Group 5: Nilo, Daniella Kim N.                      Course/Section: CSS182-2/CS-O

       Peñaflor, Rowencell A.
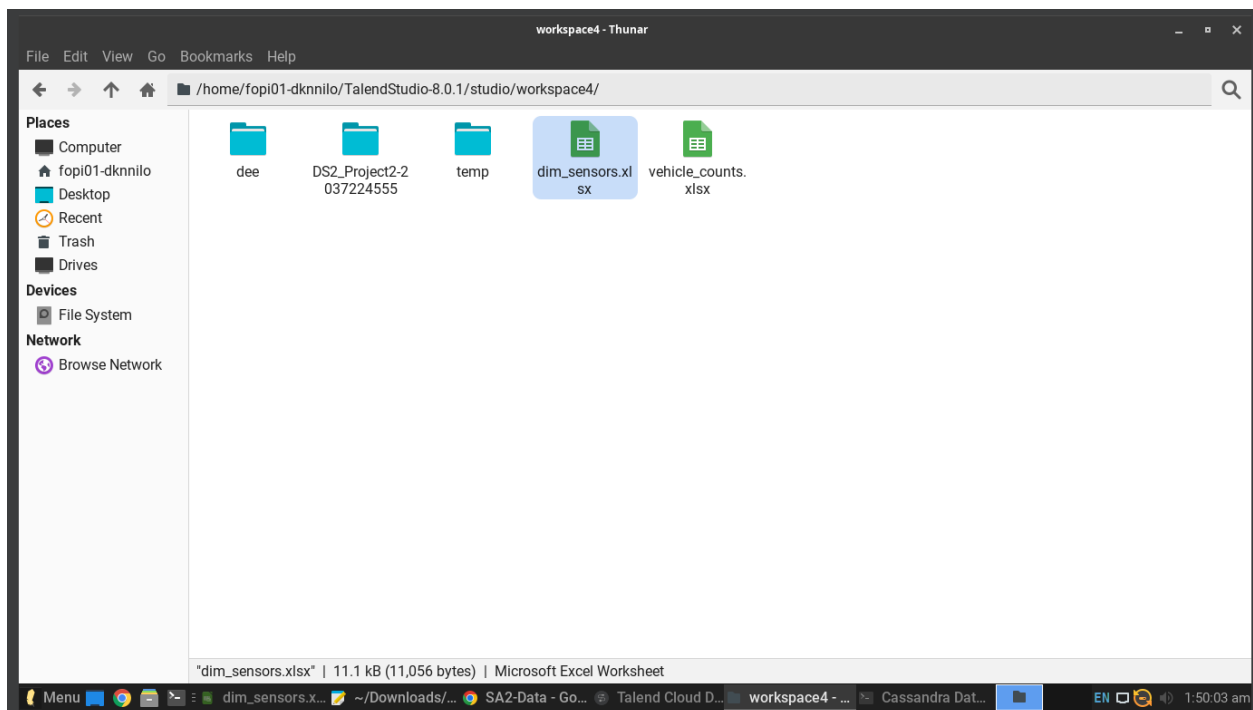
---

M2-SA2: Excel File to Cassandra Migration Process using Talend

## DOWNLOADING AND EXAMINING THE EXCEL FILE

Before doing the activity, we downloaded the excel file, dim_sensor.xls, first and then checked its contents. The Excel file contains 7 columns (`sensor_id, zone_id, node_id, sensor_type, sensor_location, survery_point, and road_monitored`) and 15 entries of sensor_id.

## CREATING CASSANDRA KEYSPACE AND TABLE

After checking the file, we utilized the `DESCRIBE KEYSPACES` command to key the available keyspaces. This will help determine if the `sensor_data` keyspace already exists. After making sure the keyspace did not exist, we created the keyspace using the command, `CREATE KEYSPACE IF NOT EXISTS sensor_data WITH replication = {'class': 'SimpleStrategy','replication_factor':1};`. We also created a table using the command, `CREATE TABLE IF NOT EXISTS sensor_data.dim_sensors (sensor_id TEXT PRIMARY KEY, zone_id INT, node_id INT, sensor_type TEXT, sensor_location TEXT, survey_point TEXT, road_monitored TEXT);`. The table columns used in the commands are taken from the columns on the Excel file.

