

4. Data Loading with Talend

Responsible & contact points:

- Alessandro Mainardi - Project Owner
- Simon Bourguignon - Delivery Manager
- Alba Carrero/ Gaetan Frenoy - Product Owner
- Rui Ferraz - Project Manager

4.1 - Talend Integration

Source data integration with Talend ETL tool

Connection details:

- **FTP Server (Meteologica):** <ftp.meteologica.com> (incremental load)
 - Talend connects to the FTP server where files containing current and future electricity prices for Spain, Italy, France, and Germany are located.
 - Talend retrieves these files from the FTP server.
 - The retrieved files are then loaded into Google Cloud Storage.
 - We need to use GCP RE on this source because in order to access Meteologica, it is required statistic IP in order to register the whitelist for data access.
- **Postgres Database (Vendohm):** ~~ses-gcp-sandbox.eu.vendohm~~ **Decommission replaced by CSV from Dataiku**
 - ~~Talend establishes a direct connection with the Vendohm PostgreSQL database that stores sensor values identified by specific curve IDs.~~
- **Dataiku files (full load)**
 - Vendohm is decommissioned and IRM energy data need to be calculated many rules. Therefore, project team decide to get the final calculation of energy from Dataiku.
 - Dataiku will generate energy files (Spot hour, Spot day and Forward combines to 1 consolidation file) as .gz file to DataOcean Industrial bucket (cs-ew1-prj-data-dm-industrial-[env]-staging/in/enr/irm-conso) twice a day
 - Dataiku will generate a file to replace IRM energy deal to DataOcean Industrial bucket (cs-ew1-prj-data-dm-industrial-[env]-staging/in/enr/irm-deals) once a day
 - Talend gets the files considering GCS as data source
 - This extracted data is loaded into Google Cloud Storage as files again in the standard way.
 - [Dataiku project](#)
- **Oracle Database (IRM):** ~~acow1pirmdb01.prod.aws.cloud.solvay.com~~
 - ~~Talend establishes a direct connection with the IRM Oracle database that stores energy deals identified by specific deal IDs.~~
 - ~~It extracts relevant data from the Oracle database.~~
 - ~~This extracted data is loaded into Google Cloud Storage as files.~~
- **Google Sheets (Hedges, wap solid fuels, CO2): (full load)**
 - Talend integrates with Google Sheets, where solid fuel wap, CO2 emissions, and hedging information are stored.
 - It retrieves this data from Google Sheets.
 - Similar to the other sources, this data is also loaded into Google Cloud Storage as files.
 - [CO2](#)
 - [WAP Solid Fuels](#)
 - [Energy Deals Hubs](#)
 - [Energy Deals Counter Parties](#)
 - [Energy Deals Sites Hedges](#)
 - [Vendohm Forwards](#)
 - [Vendohm Meteo](#)
 - [Vendohm Spot](#)

Data Transformation and Loading to Google BigQuery:

- Once data from all four sources (FTP server, PostgreSQL database, Oracle database, Google Sheets, Dataiku file) 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.

- 4.1 - Talend Integration
 - [J001_FTP_to_GCS-METEOROLOGICA_FRANCE](#)
 - [J002_FTP_to_GCS-METEOROLOGICA_ITALY](#)
 - [J003_FTP_to_GCS-METEOROLOGICA_GERMANY](#)
 - [J004_FTP_to_GCS-METEOROLOGICA_SPAIN](#)
 - [J005_FTP_to_GCS-METEOROLOGICA_SPAIN_ENS](#)
 - [J006_FTP_to_GCS-METEOROLOGICA_SPAIN_OBS](#)
 - [J011_GSheet_to_GCS_VENDOHM_FORWARDS](#)
 - [J012_GSheet_to_GCS_VENDOHM_SPOT](#)
 - [J013_GSheet_to_GCS_VENDOHM_METEO](#)
 - [J015_GSheet_to_GCS_IRM_HUBS](#)
 - [J016_GSheet_to_GCS_IRM_COUNTERPARTIES](#)
 - [J017_GSheet_to_GCS_IRM_SITES_HEDGES](#)
 - [J009_GSheet_to_GCS_CO2_EMISSIONS](#)
 - [J010_Postgres_to_GCS_VENDOHM Decommission](#)
 - [J010_FIL_Energy_FO_TO_ODS](#)
 - [J014_Oracle_to_GCS_IRM](#)
 - [J014_GCS_to_GCS_IRM_ENERGY_DEAL](#)
- 4.3 - Data Loading to Google BigQuery
 - [J001_Extraction_till_ODS_METEOROLOGICA_FRANCE](#)
 - [J002_Extraction_till_ODS_METEOROLOGICA_ITALY](#)
 - [J003_Extraction_till_ODS_METEOROLOGICA_GERMANY](#)
 - [J004_Extraction_till_ODS_METEOROLOGICA_SPAIN](#)
 - [J005_Extraction_till_ODS_METEOROLOGICA_SPAIN_ENS](#)
 - [J006_Extraction_till_ODS_METEOROLOGICA_SPAIN_OBS](#)
 - [J009_Extraction_till_ODS_GSHEET_CO2](#)
 - [J011_Extraction_till_ODS_GSHEET_VENDOHM_FORWARDS](#)
 - [J012_Extraction_till_ODS_GSHEET_VENDOHM_SPOT](#)

GCP Details:

Domain	GCP Project	Bucket
Industrial	prj-data-dm-industrial-[env]	cs-ew1-prj-data-dm-industrial-[env]-staging/ROBUSTIFY
Sustainability	prj-data-dm-sust-[env]	cs-ew1-prj-data-dm-sust-[env]-staging/ROBUSTIFY
Robustify	prj-data-robustify-[env]	N/A

Cloud Remote Engine = prj-talend-[dev] / ce-ew1-b-talend-gcp-remote-engine-[env]

- J013_Extraction_till_ODS_GSHEET_VEN DOHM_METEO
- J015_Extraction_till_ODS_GSHEET_IRM_HUBS
- J016_Extraction_till_ODS_GSHEET_IRM_COUNTER_PARTIES
- J017_Extraction_till_ODS_GSHEET_IRM_SITES_HEDGES
- J010_Extraction_till_ODS_VENDOHM
- 4.4 - Load to DM (calculations and transformations)
 - J001_ODS_TO_DM_FACT_ENERGY_PRICE_FORECAST
 - J002_ODS_TO_DM_DIM_METEO_INFO
 - J003_ODS_TO_DM_DIM_ENERGY_PRICE_INFO
 - J004_ODS_TO_DM_FACT_METEO_DATA
 - J005_ODS_TO_DM_FACT_ENERGY_PRICE replace with view in Robustify prj-data-robustify-dev.DM. V_FACT_energy_price_hourly
 - J006_ODS_TO_DM_FACT_CO2_EMISSIONS
 - J007_ODS_TO_DM_FACT_SOLID_FUEL_WAP
 - J009_ODS_TO_DM_FACT_ENERGY_DEALS
 - J009_ODS_TO_DM_FACT_IRM_DEALS
 - J001_DM_TO_DM_FACT_ENERGY_PRICE_HOURLY
- 4.5 - Scheduling and Automation
- 4.6 - Remark

4.2 - Source Data Extraction

<p>Main jobs for source extraction</p>	<ul style="list-style-type: none"> • J001_F TP_to_GCS-METEO LOGIC A_FRANCE • J002_F TP_to_GCS-METEO LOGIC A_ITALY • J003_F TP_to_GCS-METEO 	<p>--to the top --</p> <p>GCP Remote Engine(RE)</p> <p>Industrial Data Ocean(DO) project</p>
---	---	--

LOGIC
 A_GER
 MANY
 • J004_F
 TP_to_
 GCS-
 METEO
 LOGIC
 A_SPAIN
 • J005_F
 TP_to_
 GCS-
 METEO
 LOGIC
 A_SPAIN
 N_ENS
 • J006_F
 TP_to_
 GCS-
 METEO
 LOGIC
 A_SPAIN
 N_OBS

	Job description by steps	Job design
--	---------------------------------	-------------------

1. FTP Server Connection:
 - a. Establishes a secure connection to the Meteorologica FTP server.
 - b. Retrieves files from a predefined folder located on the FTP server.
2. File Consolidation and Renaming:
 - a. After fetching

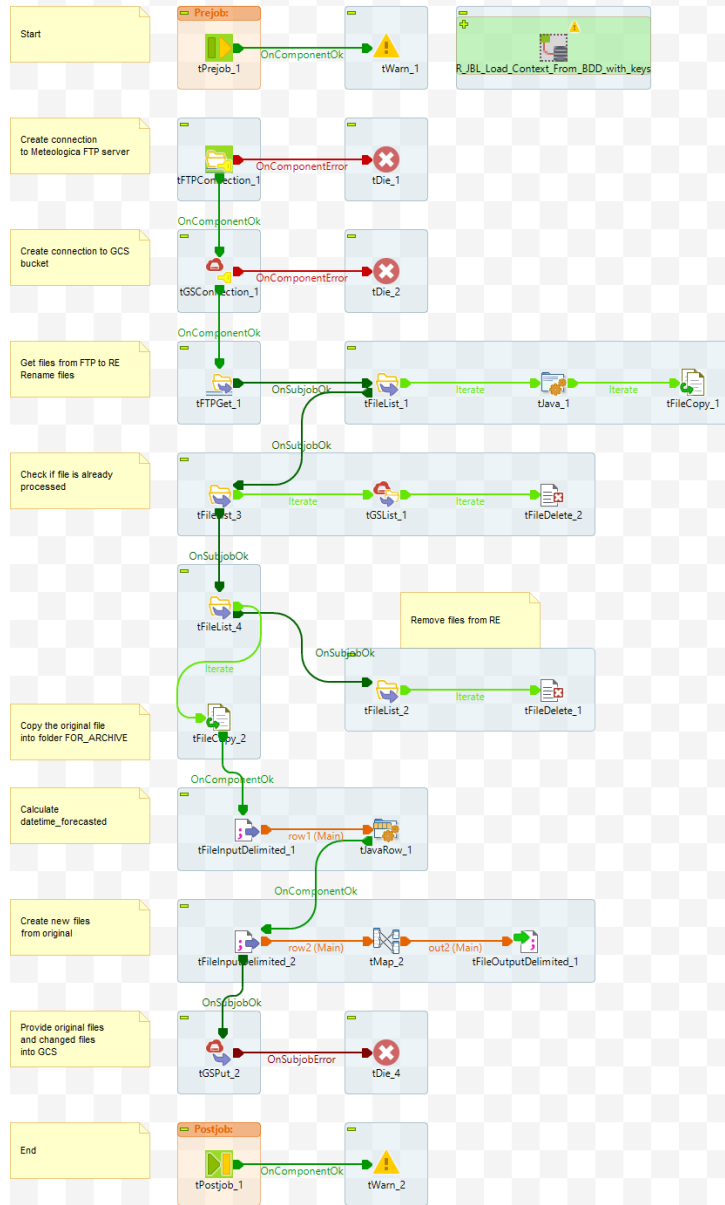
the files, the job consolidates and renames them according to predefined naming conventions.

b. This step ensures uniformity and consistency in file naming for further processing.

3. Creating File Backups:

a. As a best practice for data integrity, the job generates copies of the original files.

b. These copies are stored in an archive folder on the rem



ote
engi
ne
for
back
up
and
histo
rical
refer
ence.

