

# DATAWAREHOUSING & ETL

Hénok Agbodjogbe  
Michael Bruen  
Théo Ourvoie  
Yassine El Harrab  
Loïc Martins

# Index

- I. [Introduction](#)
- II. [Data Understanding and Exploration](#)
  - a. [Discovery](#)
  - b. [Exploration](#)
- III. [ETL Pipeline](#)
  - a. [Data Extraction](#)
  - b. [Data Transformation](#)
  - c. [Data Loading](#)
- IV. [Project Deployment](#)
- V. [Analysis](#)
  - 1. [Which product categories have the highest customer satisfaction?](#)
  - 2. [Which month do customers order the most?](#)
  - 3. [How many orders have not been delivered?](#)
  - 4. [In which of the company's cities do we have the biggest profits?](#)
  - 5. [Are top-rated products with a low average distance between customer and seller?](#)

# Introduction

The marketplace has a complex architecture because it must work with multiple sellers, multiple products, and thousands of customers. All these interactions generate a significant amount of data. To function well, the marketplace must leverage this data effectively to manage the platform and make the right decisions for future strategies.

This is why we need to create ETL processes and a data warehouse to have clean and structured data that can be used for analysis purposes.

To do it, we had a group composed of Hénok Agbodjogbe, Michael Bruen, Théo Ourvoie, Yassine El Harrab and Loic Martins.

At the organizational level, each member of the team worked on their own, and every week we had a meeting to take stock of the work done. Each member participated but Henok managed the technical part of the project.

Our project includes 4 stages:

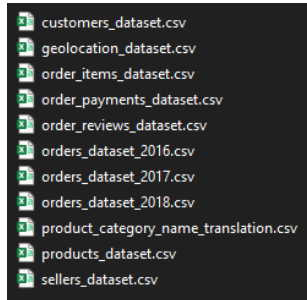
1. Data Understanding and Exploration
2. ETL Pipeline
3. Project Deployment
4. Analysis

# I. Data Understanding and Exploration

## 1. Discovery

The exploration phase is one of the most important phases because it provides the foundations of our project.

We start with a large dataset including 11 Excel files:



Without opening the tables, we can note various points:

- We have a complete view of the data related to a marketplace: orders, products, sellers, customers, payment.
- We have 3 tables concerning the orders with different dates.
- We have 3 types of Orders table: the main table with the order, the reviews, the payments and the items.
- We have 2 tables related to the products: the main table and the categories table.

## 2. Exploration

To further understand the dataset, we need to have a better representation of each tables:

product_category_name_translation	products_dataset	sellers_dataset	orders_dataset
<i>Mapping of product category names.</i>	<i>Details about the products listed in the marketplace.</i>	<i>Sellers on the marketplace.</i>	<i>Details about customer orders.</i>
product_category_name product_category_name_english	product_id product_category_name product_name_lenght product_description_lenght	seller_id seller_zip_code_prefix seller_city seller_state	order_id customer_id order_status order_purchase_timestamp order_approved_at

	product_photos_qty product_weight_g product_length_cm product_height_cm product_width_cm		order_delivered_carrier_date order_delivered_customer_date order_estimated_delivery_date
--	--	--	--

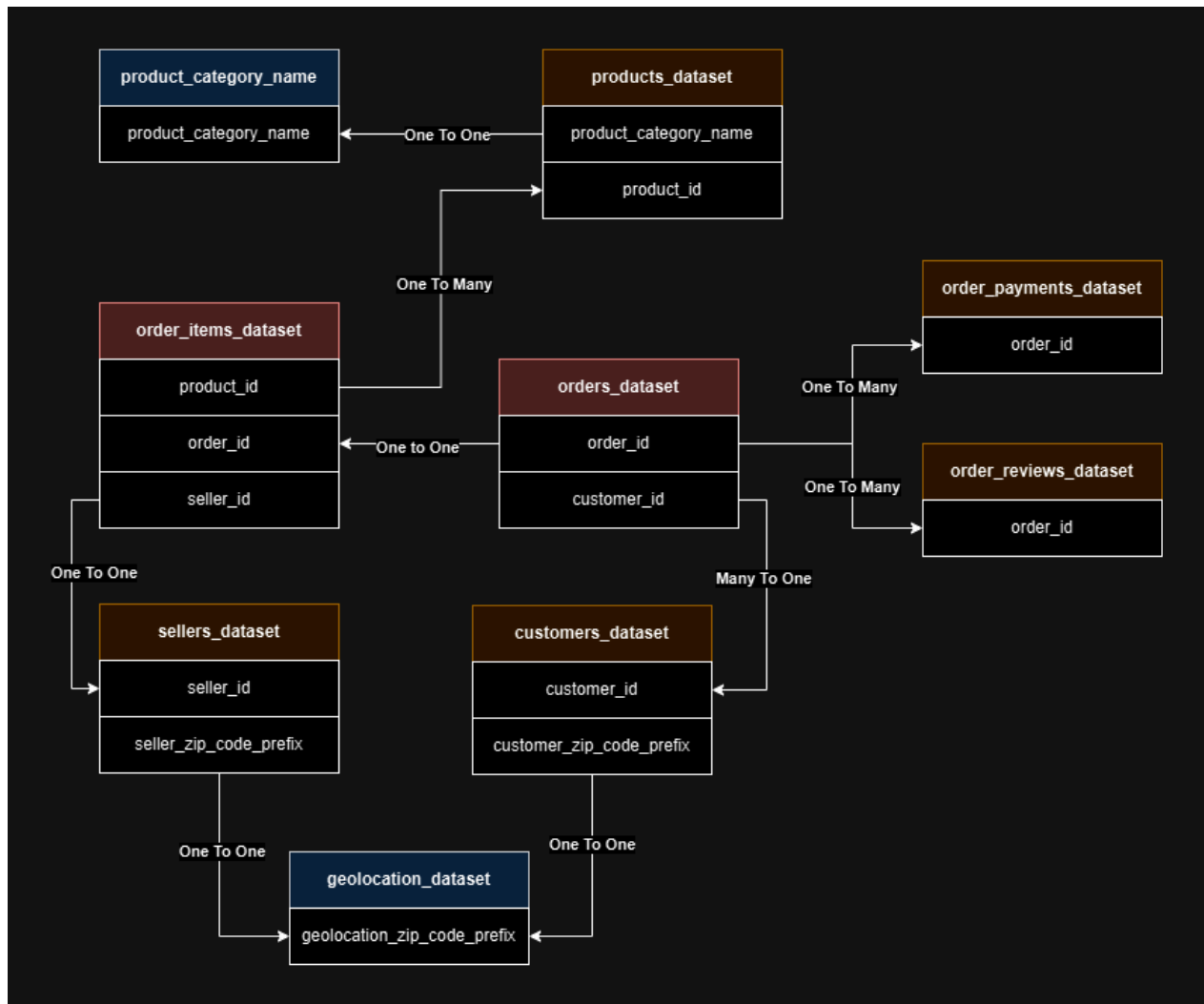
order_items_dataset	order_payments_dataset	order_reviews_dataset	geolocation_dataset
<i>Information about the individual items.</i>	<i>Details about payments made for orders.</i>	<i>Customer reviews for orders.</i>	<i>Geolocation information.</i>
order_id order_item_id product_id seller_id shipping_limit_date price freight_value	order_id payment_sequential payment_type payment_installments payment_value	review_id order_id review_score review_comment_title review_comment_message review_creation_date review_answer_timestamp	geolocation_zip_code_prefix geolocation_lat geolocation_lng geolocation_city geolocation_state

customers_dataset
<i>Customer information.</i>
customer_id customer_unique_id customer_zip_code_prefix customer_city customer_state

Looking at these different tables, we can notice that:

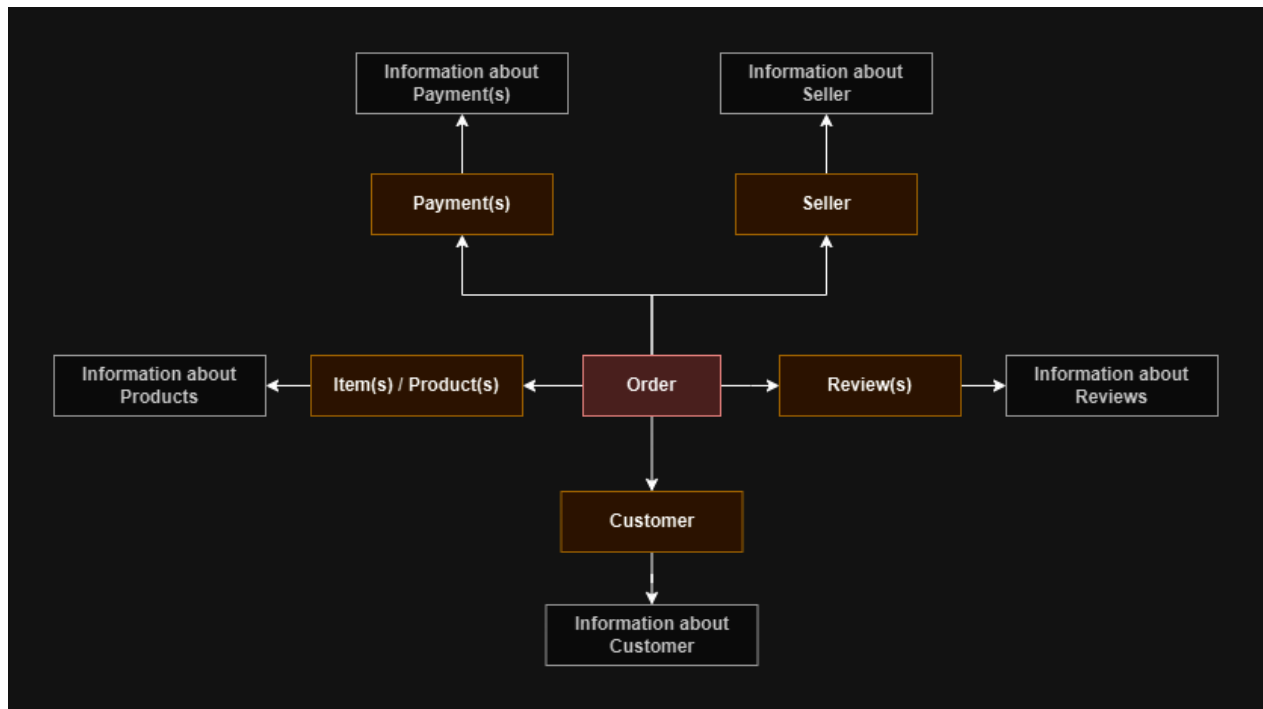
- The Order ID is our main field/variable because it is present in 4 tables.
- We can conclude that the dataset is centered around “the order”. When an order is placed, it creates an order\_id.

We now have a better understanding of what the table contains but we need to understand their relationships and how they relate to each other:



Using this diagram, we can notice that:

- The Orders table is at the center of the whole dataset because this is where all the orders are recorded first.
- The Order Items table plays an important role because thank to it we can access the Sellers and Products tables.
- Using our first diagram we can create a map around the order:



- We can notice that for one order we have:
  - Item(s) / Product(s)
  - Customer
  - Review(s)
  - Payment (s)
  - Seller

## II. ETL Pipeline

The pipeline will be done in three steps:

1. Data extraction will be done where the data will be extracted and placed in the staging area (STA). This will allow us to load the data as is, or with minimal changes.
2. The data transformation or Operational Data Store area (ODS) will allow us to clean and standardize the data.  
If the data does not pass the quality criteriums, they will be put in the “Technical\_Rejects” table as technical rejects.
3. The data loading or Data WareHouse area (DWH) will organize the data in one fact table related to multiple dimensions tables. If records can’t be integrated in the schema, they will be put in the “Functional\_Rejects” table as functional reject. Alternatively, some placeholder relations can be created.

There will be one STA and ODS package per file.

### 1. Data Extraction

Here, the role of data extraction or the staging database is to store all the data coming from the different sources. We want to accept all available data.

