



Data & AI Boot-Kon Event

Title: ELT & Lakehouse : BigQuery, Dataform, LLM for sentiment analysis

Goal of the lab

- Use Dataform for SQL Transformations in BigQuery.
- Perform sentiment analysis of customer feedback.

Author: Wissem Khlifi

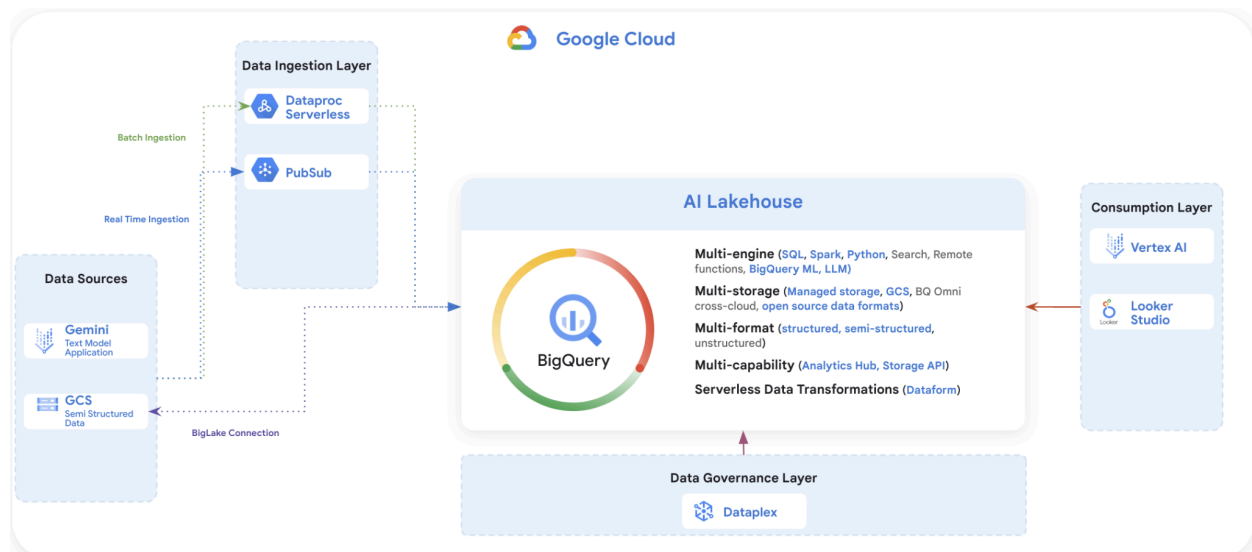
Date: 2024-04-01

Estimated Completion Time: 45 Minutes

CAUTION:

This lab is for educational purposes only and should be used with caution in production environments. Google Cloud Platform (GCP) products are changing frequently, and screenshots and instructions might become inaccurate over time. Always refer to the latest GCP documentation for the most up-to-date information.

Architecture Diagram:



ELT: Dataform & LLM for sentiment analysis from BigQuery

READING Section : Read the following Explanation of what dataform is and what is its purpose

Dataform

Dataform is a fully managed service that helps data teams build, version control, and orchestrate SQL workflows in BigQuery. It provides an end-to-end experience for data transformation, including:

- Table definition: Dataform provides a central repository for managing table definitions, column descriptions, and data quality assertions. This makes it easy to keep track of your data schema and ensure that your data is consistent and reliable.



- **Dependency management:** Dataform automatically manages the dependencies between your tables, ensuring that they are always processed in the correct order. This simplifies the development and maintenance of complex data pipelines.
- **Orchestration:** Dataform orchestrates the execution of your SQL workflows, taking care of all the operational overhead. This frees you up to focus on developing and refining your data pipelines.

Dataform is built on top of Dataform Core, an open source SQL-based language for managing data transformations.

Dataform Core provides a variety of features that make it easy to develop and maintain data pipelines, including:

- **Incremental updates:** Dataform Core can incrementally update your tables, only processing the data that has changed since the last update. This can significantly improve the performance and scalability of your data pipelines.
- **Slowly changing dimensions:** Dataform Core provides built-in support for slowly changing dimensions, which are a common type of data in data warehouses. This simplifies the development and maintenance of data pipelines that involve slowly changing dimensions.
- **Reusable code:** Dataform Core allows you to write reusable code in JavaScript, which can be used to implement complex data transformations and workflows.

