

# GSC Dashboard - 4. Data Loading with Talend

## 4.1 - Talend Integration

### Source data integration with Talend ETL tool

- **Xtract Server (WBP):**
  - Talend connects to the Xtract server and trigger extraction job using HTTP calls.
  - Xtract generates results as CSV files stored in memory.
  - Talend retrieves files, apply some transformation and then loads them into Google Cloud Storage.
- **Google Sheets (Hedges, wap solid fuels, CO2):**
  - Talend integrates with Google Sheets to extract some static files that will be use to apply mappings or lookup to BW tables.
  - It retrieves this data from Google Sheets.

### Data Transformation and Loading to Google BigQuery:

- Once data from all sources is available in Google Cloud Storage as files, Talend proceeds with data transformation and loading.
- Talend performs data transformations as needed, including cleansing, mapping, and structuring the data for consistency.
- The transformed data is loaded into various stages, operational data stores (ODS), and data mart tables within Google BigQuery.
- These tables are organized to facilitate efficient querying and reporting for energy optimization purposes.

By utilizing Talend for data extraction, transformation, and loading (ETL), the web app ensures that data from diverse sources is collected, processed, and structured for analysis and reporting within Google BigQuery, enabling users to make informed decisions based on up-to-date and accurate data.

## 4.2 - Source Data Extraction

<b>Main jobs for source extraction</b>	<ul style="list-style-type: none"> <li>• <b>F100_GSC_DATA_EXTRACT</b></li> <li>• <b>J130_mainExtract_Files_Src_to_GCS</b></li> <li>• <b>J120_GSC_mainExtract_BWQ_Data</b></li> <li>• <b>J120_GSC_Extract_BW_data_Source_14_EXT</b></li> <li>• <b>J110_GSC_mainExtract_Gsheet_data</b></li> </ul>	<p style="text-align: center;"><a href="#">--to the top --</a></p>
<b>Job description by steps</b>	<b>Job design</b>	

- [4.1 - Talend Integration](#)
- [4.2 - Source Data Extraction](#)
  - [F100\\_GSC\\_DATA\\_EXTRACT](#)
  - [J130\\_mainExtract\\_Files\\_Src\\_to\\_GCS](#)
  - [J120\\_GSC\\_mainExtract\\_BWQ\\_Data](#)
  - [J120\\_GSC\\_Extract\\_BW\\_data\\_Source\\_14\\_EXT](#)
  - [J110\\_GSC\\_mainExtract\\_Gsheet\\_data](#)
  - [J110\\_GSC\\_mainExtract\\_Gsheet\\_data](#)
  - [J130\\_GSC\\_mainExtract\\_Files\\_Src\\_to\\_GCS](#)
  - [J120\\_GSC\\_Extract\\_BW\\_data\\_Source\\_14\\_EXT](#)
  - [J120\\_GSC\\_mainExtract\\_BWQ\\_Data\\_Source](#)
- [4.2 - Load Staging and ODS](#)
  - [F300\\_GSC\\_STG\\_TO\\_ODS](#)
  - [J300\\_GSC\\_STG\\_TO\\_ODS](#)
- [4.2 - Prepare DM files](#)
  - [F400\\_GSC\\_ODS\\_to\\_Prep\\_DM\\_Files](#)
  - [F430\\_GSC\\_ODS\\_to\\_DTM](#)
  - [F420\\_GSC\\_ODS\\_to\\_DTM\\_PART\\_1](#)
  - [F420\\_GSC\\_ODS\\_to\\_DTM\\_PART\\_5](#)
- [4.3 - Compute Perfect Order Rate](#)
  - [F420\\_GSC\\_ODS\\_to\\_DTM\\_PART\\_6](#)
- [4.4 - Load to DM \(calculations and transformations\)](#)
  - [F501\\_GSC\\_Files\\_To\\_FACT](#)
  - [F502\\_GSC\\_Files\\_To\\_PIVOT](#)
- [4.5 - Scheduling and Automation](#)
- [4.6 - Data Validation](#)

### Responsible & contact points:

- Application Owner - Mathieu Pourquoi
- Delivery Manager - Donia Rachidi
- Project Manager - Vitaly Fonseca
- Data Architect - Joao Marboeuf
- Data Engineer - Virgil Lissassi ; replaced by Matteo Menghetti