#### a) Orders Table

Here is an extract of the Orders\_2016.csv file:

<u>order_id</u>	<u>customer_id</u>	<u>order_status</u>	<u>order_purchase_timestamp</u>	<u>order_approved_at</u>
d3c8851a6651eeff2f73b0e011ac45d0	957f8e082185574de25992dc659ebbc0	processing	5/10/2016 22:44	6/10/2016 15:51
cbbb524a0e1646aa6cf7a3c0bbe517ad	dacb079d55ffb1d3955c5d923df3ebb7	delivered	5/10/2016 7:31	6/10/2016 2:46
ac2b7c522d811acba0aa270ed3e112e4	ef21aebbb093a6db29ccc6aa0b89c347	delivered	5/10/2016 15:08	6/10/2016 15:44
7033745709b7cf1bac7d2533663592de	7f0ca17bb33b230b47459437cf0682c7	delivered	4/10/2016 14:13	4/10/2016 14:46
5cd498954e2b37d71b315166809b4bd7	ff1a56726b7ea149c7423865609cc0c8	delivered	7/10/2016 0:54	8/10/2016 3:56

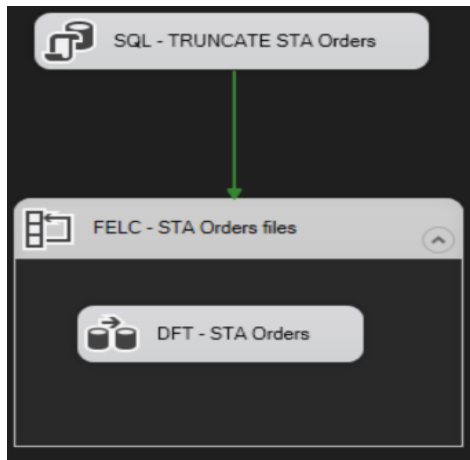
<u>order_delivered_carrier_date</u>	<u>order_delivered_customer_date</u>
10/10/2016 2:46	16/10/2016 14:36
10/10/2016 15:44	13/10/2016 15:44
8/10/2016 14:46	11/10/2016 14:46
25/10/2016 11:35	27/10/2016 17:32



There are two other files with the same columns for orders 2017 and 2018.

To be able to use the orders data, we need to put it into one table. This means that we need to extract the data from each file.

To go into all the orders files, we use a “Foreach Loop Container” where we put the extraction dataflow inside. The data flow in the container will be able to load one file at the time. Here, we also truncate the data from the previous runs.



The foreach loop is iterating over a variable “sales” that allow it to address each file separately. The variable is of the format “sales\*”, with the star meaning it will access all files that start with sales.

The screenshot shows the 'Foreach Loop Editor' configuration window. It has a table at the top with two columns: 'Enumerator' and 'Foreach File Enumerator'. Below the table is a section titled 'Foreach Loop Editor' containing 'Enumerator configuration'. This section includes a 'Folder:' field with the path 'C:\Users\HP\Desktop\Projects\ETL and DWH\Assignement\DATA' and a 'Files:' field with the pattern 'orders\_\*.csv'. At the bottom, there are radio buttons for 'Retrieve file name' with options: 'Name and extension', 'Fully qualified' (which is selected), and 'Name only'.

Enumerator	Foreach File Enumerator
Expressions	

**Foreach Loop Editor**

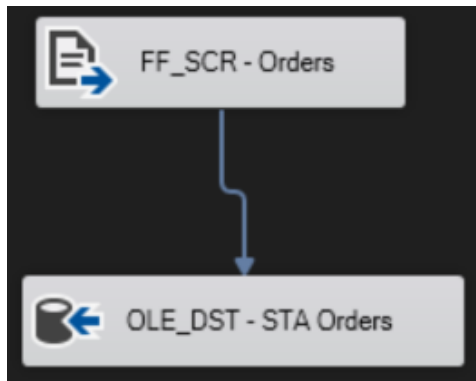
Enumerator configuration

Folder:  
C:\Users\HP\Desktop\Projects\ETL and DWH\Assignement\DATA

Files:  
orders\_\*.csv

Retrieve file name  
☐ Name and extension ☒ Fully qualified ☐ Name only

Next, the dataflow is defined as shown below:



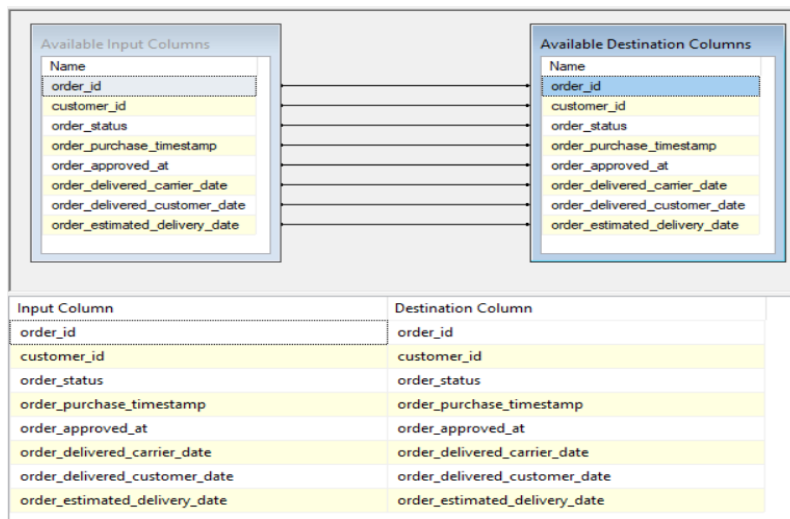
Once we have extracted the data, we can load it into our target table “Orders” in the STA database. We start this by creating a table with the following script:

```
USE PROJECT_STA
CREATE TABLE [dbo].[Orders] (
    [order_id] varchar(50),
    [customer_id] varchar(50),
    [order_status] varchar(50),
    [order_purchase_timestamp] varchar(50),
    [order_approved_at] varchar(50),
    [order_delivered_carrier_date] varchar(50),
    [order_delivered_customer_date] varchar(50),
    [order_estimated_delivery_date] varchar(50)
)
```

We then define the target destination.

The screenshot shows the 'OLE DB connection manager' dialog box. The 'Data source' dropdown is set to 'localhost.Project\_STA'. The 'Data access mode' dropdown is set to 'Table or view - fast load'. The 'Name of the table or the view' dropdown is set to '[dbo].[Orders]'. There are two 'New...' buttons. The 'Keep identity' checkbox is unchecked, and the 'Keep nulls' checkbox is unchecked. The 'Table lock' checkbox is checked, and the 'Check constraints' checkbox is checked. The 'Rows per batch' field is empty, and the 'Maximum insert commit size' field is set to '2147483647'.

And finally, we define the column mapping as follows:

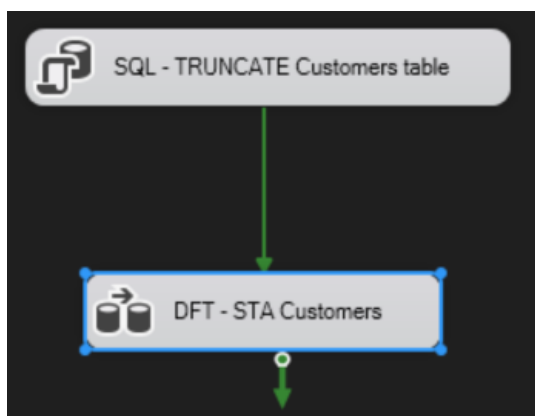


Here are the first ten lines of the results table

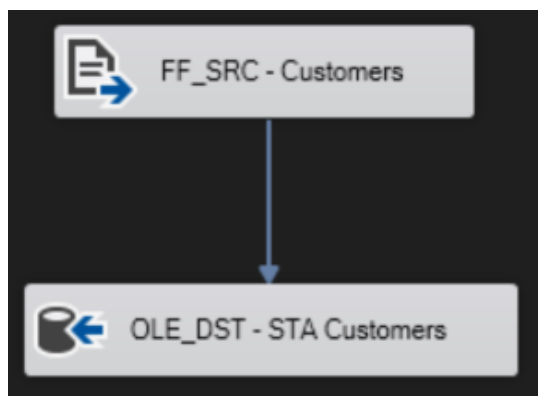
	order_id	customer_id	order_status	order_purchase_timestamp	order_approved_at	order_delivered_carrier_date	order_delivered_customer_date	order_estimated_delivery_date
1	fc1d9429da446a147d438d41350a2932	8b13d2050940cb8d6a7203e69f279254	delivered	23/10/2017 16:29	24/10/2017 7:14	06/11/2017 22:11	18/11/2017 14:31	23/11/2017 0:00
2	8a6df1083cc2efc062e037e474622de6	64f0ed97c13c1388088e93454db40488	delivered	20/03/2017 21:54	20/03/2017 21:54	25/03/2017 9:26	01/04/2017 8:47	10/04/2017 0:00
3	bcd37480e921eaf2b524c8e9724d099	fbcb39b74c648a39b750799b58a13e0f	delivered	24/11/2017 22:55	25/11/2017 1:53	28/11/2017 21:26	07/12/2017 0:24	27/12/2017 0:00
4	b951cb562f56db404cd77c81262f25	83116a97d5035c15ca0e75dc2b848bd9	delivered	26/12/2017 10:24	26/12/2017 10:36	27/12/2017 19:42	04/01/2018 21:27	24/01/2018 0:00
5	5b8e4969e32884e4f0260f1b6bd6c3d	da9f87409fcc1f75e9bce0f233ace2c	shipped	06/07/2017 14:42	06/07/2017 14:55	07/07/2017 12:04		28/07/2017 0:00
6	b825b53d8136d85eb4cbe60a1bd279f5	ec4f3f98da53dbaa4db16505b3782c0	delivered	05/09/2017 14:48	06/09/2017 2:50	06/09/2017 21:22	12/09/2017 20:53	22/09/2017 0:00
7	b1e68c1ac52a1a5cb5051d07ae91ff06	cf34464d3b5d4a957dae9c0ef4a7d27e	delivered	18/03/2017 11:36	18/03/2017 11:36	21/03/2017 8:45	27/03/2017 13:22	26/04/2017 0:00
8	4516299f280a0f30720256ae0cd592d6	582ec2f9229bb3014866b31f73def34f	delivered	13/09/2017 18:12	13/09/2017 18:25	14/09/2017 18:53	18/10/2017 13:05	09/10/2017 0:00
9	7d63cb1d349e7bd0ee1edcc61ea71077	484711e3d14ac75fb874e949b3a28395	delivered	14/08/2017 15:30	14/08/2017 15:45	16/08/2017 17:04	24/08/2017 18:05	04/09/2017 0:00
10	b11073519dc9fe626d#6281aa1965522	4f3ef3bae0d531a175d56f875c8a84d1	delivered	16/11/2017 12:29	16/11/2017 13:10	21/11/2017 1:56	27/11/2017 16:22	05/12/2017 0:00

## b) Customers Table

Next is the extraction of the data from “customers\_dataset.csv”. For this table, we don’t need to add additional data. Here we just make sure to truncate the table “Customers” before running the package:



The dataflow is an import of a flat file:



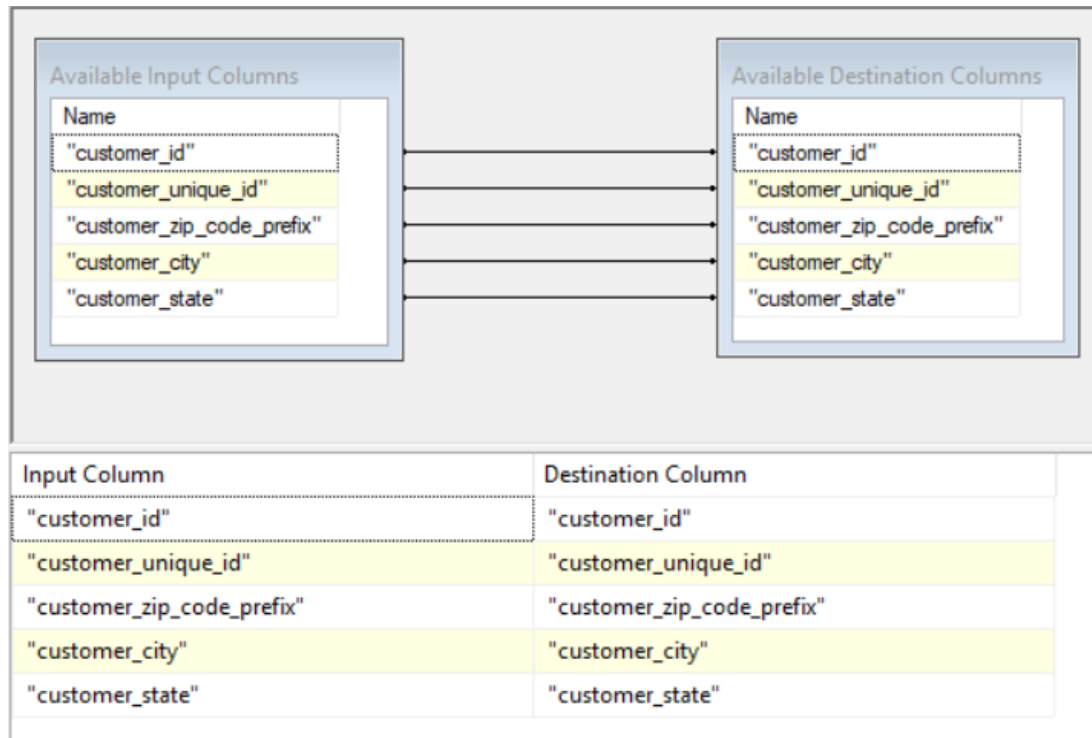
We create the destination table with the following command:

```
USE PROJECT_STA
CREATE TABLE [dbo].[Customers] (
    ["customer_id"] varchar(50),
    ["customer_unique_id"] varchar(50),
    ["customer_zip_code_prefix"] varchar(50),
    ["customer_city"] varchar(50),
    ["customer_state"] varchar(50)
)
```

The data is then loaded into the table “Customers”:

The screenshot shows the "OLE DB connection manager" dialog box. It has three main sections: "OLE DB connection manager:" with a dropdown menu showing "localhost.Project\_STA" and a "New..." button; "Data access mode:" with a dropdown menu showing "Table or view"; and "Name of the table or the view:" with a dropdown menu showing "[dbo].[Customers]" and another "New..." button.

We define the column mapping as follows:



Here is the first ten lines of the results:

	"customer_id"	"customer_unique_id"	"customer_zip_code_prefix"	"customer_city"	"customer_state"
1	a6b43413c5d37d7dc3173432f6ab2	212b273eed00f96498ccd8e91427e42	86840	faxinal	PR
2	c1bdc290c9ae68e5f0825b2dfed4d422	ae5df8e33ea496f14db277c8b8e77281	13098	campinas	SP
3	0a3e4f7efaf1bfc4a84e1adcd9b70f5	273d136f7121c493a5c994c4838fa350	59196	pedro velho	RN
4	9c64abffb2de9888aa146da8c65a547f	18e362d58b54666925aacfd86ff84213	04944	sao paulo	SP
5	04edcf9c41ea59408793171236b28e2	4e17943c8edf2f66d70f0e97b81aaf17	28085	campos dos goytacazes	RJ
6	1f954c777b0dbd306b017a3e4cd26475	d35862da9bdb38fad3994de72a77a3c2	38025	uberaba	MG
7	b769ee1c67d0e4f3c5b0abe1a1b958d4	3726ac7519981c0c5a86b98ca38297a4	28055	campos dos goytacazes	RJ
8	03bbe0ce5c28e05f22917607db798818	8f3dca4306d5a89e4ae2c65c110603a2	72465	brasilia	DF
9	cf833bbe73b3d8ddeb9889d4192b2496	03f923dbc3b3c763b22b24db6072d27	13400	piracicaba	SP
10	9716c4a296c261a120d27d058e32ab38	0647d94e80c8033ef27686bdb447e105	36015	juiz de fora	MG

The rest of the datasets are extracted the same way as the customers one above.

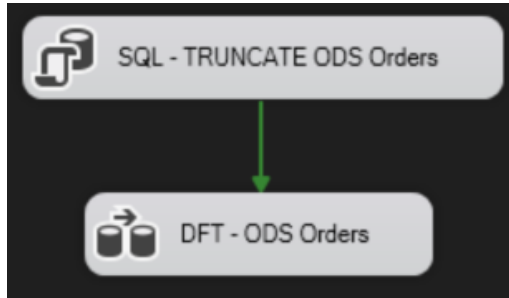
## 2. Data Transformation

The second step of the pipeline is to load usable data and transform it before putting it into the Operational Data Store (ODS). This means we need to transform the data into a usable format. We also need to clean and standardize the data. All the data that does not respect the “quality standards” will be rejected as a technical reject.

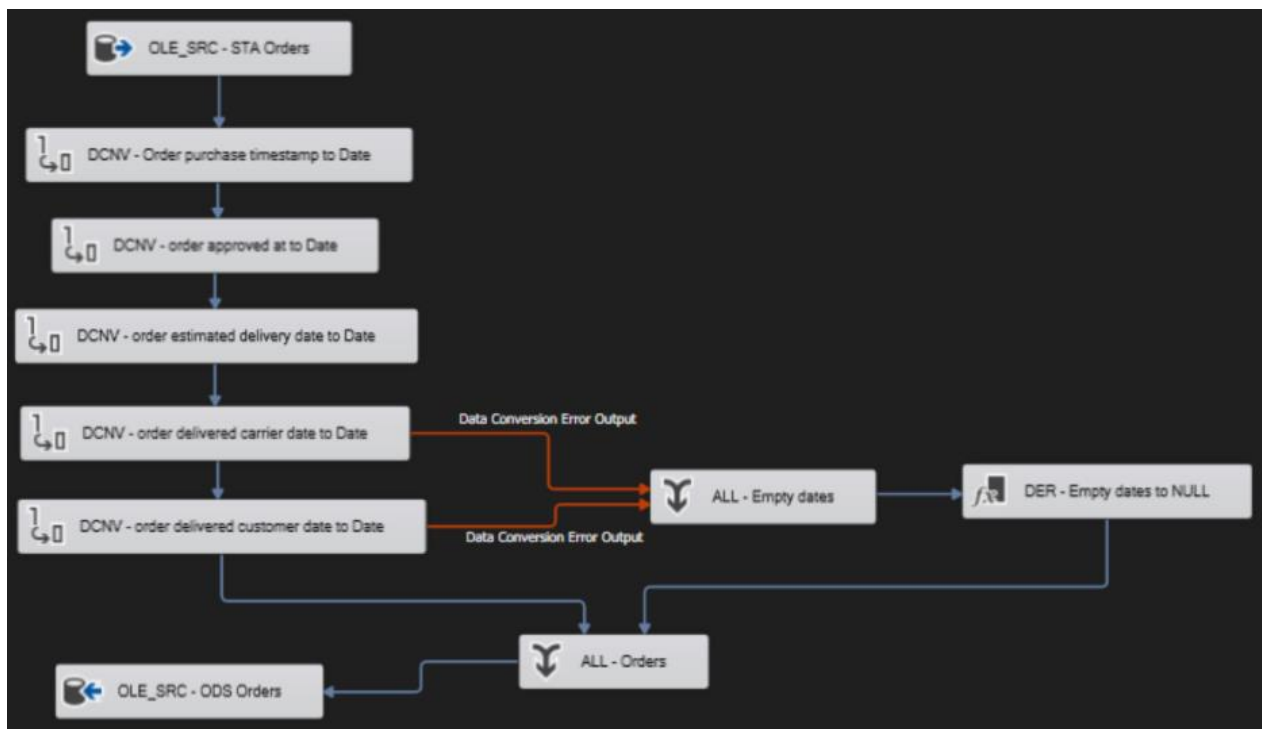
The quality standards are based on the “correctness” of the data. The output data must be consistent in data types and in values. We also need to ensure that the data can be used in queries, so we might need to reorganize and enrich the data.

## a) Orders Table

Like before, we truncate the data from the previous runs.

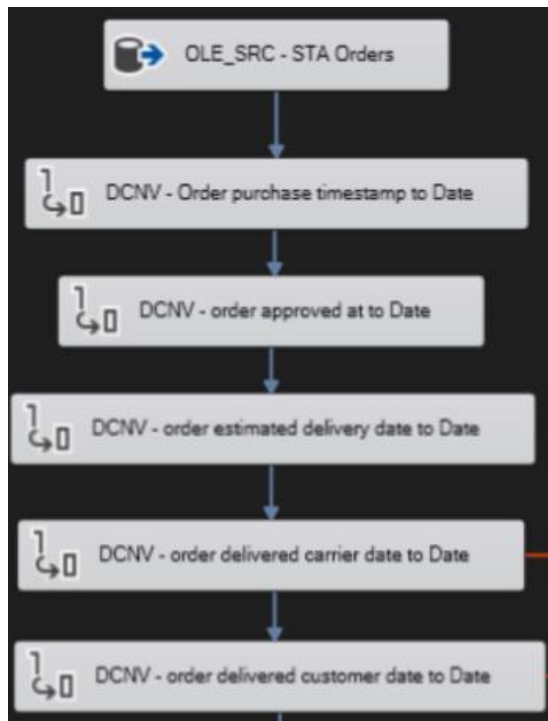


The data flow is defined below:



Below, the segments of this dataflow are explained:

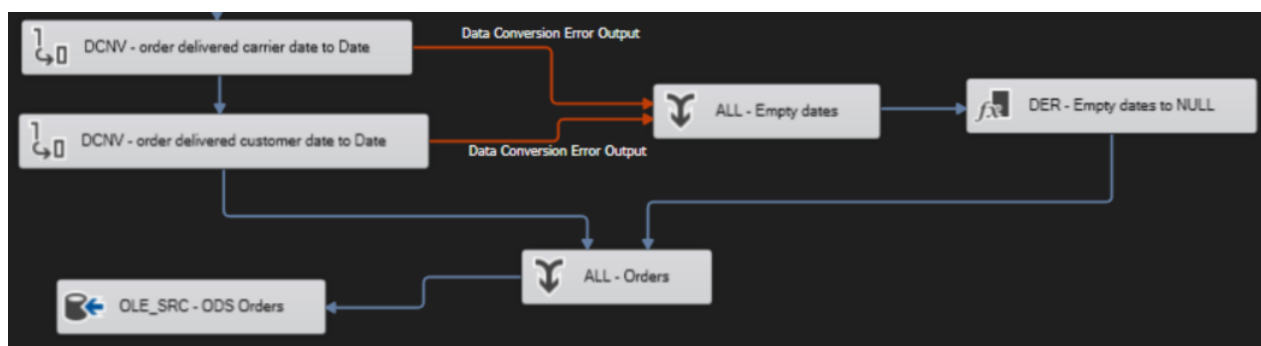
In the first step as shown below, we extract the data from STA orders. After this the columns that have date data are changed from datatype string to datatype date. Date is chosen instead of datetime because in the DWH later on we need to map this data with a column with the type date.



Within the tasks, the datatype of the column is being changed to date as shown below:

Input Column	Output Alias	Data Type
order_purchase_timestamp	order_purchase_timestamp_Date	date [DT_DATE]

Order delivered carrier date and order delivered customer date have several rows that do not contain any data at all. To handle this, the rows that produce a data conversion error have their values changed to NULL before being fed back into the dataset. This is shown in the screenshots below.