Dataform is integrated with a variety of other Google Cloud services, including GitHub, GitLab, Cloud Composer, and Workflows. This makes it easy to integrate Dataform with your existing development and orchestration workflows.

Benefits of using Dataform in Google Cloud

There are many benefits to using Dataform in Google Cloud, including:

- **Increased productivity:** Dataform can help you to increase the productivity of your data team by automating the development, testing, and execution of data pipelines.
- **Improved data quality:** Dataform can help you to improve the quality of your data by providing a central repository for managing table definitions, column descriptions, and data quality assertions.
- **Reduced costs:** Dataform can help you to reduce the costs associated with data processing by optimizing the execution of your SQL workflows.
- **Increased scalability:** Dataform can help you to scale your data pipelines to meet the needs of your growing business.

Use cases for Dataform

Dataform can be used for a variety of use cases, including:

- **Data warehousing:** Dataform can be used to build and maintain data warehouses that are scalable and reliable.
- **Data engineering:** Dataform can be used to develop and maintain data pipelines that transform and load data into data warehouses.
- **Data analytics:** Dataform can be used to develop and maintain data pipelines that prepare data for analysis.
- **Machine learning:** Dataform can be used to develop and maintain data pipelines that prepare data for machine learning models.

LAB Section : Dataform Prerequisites

1. Enable Services API (you can skip this step if you completed LAB 1)

Ensure all necessary APIs (BigQuery API, Vertex AI API, BigQuery Connection API, Dataform API, Secret Manager API) are enabled

2. Create a connection to an external data source in BigQuery (you can skip this step if you completed LAB 1)

- Create an External Connection (Enable BQ Connection API if not already done) and note down the Service Account id from the connection configuration details:
- Click the +ADD button on the BigQuery Explorer pane (in the left of the BigQuery console) and click "Connection to external data sources" in the popular sources listed





- Select Connection type as “Vertex AI remote models , remote functions and Biglake” and provide “fraud-transactions-conn” as Connection ID, select Multi Region location type.

External data source

Connection type

Vertex AI remote models, remote functions and BigLake (Cloud Resource)

- Once the connection is created, take a note of the Service Account generated from the connection configuration details

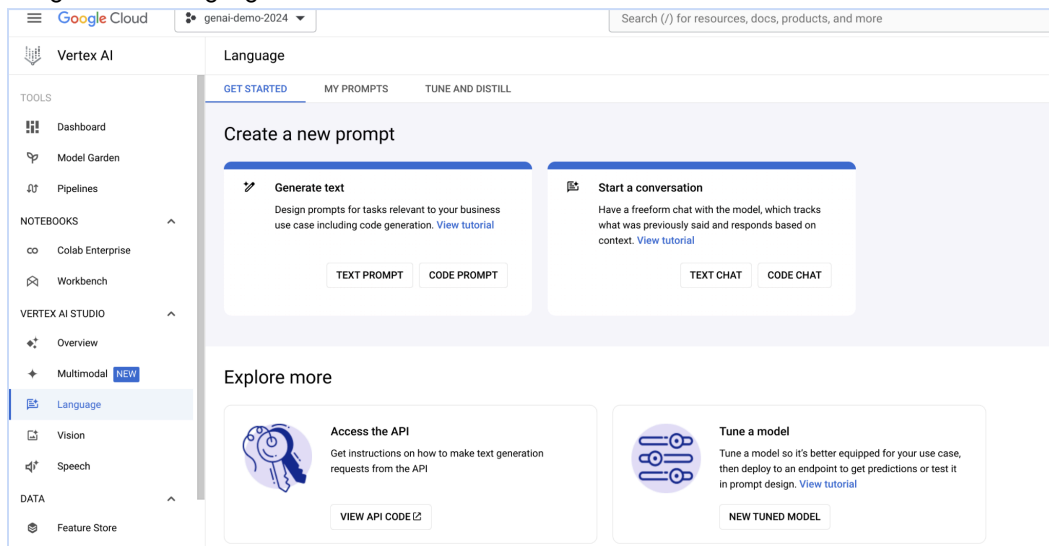
3. Grant Permissions (you can skip this step if you completed LAB 1)

In this step we will grant permissions to the Service Account to access the Vertex AI service:

Open IAM and add the Service Account you copied after creating the external connection as the Principal and select “Vertex AI User” Role.

