# ETL Project With Talend

## Introduction

ETL stands for extract, transform, and load, is the process data engineers use to extract data from different sources, transform the data into a  usable and trusted resource, and load that data into the systems end-users can access and use downstream to solve business problems.
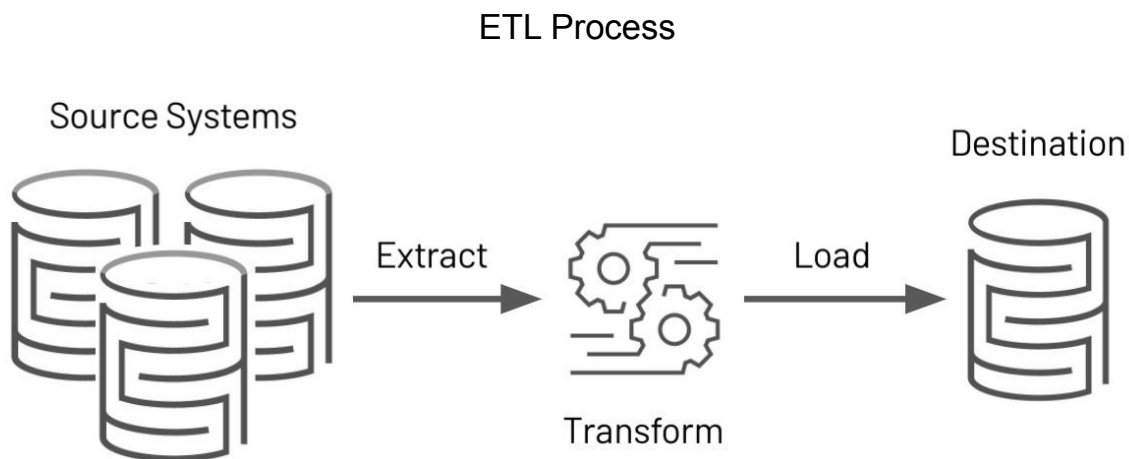
ETL Process



Image credits to: https://www.databricks.com/glossary/extract-transform-load

**Extract:**
-   Reads data from multiple data sources and extracts required set of data.
-   Recovers necessary data with optimum usage of resources.

**Transform:**
-   Filtration, cleansing, and preparation of data extracted, with lookup tables.
-   Authentication of records, refutation, and integration of data.
-   Data to be sorted, filtered, cleared, standardized, translated, or verified for consistency.

**Load:**
-   Writing data output, after transformation to a data warehouse
-   Either physical insertion of record as a new row in database table or link processes for each record from the main source

# Project Overview

In this ETL project, I extracted the data from 4 different data sources and did data transformation then uploaded it to AWS S3 storage in order to move it to Redshift later.

| Data Sources | Staging Area | Final Destination |
|---|---|---|
| 1- CSV files<br>2- Excel sheets<br>3- BigQuery<br>4- PostgreSQL database | Local PC | AWS S3<br>then AWS Redshift |

Northwind is a database created by Microsoft for training and educational purpose.

Links to the data:
https://github.com/engindemirog/Northwind-Database-Script-for-Postgre-Sql/blob/master/script.sql

https://github.com/neo4j-contrib/northwind-neo4j/tree/master/data

https://docs.google.com/spreadsheets/d/1amQgBgIaUMVEj8gYKbvmlzuoA21ABDiLe0v1orZjjkg/edit#gid=1531710140


Related article this ETL project
https://medium.com/@kmsbmadhan/dimensional-modelling-visualization-of-northwind-database-beaac7fecb20

# Objectives

There are many business drivers that can be driven from Northwind data warehouse as follows:

1. **Sales Reporting** to track sales by customer, employees, products, and suppliers to answer the following questions:

- what is our overall sales number?
- How much have we sold of each product?
- Which products are our best and worst sellers?
- Which of our clients order the most products? What do they order?
- How do our sales look when broken down by region?

2. **Request fulfillment Report** to track the order by how much time it has taken to get delivered to the customer and it can be analyzed to see if it can be improved.

3. **Employee level reporting** to track the performance of the employees and see how it can be improved by either providing rewards to the best performers or giving training to the worst performers or both.

4**. Order distribution & Product inventory analysis** to find orders distributed to customers across the world, track inventory, Order level, and Re-order level of the company for the betterment to answer important questions like

- What are the best-selling products, and do we need to store them more?
- What is the count of products left in the inventory?
- Are we going to run out of any products for delivery?
- What are the products that are going unsold and what can be done as improvement in selling or discontinuing them?
- Can we give discounts on unsold products to get attention in purchase?

Source:
https://medium.com/@kmsbmadhan/dimensional-modelling-visualization-of-northwind-database-beaac7fecb20

# Northwind schema:



# Target DWH schema:

**Solving Some Errors when splitting the data on different data sources**

Solve the problem of column header names when Upload the data to bigquery
https://medium.com/google-cloud/bigquery-create-table-from-google-sheets-causing-incorrect-column-names-string-field-0-134f6ecd3fc8

Import CSV file to PostgreSQL
https://www.neilwithdata.com/copy-permission-denied#:~:text=What%20the%20error%20means,the%20server%2C%20not%20the%20client.

Import data from Excel error solving
https://community.talend.com/s/feed/0D73p000004kIYICAM?language=en_US

**Transformation phase**

- Concatenating first name and last name for employees and customers
- Parsing date from CSV files
- Transforming date to the used format
- Append unique rows from the 4 sources together
- Create date table

**Load the data**

- Upload the new tables to AWS S3 (Simple Storage Service) which is the very popular storage service of Amazon Web Services. It is widely used by customers and Talend provides out-of-the-box connectivity with S3.

# Orders Table

# Products Table

# Employees Table

# Customers Table

# Dim time

## Upload the final tables to AWS S3

Iterate through the files in the data staging folder and upload them to AWS S3 storage in order to move them to Redshift data warehouse.



Link for the method: https://www.youtube.com/watch?v=TqJAd6RQypU&t=4s

# Courses

I have studied these data engineering courses before I made the ETL project.


University course for data warehouse
https://www.youtube.com/playlist?list=PLiJhHdYdI84DzwH47lZQWN6de1tOq8LFu

Big Data Engineering in Depth Course
https://www.youtube.com/playlist?list=PLxNoJq6k39G_m6DYjpz-V92DkaQEiXxkF

Talend course
https://www.youtube.com/playlist?list=PLOr008ImHvfan_fuDr5RVyexpeYJAp9FX

Talend ETL project (the previous project was inspired by it)
https://www.youtube.com/playlist?list=PLKdHo47jRFvf1VzST1RDYiAs8ItMFychE

data engineering podcast
https://www.youtube.com/playlist?list=PLiKvD85qG0l6pLvsQChJO6UBFLfMFfwOw



Thanks for Reading, and I will appreciate your feedback.
Kindly check my Previous projects here: https://linktr.ee/mhmod36



See you next project,
Mahmoud Sallam