Derived Column Name	Derived Column	Expression	Data Type
order_delivered_carrier_date_NULL	<add as new column>	NULL(DT_DATE)	date [DT_DATE]
order_delivered_customer_date_NULL	<add as new column>	NULL(DT_DATE)	date [DT_DATE]

After all this is done the ODS Orders table is created and the data is inserted into the database as shown below.

```

1  USE Project_ODS
2  CREATE TABLE [dbo].[Orders] (
3      [order_id] varchar(50),
4      [customer_id] varchar(50),
5      [order_status] varchar(50),
6      [order_purchase_timestamp] Date,
7      [order_approved_at] Date,
8      [order_delivered_carrier_date] Date,
9      [order_delivered_customer_date] Date,
10     [order_estimated_delivery_date] Date
11 )

```


OLE DB connection manager:

localhost.Project\_ODS New...

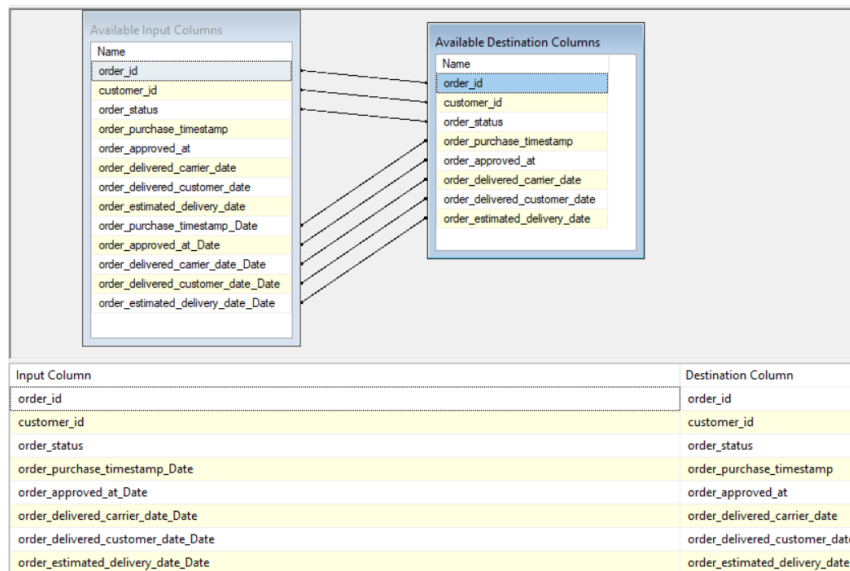
Data access mode:

Table or view - fast load ▼

Name of the table or the view:

 [dbo].[Orders] ▼ New...



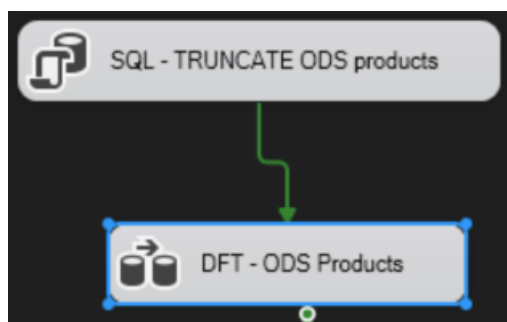


Here is the first ten lines of the results:

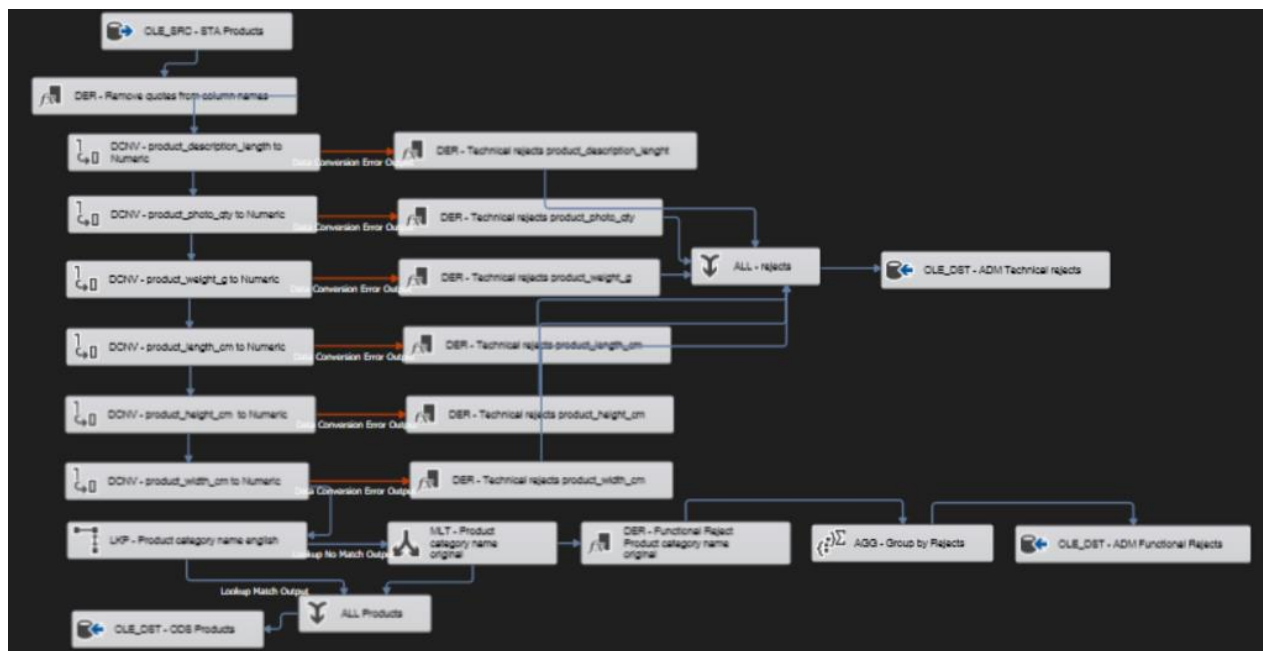
	order_id	customer_id	order_status	order_purchase_timestamp	order_approved_at	order_delivered_carrier_date	order_delivered_customer_date	order_estimated_delivery_date
1	fc1d9429da446a147d438d41350a2932	8b13d2050940cb8d6a7203e69f279254	delivered	2017-10-23	2017-10-24	2017-06-11	2017-11-18	2017-11-23
2	8a6bf1083cc2efc062e037e474622de6	64f0ed97c13c138808e93454db40488	delivered	2017-03-20	2017-03-20	2017-03-25	2017-01-04	2017-10-04
3	bcd37480e9f21eaf2b524c8e9724d099	fb0839b74c648a39b750799b58a13e0f	delivered	2017-11-24	2017-11-25	2017-11-28	2017-07-12	2017-12-27
4	b951cb6562f56db404cd77c81262f25	83116a97d5035c15ca0e75dc2b848bd9	delivered	2017-12-26	2017-12-26	2017-12-27	2018-04-01	2018-01-24
5	b825b53d8136d85eb4cbe60a1bd279f5	ec4f3f98da53dbaa4db16505b3782c0	delivered	2017-05-09	2017-06-09	2017-06-09	2017-12-09	2017-09-22
6	b1e68c1ac52a1a5cb5051d07ae91ff06	cf34464d3c5d4a957dae9c0ef4a7d27e	delivered	2017-03-18	2017-03-18	2017-03-21	2017-03-27	2017-04-26
7	4516299f280a0f30720256ae0cd592d6	582ec2f9229bb3014866b31f73def34f	delivered	2017-09-13	2017-09-13	2017-09-14	2017-10-18	2017-09-10
8	7d63cb1d349e7bd0ee1edcc61ea71077	484711e3d14ac79b874e949b3a28395	delivered	2017-08-14	2017-08-14	2017-08-16	2017-08-24	2017-04-09
9	b11073519dc9fe626df6281aa1965522	4f3ef3bae0d531a175d6f875cba84d1	delivered	2017-11-16	2017-11-16	2017-11-21	2017-11-27	2017-05-12
10	c7b440776d014b67c3164e103e3bf02	19187bc9fab7457748d4dc42a9e1c6fe	delivered	2017-11-21	2017-11-21	2017-11-23	2017-11-30	2017-01-12

## b) Products Table

Like before, we truncate the data from the previous runs.



The data flow is defined below:



Below, the segments of this dataflow are explained:

In the first step as shown below, we extract the data from STA products and remove the quotes from the column names.

Derived Column Name	Derived Column	Expression	Data Type	Length	Precision	Scale	Code Page
product_id	Replace 'product_id'	product_id	string [DT_STR]	50			1252 (ANSI - Latin I)
product_category_name	Replace 'product_category_name'	product_category_name	string [DT_STR]	50			1252 (ANSI - Latin I)
product_description_length	Replace 'product_description_length'	product_description_length	string [DT_STR]	50			1252 (ANSI - Latin I)
product_photos_qty	Replace 'product_photos_qty'	product_photos_qty	string [DT_STR]	50			1252 (ANSI - Latin I)
product_weight_g	Replace 'product_weight_g'	product_weight_g	string [DT_STR]	50			1252 (ANSI - Latin I)
product_length_cm	Replace 'product_length_cm'	product_length_cm	string [DT_STR]	50			1252 (ANSI - Latin I)
product_height_cm	Replace 'product_height_cm'	product_height_cm	string [DT_STR]	50			1252 (ANSI - Latin I)
product_width_cm	Replace 'product_width_cm'	product_width_cm	string [DT_STR]	50			1252 (ANSI - Latin I)

After this there are several columns that contain numerical data. These are changed as shown below.

Input Column	Output Alias	Data Type	Length	Precision	Scale	Code Page
product_description_length	product_description_length_Numeric	numeric [DT_NUMERIC]		18	0	

If any of the data is unable to be changed to numeric, the error is tracked and put in the “technical\_rejects” table.

Derived Column Name	Derived Column	Expression	Data Type	Length
Rejectdate	<add as new column>	GETDATE()	database timestamp [DT_DBTIMESTAMP]	
RejectPackageAndTask	<add as new column>	(DT_WSTR,100)@[System::PackageName] + " AND " + (DT_WSTR,100)@[System::TaskName]	Unicode string [DT_WSTR]	205
RejectColumn	<add as new column>	"product_description_length"	Unicode string [DT_WSTR]	26
RejectDescription	<add as new column>	"The value " + (DT_WSTR,100)product_description_length + " is not a valid integer"	Unicode string [DT_WSTR]	133

As seen above the data inserted in the technical rejects are:

4. The date of the error.
5. The column making the error.
6. An error message.
7. The package and the task causing the error.

The technical rejects table is created in the database using the following script:

```

1  USE [Project_ADM]
2  CREATE TABLE [dbo].[TechnicalRejects] (
3      [Rejectdate] datetime,
4      [RejectPackageAndTask] nvarchar(205),
5      [RejectColumn] nvarchar(29),
6      [RejectDescription] nvarchar(133)
7  )

```

This datatype conversion is performed on 6 columns. Afterwards a lookup is performed on the product category name column using the product category name translation dataset that was loaded into STA. this will change the Spanish names to English ones.

Lookup Column	Lookup Operation	Output Alias
product_category_name_english	<add as new column>	product_category_name_english

The product names that do not have a matching translation name get tracked and sent to the functional rejects table through the steps mapped out earlier and shown in more detail below.