4. Using Large Language Models from Vertex AI

Google Cloud’s language models are available within the Vertex AI Studio inside the Vertex AI service.



5. Prompt design

Prompt design is the process of creating prompts that elicit the desired response from language models. Writing well structured prompts is an essential part of ensuring accurate, high quality responses from a language model. If you need to understand this concept a bit more this is a page that introduces some basic concepts, strategies, and best practices to get you started in designing prompts

(<https://cloud.google.com/vertex-ai/docs/generative-ai/learn/introduction-prompt-design>).

The reference page above also goes into the more advanced settings you can see on the right hand side of the prompt box such as temperature, top K, top P etc.



LAB Section : Creating a Dataform Pipeline

First step in implementing a pipeline in Dataform is to set up a repository and a development environment. Detailed quickstart and instructions can be found [here](#).

1. Create a Repository in Dataform

+ CREATE REPOSITORY

< Create repository

Repositories contain a single Dataform project that can be connected to your Git provider. Within a repository you can create workspaces for development and execute your SQL workflows against BigQuery.

Repository ID *
hackaton-repository ?

Region *
us-central1 (Iowa) ▼ ?

Service account
Default Dataform service account ▼ ?

New repositories start empty and can be connected to a git provider after being created.

CREATE

2. Dataform Service Account

Take note and save somewhere the newly created service account for Dataform.

Example: service-112412469323@gcp-sa-dataform.iam.gserviceaccount.com

Repository created successfully.

Dataform will execute workflows as the service account `service-112412469323@gcp-sa-dataform.iam.gserviceaccount.com`.
Ensure that the service account has been granted the role `roles/bigquery.user` in order to create datasets, tables and jobs.
Also ensure that the service account has been granted access to read any datasets or tables that your workflow will consume. [Learn more](#)

GO TO REPOSITORIES



hackathon-repository

DEVELOPMENT WORKSPACES WORKFLOW EXECUTION LOGS RELEASES & SCHEDULING **SETTINGS**

[CONNECT WITH GIT](#) [CONFIGURE PRIVATE NPM PACKAGES](#)

Name hackathon-repository

Location us-central1

Service account Default Dataform service account

By default, dataform will use a service account derived from your project number: service-763163306074@gcp-sa-dataform.iam.gserviceaccount.com

Workspace compilation overrides [EDIT](#)

Override the target project ID, table prefix, and schema suffix settings for manual executions of all workspaces in the repository. The default settings are stored in 'workflow_settings.yaml'. Learn more about [workspace compilation overrides](#).

No workspace compilation overrides.

3. Create and initialize a Dataform development workspace

- In the Google Cloud console, go to the **Dataform** page.
Go to Dataform
- Click on the new repository you have just created.
- Click add **Create development workspace**.
- In the **Create development workspace** window, do the following:
 - In the **Workspace ID** field, enter “**hackathon-<YOURLASTNAME>-workspace**” (replace **<YOURLASTNAME>** with your name)
 - Click **Create**.
- The development workspace page appears.
- Click on the newly created development workspace
- Click **Initialize workspace**.

- You will copy the dataform files from the following repository, in the next steps.

<https://github.com/dace-de/bootkon-h2-2024/tree/main/dataform>

- Edit **workflow_settings.yaml** file :*
 - Replace **defaultDataset** value with **ml_datasets** , make sure defaultProject value should be **your project id**

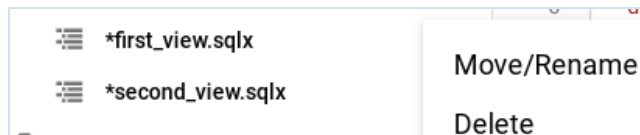
Note: Nevermind if you have a different dataform core version, just continue

```
workflow_settings.yaml
1 defaultProject: genai-demo-2024
2 defaultLocation: US
3 defaultDataset: ml_datasets
4 defaultAssertionDataset: dataform_assertions
5 dataformCoreVersion: 3.0.0-beta.4
6
```

- Click on Install Packages **Only Once**.