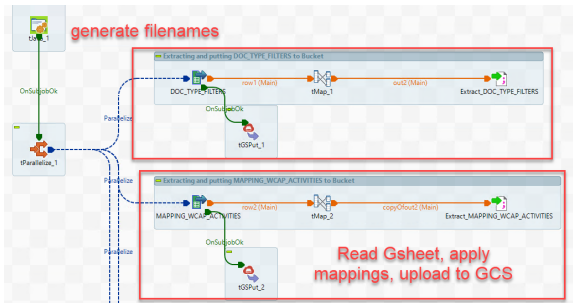
	<p>1. In the job F100_PUR_DATA_EXTRACT CT:</p> <p>a. A parallel execution of all data sources is launched</p> <p>The flow triggers the extraction of all sources needed for the Global Supply chain dashboard.</p>	
--	--	--

<p><b>Main jobs for source extraction</b></p>	<ul style="list-style-type: none"> <li>• <b>J110_GSC_mainExtract_Gsheet_data</b></li> </ul>	<p>--to the top --</p>
	<p><b>Job description by steps</b></p>	<p><b>Job design</b></p>
	<p>1. Parallel extraction of all Google Sheet sources:</p> <p>a. A parallel extraction of all Google Sheet is launched</p> <p>b. Simple transformation are done to align with specifications</p> <p>2. Upload files to GCS</p> <p>a. The files are uploaded into the GSC folder of the target bucket</p>	

<p>Main jobs for source extraction</p>	<ul style="list-style-type: none"> <li>J130_GSC_mainExtract_Files_Src_to_GCS</li> </ul>	<p>--to the top --</p>
--	---	------------------------

	<p>Job description by steps</p>	<p>Job design</p>
--	---------------------------------	-------------------

- Filename creation:
  - A specific filename following our standard naming convention is created to properly upload files into Cloud Storage

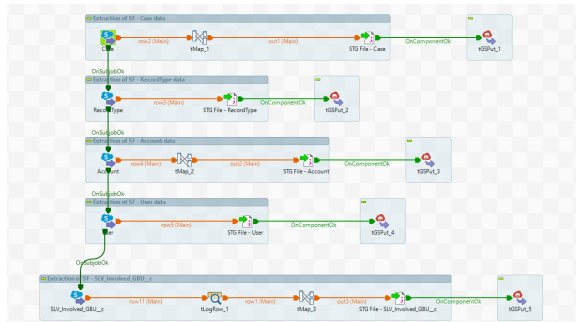
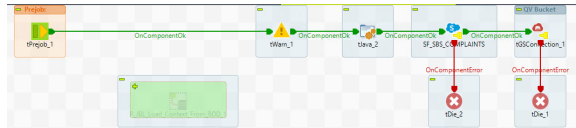


	<p>2. Parallel GSC upload</p> <ol style="list-style-type: none"> <li>a. Each Google Sheet is read</li> <li>b. Simple transformations are done to align with specifications resulting in a CSV stored in the <i>Tmp</i> folder</li> <li>c. Each file is uploaded into the GSC folder of the target bucket</li> </ol>	
--	---	--

<b>Main jobs for source extraction</b>	<ul style="list-style-type: none"> <li>• J120_GSC_Extract_BW_data_Source_14_EXT</li> </ul>	<a href="#">--to the top --</a>
	<b>Job description by steps</b>	<b>Job design</b>

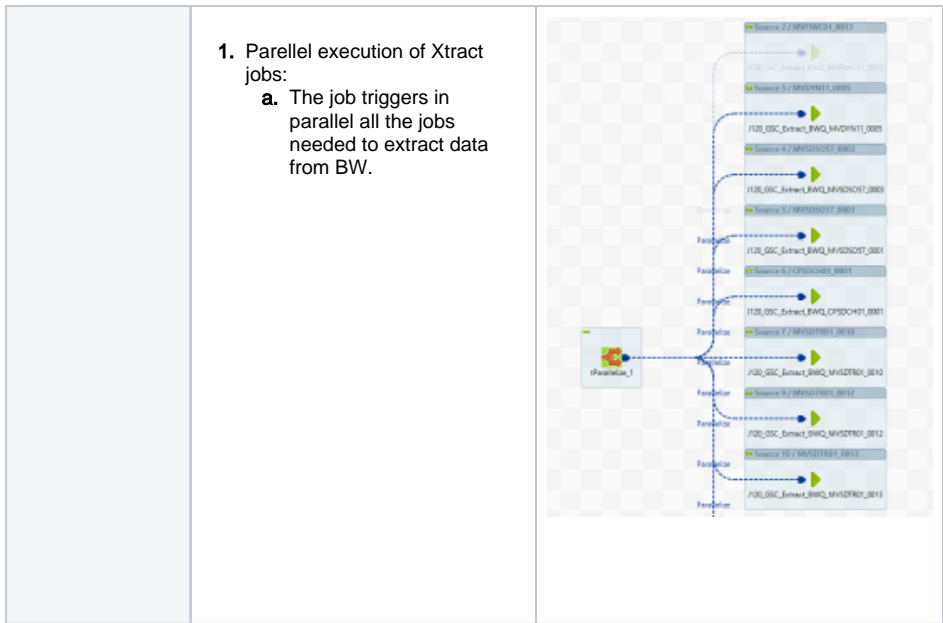
1. Login to Salesforce:

a. The job starts by establishing a connection to Salesforce using the accounts-qlikview-support.com.crm



2. Extraction of the Salesforce objects
  - a. The list of Salesforce objects is extracted in sequence
  - b. Depending on the source some mappings can be applied
  - c. Each source is uploaded to GCS as a CSV file

Main jobs for source extraction	<ul style="list-style-type: none"> <li>• J120_GSC_mainExtract_BWQ_Data_Source</li> </ul>	--to the top --
	Job description by steps	Job design



4.2 - Load Staging and ODS

<p><b>Main jobs for loading Staging</b></p>	<ul style="list-style-type: none"> <li>• F200_GS_C_Gcs_T_o_Stg</li> <li>• J200_GS_C_Gcs_T_o_Stg</li> </ul>	<p>--to the top --</p>
	<p><b>Job description by steps</b></p>	<p><b>Job design</b></p>
	<p>1. First the job retrieves the value of the Business date from the file business_date.csv file. This is to target only the files of the current execution in GCS.</p> <p>2. The flow is used to call the job J200_GS_C_Gcs_To_Stg.</p>	

• J2  
00  
\_G  
\_SC  
\_G  
\_cs  
\_T  
\_o  
\_Stg

1. The job launches a parallel load per source type into the Staging tables

The screenshot displays a data pipeline configuration. The top section shows a flow diagram with components: FixedFlowInput\_5, iFunScheme\_5, iWm\_2, iStgQueueSQLFun\_14, iSQLM\_10, and iStg\_CCS\_In\_Staging\_Local. Below this is a configuration window for 'FixedFlowInput\_5' with a 'Validation Rules' table.

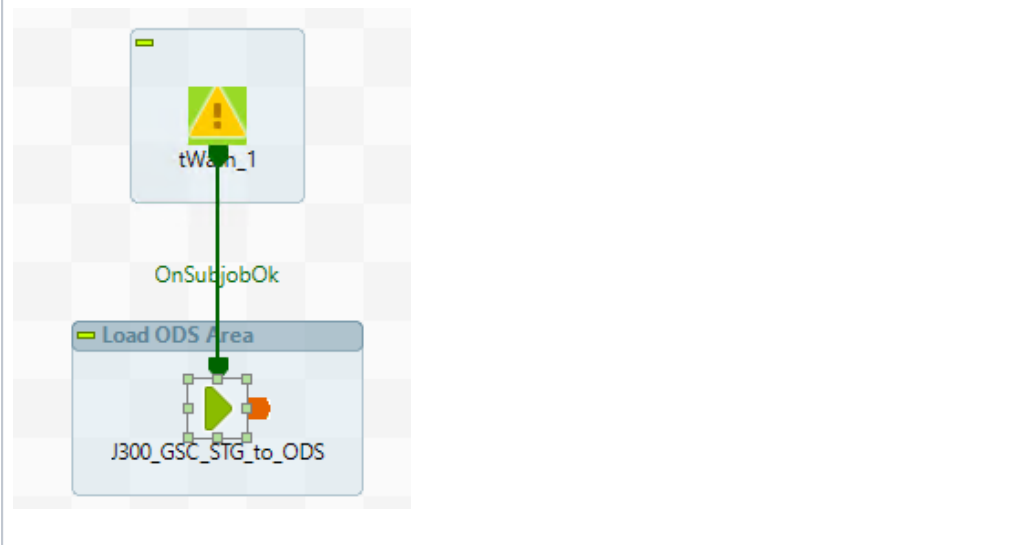
stg_table_name	file_key_profile	separator	escape_separator	width
"\$RE_0000_0007_F_W_GSC..."	"\$SC-\$SC_FT_0000_0000_0007..."	" "	"\""	1
"\$RE_0000_0000_0012_F_W_GSC..."	"\$SC-\$SC_FT_0000_0000_0012..."	" "	"\""	1
"\$RE_0000_0000_0008_F_W_GSC..."	"\$SC-\$SC_FT_0000_0000_0008..."	" "	"\""	1
"\$RE_0000_0000_0009_F_W_GSC..."	"\$SC-\$SC_FT_0000_0000_0009..."	" "	"\""	1
"\$RE_0000_0000_0012_F_W_GSC..."	"\$SC-\$SC_FT_0000_0000_0012..."	" "	"\""	1
"\$RE_0000_0000_0012_F_W_GSC..."	"\$SC-\$SC_FT_0000_0000_0012..."	" "	"\""	1
"\$RE_0000_0000_0011_F_W_GSC..."	"\$SC-\$SC_FT_0000_0000_0011..."	" "	"\""	1
"\$RE_0000_0000_0011_F_W_GSC..."	"\$SC-\$SC_FT_0000_0000_0011..."	" "	"\""	1
"\$RE_0000_0000_0009_F_W_GSC..."	"\$SC-\$SC_FT_0000_0000_0009..."	" "	"\""	1
"\$RE_0000_0000_0014_F_W_GSC..."	"\$SC-\$SC_FT_0000_0000_0014..."	" "	"\""	1
"\$RE_0000_0000_0008_F_W_GSC..."	"\$SC-\$SC_FT_0000_0000_0008..."	" "	"\""	1

