# Data Engineer Test

**Canseco García Edgar Jesús**

.                                                                    Julio de 2021

caption float

# What is data Engineering?

We are experiencing a new technological transformation, now more than ever in history, thousands and thousands of data are produced per second, and with the arrival of new ways of creating infrastructure, such as cloud computing, it is no longer crazy to store large amounts of information, now data has become a fundamental asset in the growth and operation of the vast majority of companies, and this trend is only growing with new technologies and ways of finding value and knowledge in data, as in the case of the use of interactive dashboards or deep learning models.

The problem arises that the data comes in different forms and from very different sources, the speed with which it arrives and the memory size that it will occupy when stored also varies, in addition, this data is not ready to be used for a report or ML model, this data must be cleaned, unified, transformed, etc, and finally put them somewhere where they can be available so that they can be used. This is where data engineering comes in. It is the discipline in charge of all this process described, in addition to having responsibilities such as data management, data infrastructure, and data security.

# What are the main responsabilities of a data Engineer?

The main responsibilities of a data engineer could be divided into two branches, technical and communication responsibilities. The technical responsibilities would be focused on the construction of data pipelines for the different challenges that exist today, in addition to security, infrastructure, and data management aspects that must be known for better data quality.

Due to these processes impact various areas in an organization, it is important to maintain clear communication about the various characteristics that these data pipelines must have to improve their quality.

# Explain ETL?

ETL is the acronym for Extract, Transform and Load Data, it is the process of compiling data from a diverse number of sources, its subsequent transformation, and organization to finally be stored in a repository.

In more detail on this process, the first part (E) extracts the data from the source systems, and this task causes minimal impact on the system. It is important to take into account the different formats in which the data could appear, and define whether it is a real-time or batch ingest. The second part (T) applies a series of business rules or functions to the extracted data to turn it into useful and meaningful data. In the third part (L) the data from the previous phase is loaded into the destination system. Usually, this repository is a Data Lake or Data Warehouse.

Also in the ETL, you can define the periodicity with which the data can be loaded, for example, there will be data that is required to be loaded every week or data that is very transactional and it is necessary to apply the ETL process every certain amount of hours.

## How you build a Data Pipeline?

In the projects I have worked on, it is essential to begin to listen and understand the client and thus have a better knowledge about their line of business to make a work plan that specifies how the architecture solution to be implemented will be. This work plan also defines the processes and tasks that describe the Data Pipeline, which includes how it will be executed, schedule its execution tasks, monitor the Pipeline, and make quality tests verifying the output of the Pipeline.

Once the work plan and its scope are completed, the tasks are scheduled and divided among the different members of the data team, to begin with, the planned construction of the Data Pipeline.

## In a RDBMS Joins are your friends. Explain Why

By definition on RDBMS, the data is structured in tables, where each table covers a concept or definition of the business, but if analysis or a report of the data is required, it will be necessary to work with different tables, and this split data is not very useful.

This is how joins help us in this part of data union, where selecting the correct type of join (inner, left, right, outer) and the correct union filters, depending on the use case, helps to unify that information from different tables answering various questions about the business.

# What are the main features of a production pipeline.

The production pipeline must be safe and stable mainly, in addition to having the good practices of the DevOps culture and having an environment of continuous integration and continuous delivery. They must be monitored and configured with metrics to quickly know when something is not working correctly or there are anomalies, in addition to having mapped how failures will be handled and having the necessary logs for analysis.

It must also have a reliable security infrastructure that does not jeopardize access to the information it handles and have an estimate of how long it takes to run with different workloads and thus be able to make it scalable or apply autoscaling.

# How do you monitor this data pipelines?

There are cloud services that are linked to the different services they offer to be directly monitored, such as Stackdriver, in addition to storing the logs (normally for 30 days) to be analyzed. To correctly use these services, metrics are configured to monitor and configure alarms, for example, a metric for the use of a VM CPU, and establish a threshold of use of 70 %, and if the VM exceeds that capacity of use, send a mail to administrator and data engineer.

Due the logs are also stored, these can be analyzed to find useful insights and complement the information that the metrics provide, for example, what was the workload that caused the CPU to increase up to 80 %?.

