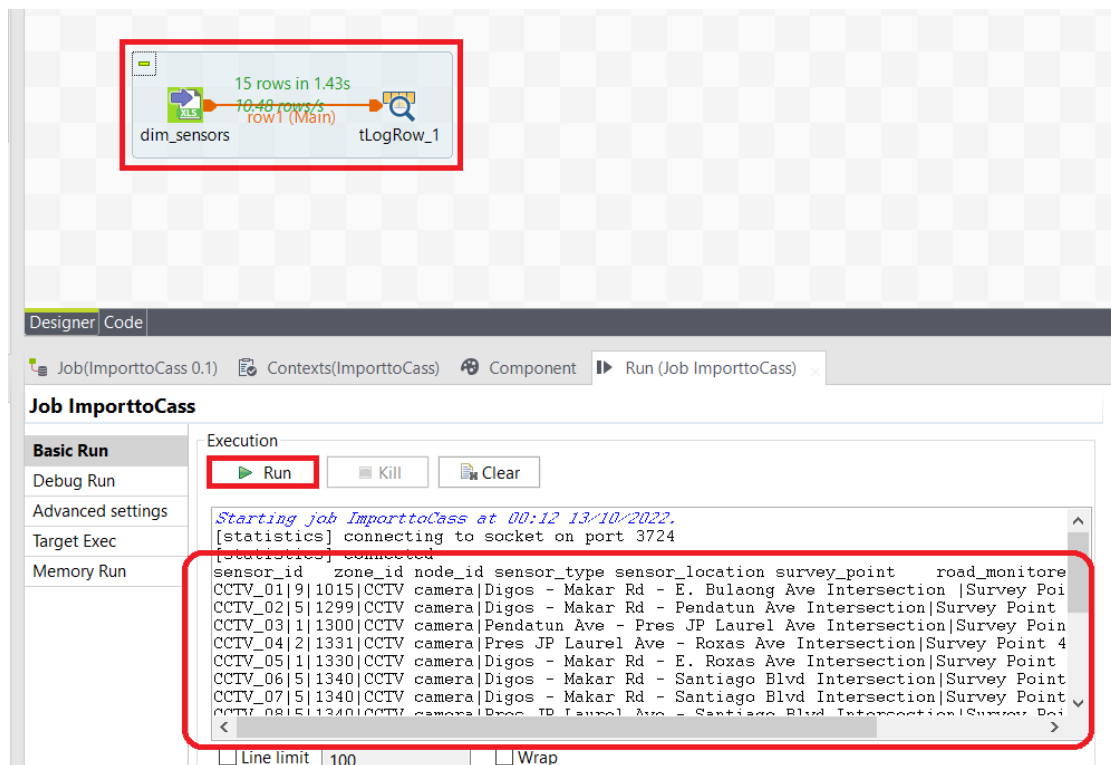There is nothing else to change in the last step so click on "Finish" to complete the process.



Next, drag the newly made File excel component *dim_sensors* to the job. Choose tFileInputExcel and click "OK".
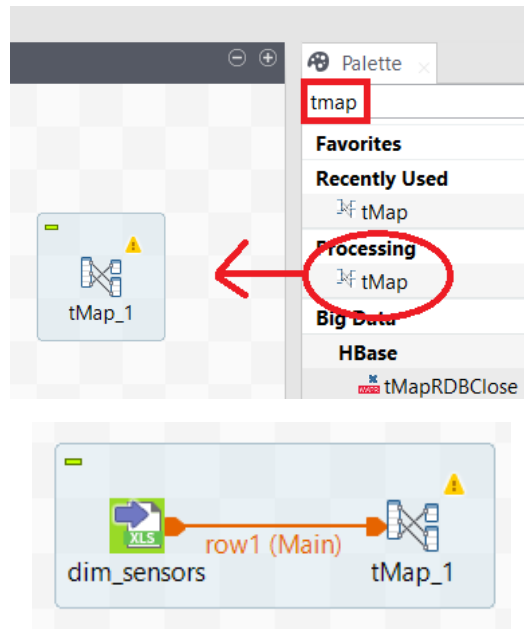
Then, in the Palette tab to your right, search for **tLog Row**, then select and drag it. This component will help verify the data. Drag the component to the workspace, as indicated.
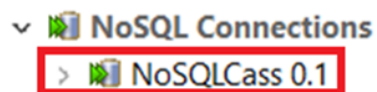