4. Checking
for
Previousl
y
Processe
d Files:

a. The
job
perf
orms
a
chec
k to
dete
rmin
e if
the
pulle
d
files
have
alrea
dy
been
proc
esse
d.

b. This
verifi
catio
n is
cruci
al to
prev
ent
the
repr
oces
sing
of
files
from
previ
ous
days
, as
older
files
may
still
exist
in
the
folde
r.

5. Data
Extraction
and
Transfor
mation:

a. For
files
that
have
not
been
proc
esse

d
earli
er,
the
job
proc
eeds
with
data
anal
ysis.

b. It
extra
cts
a
datet
ime
valu
e
from
each
file,
locat
ed
at a
spec
ified
posit
ion
withi
n
the
file.

c. Usin
g
this
datet
ime
infor
mati
on,
the
job
calc
ulate
s a
new
colu
mn,
enab
ling
preci
se
time
stam
ping
of
the
data.

6. CSV
Output
and
Google
Storage:
a. The
proc
esse
d
data,
now
enric
hed
with
the
new
time
stam
p
colu

mn,
is
conv
erte
d
into
a
CSV
form
at.

- b. The resulting CSV output files are then pushed to Google Cloud Storage for further storage and access.

7. File Deletion:

- a. Once the files have been successfully processed and their data stored in Google Cloud Storage, they are deleted from the remote engine.
- b. This deletion step helps manage

storage resources efficiently and ensures that only processed data is retained.

This Talend job ensures the efficient handling of data from the Meteologica FTP server, including consolidation, backup, and validation steps. It also extracts and enriches the data, making it ready for use in downstream processes, ultimately contributing to the accuracy and usability of the energy optimization project's data.

Main jobs for source extraction

- **J008_GS**
heet_t
GCS_WA
P_SOLID
_FUELS
- **J011_G**
Sheet_t
e_GCS
_VEND
OHM_F
ORWA
RDS
- **J012_G**
Sheet_t
e_GCS
_VEND
OHM_S
POT
- **J013_G**
Sheet_t
e_GCS
_VEND
OHM_
METEO
- **J015_G**
Sheet_t

[--to the top --](#)

GCP Remote Engine

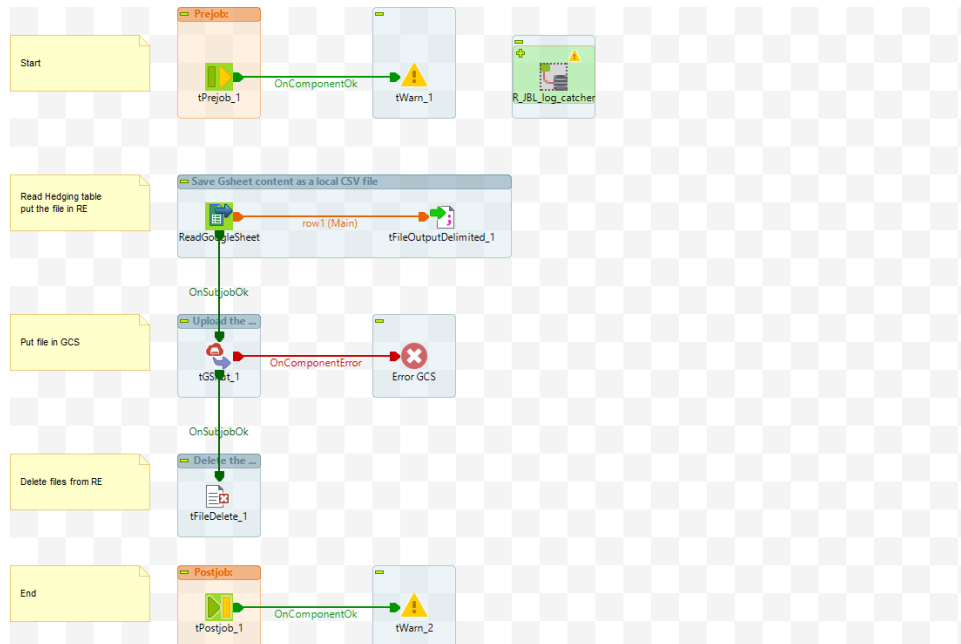
Industrial DO project

- ~~e_GGS~~
~~_IRM_H~~
~~UBS~~
- J016_G
Sheet_t
~~e_GGS~~
~~_IRM_C~~
~~OUNTE~~
~~R_PAR~~
~~TIES~~
- J017_G
Sheet_t
~~e_GGS~~
~~_IRM_S~~
~~ITES_H~~
~~EDGES~~

Job description by steps

Job design

1. Google Sheets Connection:
 - a. The job initiates a connection to a specified Google Sheets file.
 - b. It targets a specific sheet within the Google Sheets document, identified by its unique ID.
2. Data Extraction and Formatting:



- a. The job extracts data from the designated Google Sheets sheet.
- b. The extracted data is converted into a CSV format with a predefined fixed schema.
- c. The schema format is consistent and predefined to ensure data uniformity and structure.
- d. The resulting CSV file includes a predefined filename that incorporates the date of

extra
ction.

3. Copying
to
Google
Storage:

a. After
the
succ
essf
ul
extra
ction
and
form
attin
g of
the
data,
the
job
copi
es
the
CSV
file
to
Goo
gle
Clou
d
Stor
age.

b. This
step
serv
es
as a
secu
re
and
relia
ble
mea
ns
of
stori
ng
the
proc
esse
d
data
in
the
clou
d.

4. File
Deletion:

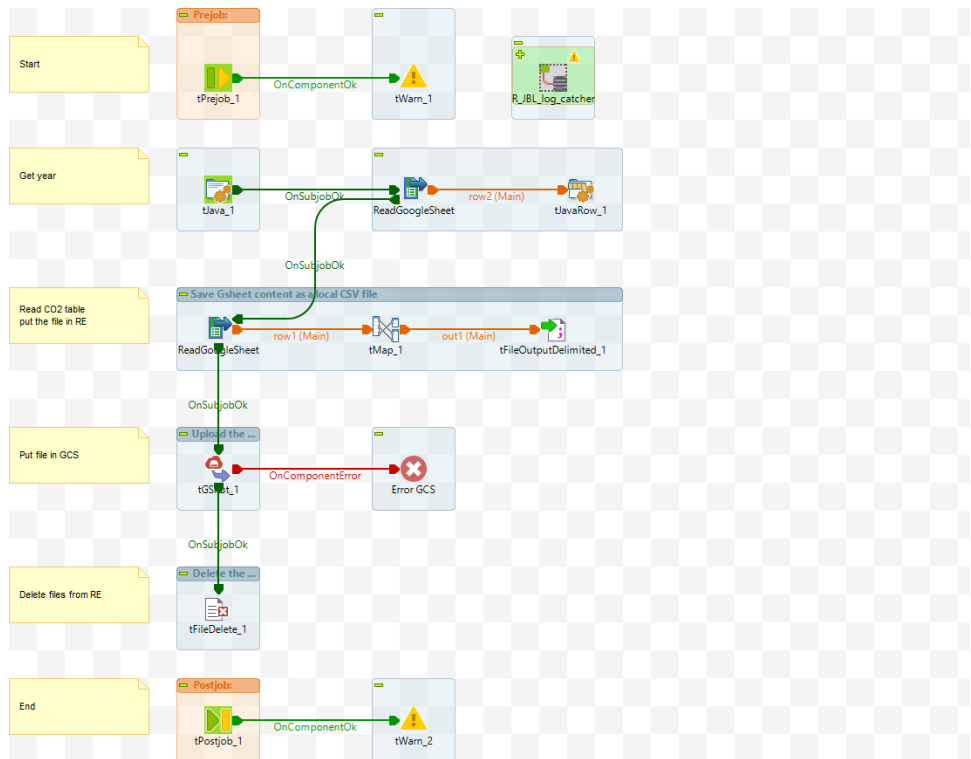
a. Follo
wing
the
succ
essf
ul
copy
ing
of
the
CSV
file
to
Goo
gle
Clou
d
Stor
age,

the original file located on the remote engine is deleted.
 b. This cleanup step helps manage storage resources efficiently and maintains data integrity.

This Talend job efficiently extracts, formats, and securely stores data from a Google Sheets document into Google Cloud Storage, ensuring that the data is readily available for further analysis and processing while adhering to a fixed schema and naming conventions.

<p>Main jobs for source extraction</p>	<ul style="list-style-type: none"> • J009_GSheet_to_GCS_CO2_EMISSIONS 	<p>--to the top -- GCP Remote Engine Sustainability DO project</p>
	<p>Job description by steps</p>	<p>Job design</p>
	<p>1. Google Sheets Connection:</p>	

- a. The job establishes a connection to a specified Google Sheet file.
- b. It focuses on a particular sheet within the Google Sheet document, identified by its unique ID.



2. Data Extraction and Formatting:

- a. The job extracts data from the designated Google Sheet.
- b. The extracted data is transformed into a CSV format, adhering to a predefined, fixed schema.
- c. The schema format is consistent and predefined.

to ensure uniformity and data structure.

d. The resulting CSV file follows a predefined naming convention, which incorporates the date of extraction.

3. Data Reading:

a. The job involves two distinct data reading steps.

i. First Reading: The job reads only a single cell within this

file, specifically the collection containing the year value. This is set to capture yearly arrears data. Second Reading

g : T h e s e c o n d r e a d o f t h e f i l e i s r e s e r v e d f o r c a p t u r i n g t h e r e m a i n i n g d a t a , d e f i n e d b y t h e s c h e m a a n d f o r

m
a
t.
T
h
i
s
s
t
e
p
c
o
m
p
r
e
h
e
n
s
i
v
e
l
y
r
e
a
d
s
t
h
e
d
a
t
a
s
p
e
c
i
f
i
e
d
i
n
t
h
e
s
c
h
e
m
a.

4. Copying to
Google
Storage:

- a. After
succe
ssfully
extract
ing
and
format
ting
the
data,
the
job
copies
the
CSV
file to
Googl
e
Cloud
Storag
e
(inside
Sustai
nabilit
y

- domain).
- b. This step serves as a secure and reliable means of storing the processed data in the cloud.

5. File

Deletion:

- a. Following the successful copying of the CSV file to Google Cloud Storage, the original file located on the remote engine is deleted.
- b. This cleanup step helps manage storage resources efficiently and maintains data integrity.

This Talend job effectively extracts, formats, and securely stores data from a Google Sheets document into Google Cloud Storage. It features specialized data reading

steps, ensuring the capture of year-related data separately, while adhering to predefined schemas and naming conventions. The deletion of the original file enhances data management and resource optimization.