## CREATING A NEW JOB FOR EXCEL TO CASSANDRA

After creating a keyspace and table, we proceeded with creating a new job for importing Excel file to Cassandra database. In this step, we first click the new job button. We input the name of the project (ImportExcelToCassandra), its purpose (Academic Project), its status (development), and then clicked the Finish button.



## CREATING A JDBC DB CONNECTION

In creating this connection, we first clicked the JDBC button for the Database Connection window and entered `CassandrSensorDB` for name and Academic Project for Purpose in the Create New database window. We clicked the next button, and it redirected us to another page. On this page, we entered the JDBC URL (`jdbc:cassandra://127.0.0.1:9042/sensor_data`) and

clicked the `cdata.jbdc.cassandra.jar` for the driver. We clicked the finish button to end the process. Another prompt showed our connection process is successful.

## IMPORTING PROCESS

Clicking `OK` on the last prompt brought us back to the Job `ImportCassandraToExcel 0.1` workspace. We dragged the `CassandraSensorDB 0.1` and `tFileInputExcel` from the left palette to the workspace. The `tFileInputExcel` is utilized because this component reads an Excel file row by row, which splits them up into fields using regular expressions. It also sends the fields as defined in the Cassandra database schema to the next component in the Talend job.

The next process is to configure the `tDBOutput_1(CassandraSensorDB)` by clicking it and adjusting the `Designer` windows upwards and write `sensor_data.dim_sensors` in the `Table` textbox. After configuring the `tDBOutput_1(CassandraSensorDB)`, the component, `tFileInputExcel_1`, is configured by typing "`dim_sensors-1`" into the `Sheet list` textbox. Then, create a schema manually in `tFileInputExcel_1` to arrange column names and click the `OK` button.

**Screenshot 1 (top):**

Talend Cloud Data Management Platform (R2024-10) | Kim, Daniella | Group5_SA1 (Connection: Signed in: Cloud (Asia Pacific on AWS))

File  Edit  View  Window  Help

Feature Manager

*Repository

*Job ExportCassandraToExcel 0.1 ×    *Job ImportExcelToCassandra 0.1 ×

Palette

Designer  Code

tFileInputExcel

GIT: Group5_SA1/main  (Local N

Job  Context (Im  Component  Run (Job Im  Cloud Artifa  Talend Git S  Modules

Favorites
Recently Used

- Metadata
  - Db Connections
    - > CassandraDB 0.1
    - > CassandraSensorDB 0.1
      - Queries
      - Synonym schemas
      - Table schemas
      - View schemas
  - File delimited
  - File positional
  - File regex
  - File XML
  - File Excel

CassandraSensorDB(tDBOutput_1)

| Basic settings | Schema | Built-In ↓ | Edit schema | Guess schema |

Advanced settings
Dynamic settings
View
Documentation

JDBC URL  "jdbc:cassandra://127.0.0.1:9042/sensor_data"

Drivers

Driver

cdta.jdbc.cassandra.jar

⊕ ✕ ↑ ↓

Driver class  "cdata.jdbc.cassandra.CassandraDriver"

User ID

Password

☐ Specify a data source alias

Table  sensor_data.dim_sensors

Data action  Insert ↓

☐ Clear data in table

☐ Die on error

Recently Used
- tFileInputExcel
- tDBOutput(JDBC)
- tDBInput
- tMap
- tFileOutputExcel

Big Data
Business Intelligence
Business
Cloud
Custom Code
Data Quality
Databases
DotNET
ELT
  Combined SQL
  Connections
  Map
  SQLTemplate

Outline ×    Code Vie

- tDBOutput_1 (CassandraSensorDB)
- tFileInputExcel_1

Menu  dim_sensors.x...  ~/Downloads/...  SA2-Data - Go...  Talend Cloud D...  workspace4 - ...  Cassandra Dat...  EN  1:31:14 am