Now, drag an arrow from the *dim_sensors* component that was created earlier to the *tLogRow_1* component. To run the execution, go to the *Run (Job ImporttoCass)* tab and select the *Run* command. The emphasized portion of the screenshot should appear; here we see the data from *dim_sensors*. Delete this *tLogRow_1* component since we do not need it for importing.
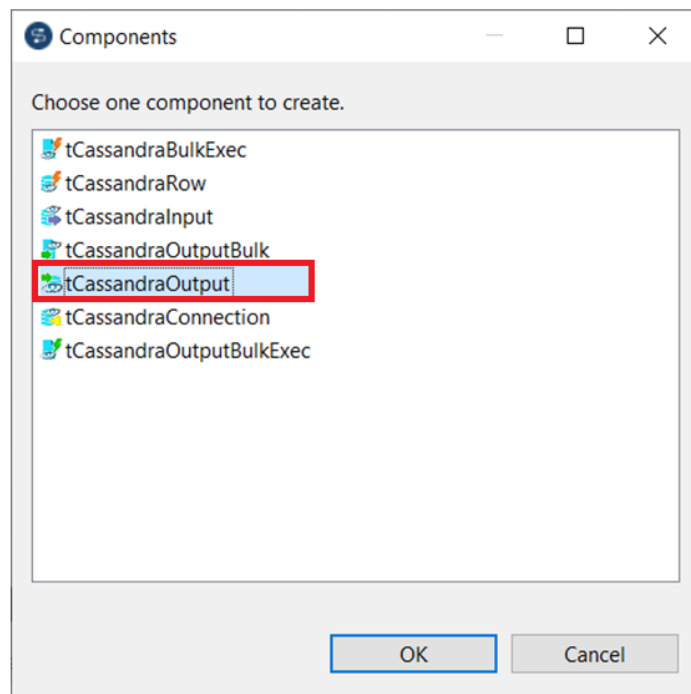
On the Palette tab located at the right-hand side area of the Talend window, search for **tMap**, then select and drag it to the workspace, as indicated. Then, drag an arrow from the *dim_sensors* component that was created earlier to the *tMap_1* job.
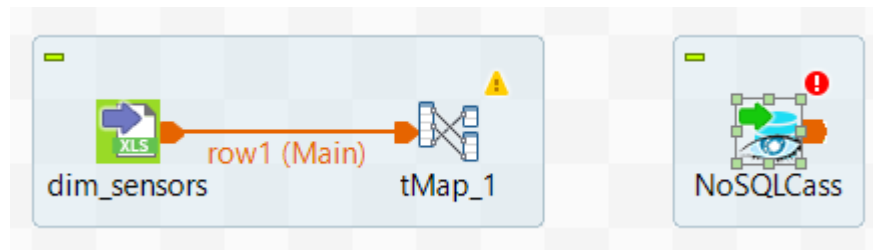