# Give us a situation where you decide to use a NoSQL database instead of a relational database. Why did you do so?

The first time I opted for a NoSQL database was in college. The prototype of a taxi services app had to be made, the way it was defined implementation logic and how to store the different reservation records with their respective information, it was very complicated to want to model it in a relational database because the data did not have a defined schema and the normalizations were not clear, so a better approach was a non-relational database and save the JSON files as documents. Additionally, the way the data was retrieving was easy to consume.

# What are the non technical soft skills that are mos valuable for data engineers?

One of them is the ability to communicate, since it not only depends on speaking appropriately, it is the ability to break down technical problems and explain them correctly and understandably to other people who often do not have that technical background.

Have the ability to work as a team and based on defined work methodologies, for a correct harmony with all the people involved in the project or job.

# Diagram of the solution on GCP Infrastructure

Supose you have to design an Anomaly Detection Solution for a client in real or near real time.

A platform for anomaly detection is about finding patterns of interest (outliers, exceptions, peculiarities, etc.) that deviate from expected behavior within dataset(s). Given this definition, it's worth noting that anomaly detection is, therefore, very similar to noise removal and novelty detection. Though patterns detected with anomaly detection are actually of interest, noise detection can be slightly different because the sole purpose of detection is removing those anomalies - or noise - from data. Which technologies do you apply for real time ingestion and stream for an anomaly detection system? Diagram the solution in AWS or GCP Infrastructure.

Solution: The architecture suggested in the diagram assumes that the data can come from on-premises sources, or the infrastructure could be in the cloud and would be deployed in a PaaS service such as App Engine or a VM or Vm´s fleets orchestrated with Kubernetes for example.

For the part of ingesting in real-time or near real-time, the Pub / Sub service is used, which is capable of ingesting millions of data with a latency of milliseconds, this is managed with a topic to which multiple instances of the system can send their messages. Cloud Dataflow integrates natively with Pub / Sub through a subscription, in addition, it contains templates for this purpose that can be adjusted if you want to do this part quickly.

The PCollection structures in Dataflow allow reading this data and making transformations distributed with the Apache Beam framework with python, so it may be able to detect anomalies and save the results in a required format to deposit them in a bucket on Cloud Storage or transfer it to a DW handled by BigQuery for analysis. There is also the possibility that for this analysis of the anomalies detected work with python and one

of its many libraries, so an instance of Jupyter Lab with Datalab is suggested.

Finally, these anomalies may be of interest to developers and that is presented in a system, in the form of a report or a future step, have an interest in making ML models and predicting scenarios, so the architecture must be capable of scaling to new services.
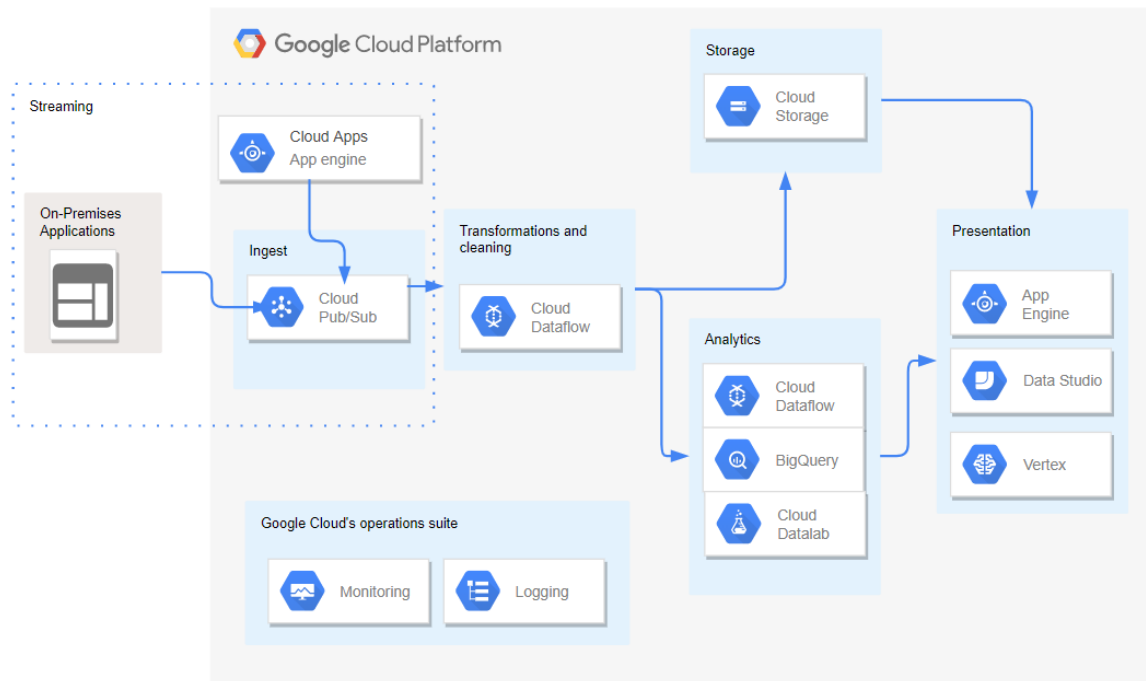


