

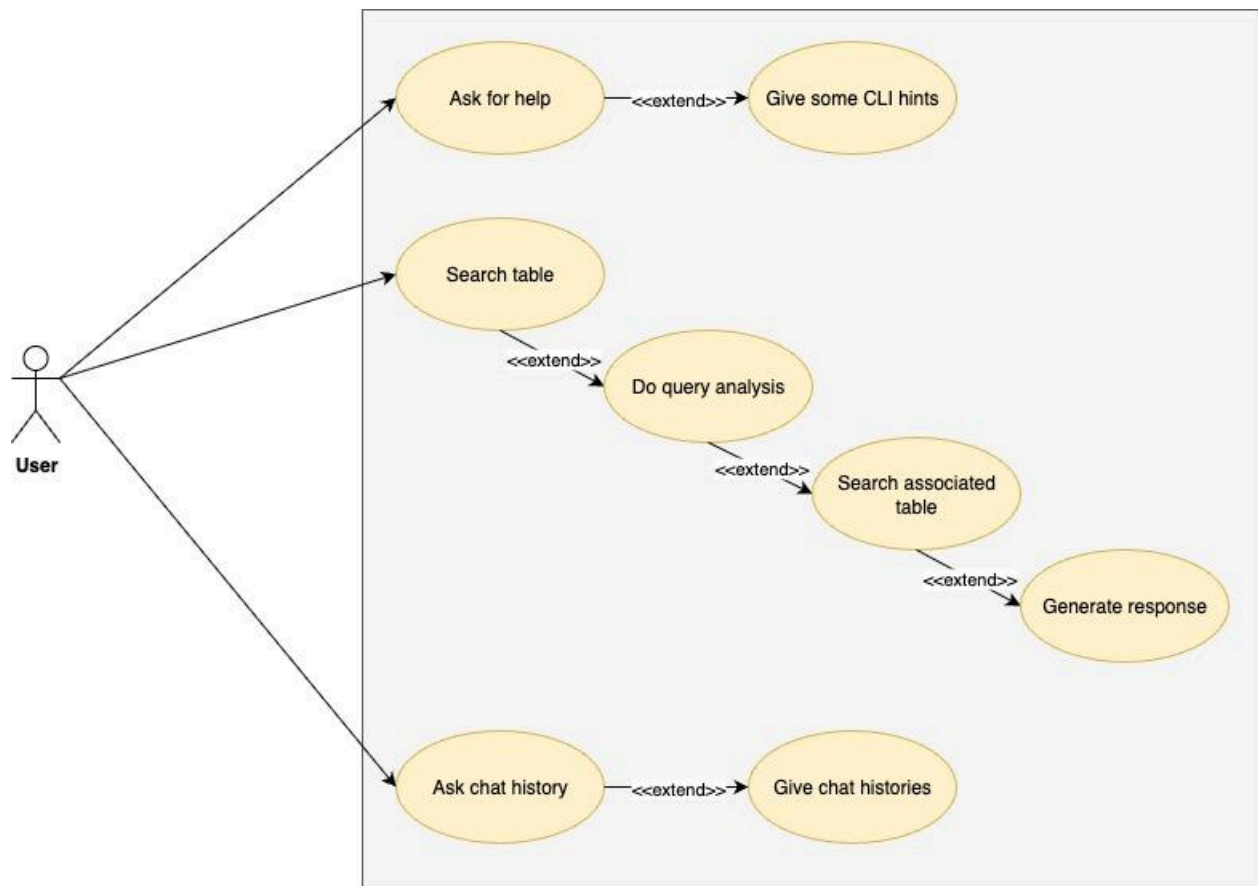
# Table Search Agent - MLOps On GCP

## Goals

The general objectives of this project are:

1. Create an interactive chatbot to search the required table after given a context words
2. The entire chat system made by some agents
3. Those agents deployed on Google Cloud Platform (GCP)
4. The overall agents' capabilities include
  - a. Give some answers for help questions
  - b. Return the chat history if asked
  - c. Search, and rank the associated table for the given chat