Main jobs for source extraction

- ~~J010_Postgres_to_GCS_VENDOHM~~ Decommission
- J010_FIL_Energy_F_O_T_O_ODS

--to the top --
 AWS Remote Engine
 Industrial DO project

Job description by steps

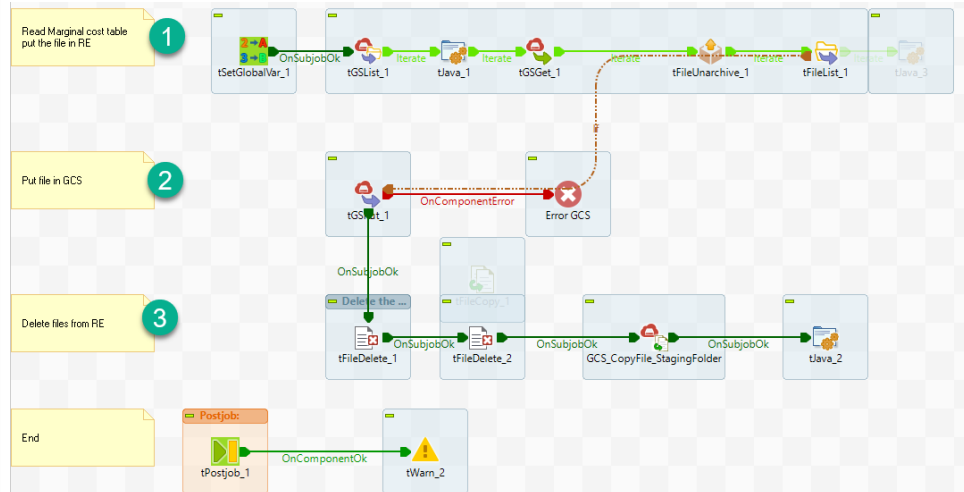
Job design

GCS Connection:

- Dataiku will generate energy a file that calculated many rules from the Energy team and save to GCS at DataOcean Industrial bucket (cs-ew1-prj-data-dm-industrial-[env]-staging/in/enr) twice a day to get the historical data to future hourly.

The mapping file and fields

1. First loading the data from cs-ew1-prj-data-dm-industrial-[env]-staging/in/enr to local drive and unzip the file
2. Upload the file to cs-ew1-prj-data-dm-industrial-[env]-



	<p>staging /ENERGY/irm- conso</p> <p>3. Storage in Google Cloud Storage (GCS) :</p> <ul style="list-style-type: none"> ○ After the successful extraction and formatting of the data, the job securely stores the CSV file in Google Cloud Storage. ○ The files are placed within a designated "ENERGY" folder in the GCS location for efficient storage and accessibilit y. ○ Also move the source file .gz to folder "Processed" ○ Following a successful data extraction and storage process, the original file located on the remote engine is deleted to manage storage resources efficiently. ○ The connection to the Oracle database is closed to ensure proper resource manageme nt and security. 	
--	---	--

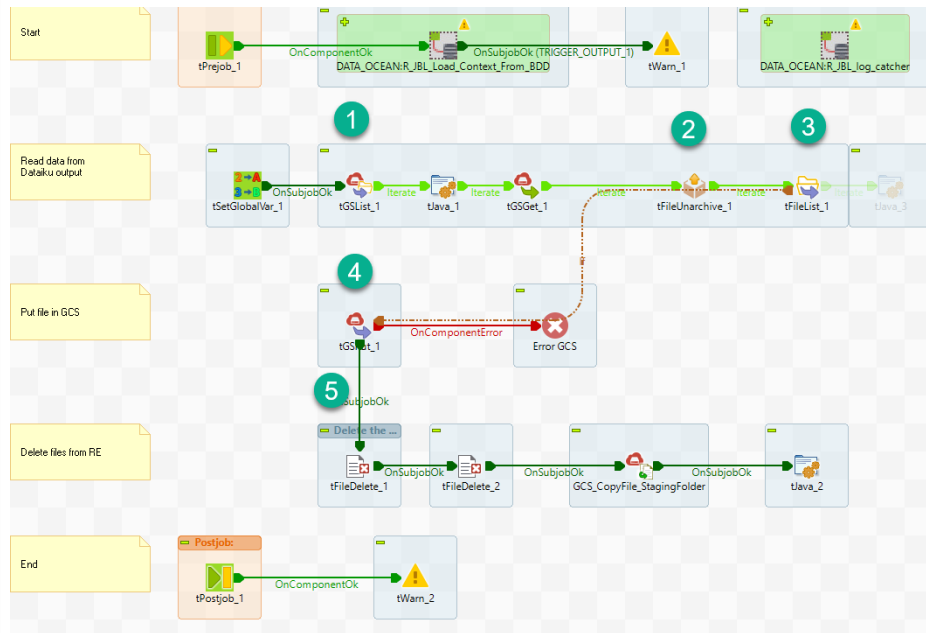
<p>Main jobs for source extraction</p>	<ul style="list-style-type: none"> • J014_Oracle to_GCS_IR M- J014_GCS_t o_GCS_IRM 	<p>--to the top --</p> <p>AWS Remote Engine</p> <p>Industrial DO project</p>
--	--	--

**_ENERGY_D
EAL**

**Job description
by steps**

1. Download file from GCS (output of Dataiku) cs-ew1-prj-data-dm-industrial-[env]-staging/in/enr/irm-deals
2. Unzip the file and keep to the GCS cs-ew1-prj-data-dm-industrial-[env]-staging /ENERGY/irm-deals
 - The data is structured according to a predefined and fixed schema to maintain data consistency.
3. Data Formatting and CSV Output:
 - The extracted data is formatted and transformed into CSV format.
 - Each resulting CSV file includes the table name and the date of extraction for reference and organization.
4. Storage in Google Cloud Storage (GCS):
 - After the successful extraction and formatting of the data, the job securely stores the CSV file in Google Cloud Storage.
 - The files are placed

Job design



	<p>within a designated "IRM" folder in the GCS location for efficient storage and accessibility.</p> <p>5. File Deletion and Connection Closure:</p> <ul style="list-style-type: none"> • Following a successful data extraction and storage process, the original file located on the remote engine is deleted to manage storage resources efficiently. • Source file from Dataiku will move to folder Process and rename to the file + timestamps 	
--	---	--

4.3 - Data Loading to Google BigQuery

<p>Main jobs for source extraction</p>	<ul style="list-style-type: none"> • J001_Extraction_timestamps_MEXICO • J002_Extraction_timestamps_ITALY • J003_Extraction_timestamps_GERMANY • J004_Extraction_timestamps 	<p>--to the top --</p> <p>GCP Remote Engine</p> <p>Industrial DO project</p>
--	---	--

n_til_O
 DS_ME
 TEOLO
 GICA_
 SPAIN
 • J005_E
 xtractie
 n_til_O
 DS_ME
 TEOLO
 GICA_
 SPAIN_
 ENS
 • J006_E
 xtractie
 n_til_O
 DS_ME
 TEOLO
 GICA_
 SPAIN_
 OBS

Job
 description
 by steps

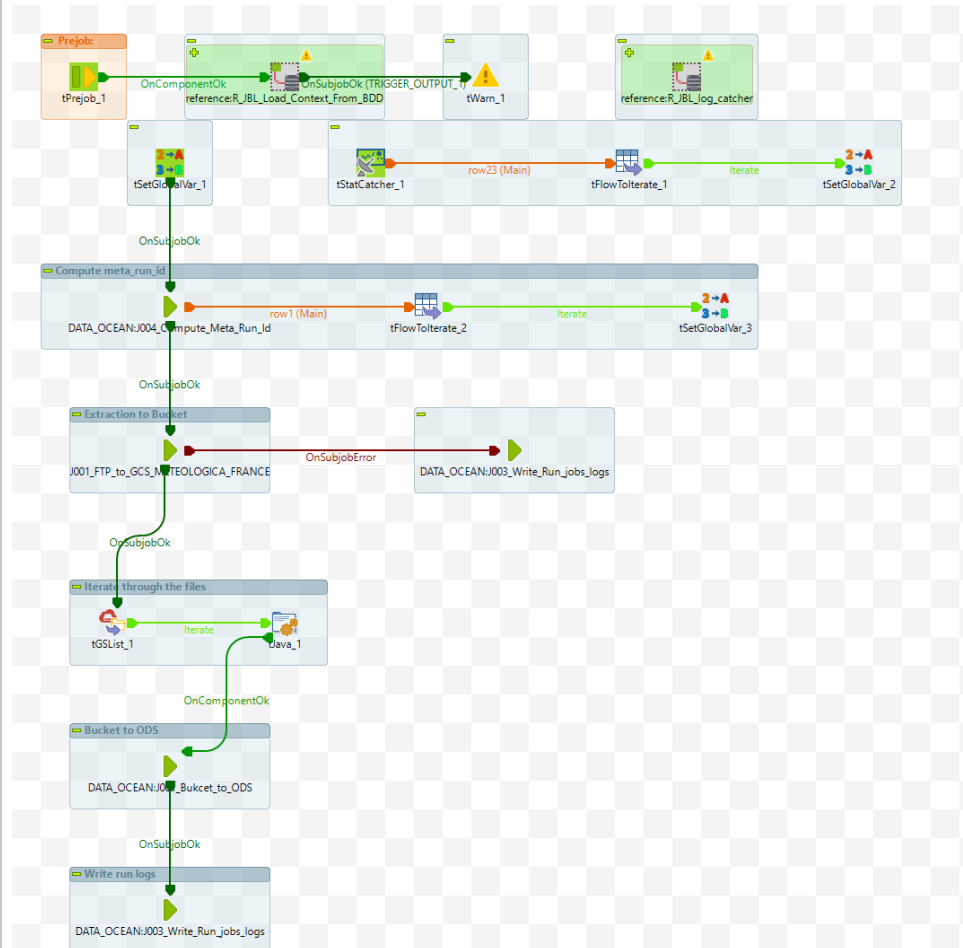
Job design

Job
 Initialization
 and Logging:

1. Job Run
 Initializati
 on
 (Subjob
 1):

a. The
 job
 starts
 with
 the
 first
 subjob,
 where
 the
 job
 ID
 is
 retrieved,
 and
 the
 extraction
 date
 is
 calculated.

b. This
 crucial
 step
 initializes
 the
 job
 and
 captures



essential
meta
data.

Data
Extraction
and
Processing:

1. Files
Extraction
from FTP
Server
(Subjob
2):

a. In
the
next
subjob,
data
files
are
extracted
from
the
Meteorological
FTP
server.

2. Error
Handling
and
Logging
(Subjob
3 - On
Failure):

a. If
the
data
extraction
subjob
(Subjob
2)
encounters
an
issue
or
does
not
finish
successfully,
the
job
transitions
to
a
subjob
for
writing
error
logs.

b. Logging is essential for tracking and diagnosing issues in the data extraction process.

3. Data Loading and Transformation (Subjob 4 - On Success):

a. When the data extraction subjob (Subjob 2) successfully completes, the job proceeds with data loading and transformation.

b. In this stage, each extracted file is processed one by one.

c. For each file, a subjob

is called to load data from the CSV file into a stage table, preparing it for further processing.

d. Subsequently, data from the stage table is loaded into an operational data store (ODS) table.

e. All necessary parameters, such as table names and connection parameters, are provided to ensure accurate data extraction and loading.

Logging and Reporting:

1. Logging (Subjob 5):
 - a. At the end of the data loading and transformation process, a subjob is called to write logs.
 - b. Logging captures critical information about the data processing, ensuring transparency and traceability.

This Talend job orchestrates the data extraction, loading, and transformation process. It begins with initialization and metadata capture, proceeds to extract data from the FTP server, handles errors if they occur, and performs data loading and transformation for energy optimization purposes. Detailed logging and

error handling mechanisms enhance job monitoring and maintain data integrity.