- a. A list of the possible CSV file with the associated Staging table name, and separator is used to create a sequence
- b. Each staging table is truncated before a new upload
- c. The job gets all the files starting with the desired prefix and for each of them calls the standard job J001\_GC\_S\_to\_ST\_AGI\_NG\_LOC\_AL

<p>Main jobs for loading the ODS</p>	<ul style="list-style-type: none"> <li>• F300_G SC_ST G_TO_ODS</li> <li>• J300_G SC_ST G_TO_ODS</li> </ul>	<p>--to the top --</p>
--------------------------------------	--	------------------------

	<p>Job description by steps</p>	<p>Job design</p>
--	---------------------------------	-------------------

1. The flow is used to call the job J300\_G SC\_ST G\_TO\_ODS

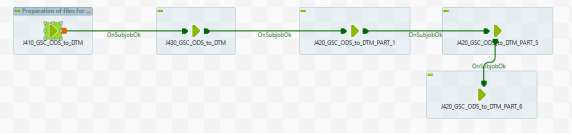


1. Parallel loading of the ODS table by source type
  - a. A parallel loading of ODS table is done by source type (BW, Google Sheet, Salesforce)
2. Load in sequence all ODS tables of a source
  - a. Give a source type the list of staging and ODS tables is hardcoded.
  - b. ODS table is truncated
  - c. The standard job J001\_ST\_AGING\_TO\_ODS is called to load the ODS table

	<p>Given the fact that staging and ods are identical and tables are always truncated, this step is unnecessary.</p>	
--	---	--

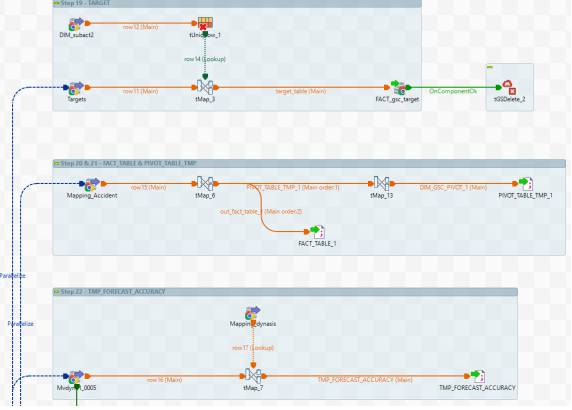
4.2 - Prepare DM files

<p>Main jobs for building DM</p>	<ul style="list-style-type: none"> <li>• F400_G SC_OD S_to_P rep_DM _Files F400_PU RCHASI NG_STEP <ul style="list-style-type: none"> <li>◦ F41 0_G SC_ods _to_ DTM</li> <li>◦ F42 0_G SC_ods _to_ DTM _Par t_1</li> <li>◦ F42 0_G SC_ods _to_ DTM _Par t_5</li> <li>◦ F42 0_G SC_ods _to_ DTM _Par t_6</li> <li>◦ F43 0_G SC_ods _to_ DTM</li> </ul> </li> </ul>	<p style="color: blue;">--to the top --</p>
	<p>Job description by steps</p>	<p>Job design</p>