## Use-Case



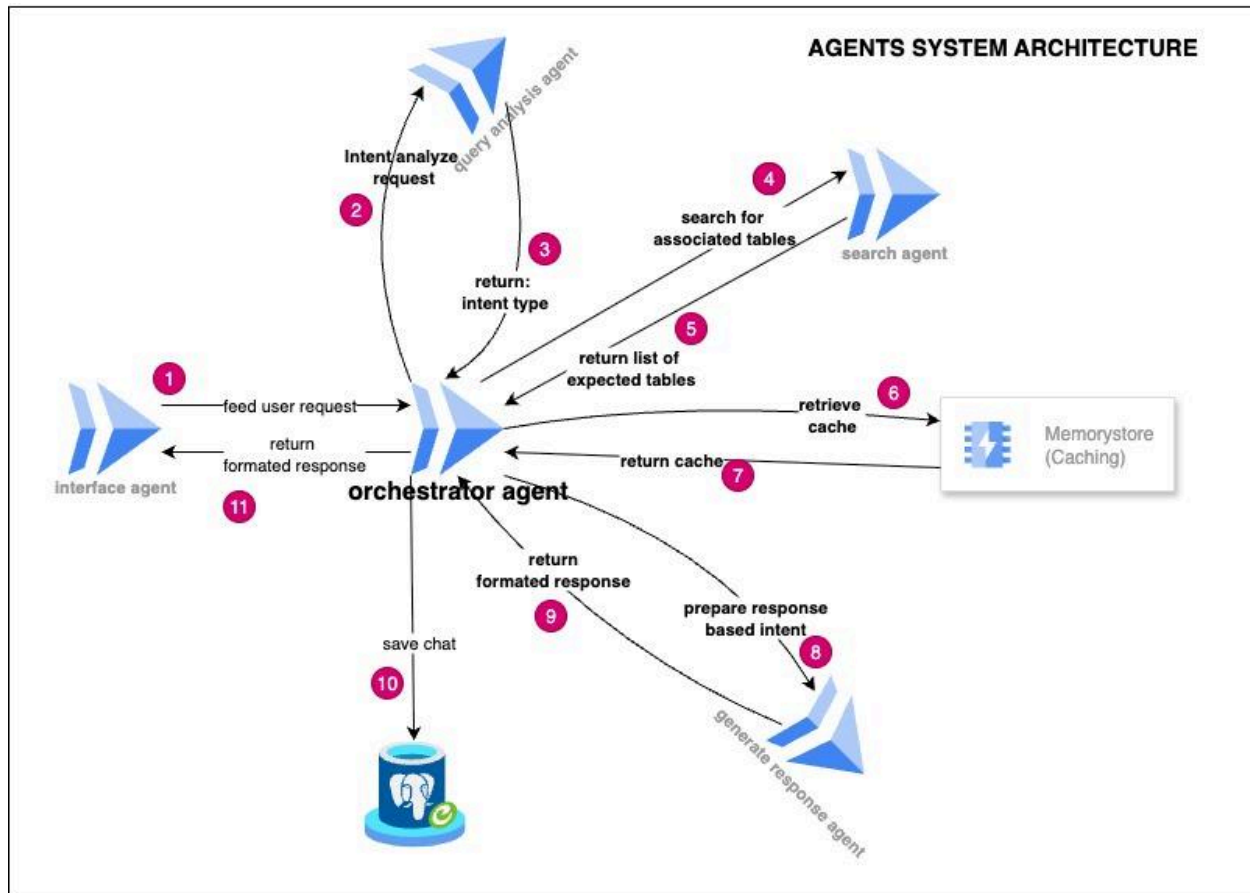
As we mentioned before, the main capabilities of our agents are response the user help request by giving list of required CLI, response user request for BigQuery tables searching, and give chat history list if requested.

Since the system merely depend on classic and simplest NLP algorithm (TF-IDF), we are set some question boundaries, described by the following table

Chat context	Format style	Example
Table search	<b>Flexible chat</b>	Help me find the marketing tables  Can you search for sales tables  Mention user related tables  Where products are stored
Ask for help	<b>FIXED chat</b>	Help
Ask for chat history	<b>FIXED chat</b>	history

# System Architecture

Since we are only utilizing classic NLP algorithms like TF-IDF and word matching, so there is no AI model included in this architecture. For further considerations, like enhancing the system capability to be more smart, we require to integrate with AI model like ChatGPT or Claude AI.