---

**Screenshot 2 (bottom):**

Talend Cloud Data Management Platform (R2024-10) | Kim, Daniella | Group5_SA1 (Connection: Signed in: Cloud (Asia Pacific on AWS))

File  Edit  View  Window  Help

Feature Manager

*Repository

*Job ExportCassandraToExcel 0.1 ×    *Job ImportExcelToCassandra 0.1 ×

Palette

Designer  Code

tFileInputExcel

GIT: Group5_SA1/main  (Local N

Job  Context (Im  Component  Run (Job Im  Cloud Artifa  Talend Git S  Modules

Favorites
Recently Used

- Metadata
  - Db Connections
    - > CassandraDB 0.1
    - > CassandraSensorDB 0.1
      - Queries
      - Synonym schemas
      - Table schemas
      - View schemas
  - File delimited
  - File positional
  - File regex
  - File XML
  - File Excel

tFileInputExcel_1

| Basic settings | Property Type | Built-In ↓ |

Advanced settings
Dynamic settings
View
Documentation

☐ Read excel2007 file format (xlsx)

File name/Stream  "/home/fopi01-dknnilo/TalendStudio-8.0.1/studio/workspace4/dim_se ... *

☐ All sheets

Sheet list

| Sheet (name or position) | ☐ Use Regex |
|---|---|
| "dim_sensors-1" | ☐ |

⊕ ✕ ↑ ↓

Header  1    Footer  0    Limit

☐ Affect each sheet(header&footer)

☐ Die on error

First column  1    Last column

The variable attached to this parameter is:__LAST_COLUMN_

Recently Used
- tFileInputExcel
- tDBOutput(JDBC)
- tDBInput
- tMap
- tFileOutputExcel

Big Data
Business Intelligence
Business
Cloud
Custom Code
Data Quality
Databases
DotNET
ELT
  Combined SQL
  Connections
  Map
  SQLTemplate

Outline ×    Code Vie

- tDBOutput_1 (CassandraSensorDB)
- tFileInputExcel_1

Menu  dim_sensors.x...  ~/Downloads/...  SA2-Data - Go...  Talend Cloud D...  workspace4 - ...  Cassandra Dat...  EN  1:35:30 am

The next step is to connect the `tFileInputExcel_1` to `CassandraSensorDB(tDBOuput_1)` by manually establishing a connection line between the two components. Run the job and verify if the data is imported by querying in Cassandra. Use the command `DESCRIBE_KEYSPACE` to determine if the `sensor_data` keyspace is present. After confirming it is present, use the command `DESCRIBE KEYSPACE sensor_data;` to view the necessary information about the keyspace and the table present in the keyspace.

Afterwards, display the data using the commands `SELECT  sensor_id, zone_id, node_id, sensor_type FROM sensor_data.dim_sensors;` to view the first 4 columns of the table and `SELECT    sensor_location,   survey_point,   road_monitored   FROM sensor_data.dim_sensors;` to view the location of the sensor, their survey point, and the road assigned to them for monitoring.

The order of the queries in the image is unsorted since Cassandra is designed for rapid distributed data access rather than preserving the order of the responses. The clustering columns sort the data in Cassandra. If the clustering columns are not defined, Cassandra is unable to sort rows within a partition. More specifically, if the table simply has a partition key, such as a PRIMARY KEY, the rows within the partition are not sorted in a certain order by default. Since the sorting between partitions is not supported, the queried data from different partitions appears to be unsorted.

**REFERENCES:**

Qlik Talend. (n.d.). *tFileInputExcel*. tFileInputExcel | Talend Components for Jobs Help. https://help.qlik.com/talend/en-US/components/8.0/excel/tfileinputexcel

Rowe, W. (2019, January 21). *Partition key vs composite key vs clustering columns in cassandra*. BMC Blogs. https://www.bmc.com/blogs/cassandra-clustering-columns-partition-composite-key/