	<p>1. This job is used to call a series of subjob that creates files needed to create the dimensions and fact tables</p>	
--	--	--

<p>Prepare DM file</p>	<ul style="list-style-type: none"> <li>F410_GCS_ODS_t_o_DTM</li> </ul>	<p>--to the top --</p>
------------------------	--	------------------------

	<p>Job description by steps</p>	<p>Job design</p>
--	---------------------------------	-------------------

<p>The job computes steps 19-28 of the <a href="#">specification document</a></p> <ol style="list-style-type: none"> <li>The job reads data from the ODS table and apply the necessary mappings;</li> <li>The job output PIVOT and FACT files which are stored on the remote engine and will be uploaded to the DM in the J500 jobs.</li> <li>The jobs creates the table <b>FACT_gsc_target</b></li> </ol>		
--	--	---

<b>Main jobs for source extraction</b>	<ul style="list-style-type: none"> <li>• F430_G CS_OD S_to_D TM</li> </ul>	--to the top --
--	--	-----------------

	<b>Job description by steps</b>	<b>Job design</b>
--	---------------------------------	-------------------

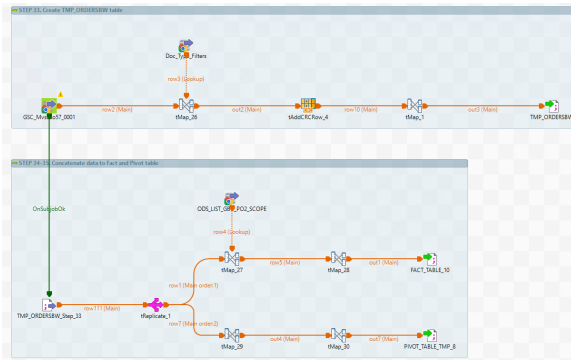
<p>The job computes steps 41-49 of the <a href="#">specification document</a></p> <ol style="list-style-type: none"> <li>1. The job reads data from the ODS table and apply the necessary mappings;</li> <li>2. The job output PIVOT and FACT files which are stored on the remote engine and will be uploaded to the DM in the J500 jobs.</li> <li>3. The job also builds the files necessary to create the airfreight KPI following <a href="#">these specifications</a></li> </ol>		
---	--	--

<b>Main jobs for source extraction</b>	<ul style="list-style-type: none"> <li>• F420_G SC_OD S_to_D TM_PA RT_1</li> </ul>	--to the top --
--	--	-----------------

	<b>Job description by steps</b>	<b>Job design</b>
--	---------------------------------	-------------------

The job computes steps 33,34, and 35-49 of the [specification document](#)

1. The job reads data from the ODS table and apply the necessary mappings;
2. The job output PIVOT and FACT files which are stored on the remote engine and will be uploaded to the DM in the J500 jobs.



<p><b>Main jobs for source extraction</b></p>	<ul style="list-style-type: none"> <li>• F420_G SC_OD S_to_D TM_PA RT_5</li> </ul>	<p>--to the top --</p>
	<p><b>Job description by steps</b></p>	<p><b>Job design</b></p>

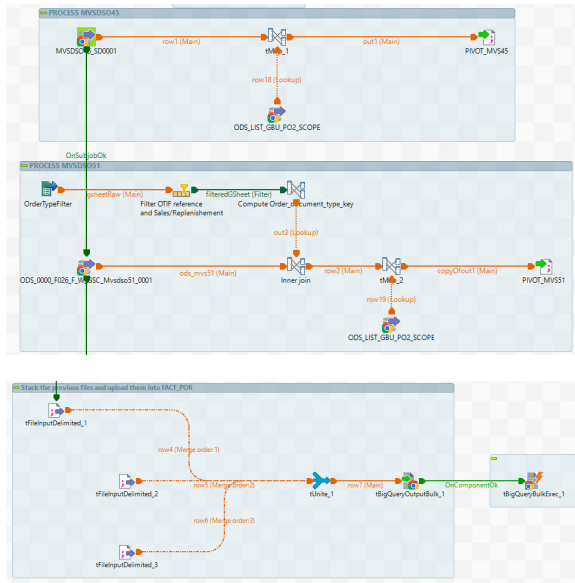


The job creates the table *FACT\_POR* which is displayed by the Perfect Order rate dashboard.