Package installation succeeded

- Remove from your local dataform repository the following SQLX files; Delete the following files;
 - first_view.sqlx**
 - second_view.sqlx**



- k. Click on definitions and create a new directory called “models”:

Create new directory

Add a directory path *

definitions/models/

+ CREATE DIRECTORY

CANCEL

- l. Click on models directory and create 2 new files ; (make sure all file names are in lowercase)

- create_dataset.sqlx
- llm_model_connection.sqlx

Those files should be created under **definitions/models** directory

Example:

Create new file

Dataform only compiles .sqlx and .js files in the c in the includes / directory. Files in other root direct ignored unless explicitly referenced.

Add a file path *

definitions/models/create_dataset.sqlx

+ CREATE FILE

CANCEL

- m. Copy the contents from <https://github.com/dace-de/bootkon-h2-2024/tree/main/dataform/definitions/models> to each of those files.
- n. Click on definitions and create 3 new files: (make sure all file names are in lowercase)
- mview_ulb_fraud_detection.sqlx
 - sentiment_inference.sqlx
 - ulb_fraud_detection.sqlx

Those files should be created under **definitions** directory

Example:



Create new file

Dataform only compiles .sqlx and .js files in the definitions/ directory and .js files in the includes/ directory. Files in other root directories and other file types will be ignored unless explicitly referenced.

Add a file path *

definitions/mview_ulb_fraud_detection.sqlx

?

+ CREATE FILE

CANCEL

- o. Copy the contents from <https://github.com/dace-de/bootkon-h2-2024/tree/main/dataform/definitions> to each of those files.
- p. Make sure to replace database by your project ID value in ulb_fraud_detection.sqlx file

definitions/ulb_fraud_detection.sqlx FORMAT ✓

Press Alt+F1 for Accessibility Options.

```
1 config
2   type: "declaration",
3   database: "bootkon-2024",
4   schema: "ml_datasets",
5   name: "ulb_fraud_detection"
6 }
```

- q. In llm_model_connection.sqlx, replace the 'us.llm-connection' connection with the connection name you have created in LAB 2 during the BigLake section. If you have followed the steps in LAB 2, the connected name should be "us.fraud-transactions-conn"

Notice the usage of \$ref in line 10, of **definitions/mview_ulb_fraud_detection.sqlx** "sqlx" file. The advantages of using \$ref in Dataform are

- Automatic Reference Management: Ensures correct fully-qualified names for tables and views, avoiding hardcoding and simplifying environment configuration.
 - Dependency Tracking: Builds a dependency graph, ensuring correct creation order and automatic updates when referenced tables change.
 - Enhanced Maintainability: Supports modular and reusable SQL scripts, making the codebase easier to maintain and less error-prone.
5. Run the dataset creation by **TAG**. TAG allows you to just execute parts of the workflows and not the entire workflow. Click on **Start Execution > Tags > "dataset_ulb_fraud_detection_llm" > Start Execution**



Execute

Execute all actions, or select a subset of actions. Service account **service-112412469323@gcp-sa-dataform.iam.gserviceaccount.com** will be used.
[Learn more](#)

[ALL ACTIONS](#) [SELECTION OF ACTIONS](#) [SELECTION OF TAGS](#)

Select tags to execute
dataset_ulb_fraud_detection_llm

Execution options

☐ Include dependencies

☐ Include dependents

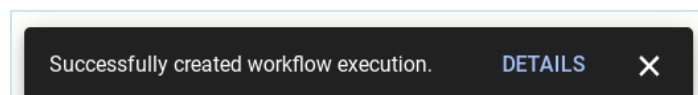
☐ Run with full refresh

1 action selected for execution

Destination	Type	File
bootkon-2024.ml_datasets.create_dataset	Operations	definitions/models/create_dataset.sqlx

[START EXECUTION](#) [CANCEL](#)

6. Click on Details;



7. Notice the Access Denied error on BigQuery for the dataform service account XXX@gcp-sa-dataform.iam.gserviceaccount.com;

Failure reason	Access Denied: Project 112412469323: User does not have bigquery.jobs.create permission in project 112412469323.
Action destination	bootkon-2024.ml_datasets.create_dataset

8. Go to IAM & Admin > Grant access and grant **BigQuery Data Editor** , **BigQuery Job User** and **BigQuery Connection User** to the data from the service account. Click on Save.