Main jobs for source extraction

- J008_Ext
raction_t
il_ODS_
GSHEET
_WAP_S
OLID_FU
ELS
- J009_E
xtractio
n_till_O
DS_GS
HEET_
CO2
- J011_E
xtractio
n_till_O
DS_GS
HEET_
VENDO
HM_FO
RWAR
DS
- J012_E
xtractio
n_till_O
DS_GS
HEET_
VENDO
HM_SP
OT
- J013_E
xtractio
n_till_O
DS_GS
HEET_
VENDO
HM_ME
TEO
- J014_Ext
raction_#
l_ODS_O
RACLE_
RM J01
4_FIL_IR
M_ENER
GY_DEA
L_F_O_T
O_ODS
- J015_E
xtractio
n_till_O
DS_GS
HEET_
RM_HU
BS
- J016_E
xtractio
n_till_O
DS_GS
HEET_
I

--to the top --

GCP Remote Engine

Industrial DO project

Exception :

- J009_Extraction_til_ODS_GSHEET_CO2 on GCP RE and on Sustainability DO project
- J014_FIL_IRM_ENERGY_DEAL_F_O_TO_ODS on AWS RE and on Industrial DO project

RM_CO
 UNTER
 PARTI
 ES
 • J017_E
 xtractie
 n_till_O
 DS_CS
 HEET_I
 RM_SIT
 ES_HE
 DGES

Job
 description
 by steps

Job design

Job
 Initialization
 and Logging:

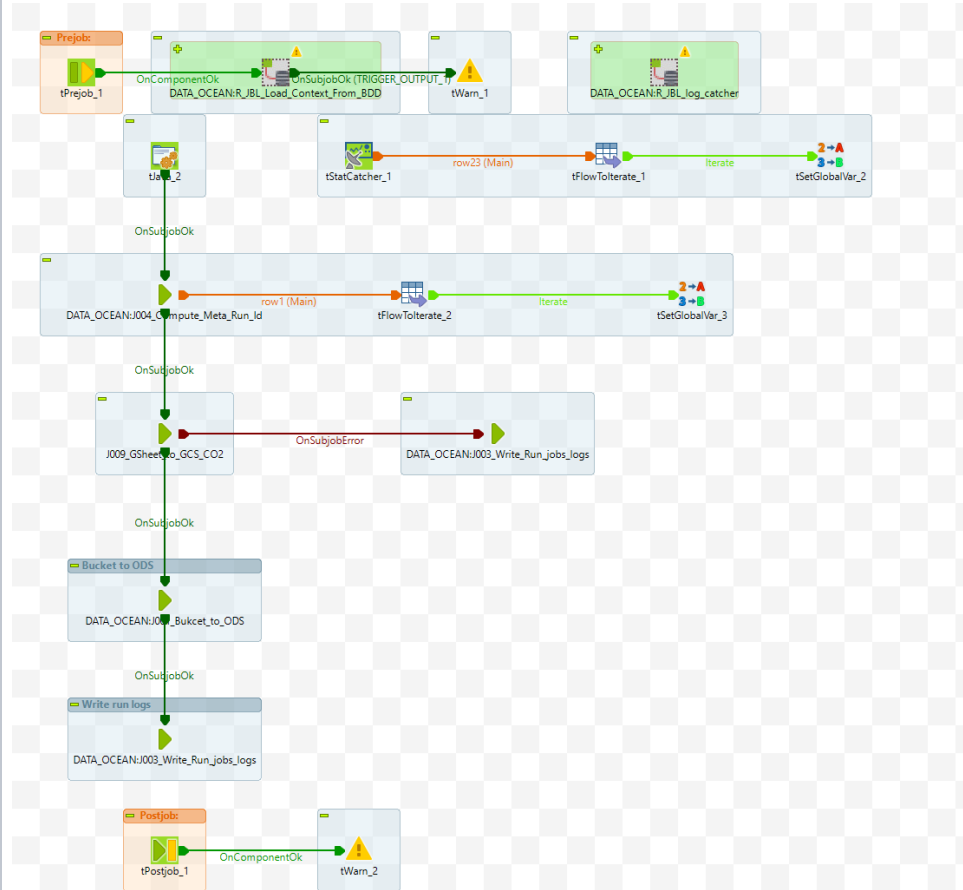
1. File
 Name
 Configura
 tion
 (Global
 Variable):

a. The
 job
 begi
 ns
 by
 confi
 gurin
 g
 the
 file
 nam
 e
 as
 a
 glob
 al
 varia
 ble.

b. This
 varia
 ble
 is
 used
 thro
 ugho
 ut
 the
 job
 to
 ensu
 re
 cons
 isten
 cy
 in
 file
 hand
 ling.

2. Job Run
 Initializa
 tion
 (Subjob
 1):

a. In
 the
 first
 subj
 ob,



the job ID is retrieved, and the extraction date is calculated.

b. This step sets the stage for the data processing and captures essential meta data.

Data Extraction from Google Sheets / IRM database:

1. Files Extraction from Google Sheets / IRM database (Subjob 2):

a. The next subjob is responsible for extracting data from Google Sheets / IRM database.

b. It targets the specified file base

d on
the
glob
al
varia
ble
(file
nam
e)
confi
gure
d
earli
er.

2. Error
Handling
and
Logging
(Subjob
3 - On
Failure):

a. In
case
the
data
extra
ction
subj
ob
(Sub
job
2)
enco
unte
rs
any
issu
es
or
does
n't
finis
h
succ
essf
ully,
the
job
trans
ition
s to
a
subj
ob
for
writi
ng
error
logs.

b. Logg
ing
is
cruci
al
for
track
ing
and
diag
nosi
ng
probl
ems
durin
g
the
data
extra
ction

process.

Data Loading and Transformation:

1. Data Loading and Transformation (Subjob 4 - On Success):
 - a. When the data extraction subjob (Subjob 2) successfully completes for the provided filename, the job proceeds with data loading and transformation.
 - b. A subjob is called to load data from the CSV file into a stage table, preparing it for further processing.
 - c. Subsequently

, data from the stage table is loaded into an operational data store (ODS) table.
d. All necessary parameters, such as table names and connection parameters, are provided to ensure accurate data extraction and loading.

Logging and Reporting:

1. Logging (Subjob 5):
 - a. At the end of the data loading and transformation process, a subjob is called to

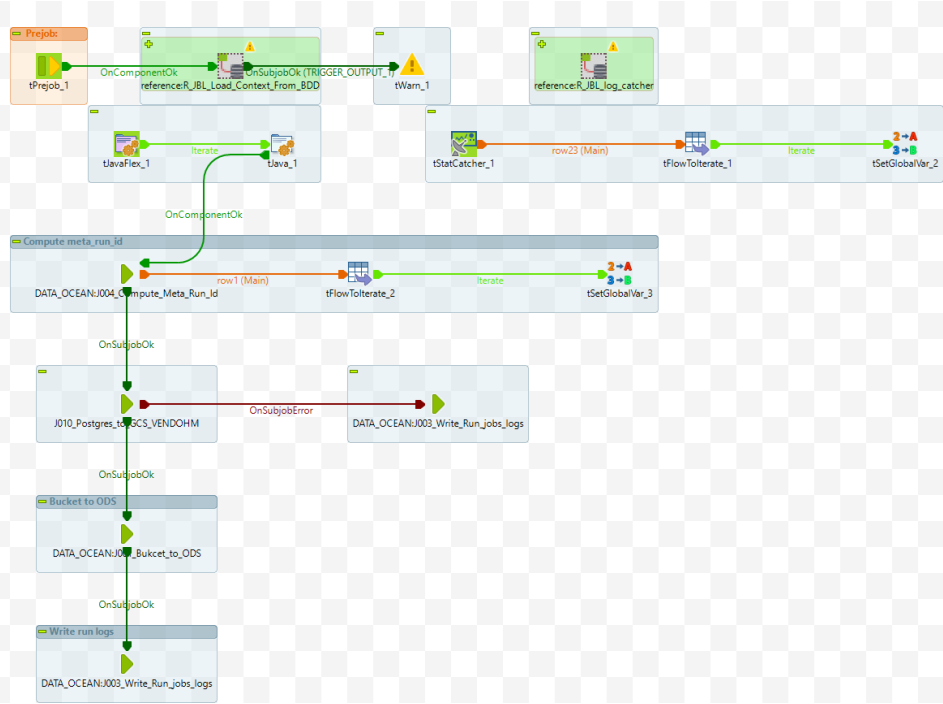
write logs.
 b. Logging captures crucial information about the data processing, enhancing transparency and traceability.

This Talend job streamlines the data extraction, loading, and transformation process. It begins with file name configuration, proceeds with metadata capture, extracts data from Google Sheets, handles errors if they occur, and performs data loading and transformation, all while maintaining detailed logs for monitoring and maintaining data integrity.

<p>Main jobs for source extraction</p>	<ul style="list-style-type: none"> • J010-Extraction-til-ODS-VE-NDOHM 	<p>--to the top -- GCP Remote Engine Industrial DO project</p>
	<p>Job description by steps</p>	<p>Job design</p>
	<p>Global Variable Configuration: 1. File Names</p>	

and Table Names Configuration (Global Variables postgres table names):

- The job starts by configuring global variables for file names and table names.
- These variables are populated by reading a list of all tables and files that subjects need to load.
- This dynamic configuration ensures that the job adapts to the specific files and tables required for processing.



Data-
Extraction-
and-
Processing in-
Iterations:

1. Iteration-
Loop-
(One-by-
One):

a. The-
job-
pre-
cesse-
s-
table-
s-
and-
files-
iterat-
ively-
-
one-
by-
one-
base-
d-on-
the-
confi-
gure-
d-
glob-
al-
varia-
bles.

2. Job Run-
Initializati-
on-
(Subjob-
4):

a. Withi-
n-
each-
iterat-
ion-
the-
first-
subj-
ob-
is-
calle-
d-to-
retri-
eve-
the-
job-
ID-
and-
calc-
ulate-
the-
extra-
ction-
date.

b. This-
initia-
lizati-
on-
step-
capt-
ures-
vital-
meta-
data-
for-
the-
ongo-

ing-
data-
proc-
essi-
ng-

3. Files-
Extraction
from-
PostgreS
QL-
Database
(Subjob-
2):

a. In-
the-
next-
subj-
ob-
data-
files-
are-
extra-
cted-
from-
the-
Post-
greS-
QL-
data-
base-

b. Eac-
h-
iterat-
ion-
targ-
ets-
a-
spee-
ific-
file-
base-
d on-
the-
glob-
al-
varia-
bles-
(file-
nam-
es)-
confi-
gure-
d-
earli-
er-

4. Error-
Handling-
and-
Logging-
(Subjob-
3-On-
Failure):

a. If-
the-
data-
extra-
ction-
subj-
ob-
(Sub-
job-
2)-
encoe-
unte-
rs-
any-
issu-
es-

of
does
n't
finis
h
succe
essf
ully,
the
job
trans
ition
s-to
a
subj
ob-
for
writi
ng-
error
logs.

b. Logg
ing-
is-
cruci
al
for
track
ing-
and
diag
nosi
ng-
probl
ems-
durin
g-
the
data-
extra
ction
proce
ss.