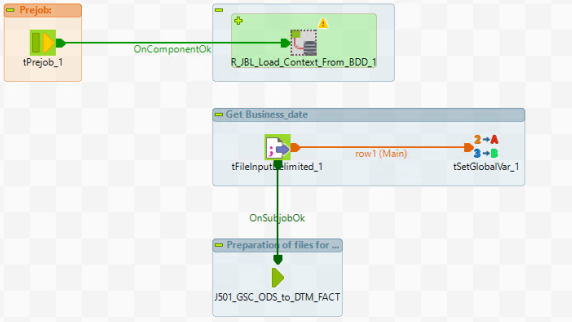
1. The job reads data from the ODS table and apply the necessary mappings. The main ODS tables are:

- a. ODS\_000\_0\_F\_025\_F\_W\_GS\_C\_M\_VSD\_SO4\_5\_S\_D00\_01
- b. ODS\_000\_0\_F\_026\_F\_W\_GS\_C\_M\_vsdso51\_0001
- c. ODS\_000\_0\_F\_020\_F\_W\_GS\_C\_Sf\_Co\_mplaint\_Case

2. Each step outputs a FACT file which is stored in the *Tmp* folder
3. The files are stacked to create a unique output bulk file which is used to create the *FACT\_POR* table.



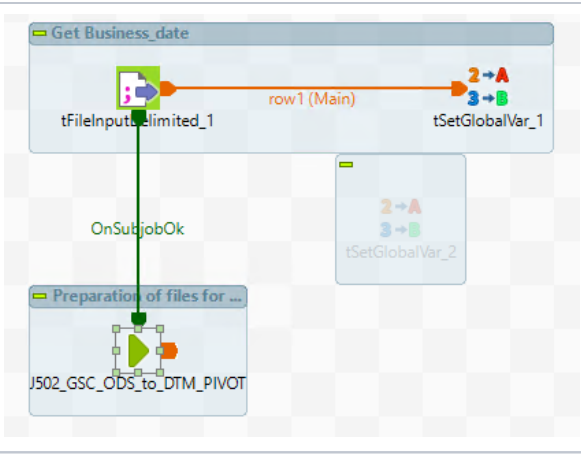
#### 4.4 - Load to DM (calculations and transformations)

<p><b>Main jobs for source extraction</b></p>	<ul style="list-style-type: none"> <li>• F501_GSC_File_s_To_FACT</li> <li>• 501_GSC_ODS_to_DTM_FACT</li> </ul>	<p>--to the top --</p>
	<p><b>Job description by steps</b></p>	<p><b>Job design</b></p>
	<p>This jobs is used to call the job J501_GSC_ODS_to_DTM_FACT</p>	 <p>The diagram illustrates the job design for 'Prejob_1'. It consists of the following components and connections:</p> <ul style="list-style-type: none"> <li><b>Prejob_1</b> (Prejob task) is connected to <b>RuBl_Load_Context_From_BDD_1</b> (Data Flow Task) via the <b>OnComponentOk</b> event.</li> <li><b>RuBl_Load_Context_From_BDD_1</b> is connected to <b>Get Business_date</b> (Data Flow Task) via the <b>row1 (Main)</b> connection.</li> <li>Inside the <b>Get Business_date</b> task, <b>tFileinputLimited_1</b> (Source) is connected to <b>tSetGlobalVar_1</b> (Destination).</li> <li><b>Get Business_date</b> is connected to <b>Preparation of files for...</b> (Data Flow Task) via the <b>OnSubjobOk</b> event.</li> <li><b>Preparation of files for...</b> is connected to <b>J501_GSC_ODS_to_DTM_FACT</b> (Data Flow Task).</li> </ul>