Derived Column Name	Derived Column	Expression	Data Type	Length
RejectDate	<add as new column>	GETDATE()	database timestamp [DT_DBTIMESTAMP]	
RejectPackageAndTask	<add as new column>	(DT_WSTR,100)@[System::PackageName] + " AND " + (DT_WSTR,100)@[System::TaskName]	Unicode string [DT_WSTR]	205
RejectColumn	<add as new column>	"Product_category_name"	Unicode string [DT_WSTR]	21
RejectDescription	<add as new column>	(DT_WSTR,100)product_category_name + " has not translation in english"	Unicode string [DT_WSTR]	131

Available Input Columns		
<input checked="" type="checkbox"/>	Name	
<input type="checkbox"/>	product_photos_qty	
<input type="checkbox"/>	product_weight_g	
<input type="checkbox"/>	product_length_cm	
<input type="checkbox"/>	product_height_cm	
<input type="checkbox"/>	product_width_cm	
<input type="checkbox"/>	product_description_leng...	
<input type="checkbox"/>	product_photos_qty_Nu...	
<input type="checkbox"/>	product_weight_g_Numeric	
<input type="checkbox"/>	product_length_cm_Num...	
<input type="checkbox"/>	product_height_cm_Num...	
<input type="checkbox"/>	product_width_cm_Nume...	
<input checked="" type="checkbox"/>	RejectDate	
<input checked="" type="checkbox"/>	RejectPackageAndTask	
<input checked="" type="checkbox"/>	RejectColumn	
<input checked="" type="checkbox"/>	RejectDescription	

Input Column	Output Alias	Operation
(*)	Count all	Count all
RejectDate	RejectDate	Group by
RejectPackageAndTask	RejectPackageAndTask	Group by
RejectColumn	RejectColumn	Group by
RejectDescription	RejectDescription	Group by

```

1  USE [Project_ADM]
2  CREATE TABLE [dbo].[FunctionalRejects] (
3      [RejectDate] datetime,
4      [RejectPackageAndTask] nvarchar(205),
5      [RejectColumn] nvarchar(29),
6      [RejectDescription] nvarchar(131),
7      [Count] int
8  )

```

These functional rejects are also fed back into the dataset with the use of a multi. This is done for potential analysis later on translation and sales. After all this is done the data is loaded into an ODS products table as shown in the steps below.

OLE DB connection manager:

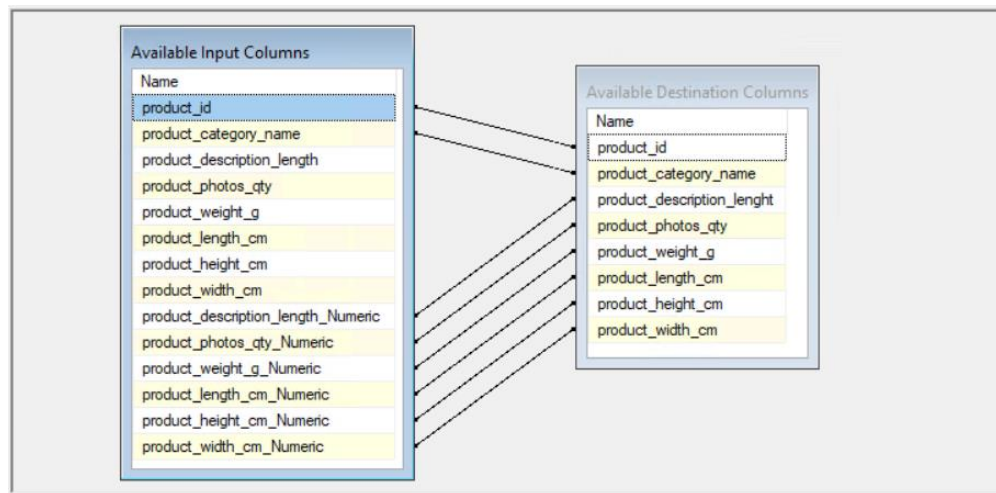
localhost.Project\_ODS ▼ New...

Data access mode:

Table or view - fast load ▼

Name of the table or the view:

[dbo].[Products] ▼ New...

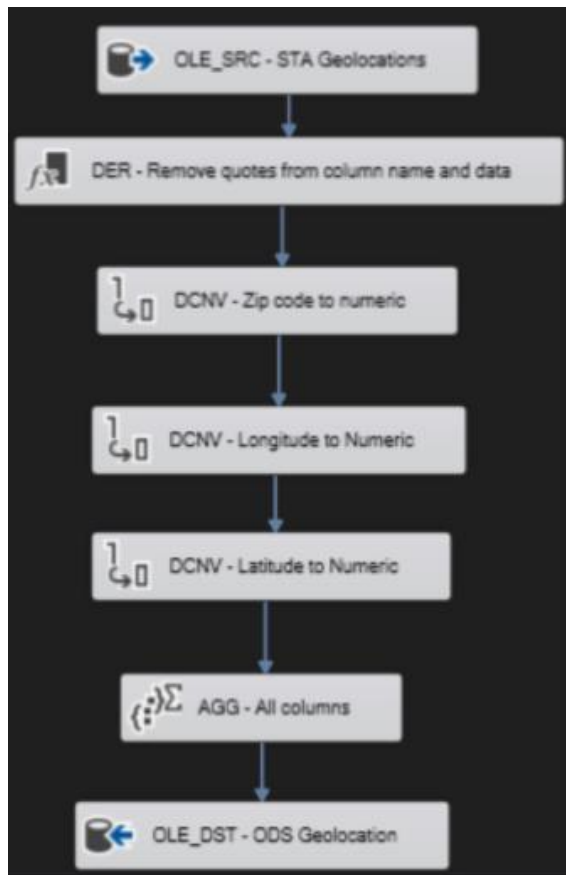


Input Column	Destination Column
product_id	product_id
product_category_name	product_category_name
product_description_length_Numeric	product_description_lenght
product_photos_qty_Numeric	product_photos_qty
product_weight_g_Numeric	product_weight_g
product_length_cm_Numeric	product_length_cm
product_height_cm_Numeric	product_height_cm
product_width_cm_Numeric	product_width_cm

	product_id	product_category_name	product_description_lenght	product_photos_qty	product_weight_g	product_length_cm	product_height_cm	product_width_cm
1	b5694f734bee2683d25bb188a745b40	stationery	684	2	100	16	2	11
2	41b24e4f49504cc023666ed7cde3cf16	sports_leisure	628	7	1800	25	30	40
3	e987c6af8290cf1db0510d52b65795	health_beauty	99	2	250	20	10	15
4	6e2d41323c4b6d240b530b9e68e83de1	watches_gifts	626	1	400	23	15	19
5	e76572c2ae4cede1e3b2ead512423efe	bed_bath_table	352	2	450	45	11	11
6	e97ee7673459379fac14fe9a01a64e6c	toys	515	1	289	17	15	17
7	d212ad4ac52091c8d50442a89fbbc7af	bed_bath_table	160	1	350	16	16	16
8	99f83ce7f14247f186cb52d73bdd6e79	construction_tools_construction	453	1	100	16	7	12
9	9539cddd99ab3f8ca77b6cad0a0899ea	sports_leisure	1067	4	433	23	12	17
10	9fae04deab3bb1c353c27c50cbf3b42b	toys	927	4	450	21	36	11

### c) Geolocation Table

This package starts with a truncate the same as the previous packages. The DFT is shown below:



In the first step, we extract the data from STA Geolocation and remove the quotes from the column names the same as previous packages. After this there are several columns that contain numerical data. These are changed the same way as the previous package without errors. Once this is completed, a group by is used to aggregate all the data and remove duplicates.

Available Input Columns		
<input checked="" type="checkbox"/>	Name	
<input type="checkbox"/>	(*)	
<input checked="" type="checkbox"/>	geolocation_zip_code_prefix	
<input checked="" type="checkbox"/>	geolocation_lat	
<input checked="" type="checkbox"/>	geolocation_lng	
<input checked="" type="checkbox"/>	geolocation_city	
<input checked="" type="checkbox"/>	geolocation_state	
<input checked="" type="checkbox"/>	geolocation_zip_code_prefix_Numeric	
<input checked="" type="checkbox"/>	geolocation_lng_Numeric	
<input checked="" type="checkbox"/>	geolocation_lat_Numeric	

Input Column	Output Alias	Operation
geolocation_zip_code_prefix_Numeric	geolocation_zip_code_prefix_Numeric	Group by
geolocation_lat	geolocation_lat	Group by
geolocation_lng	geolocation_lng	Group by
geolocation_city	geolocation_city	Group by
geolocation_lng_Numeric	geolocation_lng_Numeric	Group by
geolocation_lat_Numeric	geolocation_lat_Numeric	Group by
geolocation_state	geolocation_state	Group by
geolocation_zip_code_prefix	geolocation_zip_code_prefix	Group by

After all is completed, the data is loaded into an ODS geolocation table in the same steps as shown previously. The script for the table is shown below.

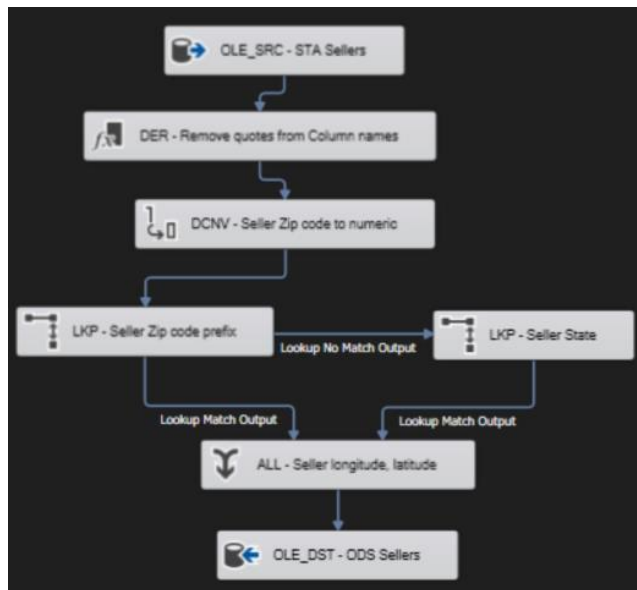
```

1      USE [Project_ODS]
2      CREATE TABLE dbo.Geolocation (
3          [geolocation_zip_code] numeric(18,0),
4          [geolocation_lng] decimal(18,10),
5          [geolocation_lat] decimal(18,10),
6          [geolocation_city] varchar(50),
7          [geolocation_state] varchar(50),
8      )

```

#### d) Sellers Table

This package starts with a truncate the same as the previous packages. The DFT is shown below:



In the first step, we extract the data from STA Geolocation and remove the quotes from the column names the same as previous packages. After this the zip code column is changed to numerical data using a data conversion as shown in previous packages. After this a lookup is performed on the seller zip code column using the geolocation dataset that was loaded into ODS. The data that had no matching geolocation data for zip code are then matched with state to get a geolocation.

Available Input Columns		Available Lookup Columns	
Name		Name	Index
seller_id		<input type="checkbox"/> geolocation_zip_code	
seller_zip_code_prefix		<input checked="" type="checkbox"/> geolocation_lng	
seller_city		<input checked="" type="checkbox"/> geolocation_lat	
seller_state		<input type="checkbox"/> geolocation_city	
seller_zip_code_prefix_Numeric		<input type="checkbox"/> geolocation_state	

Lookup Column	Lookup Operation	Output Alias
geolocation_lng	<add as new column>	geolocation_lng
geolocation_lat	<add as new column>	geolocation_lat



Available Input Columns

Name
seller_id
seller_zip_code_prefix
seller_city
seller_state
seller_zip_code_prefix_Numeric

Available Lookup Columns

<input checked="" type="checkbox"/>	Name	Index
<input type="checkbox"/>	geolocation_zip_code	
<input checked="" type="checkbox"/>	geolocation_lng	
<input checked="" type="checkbox"/>	geolocation_lat	
<input type="checkbox"/>	geolocation_city	
<input type="checkbox"/>	geolocation_state	

Lookup Column	Lookup Operation	Output Alias
geolocation_lng	<add as new column>	geolocation_lng
geolocation_lat	<add as new column>	geolocation_lat

Output Column Name	Union All Input 1	Union All Input 2
seller_id	seller_id	seller_id
seller_zip_code_prefix	seller_zip_code_prefix	seller_zip_code_prefix
seller_city	seller_city	seller_city
seller_state	seller_state	seller_state
seller_zip_code_prefix_Numeric	seller_zip_code_prefix_Numeric	seller_zip_code_prefix_Numeric
geolocation_lng	geolocation_lng	geolocation_lng
geolocation_lat	geolocation_lat	geolocation_lat

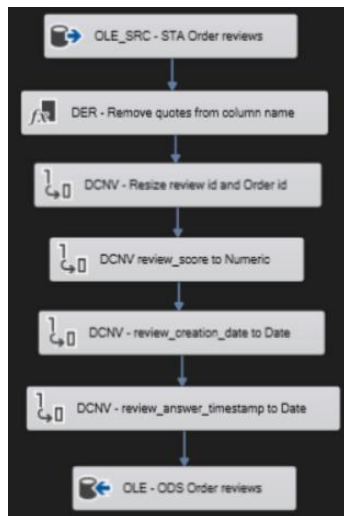
At the end of this process, all the sellers will have geolocation data attached to them. The reason for this is to have the potential to evaluate the distance between sellers and customers in analysis. This also means you can create maps with your data during analysis. After all this the data is loaded into a table in ODS the same way as the previous packages.