## Interface Agent:

**Tools** : GCP Cloud run  
**Capacity** : 256 MB memory  
**Language** : Golang  
**CI/CD** : Github action, Gitlab CI/CD, GCP Cloud Run  
**Docker Image** : Stored on GCP Artifact registry  
**Service Type** : HTTP  
**Job Scopes** : As a frontliner for incoming request, directly integrated with the front end  
**Further improvement** : - Add capability as rate limiter  
- As a bridge for all data engineering HTTP services

## Orchestrator Agent:

**Tools** : GCP Cloud run

*Capacity* : 256 MB memory  
*Languange* : Golang  
*CI/CD* : Github action or Gitlab CI/CD or GCP Cloud Run  
*Docker Image* : Stored on GCP Artifact registry  
*Service Type* : HTTP  
*Job Scopes* : As agent orchestrator, handle database (SQL and NoSQL) orchestration as well  
*Further improvement* : Consume event via Pub/Sub for asynchronous working capability

#### **Query Analysis Agent:**

*Tools* : GCP Cloud run  
*Capacity* : 256 MB memory  
*Languange* : Python  
*CI/CD* : Github action or Gitlab CI/CD or GCP Cloud Run  
*Docker Image* : Stored on GCP Artifact registry  
*Service Type* : HTTP  
*Job Scopes* : Analysis intent, refers to **intent enumeration**  
*Further improvement* : Integrate with AI Models like ChatGP, Claude AI

#### **Search Agent:**

*Tools* : GCP Cloud run  
*Capacity* : 256 MB memory  
*Languange* : Python  
*CI/CD* : Github action or Gitlab CI/CD or GCP Cloud Run  
*Docker Image* : Stored on GCP Artifact registry  
*Service Type* : HTTP  
*Job Scopes* : Search associated tables, rank them, based on TF-IDF and keyword mach  
*Further improvement* : Integrate with AI Models like ChatGP, Claude AI

#### **Response Generation Agent:**

*Tools* : GCP Cloud run  
*Capacity* : 256 MB memory  
*Languange* : Python  
*CI/CD* : Github action or Gitlab CI/CD or GCP Cloud Run  
*Docker Image* : Stored on GCP Artifact registry  
*Service Type* : HTTP  
*Job Scopes* : Generate response wording based on intent (intent: refers to **intent enumeration**)

#### **Response Caching:**

*Tools* : GCP Memory store (or replaced by Redis)  
*Job Scopes* : Cache agent response based on **intent type** and **table name**.  
*Expiry time* : 7 days

Job Scopes : Reduce latency, since the highly requested table's responses already available

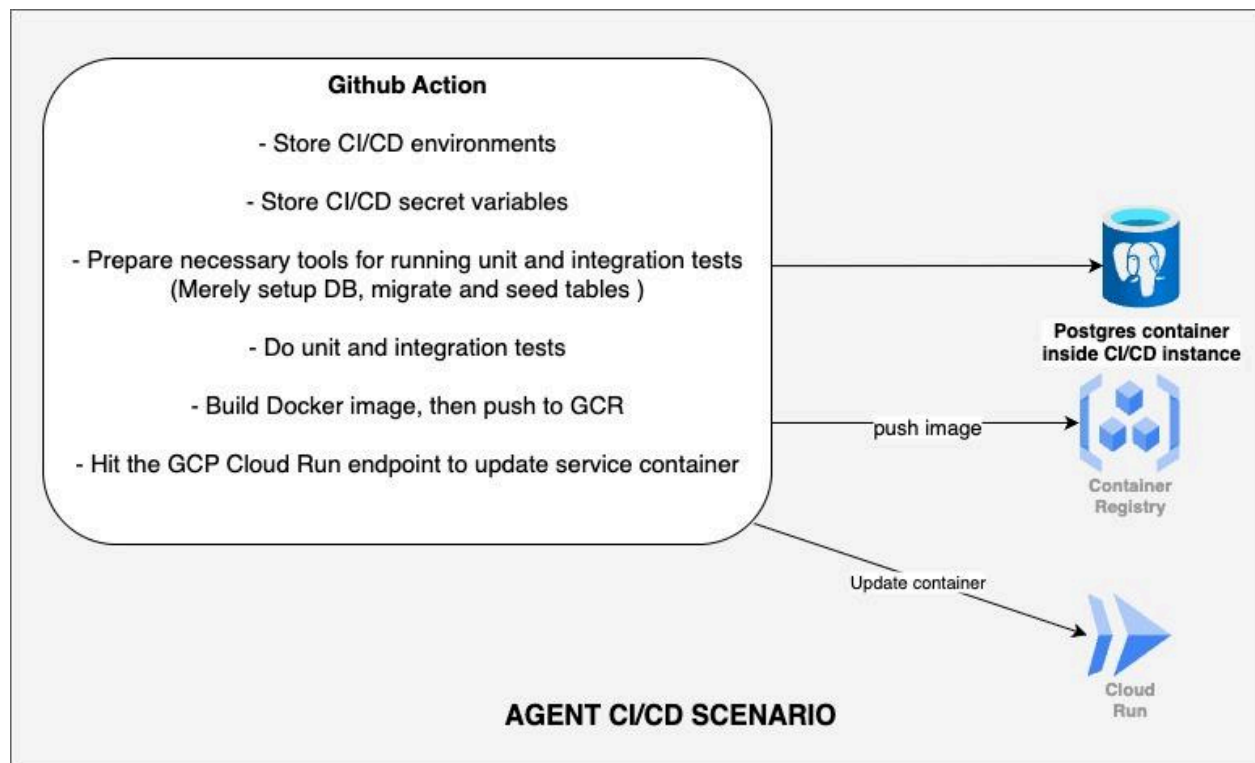
#### System Database:

Tools : Postgres SQL  
Job Scopes : System database  
E/R Data Model : Will be provided in the next section

#### Other Specifications:

IaC : Terraform

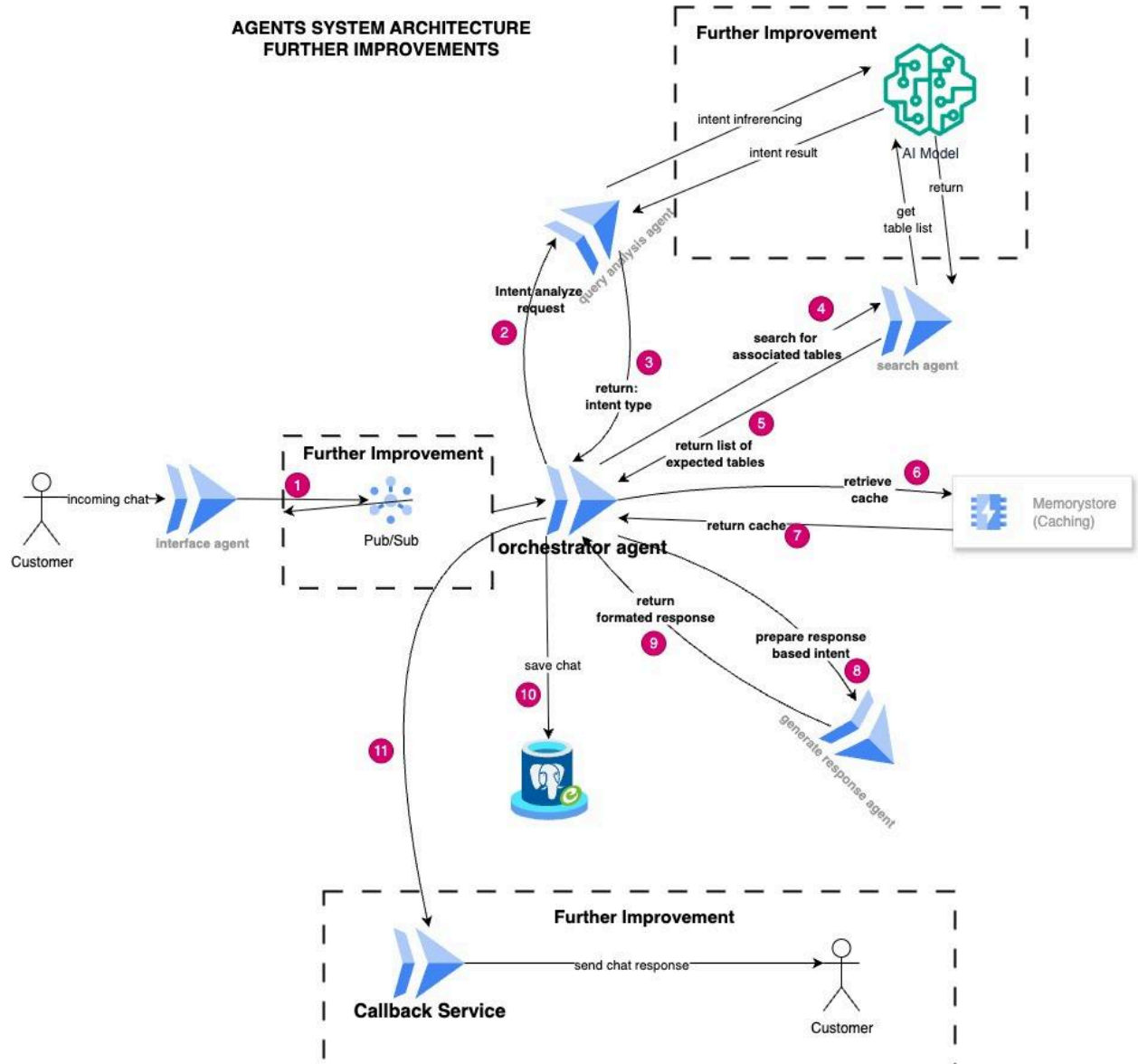
#### CI/CD Pipeline :