<p><b>501_GSC_ODS_to_DTM_FACT</b></p>	<ol style="list-style-type: none"> <li>1. First the job truncates the table <i>DM.FACT_gsc_c_tmp</i></li> <li>2. The job reads all files whose name starts with <i>FACT_FILE</i> inside the <i>InOut</i> folder and uploads them into the <i>FACT_gsc_c_tmp</i> table. If files is successfully inserted the local file is deleted from the remote engine</li> <li>3. The job truncates the table <i>DM.FACT_gsc</i></li> <li>4. The job copies data from <i>FAC T_gsc_tm p</i> into <i>FA CT_gsc</i></li> <li>5. The job truncates the table <i>DM.FACT_gsc_c_tmp</i></li> </ol>	
---------------------------------------	---	--

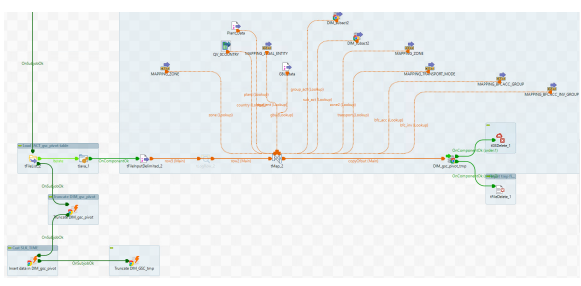
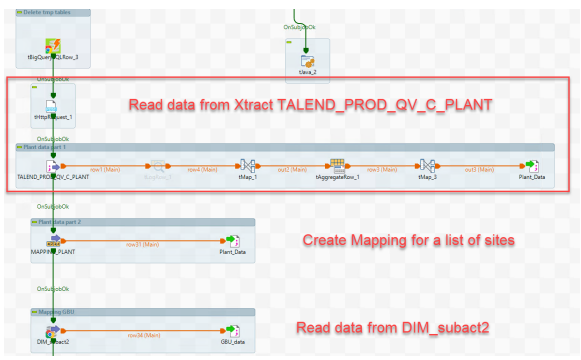
<p><b>Main jobs for source extraction</b></p>	<ul style="list-style-type: none"> <li>• F502_GSC_Files_To_PIVOT</li> <li>• J502_GSC_ODS_to_DTM_PIVOT</li> </ul>	<p>--to the top --</p>
	<p><b>Job description by steps</b></p>	<p><b>Job design</b></p>

1. This jobs is used to call the job J502\_GSC\_ODS\_to\_DTM\_PIVOT



- J502\_GSC\_ODS\_to\_DTM\_PIVOT

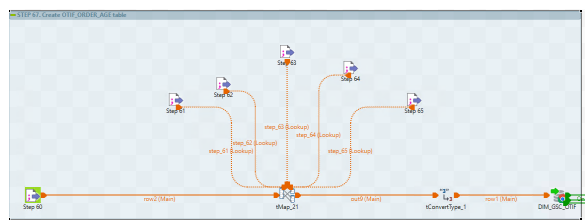
1. First the job truncates the table DM.DIM\_gsc\_pivot\_tmp
2. The job gets data used to perform lookups
  - a. The job executes the TALEND\_PROD\_QV\_C\_PLANT Xtract job
  - b. Creates mapping for a list of sites
  - c. Read data from DIM\_subact2
  - d. Creates a mapping for some old GBU codes



3. The job reads all files whose name starts with *PIVO T\_* and *T PIVOT\_* inside the InOut folder and uploads them into the *DIM\_gsc\_tmp* table. If files is successfully inserted the local file is deleted from the remote engine. Some joins are done using the sources previously extracted
4. The job truncates the table *DM.DIM\_gsc\_pivot*
5. The job copies data from *DIM\_gsc\_pivot\_tmp* into *DIM\_gsc\_pivot*
6. The job truncates the table *DM.DIM\_gsc\_pivot\_tmp*

<b>Main jobs for source extraction</b>	<ul style="list-style-type: none"> <li>• F503_GSC_Files_To_OTIF</li> <li>• J503_GSC_ODS_</li> </ul>	<a href="#">--to the top --</a>
	<b>Job description by steps</b>	<b>Job design</b>

1. The flow is used to build the table *DIM\_gsc\_otif*
2. The job reads data from the Step\_60.csv file computed in the job and joins it with the fields Step\_{61..65}. These are the output of job J420\_GSC\_O DS\_to\_DTM\_P ART\_5
3. The combined file is used to overwrite the content of the table *DIM\_gsc\_otif*



## 4.5 - Scheduling and Automation

TMC - PL\_QV\_TO\_TABLEAU\_GSC - Daily at 9 PM.

## 4.6 - Data Validation