Figura 1- Architecture on GCP

## Differences between OLAP and OLTP Systems.

In its essence, OLTP is a database system oriented to transaction processing, so access to data is optimized for frequent reading and writing tasks since its transactionality is high and queries are not usually complex, while OLAP systems are databases oriented to analytical processing, access to the data is usually read-only. The most common action is retrieving data, with additions and joins to consult historical information and obtain useful information in its analysis.

## Create the following tables in a SQL Database

For the challenges, a temporary test laboratory instance was created to do the challenge.

A project on GCP was used to configure a Cloud SQL instance with postgresql.



Figura 2- Cloud SQL

It was configured to be able to work in Dbeaver.



Figura 3- Dbeaver configured

Finally the tables were created and the data was inserted.

Figura 4- Data inserted

Now, for the NoSQL database I used mongodb, using the Atlas service configured on a GCP instance



Figura 5- MongoDB cluster

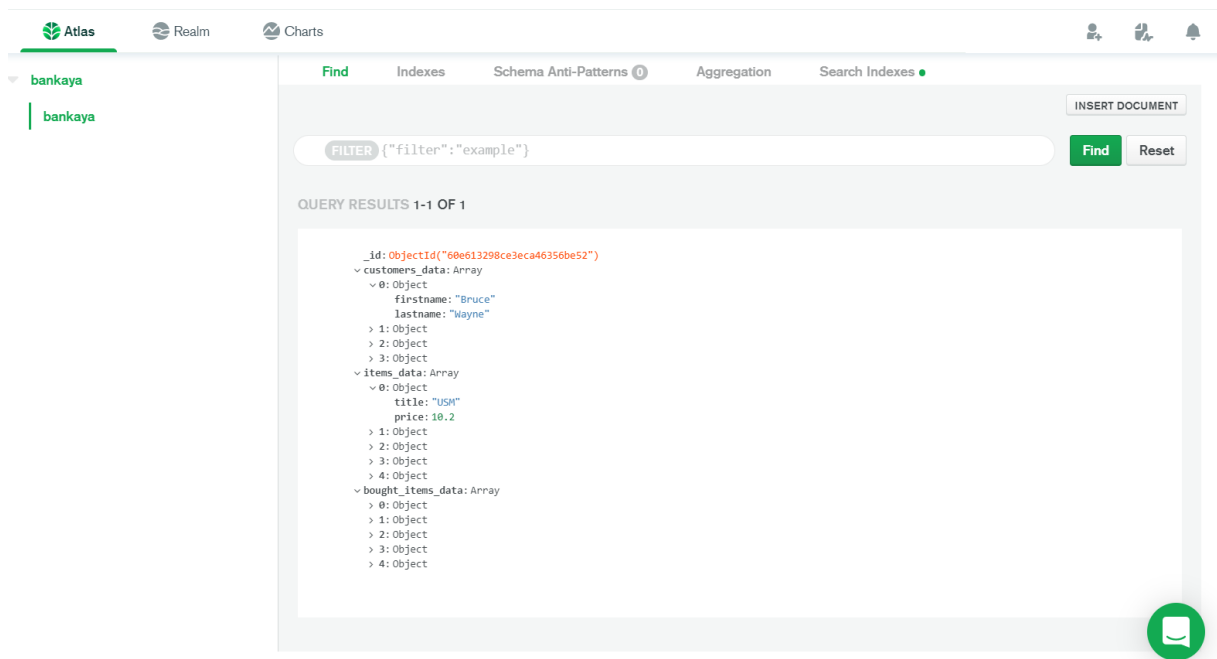The data was inserted into a document as proposed in the challenge.