### e) Customers Table

The customers table follows the same steps as the sellers table above with its own respective data. Zip code is changed to numeric, geolocation data is added for the customers using a lookup from the geolocation dataset, and the data is loaded into a table in ODS.

### f) Order Reviews Table

This package starts with a truncate the same as the previous packages. The DFT is shown below:



In the first step, we extract the data from STA Customer\_reviews and remove the quotes from the column names the same as previous packages. After this, data conversion is done on various columns in the same steps as done in previous packages. Review ID and order ID have their length changed to 50, review score is changed to numeric, and review creation date and review answer timestamp is changed to date. After all this is done the data is loaded into a table in ODS the same way as previous packages.

Order payments and order items follow the same steps as order reviews above. Data is loaded from STA, quotes are removed from the column names, datatypes are changed on the columns required, and the data is loaded into ODS.

### 3. Data Loading

The last step in the data pipeline is to load or integrate the data into the Data Warehouse (DWH). One common database schema that is used is the Star schema. In this schema, we have one main table called the “Fact Table” That is surrounded by “Dimension Tables”. The fact table contains the most important data called “facts”, whereas the dimensions tables give addition descriptive information. We will design our database around this schema. Our design will have orders as our fact table with everything else being a dimension table. Order items, order reviews, and order payments are part of our orders fact table however due to limitations with the data, we decided to leave them as separate tables. They effectively work as 1 large fact table broken up into different sections.

A common dimension table that we added into the data is the “Date” dimension. It allows you to describe the dates with multiple temporal aggregate categories (year, month, quarter). The

table will contain multiple variations of those three data to be able to adapt to various styles of queries. We define the relation to the fact table with a technical key with the format “YYYYMM”.

## a) Integration of the Date Dimension

This dimension is describing dates, in particular months. Since we don’t need external data to build this dimension, we use a SQL script to make it.

First, we create an empty

table:

```

1  USE [Project_DWH]
2  CREATE TABLE [dbo].[DimDate] (
3      DateKey INT NOT NULL PRIMARY KEY,
4      [Date] DATE NULL,
5      [Day] TINYINT NOT NULL,
6      [DaySuffix] CHAR(2) NOT NULL,
7      [Weekday] TINYINT NOT NULL,
8      [WeekDayName] VARCHAR(10) NOT NULL,
9      [WeekDayName_Short] CHAR(3) NOT NULL,
10     [WeekDayName_FirstLetter] CHAR(1) NOT NULL,
11     [DOWInMonth] TINYINT NOT NULL,
12     [DayOfYear] SMALLINT NOT NULL,
13     [WeekOfMonth] TINYINT NOT NULL,
14     [WeekOfYear] TINYINT NOT NULL,
15     [Month] TINYINT NOT NULL,
16     [MonthName] VARCHAR(10) NOT NULL,
17     [MonthName_Short] CHAR(3) NOT NULL,
18     [MonthName_FirstLetter] CHAR(1) NOT NULL,
19     [Quarter] TINYINT NOT NULL,
20     [QuarterName] VARCHAR(6) NOT NULL,
21     [Year] INT NOT NULL,
22     [MYYYYY] CHAR(6) NOT NULL,
23     [MonthYear] CHAR(7) NOT NULL,
24     IsWeekend BIT NOT NULL,
25 )
26
27
28  GO

```

Then to build the data, we use the following script:

```

31  SET NOCOUNT ON
32
33  TRUNCATE TABLE DimDate
34
35  DECLARE @CurrentDate DATE = '2016-01-01'
36  DECLARE @EndDate DATE = '2021-12-31'
37
38
39  --Insertion of a row to get a DateKey for the NULL Dates in Orders (Not delivered yet orders)
40  INSERT INTO [dbo].[DimDate] (
41      [DateKey],
42      [Date],
43      [Day],
44      [DaySuffix],
45      [Weekday],
46      [WeekDayName],
47      [WeekDayName_Short],
48      [WeekDayName_FirstLetter],
49      [DOWInMonth],
50      [DayOfYear],
51      [WeekOfMonth],
52      [WeekOfYear],
53      [Month],
54      [MonthName],
55      [MonthName_Short],
56      [MonthName_FirstLetter],
57      [Quarter],
58      [QuarterName],
59      [Year],
60      [MYYYYY],
61      [MonthYear],
62      [IsWeekend]
63
64      VALUES (
65          -1,
66          NULL,
67          0,
68          '',
69          0,
70          '',
71          ,
72          '',
73          '',
74          0,
75          0,
76          0,
77          0,
78          0,
79          '',
80          '',
81          '',
82          0,
83          '',
84          0,
85          '',
86          '',
87          0
88      )
89
90  WHILE @CurrentDate < @EndDate
91  BEGIN

```

```

91 BEGIN
92 INSERT INTO [dbo].[DimDate] (
93     [DateKey],
94     [Date],
95     [Day],
96     [DaySuffix],
97     [Weekday],
98     [WeekDayName],
99     [WeekDayName_Short],
100    [WeekDayName_FirstLetter],
101    [DOWInMonth],
102    [DayOfYear],
103    [WeekOfMonth],
104    [WeekOfYear],
105    [Month],
106    [MonthName],
107    [MonthName_Short],
108    [MonthName_FirstLetter],
109    [Quarter],
110    [QuarterName],
111    [Year],
112    [MMYYYY],
113    [MonthYear],
114    [IsWeekend]
115 )
116 SELECT DateKey = YEAR(@CurrentDate) * 10000 + MONTH(@CurrentDate) * 100 + DAY(@CurrentDate),
117        DATE = @CurrentDate,
118        Day = DAY(@CurrentDate),
119        [DaySuffix] = CASE
120            WHEN DAY(@CurrentDate) = 1
121            OR DAY(@CurrentDate) = 21
122            OR DAY(@CurrentDate) = 31
123            THEN 'st'
124            WHEN DAY(@CurrentDate) = 2
125            OR DAY(@CurrentDate) = 22
126            THEN 'nd'
127            WHEN DAY(@CurrentDate) = 3
128            OR DAY(@CurrentDate) = 23
129            THEN 'rd'
130            ELSE 'th'
131        END,
132        WEEKDAY = DATEPART(dw, @CurrentDate),
133        WeekDayName = DATENAME(dw, @CurrentDate),
134        WeekDayName_Short = UPPER(LEFT(DATENAME(dw, @CurrentDate), 3)),
135        WeekDayName_FirstLetter = LEFT(DATENAME(dw, @CurrentDate), 1),
136        [DOWInMonth] = DAY(@CurrentDate),
137        [DayOfYear] = DATENAME(dy, @CurrentDate),
138        [WeekOfMonth] = DATEPART(WEEK, @CurrentDate) - DATEPART(WEEK, DATEADD(WH, 0, @CurrentDate), 0) + 1,
139        [WeekOfYear] = DATEPART(WK, @CurrentDate),
140        [Month] = MONTH(@CurrentDate),
141        [MonthName] = DATENAME(mm, @CurrentDate),
142        [MonthName_Short] = UPPER(LEFT(DATENAME(mm, @CurrentDate), 3)),
143        [MonthName_FirstLetter] = LEFT(DATENAME(mm, @CurrentDate), 1),
144        [Quarter] = DATEPART(q, @CurrentDate),
145        [QuarterName] = CASE
146            WHEN DATENAME(qq, @CurrentDate) = 1
147            THEN 'First'
148            WHEN DATENAME(qq, @CurrentDate) = 2
149            THEN 'Second'
150            WHEN DATENAME(qq, @CurrentDate) = 3
151            THEN 'Third'
152            WHEN DATENAME(qq, @CurrentDate) = 4
153            THEN 'Fourth'
154            ELSE 'Fourth'
155        END,
156        [Year] = YEAR(@CurrentDate),
157        [MMYYYY] = RIGHT('0' + CAST(MONTH(@CurrentDate) AS VARCHAR(2)), 2) + CAST(YEAR(@CurrentDate) AS VARCHAR(4)),
158        [MonthYear] = CAST(YEAR(@CurrentDate) AS VARCHAR(4)) + UPPER(LEFT(DATENAME(mm, @CurrentDate), 3)),
159        [IsWeekend] = CASE
160            WHEN DATENAME(dw, @CurrentDate) = 'Sunday'
161            OR DATENAME(dw, @CurrentDate) = 'Saturday'
162            THEN 1
163            ELSE 0
164        END
165 SET @CurrentDate = DATEADD(DD, 1, @CurrentDate)
166 END
167
168

```

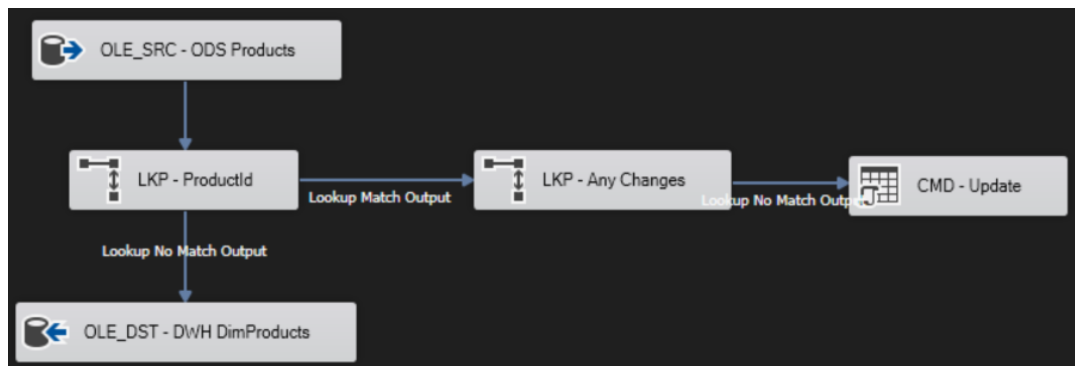
Here is the first ten lines of the results:

	DateKey	Date	Day	DaySuffix	Weekday	WeekDayName	WeekDayName_Short	WeekDayName_FirstLetter	DOWInMonth	DayOfYear	WeekOfMonth	WeekOfYear	Month	MonthName	MonthName_Short
1	-1	NULL	0		0				0	0	0	0	0		
2	20160101	2016-01-01	1	st	6	Friday	FRI	F	1	1	1	1	1	January	JAN
3	20160102	2016-01-02	2	nd	7	Saturday	SAT	S	2	2	1	1	1	January	JAN
4	20160103	2016-01-03	3	rd	1	Sunday	SUN	S	3	3	2	2	1	January	JAN
5	20160104	2016-01-04	4	th	2	Monday	MON	M	4	4	2	2	1	January	JAN
6	20160105	2016-01-05	5	th	3	Tuesday	TUE	T	5	5	2	2	1	January	JAN
7	20160106	2016-01-06	6	th	4	Wednesday	WED	W	6	6	2	2	1	January	JAN
8	20160107	2016-01-07	7	th	5	Thursday	THU	T	7	7	2	2	1	January	JAN
9	20160108	2016-01-08	8	th	6	Friday	FRI	F	8	8	2	2	1	January	JAN
10	20160109	2016-01-09	9	th	7	Saturday	SAT	S	9	9	2	2	1	January	JAN

To accept new data in the data warehouse, this dimension can be updated. We would need to change the “@EndDate” variable to include all the dates required.

## b) Integration of the Products Dimension

We load this data with the process below:



We need to make sure that we can make joins between the fact table and the dimension table, so we check the link with "ProductID".