5. Data-
Loading-
and-
Transfer
mation
(Subjob-
4-On-
Success):

a. Whe
n-
the
data-
extra
ction
subj
ob-
(Sub
job-
2)-
succe
essf
ully-
com
plete
s-for
the-
provi
ded-
file-
nam
e,-
the
job-
proce
ceeds
with

data-
loadi-
ng-
and-
trans-
form-
ation-

- b. A-
subj-
ob-
is-
calle-
d-to-
load-
data-
from-
the-
CSV-
file-
into-
a-
stag-
e-
table-
-
prep-
aring-
it-for-
furth-
er-
pro-
cessi-
ng-
- c. Sub-
sequ-
ently-
-
data-
from-
the-
stag-
e-
table-
is-
load-
ed-
into-
an-
oper-
ation-
at-
data-
store-
(OD-
S)-
table-
- d. All-
nece-
ssar-
y-
para-
mete-
rs,-
such-
as-
table-
nam-
es-
and-
conn-
ectio-
n-
para-
mete-
rs,-
are-
provi-
ded-
to-

ensu-
re-
accu-
rate-
data-
extra-
ction
and-
loadi-
ng-

Logging and-
Reporting:

1. Logging-
(Subjob-
5):
 - a. At-
the-
end-
of-
each
iterat-
ion,-
a-
subj-
ob-
is-
calle-
d-to-
write
logs.
 - b. Logg-
ing-
capt-
ures-
critic-
al-
infor-
mati-
on-
about
the-
data-
proce-
ssi-
ng,-
ensu-
ring-
trans-
pare-
ncy-
and-
trace-
abilit-
y-for-
each
proce-
sse-
d-
file-
and-
table.

This Talend-
job-is-
designed-for-
dynamic-data-
extraction-and-
processing,-
adapting-to-
different-
tables-and-
files.-It-
initializes-
each iteration-
with metadata-
capture,-
extracts data-
from the-

	<p>PostgreSQL database, handles errors when necessary, performs data loading and transformation with parameterized configurations, and maintains detailed logs for monitoring and data integrity across multiple files and tables.</p>
--	--

4.4 - Load to DM (calculations and transformations)

<p>Main jobs for source extraction</p>	<ul style="list-style-type: none"> • J001_ODS_TO_DM_FACT_ENERGY_PRICE_FORECAST 	<p>--to the top --</p> <p>GCP Remote Engine</p> <p>Industrial DO project</p>
	<p>Job description by steps</p>	<p>Job design</p>
	<p>Metadata Calculation and Cache Initialization:</p> <ol style="list-style-type: none"> 1. Job Metadata Calculation: <ul style="list-style-type: none"> • The job begins by calculating essential job metadata, which include information such as job identifiers, time 	

Data

Extraction:

- The job extracts energy forecast data from four Operational Data Store (ODS) tables, each containing information for a specific region (Italy, Germany, Spain, and France).

2. Data

Comparison and

Import:

- The extracted energy forecast data is compared with the cached keys from the Fact table.
- The primary purpose of this step

is to identify and import only new rows of data for the latest forecasted date.

- By importing only new data, the job ensures that the dataset remains up to date and avoids unnecessary duplication.

Logging and Cleanup:

1. Log Generation:

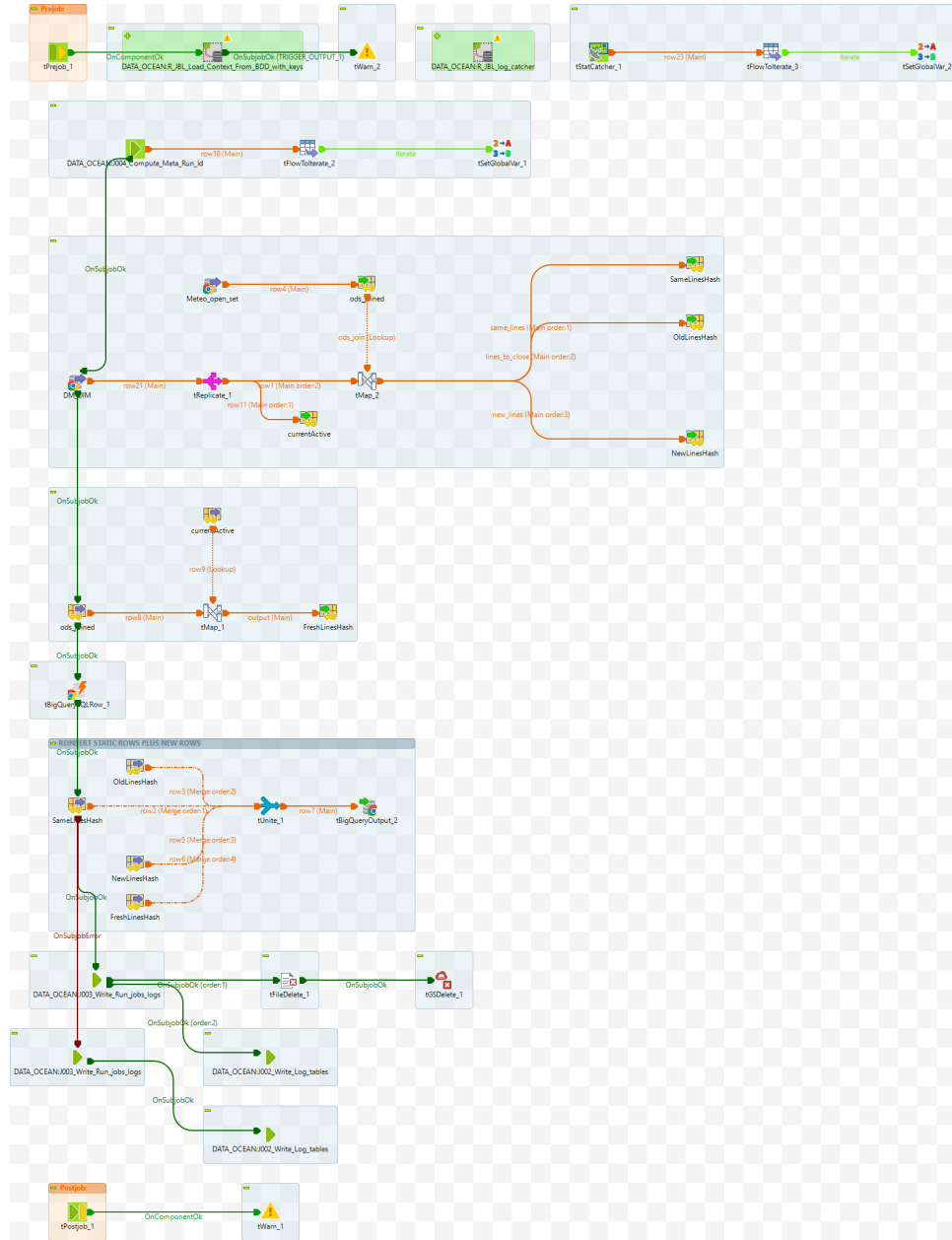
- At the conclusion of data import and processing, the job generates logs.
- These logs provide a record of the

	<p>job's activities, facilitating monitoring and troubleshooting.</p> <p>2. Temporary File Deletion:</p> <ul style="list-style-type: none"> As a final step, the job deletes any temporary files created during the data transfer process. This cleanup helps manage storage resources and maintain data integrity. 	
--	---	--

<p>Main jobs for source extraction</p>	<ul style="list-style-type: none"> J002_0 DS_TO DM_DI M_MET EO_INFO 	<p>--to the top --</p> <p>GCP Remote Engine</p> <p>Industrial DO project</p>
	<p>Job-description-by-steps</p>	<p>Job design</p>
	<p>Metadata-Calculation:</p> <p>1. Job-Metadata</p>	

Calculatio

- The job initiates by calculating essential job meta data. This meta data includes information on job identifiers, timestamps, or job specific configuration details.
- This step sets the foundation for the subsequent data processing.



Slowly Changing Dimension (SCD2) Data Loading:

1. SCD2 Logic for Dimension Data Loading:
 - In this step, the job implements Slowly Changing

g-
Dim
ensi
on-
Type
2-
(SC
D2)-
logie.

- The
purp
ose
of
this
logie
is to
load
data
into
a
dime
nsio
n-
table
-
spee
ificall
y
from
the
Oper
ation
at
Data
Stor
e
(OD
S)-
mete
o-
table
into
the
Data
Mart
(DM)
mete
o-
table.
- The
SCD
2-
logie
is
appli
ed
by
com
parin
g
multi
ple
field
s-
acro
ss
the
two
table
s,
such
as
nam
e,
data
base
nam
e, le
grou

p-
asse
t
clas
s-
asse
t
sub-
clas
s-
asse
t
type,
capa
city-
unit,
som
modi
ty,
curv
e-
publi
shin
g-
curv
e-
type,
hub,
mark
et,
obje
ct
type,
own
er,
regio
n-
unit,
and
volu
me-
unit.

Logging and Cleanup:

1. Log-
Generatio
n:
 - Upe
n-
com
pleti
on-
of
the
data-
loadi
ng-
pro
cess,
the
job-
gene
rates
logs.
Thes
e-
logs-
serv
e as
a
rece
rd
of
the
job's
activ
ities,
provi

ding-
visibi-
lity-
into-
the-
data-
loadi-
ng-
proce-
sses.

2. Temporary File Deletion:

- As a final step, the job takes care of deleting any temporary files created during the data transfer process. This cleanup ensures the efficient management of storage resources.

3. In case of error to save data in BQ, it the last data to save to BQ will save to /DATA /DEV/IND /ROBUST IFY/nOut /DM/ [meta_run_id].csv

Main jobs for source

- **J003_Q DS_TO**

[to the top](#)

GCP Remote Engine

extraction	<p>DM-DM M-ENE RGY-P RICE-I NFO</p>	Industrial-DO project
------------	--	-----------------------