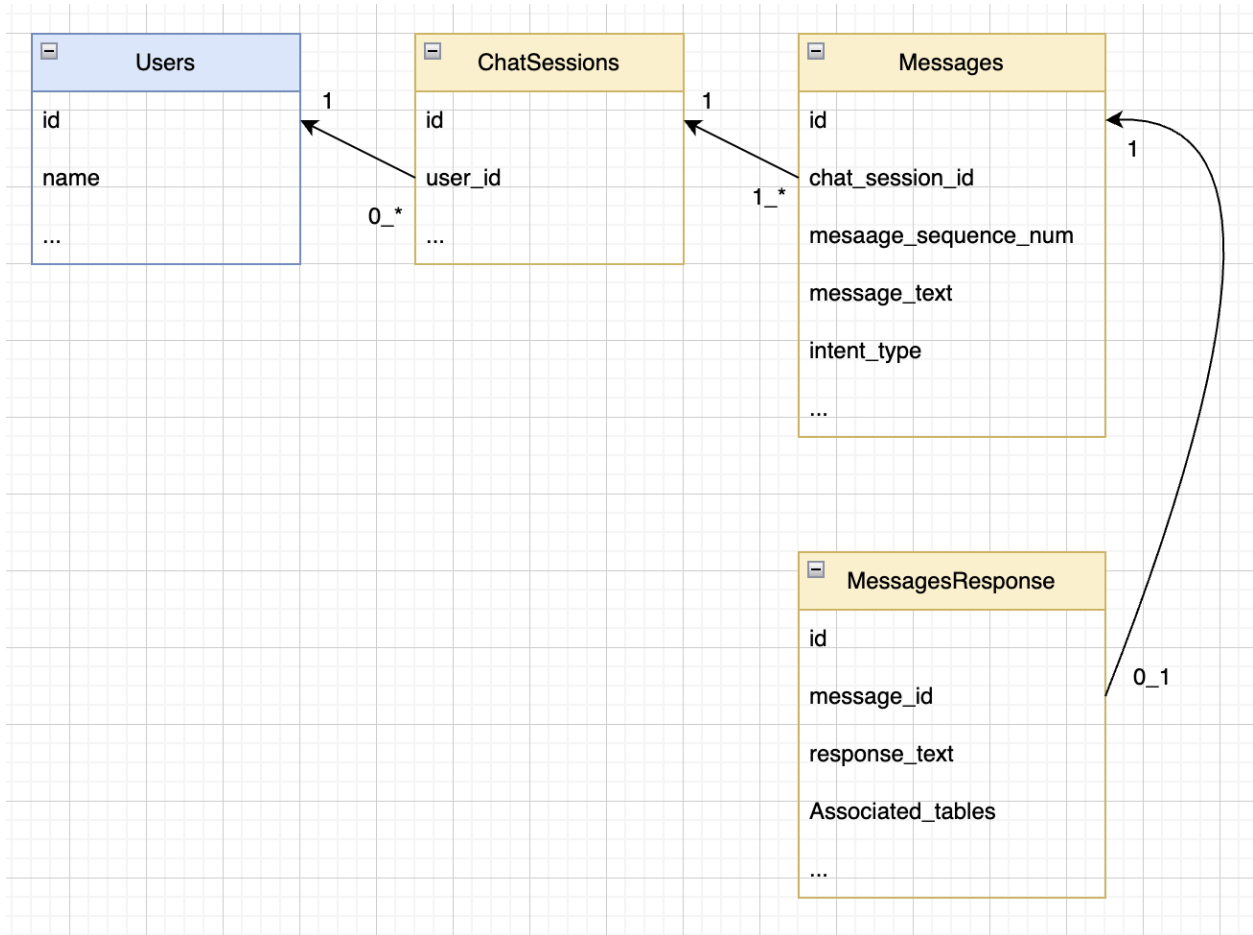
#### FURTHER IMPROVEMENT

- **Reduce Latency:** add event driven mechanisme between **inteface agent** and **orchestrator agent**, Add callback mechanisme when the process finished
- **Increase capability:** Integrate Query Analysis Agent and Table Search Agent with AI models. It can be in house AI model, or ChatGPT, CLaude AI.