Grant access to "genai-demo-2024"

Grant principals access to this resource and add roles to specify what actions the principals can take. Optionally, add conditions to grant access to principals only when a specific criteria is met. [Learn more about IAM conditions](#)

Resource
genai-demo-2024

Add principals
Principals are users, groups, domains, or service accounts. [Learn more about principals in IAM](#)

New principals *
service-292219499736@gcp-sa-dataform.iam.gserviceaccount.com

Assign roles
Roles are composed of sets of permissions and determine what the principal can do with this resource. [Learn more](#)

Role *	IAM condition (optional)	
BigQuery Data Editor Access to edit all the contents of datasets	+ ADD IAM CONDITION	
BigQuery Job User Access to run jobs	+ ADD IAM CONDITION	
BigQuery Connection User	+ ADD IAM CONDITION	

+ ADD ANOTHER ROLE

SAVE CANCEL

Note: If you encounter the following policy update screen, just click on update.

Policy was out of date

The version of the policy you are currently viewing is out of date. The following changes have been committed since this page was loaded - these changes may have been committed by another user, in another tab, or the result of a user accepting a project invite.

Changes committed since initial load

Action	Role	Principal	Condition
Add	Dataform Service Agent	service-904843752328@gcp-sa-dataform.iam.gserviceaccount.com	
Add	Dataproc Service Agent	service-904843752328@dataproc-accounts.iam.gserviceaccount.com	
Add	Cloud Pub/Sub Service Agent	service-904843752328@gcp-sa-pubsub.iam.gserviceaccount.com	

Proposed changes

These are the proposed changes that will be applied on the new version of the policy if you continue with your update. Do you wish to continue?

Action	Role	Principal	Condition
Add	BigQuery Data Editor	service-904843752328@gcp-sa-dataform.iam.gserviceaccount.com	
Add	BigQuery Job User	service-904843752328@gcp-sa-dataform.iam.gserviceaccount.com	
Add	BigQuery Connection User	service-904843752328@gcp-sa-dataform.iam.gserviceaccount.com	

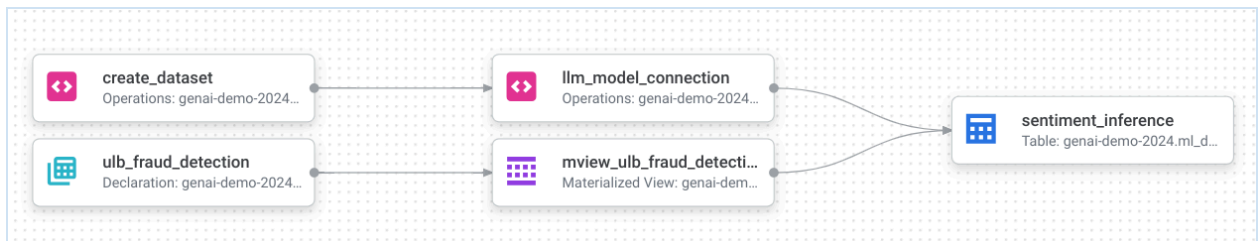
CANCEL UPDATE

9. Go back to dataform from the BigQuery console, and retry step 5. Notice the execution status. It should be a success.

Status Success

10. Click on Compiled graph and explore it;

Go to **Dataform > Compiled Graph**



LAB Section : Execute the workspace workflow

1. For the sentiment inference step to succeed . You need to grant the external connection service account the Vertex AI user privilege. More details can be found in this [link](#).

🔗 fraud-transactions-conn

Connection info

Connection ID	projects/genai-demo-2024/locations/us/connections/fraud-transactions-conn
Friendly name	
Created	Mar 28, 2024, 10:34:35 AM UTC+1
Last modified	Mar 28, 2024, 10:34:35 AM UTC+1
Data location	us
Description	
Connection type	Vertex AI remote models, remote functions and BigLake (Cloud Resource)
Service account id	<u>bqcx-292219499736-a17a@gcp-sa-bigquery-condel.iam.gserviceaccount.com</u>

2. Take note of the service account and grant it the **Vertex AI User** role.

Grant access to "bootkon-2024"

Grant principals access to this resource and add roles to specify what actions the principals can take. Optionally, add conditions to grant access to principals only when a specific criteria is met. [Learn more about IAM conditions](#)