Then, under the Repository tab click on NoSQL Connections and drag the NoSQLCass component we made earlier to the job. Select *tCassandraOutput* and click "OK."
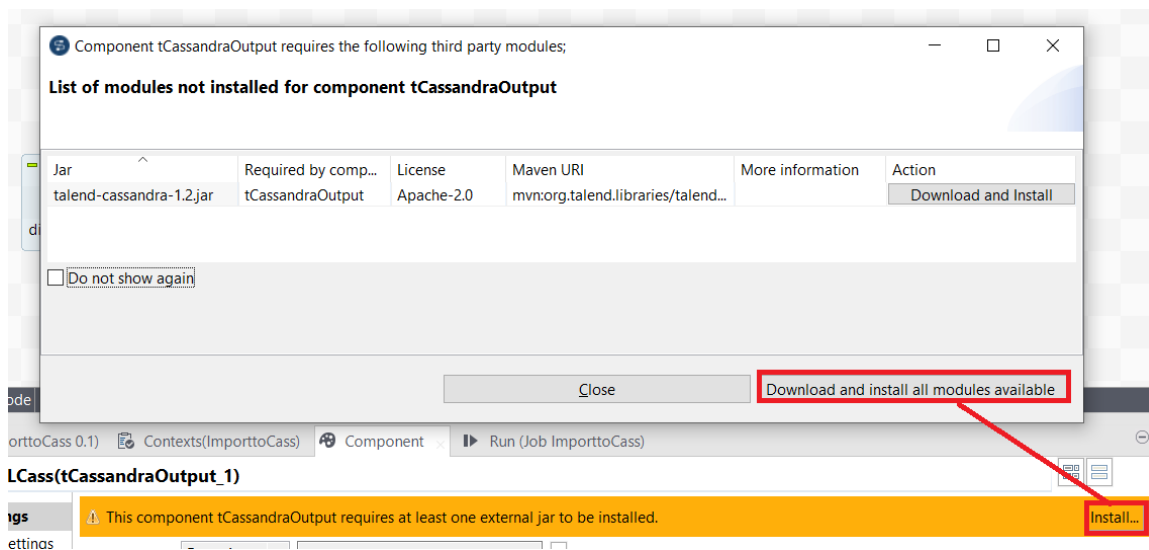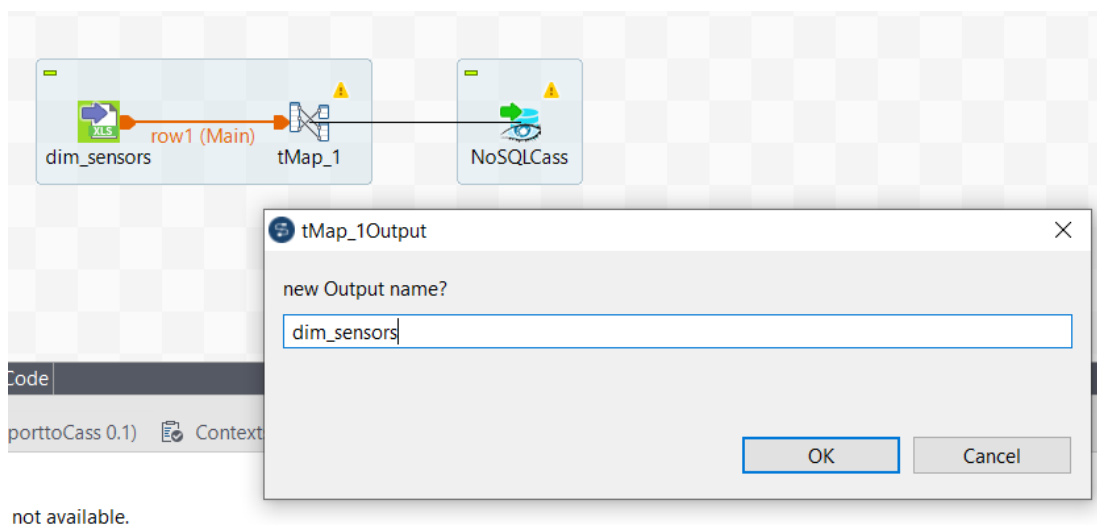
The workspace should appear like the figure below.



Since the *tCassandraOutput* component requires at least one external jar file to be installed, select *Install…* then click on the *Download and install all modules available* button to download the required modules.



Next, drag an arow from the *tMap* component to the *tCassandraOutput*. There will be a prompt to name the tMap output. We named this output as *dim_sensors* but you may name it whatever you want.

Next, double click on the *tCassandraOutput* component to check its basic settings. Make sure you selected "Built-In" for **PROPERTY**. After, under Keyspace configuration, properly indicate the **Keyspace** we made earlier which would be *sensors*. Then, under Column family configuration, indicate the **Column family** which would be *dim_sensors* and select "Create column family if does not exist" for **Action on column family**. Lastly, click on the button with the three dots to edit the table schema.