# AGENTS SYSTEM ARCHITECTURE FURTHER IMPROVEMENTS

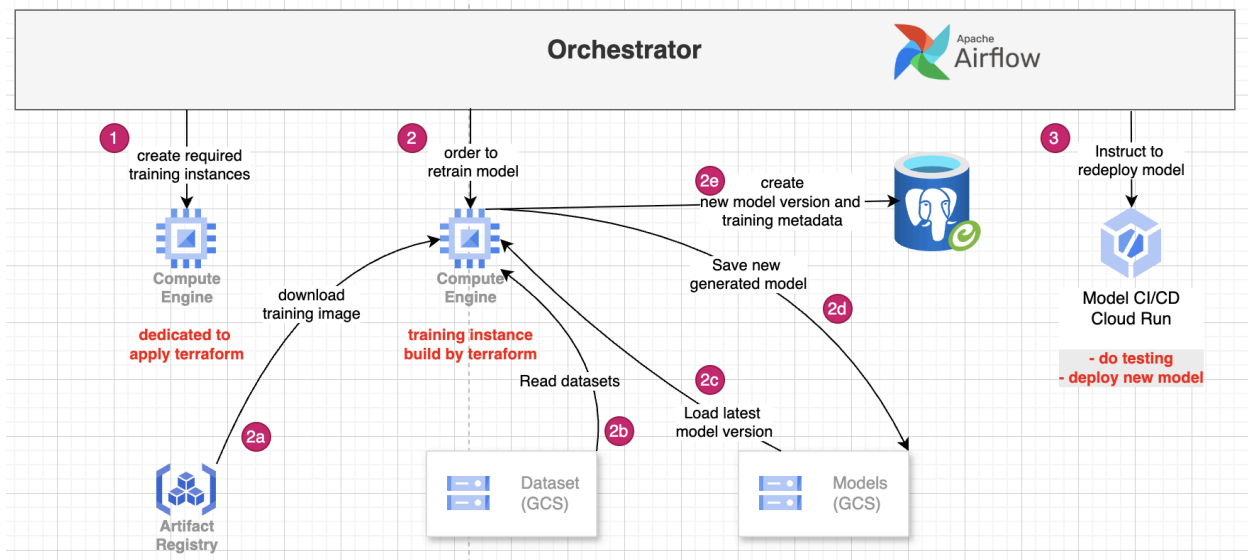


## E/R Data Model



# MODEL LIFECYCLE

(Case specific for inhouse models developments)



This section demonstrate model retraining and redeployment in GCP. We put Apache airflow as orchestrator. Its main tasks including

1. Orchestrate a dedicated terraform apply instance to create required training infrastructure
2. Once training instance created, Apache Airflow will manage to instruct start training. Then, the insantance will pull the training image from GCR, and run training container. Then, it manages to load dataset from GCS (all or batch based on engineering SLA). After that, the model will be loaded for retraining. Once the retraining process was finished, the newly created model will be saved in GCS. The training instance will save new model version, performances, and metadata in model DB.
3. After that, Apache Airflow will manage to orchestrate GCP cloud run for model CI/CD.
4. Finally, Airflow managed to instruct terraform instance to destroy the training instance.

## Terraform instance:

Tools	: GCP Compute Engine
Capacity	: 256 - 1GB MB memmory
Standby	: 1 (dedicated)

## Training instance:



<i>Tools</i>	: GCP Compute Engine
<i>Capacity</i>	: depend on the model size and dataset
<i>Use GPU</i>	: depend on engineering and model SLA
<i>Status</i>	: temporary instance, destroyed when training pipeline was finished

**Why don't we managed to utilize GCP vertex or AWS sagemaker instead ?**

They are absolutely expensive, by designing our training instance orchestrations, we manage to reduce huge costs, and increase flexibility