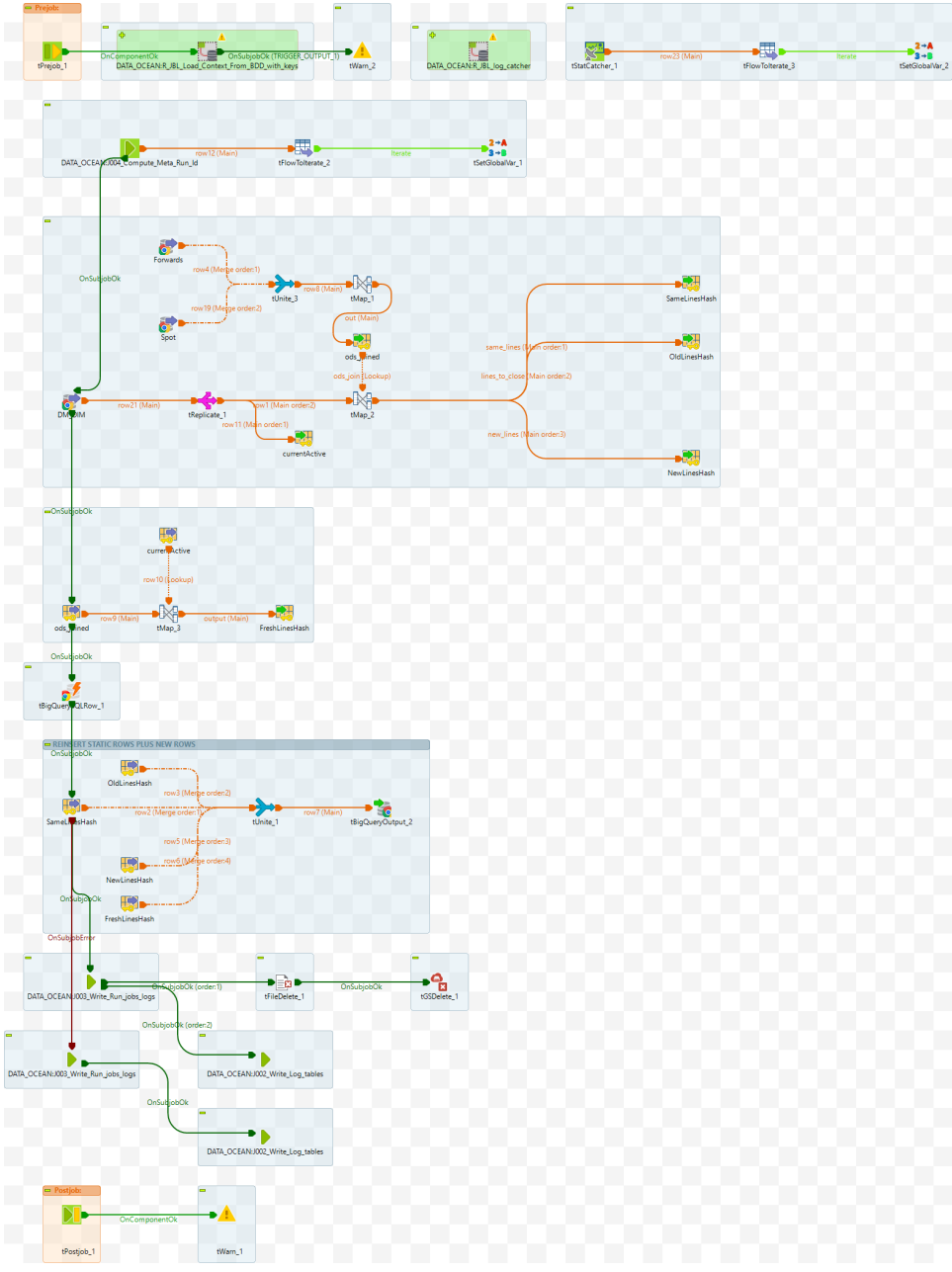
	<p>Job description by steps</p>	<p>Job design</p>
--	--	--------------------------

Metadata Calculation:

1. Job-Metadati Calculation:

- The job initiates by calculating essential job metadata. This metadata includes information on like job identifiers, timestamps, or job specific configuration details.
- This step sets the foundation for the subsequent data processing.

Slowly Changing Dimension (SCD2) Data Loading:



1. SCD2-
Logic for
Dimension
Data
Loading:

- In this step, the job applies Slowly Changing Dimension Type 2 (SCD2) logic to load data into a dimension table.
- The source of the data is the Operational Data Store (ODS) energy price metadata table - which is a union of tables named "forwards" and "spot". The data is loaded into the

Data Mart (DM) energy price table using the SCD 2 methodology.

- SCD 2 logic involves comparing specific fields, such as name, commodity, contract, proprietary, time span, rolling, and time zone, to identify changes and maintain historical records in the DM table.

Logging and Cleanup:

1. Log Generation:
 - Upon the completion of the data load

ng-
proe
ess,
the
job-
gene
rates
logs.
Thes
e-
logs-
serv
e as
a
reco
rd
of
the
job's
activ
ities
and
provi
de
visibi
lity
into
the
data-
loadi
ng-
proe
ess.

2. Temporary File Deletion:

- As
a
final
step,
the
job-
take
s
care-
of
delet
ing
any
temp
orar
y-
files-
creat
ed
durin
g
the
data-
trans
fer-
proe
ess.
This
clea
nup
ensu
res
the
effici
ent
man
age
ment
of
stora
ge-

rese
uree
s-
3. In case
of error
to save
data in
BQ, if the
last data
to save
to BQ
will save
to /DATA
/DEV/IND
/ROBUST
IFY/inOut
/DM/
[meta_ru
n_id].csv

Main jobs for source extraction

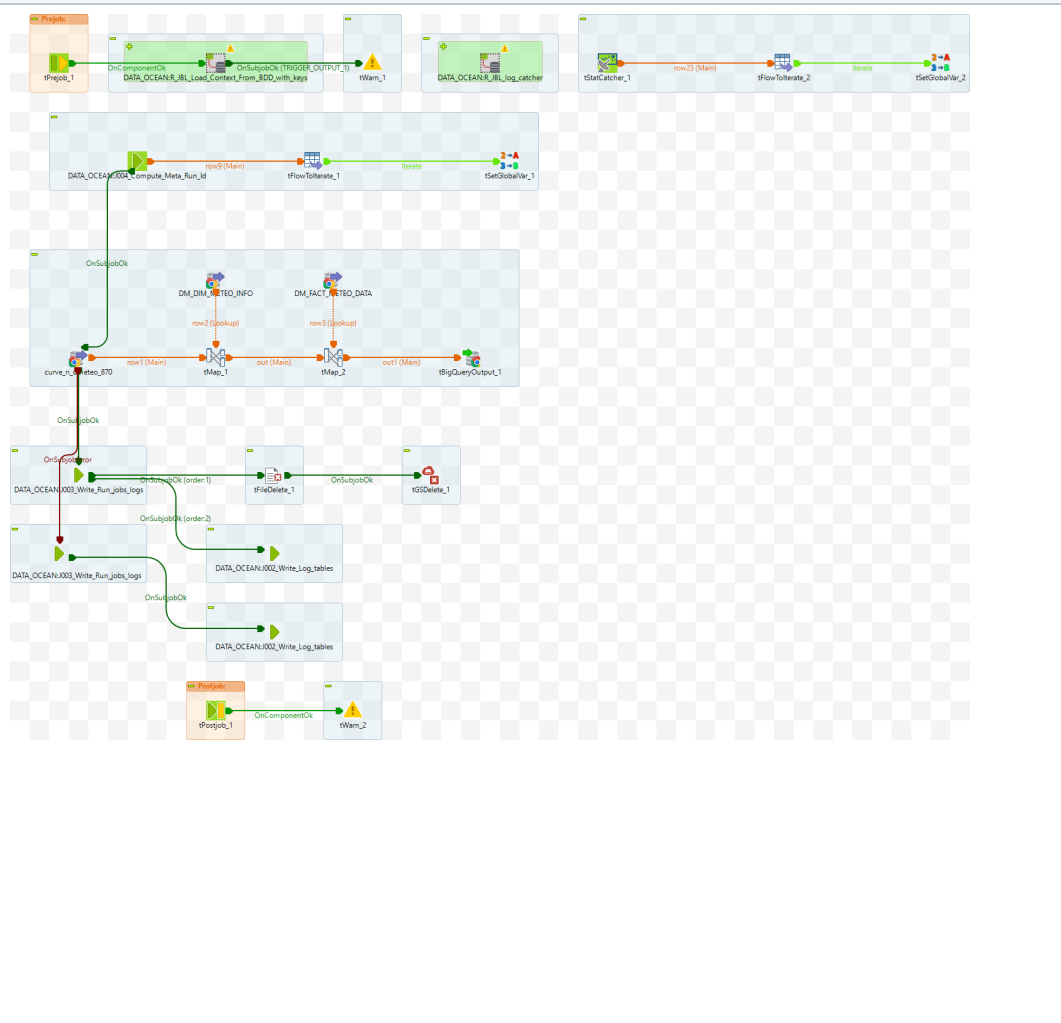
- J004_0 DS_TO DM_FACT_M ETEO_DATA

[to the top](#)
GCP Remote Engine
Industrial DO project

Job description by steps

Job design

Metadata Calculation:
1. Job Metadata Calculation:
• The job initiates by calculating essential job metadata. This metadata includes information on job identifiers, time stamps, or job specific configuration



details:
• This step establishes the foundation for the subsequent data processing.

Data Enrichment and Loading:

1. Data Enrichment from DIM_metadata_info:
 - In this step, the job enriches data from the Vendor metadata table with attributes obtained from the DIM_metadata_info table.
2. Loading into FACT_metadata_data:
 - The enriched data is then loaded into the FACT_metadata_data table.

- This table serves as a fact table containing data relevant to meteorological information.

Logging and Cleanup:

1. Log Generation:

- Upon the completion of the data entry and loading process, the job generates logs. These logs serve as a record of the job's activities and provide visibility into the data processing steps.

2. Temporary File Deletion:

- As a

final step, the job takes care of deleting any temporary files created during the data transfer and ensure that processes. This ensures efficient storage resource management.

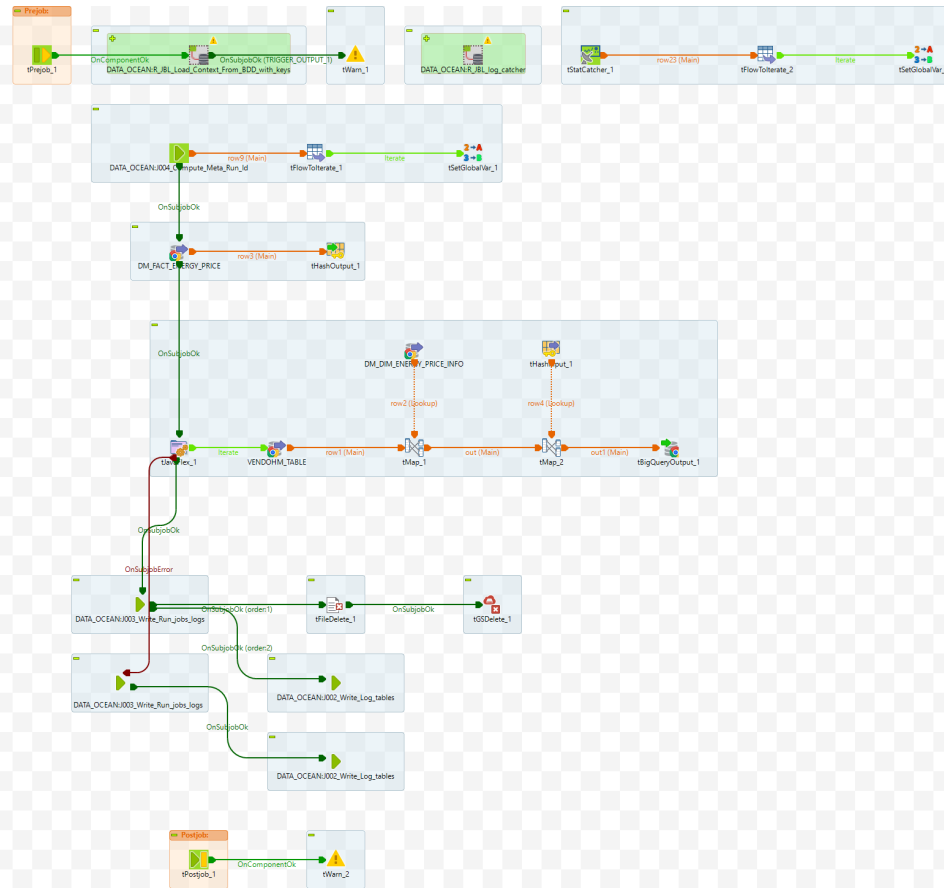
<p>Main jobs for source extraction</p>	<ul style="list-style-type: none"> • J005_0 DS_TO _DM_F ACT_E ENERGY _PRICE <p>replace with view in Robustify project-data-robustify-dev. DM. V_FACT_energy_price_hourly</p>	<p>--to the top --</p> <p>GCP Remote Engine Industrial DO project</p>
	<p>Job description by steps</p>	<p>Job design</p>
	<p>Metadata Calculation:</p>	

1. Job Metadata Calculation

- The job comes with the calculation of essential job metadata. This metadata includes information like job identifiers, timestamps, or specific configuration details.
- This initial step provides a foundation for subsequent data processing.