In the schema, we need to make sure that sensor_id is our primary key. Click on "OK" once finished.

Go back to your job and double-click on the tMap component to open the Schema editor. Drag every component from row1 to the table (dim_sensors) we have for Cassandra.



Now, we will run the job. The emphasized portion of the screenshot below indicates that all 15 rows from the excel file were imported successfully to our Cassandra database.

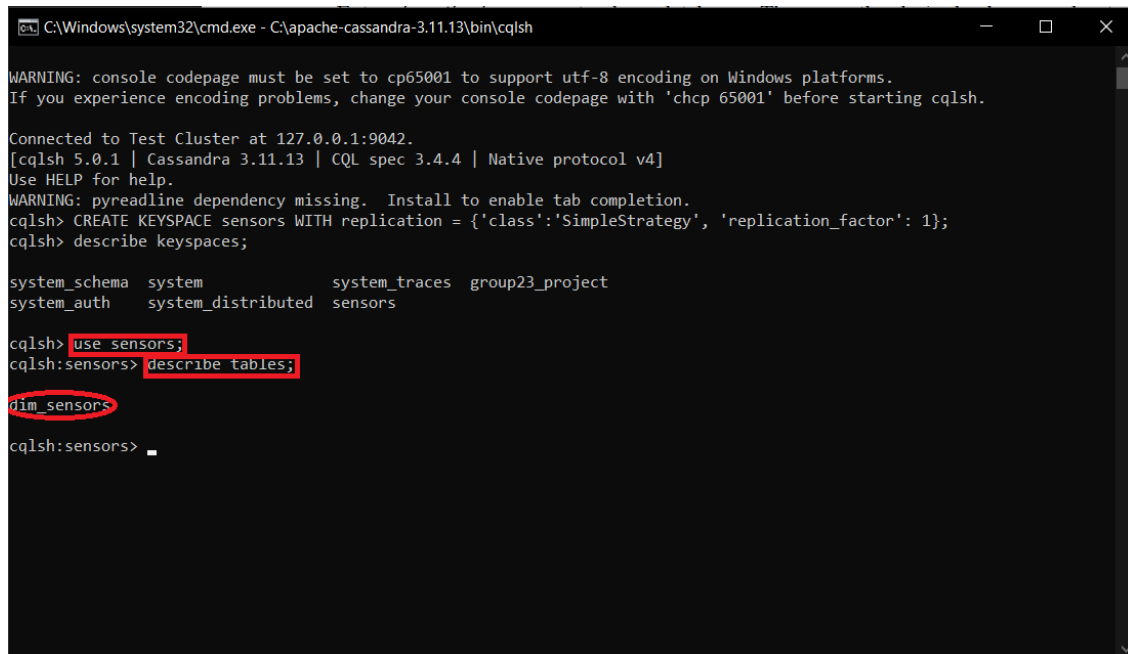## V.     Verifying the Importing of data to Cassandra

To verify that the import from the dataset to a Cassandra table via Talend was successful, run the following commands on Cassandra's CQL shell:

`use sensors;`

*This command opens or uses the database in which the table was saved.*

`describe tables;`

*This command shows or describes the tables saved in the database sensors.*



Note: If all the configurations made previously were correct, the table `dim_sensors` should appear after executing the `describe tables;` command.

Now, run the following command to see if table dim_sensors has been populated:

`select * from dim_sensors;`



This proves that our Talend project was a success, and we were able to meet the objectives of the activity.

# REFERENCES

*CREATE KEYSPACE*. (n.d.). DataStax. Retrieved from *https://docs.datastax.com/en/cql-oss/3.3/cql/cql_reference/cqlCreateKeyspace.html*

*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

*Read Excel File in Talend*. (n.d.). Tutorial Gateway. Retrieved from https://www.tutorialgateway.org/read-excel-file-in-talend/