OLE DB connection manager:

localhost.Project\_DWH

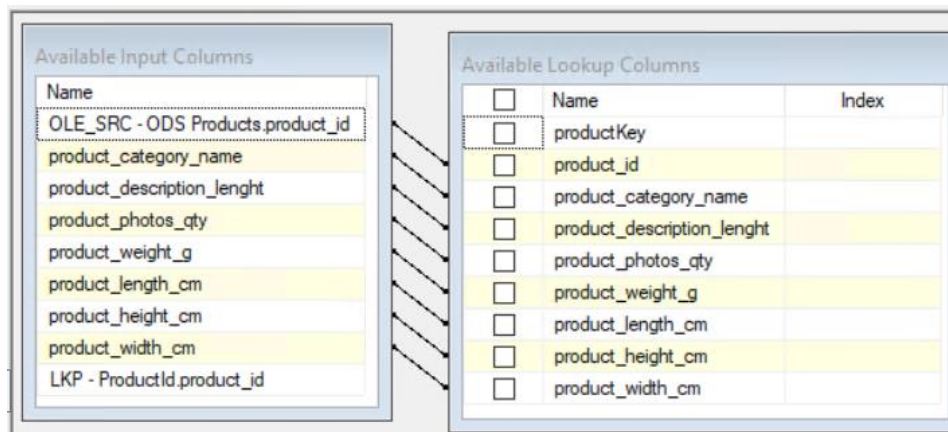
☒ Use a table or a view:

[dbo].[DimProducts]

Available Input Columns	
Name	
product_id	<input checked="" type="checkbox"/>
product_category_name	<input type="checkbox"/>
product_description_lenght	<input type="checkbox"/>
product_photos_qty	<input type="checkbox"/>
product_weight_g	<input type="checkbox"/>
product_length_cm	<input type="checkbox"/>
product_height_cm	<input type="checkbox"/>
product_width_cm	<input type="checkbox"/>

Available Lookup Columns		
	Name	Index
<input type="checkbox"/>	productKey	
<input checked="" type="checkbox"/>	product_id	
<input type="checkbox"/>	product_category_name	
<input type="checkbox"/>	product_description_lenght	
<input type="checkbox"/>	product_photos_qty	
<input type="checkbox"/>	product_weight_g	
<input type="checkbox"/>	product_length_cm	
<input type="checkbox"/>	product_height_cm	
<input type="checkbox"/>	product_width_cm	

For the dimensions tables, we also need to have an update policy. We choose the SCD1 strategy, implemented by checking is there is any change and updating the table if necessary.



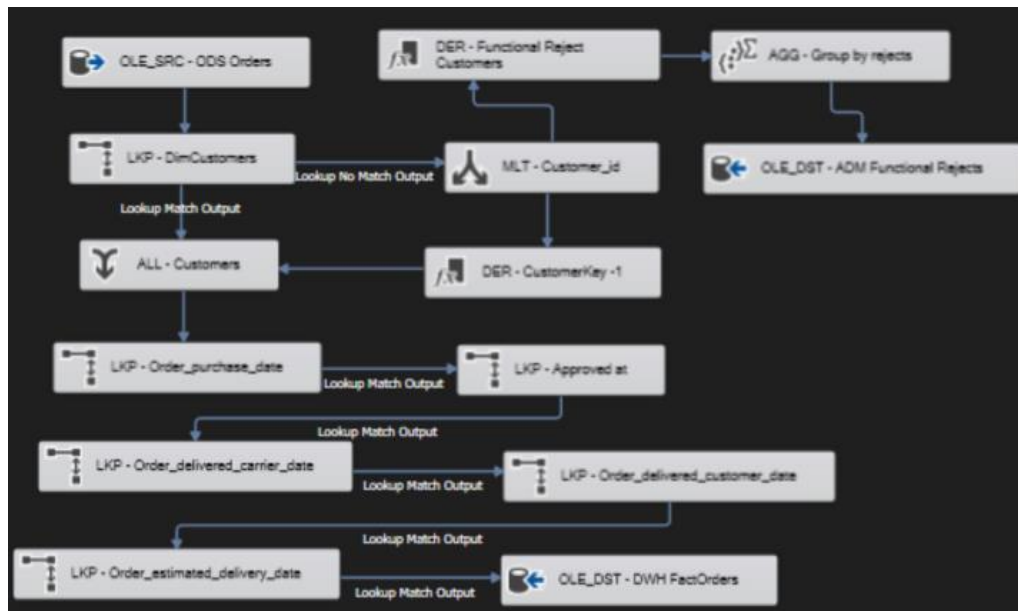
Here is the first ten lines of the results:

	productKey	product_id	product_category_name	product_description_lenght	product_photos_qty	product_weight_g	product_length_cm	product_height_cm	product_width_cm
1	1	b5694f734bee2683d25bb188a745b40	stationery	684	2	100	16	2	11
2	2	41b24e4f49504cc023666ed7cde3cf16	sports_leisure	628	7	1800	25	30	40
3	3	e987fc6af8290cf1db0510d52b65795	health_beauty	99	2	250	20	10	15
4	4	6e2d41323c4b6d240b530b9e68e83de1	watches_gifts	626	1	400	23	15	19
5	5	e76572c2ae4cede1e3b2ead512423efe	bed_bath_table	352	2	450	45	11	11
6	6	e97ee7673459375fac14fe9a01a64e6c	toys	515	1	289	17	15	17
7	7	d212ad4ac52091c8d50442a89fbbc7af	bed_bath_table	160	1	350	16	16	16
8	8	99f83ce7f14247f186cb52d73bdd6e79	construction_tools_construction	453	1	100	16	7	12
9	9	9539cdd99ab3f8ca77b6cad0a0899ea	sports_leisure	1067	4	433	23	12	17
10	10	9fae04deab3bb1c353c27c50cbf3b42b	toys	927	4	450	21	36	11

The other Dimension tables follow the same structure as the one above.

### c) Integration of the orders fact table

Now that we finally have our dimensions tables, we can build our fact table while checking valid relations with the dimensions. We verify and load the data as follows:



In case of an error in the processing, we generate a functional reject and insert it into a table “Functional\_Rejects”.

The first dimension we check is “DimCustomers”. We check it with the key between “customer\_id” and “customer\_id”. We also add a technical key “customerKey\_FK”.

Available Input Columns		Available Lookup Columns	
Name		<input type="checkbox"/> Name	Index
order_id		<input checked="" type="checkbox"/> customerKey	
customer_id		<input type="checkbox"/> customer_id	
order_status		<input type="checkbox"/> customer_city	
order_purchase_timestamp		<input type="checkbox"/> customer_state	
order_approved_at		<input type="checkbox"/> customer_zip_code_prefix	
order_delivered_carrier_date		<input type="checkbox"/> geolocation_lng	
order_delivered_customer_date		<input type="checkbox"/> geolocation_lat	
order_estimated_delivery_date			

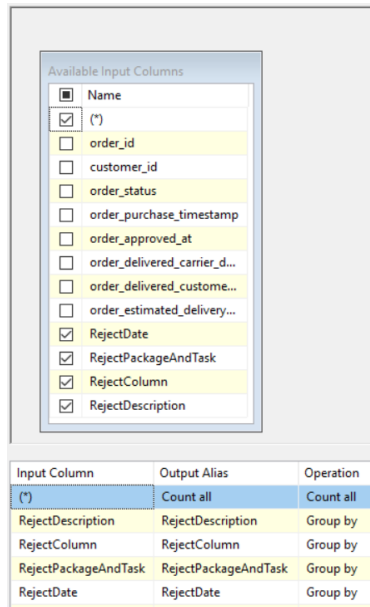
  

Lookup Column	Lookup Operation	Output Alias
customerKey	<add as new column>	customerKey

We tack the errors with the same method that in ODS. For a missing relation with “DimCustomers”, we generate the following data:

Derived Column Name	Derived Column	Expression	Data Type	Length
RejectDate	<add as new column>	GETDATE()	database timestamp [DT_DBTIMESTAMP]	
RejectPackageAndTask	<add as new column>	(DT_WSTR,100)@[System::PackageName] + " AND " + (DT_WSTR,100)@[System::TaskName]	Unicode string [DT_WSTR]	205
RejectColumn	<add as new column>	"customer_id"	Unicode string [DT_WSTR]	11
RejectDescription	<add as new column>	"This customer_id" + (DT_WSTR,50)customer_id + " doesn't exist in our DimCustomer table"	Unicode string [DT_WSTR]	106

We then group by the rejects:



The screenshot shows a software interface for configuring a data transformation. It features a list of 'Available Input Columns' on the left, including Name, (\*), order\_id, customer\_id, order\_status, order\_purchase\_timestamp, order\_approved\_at, order\_delivered\_carrier\_d..., order\_delivered\_custome..., order\_estimated\_delivery..., RejectDate, RejectPackageAndTask, RejectColumn, and RejectDescription. The 'Reject' columns are checked. Below this is a table mapping input columns to output aliases and operations.

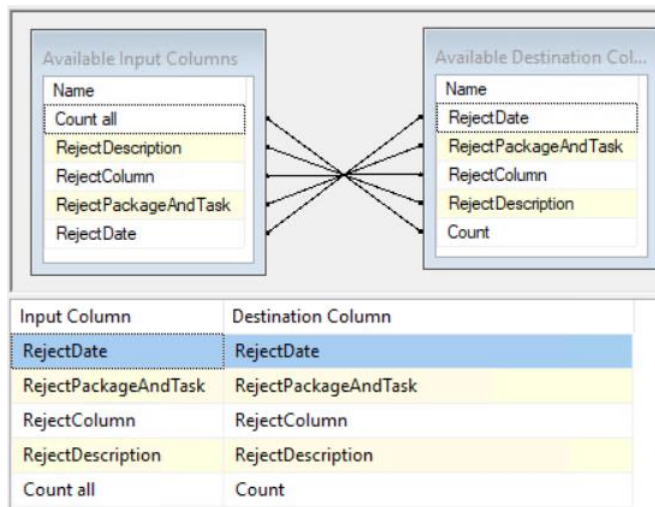
Input Column	Output Alias	Operation
(*)	Count all	Count all
RejectDescription	RejectDescription	Group by
RejectColumn	RejectColumn	Group by
RejectPackageAndTask	RejectPackageAndTask	Group by
RejectDate	RejectDate	Group by

In this case, we will also integrate the record with the default value (-1) for the Customer:

Derived Column Name	Derived Column	Expression	Data Type
customerKey	<add as new column>	-1	four-byte signed integer [DT_I4]

As shown here, we are keeping track of the rejects as they are sent into functional rejects:





Next for all the columns containing dates, a lookup is performed using the Dimdate table.