Data Enrichment and Loading:

- Data Enrichment from DIM_energypri_inf:
 - In this step, the job enrie



hed-
data-
from
the
ener-
gy-
price
table
by
incor-
pora-
ting
attrib-
utes-
sour-
ced
from
the
DIM
_ene-
rgy-
price
_info
table.

2. Loading-
into
FACT_ene-
rgy_pri-
ce_data:

- The
entire
hed-
data-
is
subse-
quently
load-
ed
into
the
FACT_ene-
rgy_pri-
ce_data
table.
- The
FACT_ene-
rgy_pri-
ce_data
table
is
desi-
gned
to
store
data
relev-
ant
to
ener-
gy-
price
s.

Logging and
Cleanup:

1. Log-
Generatio-
n:
- Follo-
wing
the
com-

pletion of the data enrichment and loading process, the job generates logs. These logs serve as a record of the job's activities and provide visibility into the processing steps.

2. Temporary File Deletion:

- As a final step, the job manages the deletion of any temporary files that were created during the data transfer and enrichment process.

ess-
This
clea
rup
ensu
res
effici
ent
stora
ge-
rese
urces
man
age
ment

Main
jobs
for
source
extraction

- J006_O
DS_TO
_DM_F
ACT_C
O2_EMI
SSIONS

--to the top --

GCP Remote Engine

Industrial DO project

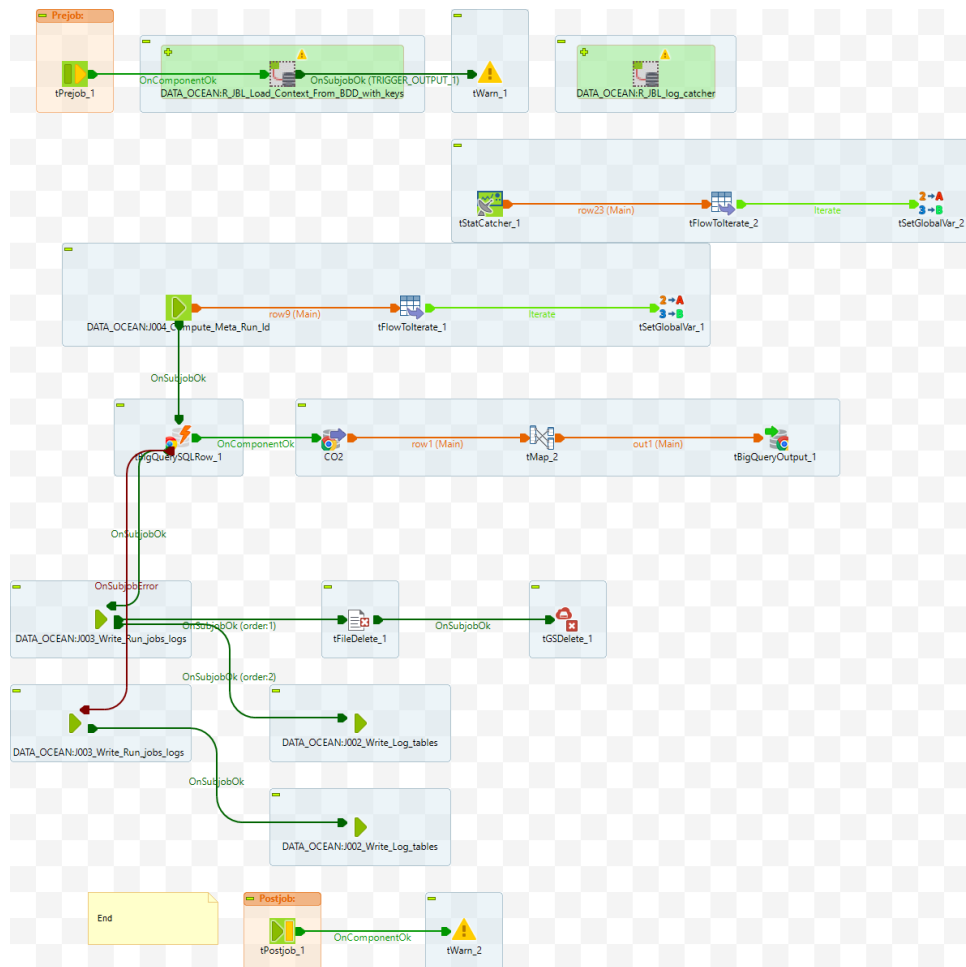
Job
description
by steps

Job design

Metadata
Calculation:

1. Job
Metadata
Calculatio
n:

- The job commences by calculating essential job metadata. This metadata includes information such as job identifiers, timestamps, or specific configuration details.



- This initial step provides a foundation for subsequent data processing.

Latest CO2
Data Refresh:

1. Refreshing/Providing Latest CO2 Data:

- In this step, the job focuses on refreshing or providing the latest CO2 emissions data.
- The source of this data is the Operational Data Store (ODS) table, which serves as a central repository for CO2 - relat

ed
infor
mati
on.

Logging and
Cleanup:

1. Log
Generatio
n:

- Follo
wing
the
com
pleti
on
of
the
CO2
data
refre
sh,
the
job
gene
rates
logs. Thes
e
logs
serv
e as
a
reco
rd
of
the
job's
activ
ities
and
provi
de
visibi
lity
into
the
data
refre
sh
proc
ess.

2. Temporar
y File
Deletion

- As
a
final
step,
the
job
man
ages
the
delet
ion
of
any
temp
orar
y
files
that
were
creat
ed
durin
g
the

data transfer process. This cleanup ensures efficient storage resource management.

Main jobs for source extraction

- J007_0DS_TO_DM_FACT_SOLID_FUEL_WAP

--to the top --

GCP Remote Engine

Industrial DO project

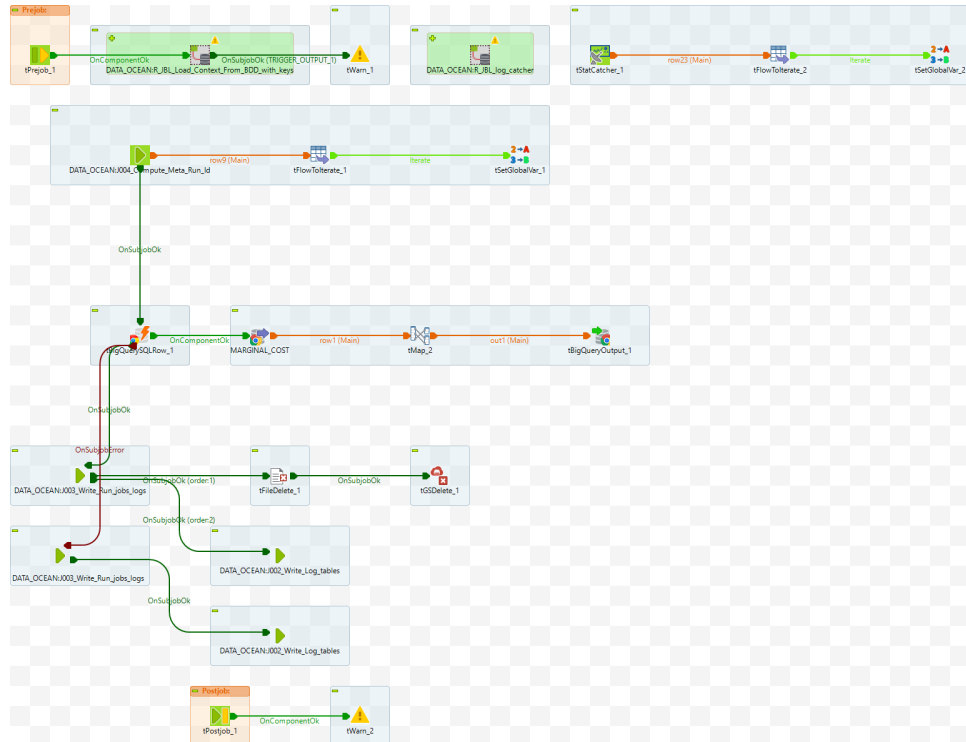
Job description by steps

Job design

Metadata Calculation:

1. Job Metadata Calculation:

- The job commences by calculating essential job metadata. This metadata includes information such as job identifiers, time stamps, or specific confi



guration details.

- This initial step provides a foundation for subsequent data processing.

Latest wap solid fuel data
Refresh:

1. Refreshing /Providing Latest Cost Data:

- In this step, the job focuses on refreshing or providing the latest weighted average fuel prices data.
- The source of this data is the Operational Data Store (ODS) table, which serves as a

central repository for wapsolid fuel-related information.

Logging and Cleanup:

1. Log Generation:

- Following the completion of the Wapsolid fuel data refresh, the job generates logs. These logs serve as a record of the job's activities and provide visibility into the data refresh process.

2. Temporary File Deletion

- As a final step, the job manages the deletion of any

	<p>temporarily files that were created during the data transfer process. This cleanup ensures efficient storage resource management.</p>	
--	--	--

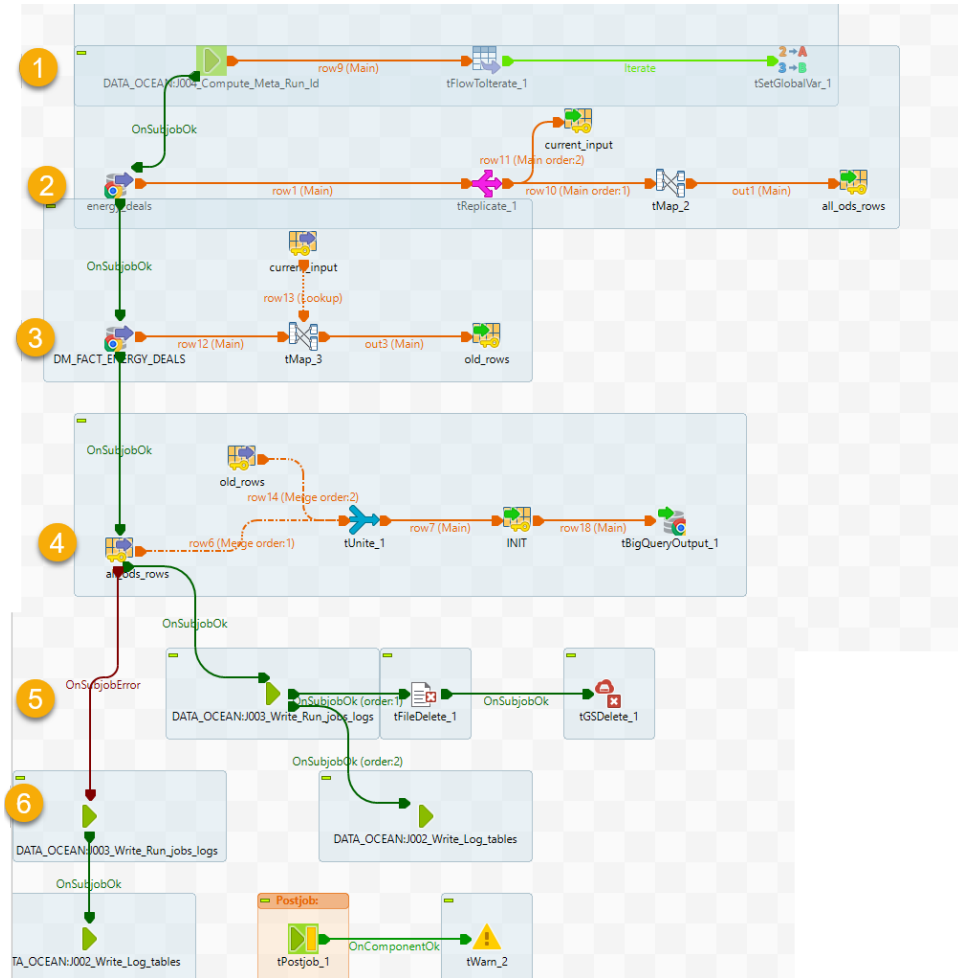
<p>Main jobs for source extraction</p>	<ul style="list-style-type: none"> • J009_0_DS_TO_DM_FACT_ENERGY_DEALS J009_ODS_TO_DM_FACT_IRM_DEALS <p>(DATA_OCEAN_DOMAIN_SUSTAINABILITY)</p>	<p>--to the top --</p> <p>GCP Remote Engine</p> <p>Industrial DO project</p>
--	--	--