Figura 6- MongoDB data

And this part was done through a python script



Figura 7- script in python

# Create a ETL Script pipeline

Dataprep is an interesting service that fits very well with the proposed task, it is capable of extracting various sources of information, combining them into elements called -ecipe", cleaning with a very complete suite of tools, selecting different loading destinations,

among They BigQuery, and above all, it is a managed and scalable service, since in the end all the transformation flows are transformed into a dataflow pipeline, to distribute the workload, and finally, it can schedule the frequency of execution of the flows, making a merge of the records, to only update the most current data once the history is loaded.

The first part consists of adding the data sources, which in this case is Mongo and Postgresql database.



Figura 8- adding sources

You can select the preprocessing and cleaning tasks in steps.



Figura 9- transformations

Finally the load will be to the BigQuery data warehouse where a dataset called Bankaya was created

Figura 10- pipeline

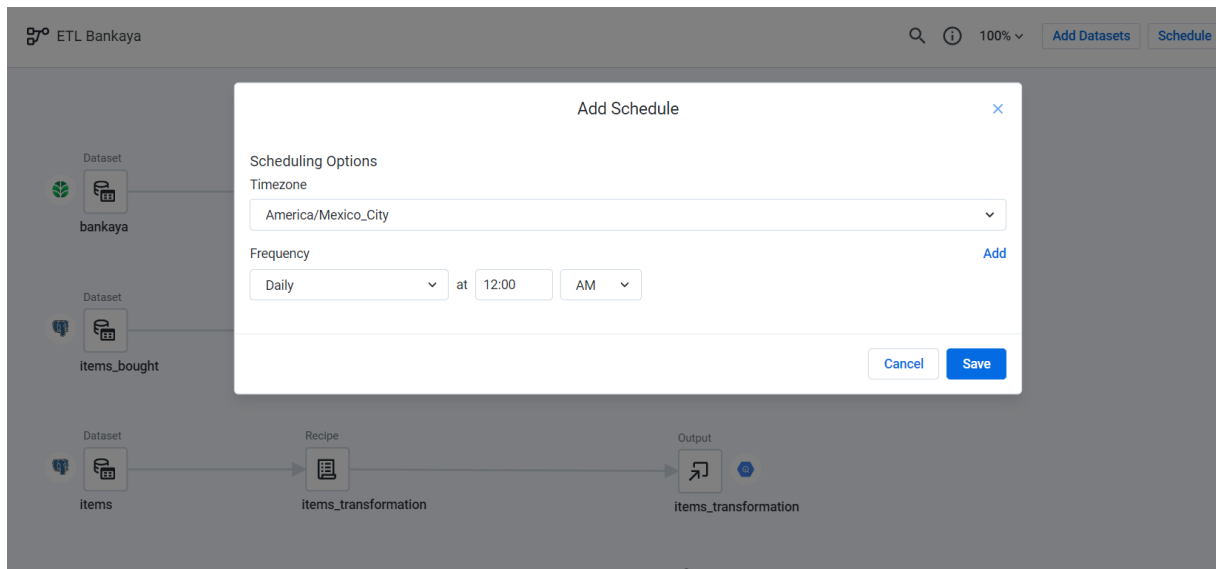It is scheduled so that the load is daily at 12:00 am with respect to the Mexico area (GMT -5)



Figura 11- scheduling

This is the result of the pipeline with the tables in Bigquery

Figura 12- result

And this is what a workflow launched in dataflow looks like for data transformation.



Figura 13- dataflow