Resource

bootkon-2024

Add principals

Principals are users, groups, domains, or service accounts. [Learn more about principals in IAM](#)

New principals *

bqcx-112412469323-7jrd@gcp-sa-bigquery-condel.iam.gserviceaccount.com

Assign roles

Roles are composed of sets of permissions and determine what the principal can do with this resource. [Learn more](#)

Role *
Vertex AI User

IAM condition (optional) ?
[+ ADD IAM CONDITION](#)

Grants access to use all resource in Vertex AI

[+ ADD ANOTHER ROLE](#)

[SAVE](#) [CANCEL](#)

3. Click **START EXECUTION** from the top menu, then **"Execute actions"**.



START EXECUTION ▾

Execute actions

- Click on **ALL ACTIONS** Tab then Click on **START EXECUTION**

Execute

Execute all actions, or select a subset of actions. Service account **service-112412469323@gcp-sa-dataform.iam.gserviceaccount.com** will be used.
[Learn more](#)

ALL ACTIONS

SELECTION OF ACTIONS

SELECTION OF TAGS

Execution options

☐ Run with full refresh ⓘ

5 actions selected for execution

Destination	Type	File
bootkon-2024.ml_datasets.create_dataset	Operations	definitions/models/create_dataset...
bootkon-2024.ml_datasets.llm_model_co...	Operations	definitions/models/llm_model_co...
bootkon-2024.ml_datasets.mview_ulb_fra...	Materialized View	definitions/mview_ulb_fraud_detec...
bootkon-2024.ml_datasets.sentiment_infe...	Table	definitions/sentiment_inference.sqlx
bootkon-2024.ml_datasets.ulb_fraud_dete...	Declaration	definitions/ulb_fraud_detection.sqlx

START EXECUTION

CANCEL

- Check the execution status. It should be a success.
- Verify the new table **sentiment_inference** in the ml_datasets dataset in BigQuery.
- Query the BigQuery table content (At this point you should be familiar with running BigQuery SQL)

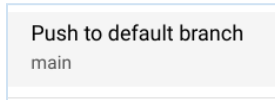
BigQuery SQL : Check few rows of sentiment_inference table
Replace the project-id by your project id value.

```
SELECT distinct ml_generate_text_llm_result,  
prompt,  
Feedback  
FROM `project-id.ml_datasets.sentiment_inference` LIMIT 10;
```

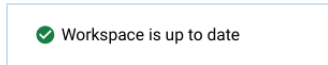
Google Cloud



8. **[Max 2 minutes]** Discuss the table results within your team group.
9. Before moving to the challenge section of the Lab, make sure to “**Commit X Changes**” (X should be about 7) and then “**Push to Default Branch**”



You should now have the message



CHALLENGE Section : Production, Scheduling and Automation

Automate and schedule the compilation and execution of the pipeline. This is done using release configurations and workflow configurations.

Release Configurations:

Release configurations allow you to compile your pipeline code at specific intervals that suit your use case. You can define:

- Branch, Tag, or Commit SHA: Specify which version of your code to use.
- Frequency: Set how often the compilation should occur, such as daily or weekly.
- Compilation Overrides: Use settings for testing and development, such as running the pipeline in an isolated project or dataset/table.

Common practice includes setting up release configurations for both test and production environments. For more information, refer to the [release configuration documentation](#).

Workflow Configurations

To execute a pipeline based on your specifications and code structure, you need to set up a workflow configuration. This acts as a scheduler where you define:

- Release Configuration: Choose the release configuration to use.
- Frequency: Set how often the pipeline should run.
- Actions to Execute: Specify what actions to perform during each run.

The pipeline will run at the defined frequency using the compiled code from the specified release configuration. For more information, refer to the [workflow configurations documentation](#).

[TASK] Challenge : Take up to 10 minutes to Setup a Daily Frequency Execution of the Workflow

Goal: Set up a daily schedule to automate and execute the workflow you created.

1. Automate and schedule the pipeline’s compilation and execution.
2. Set up a daily frequency execution of the workflow you have created.
3. Define release configurations for different environments.
4. Set up workflow configurations to schedule pipeline execution.



Note: If you are stuck and cannot figure out how to proceed after a few minutes, ask the event moderator for help.



Congratulations on completing Lab 3!
You can now move on to Lab 4 for further practice.