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

	<p>1. Job Metadata Calculation:</p> <ul style="list-style-type: none"> • The job initiates by calculating essential job metadata. This metadata 	
--	--	--

include information like job identifiers, timestamps, or job-specific configuration details.

- This step establishes the foundation for the subsequent data processing.
- 2. Get full load from ODS
- 3. Get data from DM. FACT_irm_energy_deals_daily in case there is some data, which not exist in ODS
- 4. Combine data from 2 and 3 to DM. FACT_irm_energy_deals_daily with truncate option
- 5. Log Generation:
 - Upon the completion of the data



enrichment and loading process, the job generates logs. These logs serve as a record of the job's activities and provide visibility into the data processing steps.

6. Temporary File Deletion:

- As a final step, the jobs take care of deleting any temporary files created during the data transfer and enrichment process. This cleanup ensures efficiency

ent
stora
ge
reso
urce
man
age
ment.

Main jobs for source extraction

- **J001_DM_TO_DM_FACT_ENERGY_PRICE_HOURLY**
(ROBUSTIFY)

Job description by steps

Data Extraction and Transformation:

1. The code extracts data related to energy prices and contracts from a database.
 - The source table is DataOcean.V_FACT_energy_irm
 - To get the raw data from IRM in ODS and transform to the format that Dataiku can consume
 - Split data into hourly
 - The past to current date will get from Spot
 - The current month will get from the first date of the next month
 - Next month data will be all the same within the month by the first date of each month

Explanation the job:

1. Setup global variable on output file / source table and define the value of gas = 0.0
2. Get data from the source with format below

delivery_date	elec_DE_EPEX	elec_FR_EPEX	elec_IT_PUN	elec_SP_OMEI	gas_DE_THE	gas_FR_PEG	gas_IT_PSV	gas_SP_PVB	gas_TTF
2024-06-21 21:00:00 UTC	105.26	105.14	135.3151	114.86	0.0	0.0	0.0	0.0	0.0
2024-06-21 22:00:00 UTC	95.99	96.49	115.0116	115.0	0.0	0.0	0.0	0.0	0.0
2024-06-21 23:00:00 UTC	82.14	83.26	112.0	105.74	0.0	0.0	0.0	0.0	0.0
2024-07-01 00:00:00 UTC	73.15	50.94	106.9	74.56	34.5	34.31	36.15	34.45	34.48
2024-08-01 00:00:00 UTC	76.32	50.1	106.04	77.89	34.92	34.64	35.44	34.43	34.87

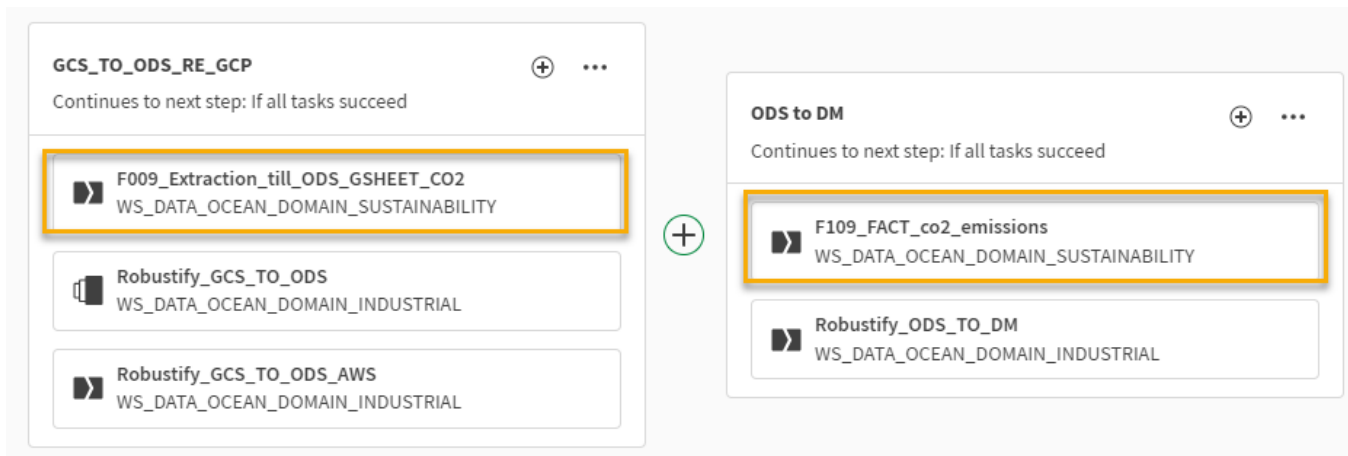
3. Apply the business rule and save into the output file by JavaFlex
4. Transform to the Dataiku format

date	elec_DE_EPEX	elec_FR_EPEX	elec_IT_PUN	elec_SP_OMEI	gas_DE_THE	gas_FR_PEG	gas_IT_PSV	gas_SP_PVB	gas_TTF	elec_IT_MGP	inserted_date
2024-06-21 21:00:00 UTC	105.26	105.14	135.3151	114.86	33.963	34.028	36.875	34.49	33.867	null	2024-06-21 10:09:55 UTC
2024-06-21 22:00:00 UTC	95.99	96.49	115.0116	115.0	33.963	34.028	36.875	34.49	33.867	null	2024-06-21 10:09:55 UTC
2024-06-21 23:00:00 UTC	82.14	83.26	112.0	105.74	33.963	34.028	36.875	34.49	33.867	null	2024-06-21 10:09:55 UTC
2024-06-22 00:00:00 UTC	73.15	50.94	106.9	74.56	34.5	34.31	36.15	34.45	34.48	null	2024-06-21 10:09:55 UTC
2024-06-22 01:00:00 UTC	73.15	50.94	106.9	74.56	34.5	34.31	36.15	34.45	34.48	null	2024-06-21 10:09:55 UTC
2024-06-22 02:00:00 UTC	73.15	50.94	106.9	74.56	34.5	34.31	36.15	34.45	34.48	null	2024-06-21 10:09:55 UTC
2024-06-22 03:00:00 UTC	73.15	50.94	106.9	74.56	34.5	34.31	36.15	34.45	34.48	null	2024-06-21 10:09:55 UTC
2024-06-22 04:00:00 UTC	73.15	50.94	106.9	74.56	34.5	34.31	36.15	34.45	34.48	null	2024-06-21 10:09:55 UTC
2024-06-22 05:00:00 UTC	73.15	50.94	106.9	74.56	34.5	34.31	36.15	34.45	34.48	null	2024-06-21 10:09:55 UTC

5. Delete temporary file

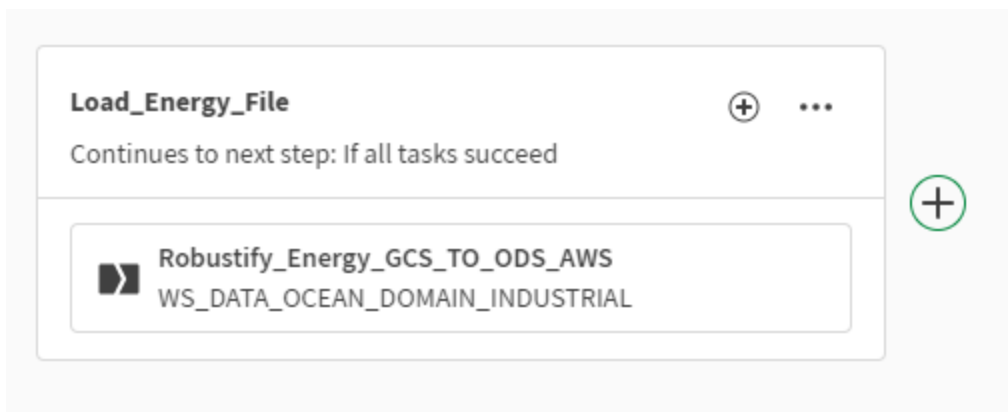
4.5 - Scheduling and Automation

Although there are 3 projects involve but main plan will be in TMC on WS_DATA_OCEAN_DOMAIN INDUSTRIAL with plan PL_INDUS_ROBUSTIFY_LOAD load - Daily at 12:15 PM.



**Yellow box is the job in Sustainability domain.

PL_INDUS_ROBUSTIFY_ENERGY load hourly at xx:30



Monitor the loading in the log tables on prj-data-dm-industrial-[env] by

```
select job.job_name , job.meta_start_date , logs.meta_run_id , logs.meta_source_system , logs.meta_step , logs.meta_status , logs.meta_num_lines ,
logs.meta_error_lines from STG.log_tables logs join STG.run_jobs job on logs.meta_run_id = job.meta_run_id
where logs.meta_run_id in ( SELECT meta_run_id FROM STG.run_jobs order by meta_start_date desc limit 1000 )
and (
job_name like '%METEOROLOGICA%' or
upper(job_name) like '%ENERGY%' or
job_name like '%WAP%' or
job_name like '%IRM%' )
and meta_start_date > DATE_SUB ( CURRENT_TIMESTAMP () , INTERVAL 1 DAY )
order by job.meta_start_date desc
```

4.6 - Remark

1. Most of Talend jobs will be on project DATA_OCEAN_DOMAIN_INDUSTRIAL except:

- Gsheet Co2 (J009_ODS_TO_DM_FACT_ENERGY_DEALS) will be in DATA_OCEAN_DOMAIN_SUSTAINABILITY prj-data-dm-sust-[env].STG.
STG_FIL_0000_0000_F001_F_D_co2_emissions

- DM to DM (J001_DM_TO_DM_FACT_ENERGY_PRICE_HOURLY) will be in ROBUSTIFY project prj-data-robustify-[env].DM.
FACT_energy_price_hourly

This is following [data architecture](#)

2. Most of Talend jobs are required to use remote engine on Cloud except J014_Extraction_till_ODS_ORACLE_IRM is required AWS. This is because of security reason