OLE DB connection manager:

localhost.Project\_DWH

☒ Use a table or a view:

[dbo].[DimDate]

Lookup Column	Lookup Operation	Output Alias
DateKey	<add as new column>	DateKey

Finally we can create and load the fact table into the data warehouse.

```

1  USE [Project_DWH]
2  CREATE TABLE [dbo].[FactOrders] (
3      [order_id] varchar(50),
4      [customer_id] varchar(50),
5      [order_status] varchar(50),
6      [customerKey] int,
7      [order_purchase_DateKey] int,
8      [approved at_DateKey] int,
9      [order_delivered_carrier_DateKey] int,
10     [order_delivered_customer_DateKey] int,
11     [order_estimated_delivery_DateKey] int
12 )

```

OLE DB connection manager:

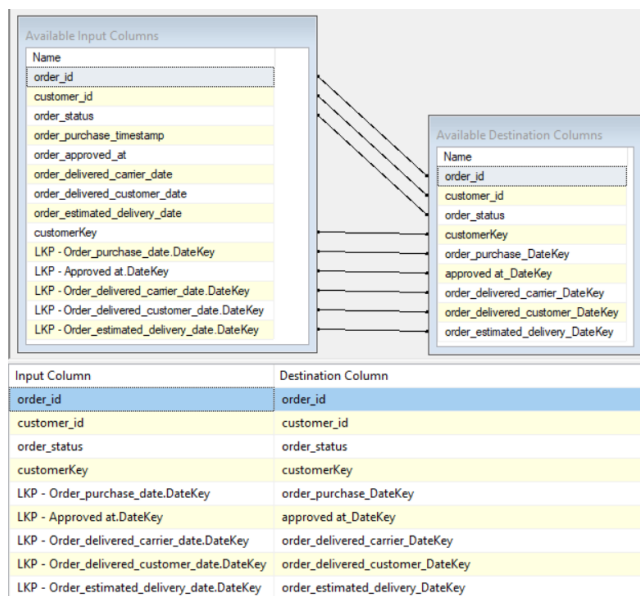
localhost.Project\_DWH

Data access mode:

Table or view - fast load

Name of the table or the view:

[dbo].[FactOrders]



Here is the first ten lines of the results of the facts table “orders”:

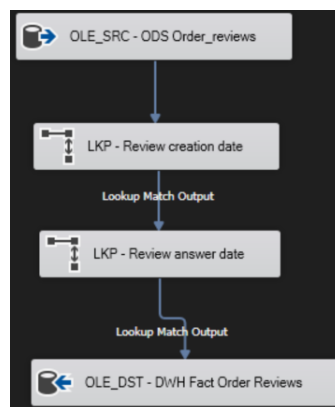
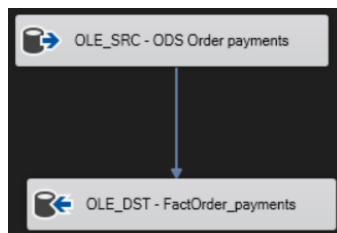
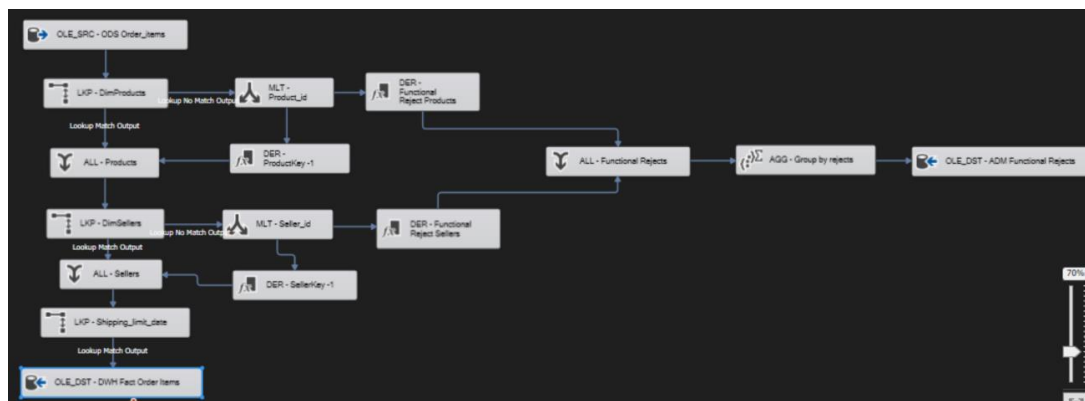
	order_id	customer_id	order_status	customerKey	order_purchase_DateKey	approved at_DateKey	order_delivered_carrier_DateKey	order_delivered_customer_DateKey
1	fc1d9429da446a147d438d41350a2932	8b13d2050940cb8d6a7203e69f279254	delivered	3566	20171023	20171024	20170611	20171118
2	8a6bf1063cc2efc062e037e474622de6	64f0ed97c13c1389088e93454db40488	delivered	16866	20170320	20170320	20170325	20170104
3	bcd37480e9f21eaf2b524c9e9724d099	fbcb839b74c648a39b750799b58a13e0f	delivered	47155	20171124	20171125	20171128	20170712
4	b951cb656256cfb404cd77c81262f25	83116a97d503c15ca0e75dc2b848bd9	delivered	64538	20171226	20171226	20171227	20180401
5	b829b53d8136d85eb4cbe60a1bd279f5	ec4f3f98da53bbaa4db16509b3782c0	delivered	79211	20170509	20170609	20170609	20171209
6	b1e68c1ac52a1a5cb5051d07ae91ff06	cf34464d3b5d4a957dae3c0ef4a7d27e	delivered	32037	20170318	20170318	20170321	20170327
7	4516299f280a0f30720256ae0cd592d6	582ec29229bb3014866b31f73def34f	delivered	28312	20170913	20170913	20170914	20171018
8	7d63cb1d349e7bd0ee1edcc61ea71077	484711e3d14ac75fb874e949b3a28395	delivered	72746	20170814	20170814	20170816	20170824
9	b11073519dc9fe626df6281aa1965522	4f3ef3bae0d531a175d66f875cba84d1	delivered	22825	20171116	20171116	20171121	20171127
10	c7b440776d014b67c3164e103e3bf02	19187bc9fab7457748d4dc42a9e1cffe	delivered	37298	20171121	20171121	20171123	20171130

The rest of the fact tables follow the same structure as above. This structure in summary is:

For the dimension(s) it contains:

- Lookup with the Dimension primary key and grab the surrogate key. In case of a no match, redirect the row to functional error and use a multicast to, at the same time follow the row in the main stream but the surrogate key is -1
- Lookup with Dimdate for each date it contains.
- In what is inserted in the DWH Fact table, the primary keys of the dimensional columns are replaced with their surrogate key. The dates are replaced with datekey.

Their dataflows are shown below:



# III. Project Deployment

We have defined all required steps to deploy our data warehouse with the available data. To ensure that the pipeline is executed in a reproducible manner, we define an additional package “Project\_ETL” with containers of all the tasks.



Using this package allows us to use one interface for the entire pipeline. The containers make sure everything is launched in the correct order.

# VI. Analysis

In this part, we decided to choose 4 business questions:

1. Which product categories have the highest customer satisfaction?
2. Which month do customers order the most?
3. How many orders have not been delivered?
4. In which of the company's cities do we have the biggest profits?
- 5.

To answer each question, we built a specific query.

## 1. Which product categories have the highest customer satisfaction?

It's always important to know which products are popular with customers. This data can be used by the marketing or business team to make decisions:

### a) Query

```
SELECT TOP(5) prod.product_category_name, rev.review_score, COUNT(rev.review_score) as  
number_of_reviews  
FROM  
DimProducts as prod  
JOIN FactOrder_items as ord ON prod.productKey = ord.productKey  
JOIN FactOrder_reviews as rev on ord.order_id = rev.order_id  
WHERE rev.review_score = 5  
GROUP BY prod.product_category_name, rev.review_score  
ORDER BY number_of_reviews DESC
```

### b) Result

	product_category_name	review_score	number_of_reviews
1	health_beauty	5	11716
2	bed_bath_table	5	11570
3	sports_leisure	5	10242
4	furniture_decor	5	8904
5	computers_accessories	5	8400

As far as our results are concerned, we can see that the health and beauty category is the best-selling, which seems logical. But the second category seems more surprising and unclear. This category may be linked to a famous company that sells products for the home.

## 2. Which month do customers order the most?

Knowing the months in which customers buy the most enables the company to make strategic and financial decisions. For example, to advertise or restock certain products.

### a) Query

```
SELECT COUNT(order_id) as nb_orders, MonthName
FROM FactOrders as ord
JOIN DimDate as date_ ON ord.[approved at _DateKey] = date_.DateKey
GROUP BY MonthName
ORDER BY COUNT(order_id) DESC
```

### b) Result

	nb_orders	MonthName
1	10450	July
2	9857	May
3	9239	March
4	9073	August
5	8992	June
6	8777	November
7	8572	April
8	8388	January
9	8009	February
10	6140	October
11	5903	December
12	5881	September

Surprisingly, July is the month when customers buy the most. We can see that the top 5 months are in summer or spring.

## 3. How many orders have not been delivered?

One of the most important parts of a marketplace is delivery. This means that the order has been completed. When the order is not fulfilled, for whatever reason, it means that the marketplace loses a sale. So it's important to know how many orders are not delivered.

### a) Query

```
SELECT
COUNT(order_id) AS total_orders,
COUNT(CASE WHEN order_status != 'delivered' THEN 1 END) AS not_delivered_orders
FROM
FactOrders;
```

*\* this query was created with the help of ChatGPT. After 1 hour of trying to figure out why the query wasn't working, we asked ChatGPT to help us understand our mistake: we hadn't put the CASE WHEN in the right place.*

### b) Result

	total_orders	not_delivered_orders
1	99281	2817

The number of orders approved but not delivered represents around 3%. Depending on the company's wishes, this figure may seem high.

## 4. In which of the company's cities do we have the biggest profits?

It will be interesting to find out which city sellers make the most profit. We could find correlations between the city and profits.

### a) Query

```
SELECT seller.seller_city, SUM(ord_items.price) as benef
FROM DimSellers as seller
JOIN FactOrder_items as ord_items ON seller.sellerKey = ord_items.sellerKey
GROUP BY seller_city
ORDER BY benef DESC FactOrders;
```

### b) Result

	seller_city	benef
1	sao paulo	5405756.28
2	ibitinga	1249185.88
3	curitiba	941519.64
4	rio de janeiro	716827.18
5	guarulhos	658988.76
6	ribeirao preto	551952.88
7	itaquaquecetuba	461136.24
8	guariba	458945.26
9	santo andre	457123.20
10	lauro de freitas	451050.10
11	piracicaba	425355.64
12	belo horizonte	415342.46

Looking at the data, we can see that the sellers are located in Brazil and that the first city, Sao Paulo, has a significant lead over the other cities. We can assume that there are more companies located in Sao Paulo. It will be interesting to identify this number.

## 5. Are top-rated products with a low average distance between customer and seller?

The result of this query can help us understand whether the distance between the customer and the salesperson can have an impact on the price of the product. Distance is important because it can impact the delivery delay or the condition of the parcel on arrival.

To calculate the distance we use this formula:

### Distance

This uses the '**haversine**' formula to calculate the great-circle distance between two points – that is, the shortest distance over the earth's surface – giving an 'as-the-crow-flies' distance between the points (ignoring any hills they fly over, of course!).

*Haversine*  $a = \sin^2(\Delta\phi/2) + \cos \phi_1 \cdot \cos \phi_2 \cdot \sin^2(\Delta\lambda/2)$

*formula:*  $c = 2 \cdot \text{atan2}(\sqrt{a}, \sqrt{1-a})$

$d = R \cdot c$

*where:  $\phi$  is latitude,  $\lambda$  is longitude, R is earth's radius (mean radius = 6,371km);*

*note that angles need to be in radians to pass to trig functions!*

### a) Query

WITH lat\_lng\_radians AS

(

SELECT RADIANS(sellers.geolocation\_lat) as seller\_lat, RADIANS(sellers.geolocation\_lng)  
as seller\_lng,



```

        RADIANS(customers.geolocation_lat) as customer_lat,
        RADIANS(customers.geolocation_lng) as customer_lng, order_reviews.review_score,
        customers.customer_id

    FROM dbo.FactOrders as orders

        JOIN dbo.DimCustomers as customers ON orders.customerKey =
        customers.customerKey

        JOIN dbo.FactOrder_items as order_items ON order_items.order_id =
        orders.order_id

        JOIN dbo.DimSellers as sellers ON sellers.sellerKey = order_items.sellerKey

        JOIN dbo.FactOrder_reviews as order_reviews ON order_reviews.order_id =
        orders.order_id

    WHERE order_status = 'delivered'

)

SELECT AVG(

    ROUND(2 * 6371 * ASIN(

        SQRT(

            POWER(SIN((customer_lat - seller_lat)/2), 2) + COS(seller_lat) *
            COS(customer_lat) * POWER(SIN((customer_lng - seller_lng)/2), 2)

        )

    ), 1)

    ) as mean_distance_km, review_score

FROM lat_lng_radians

GROUP BY review_score

ORDER BY review_score DESC

```

## b) Result

	mean_distance_km	review_score
1	569.374808858838	5
2	620.466844515795	4
3	629.699567146411	3
4	645.225628548257	2
5	659.915650099402	1

We can see that the highest score, 5, is linked to the short distance between customer and salesperson. Using other variables, we can deduce that one of the most important factors impacting evaluation is delivery time.