# Data Engineering project

ZoomCamp - 2024
Maria Fisher

## Introduction

The project is built on Azure with Terraform, using the Medallion architecture for data transformation with three layers: bronze, silver, and gold. Key Azure services like Azure Data Factory, Azure Databricks, and Azure Data Lake Storage Gen2.

Synthetic datasets were used to simulate food storage, environmental data scenarios and evaluate fungus and consequently mycotoxin contamination risks. A pipeline was created to transfer data from blob storage to an SQL database, and another pipeline was established to collect data from the database for processing in the medallion layers.
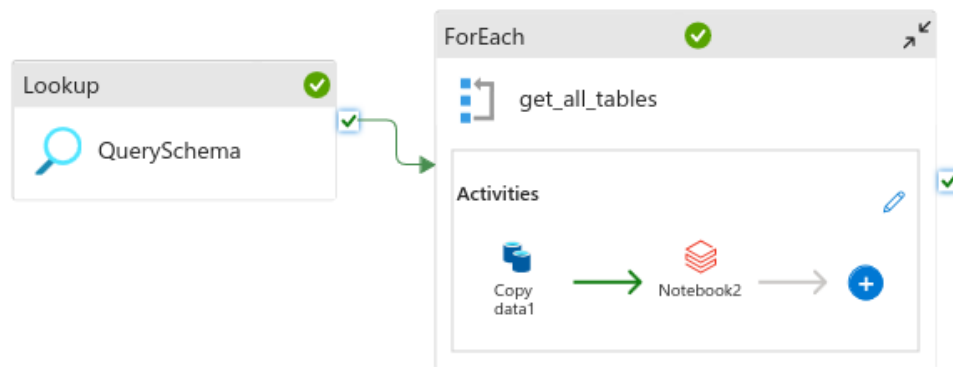
In Databricks a workspace was created within Azure, and a secret key was set up using Azure Key Vault for secure access.

The pipeline was structured to adhere to the Medallion architecture's Bronze-Silver-Gold layers, ensuring accurate data copying and processing.

For data transformation dbt-Databricks was used and SQL files for snapshots of tables were created. The transformed data was saved in the silver layers, and the final output in the gold layer.

Dashboard was created with Databricks and connected with Looker.

**Project pipeline:**



This pipeline extracts data from an SQL database and loads it into a Databricks notebook for analysis and storage in a DataLake for further use.

# Steps:

Firstly, ensure that your SQL database contains the necessary data. You have two options: populate your database from the Azure marketplace or use your own data. If you prefer the latter, I have prepared a detailed document outlining the steps to create your database using your own data. You can find the instructions in the file named ''Blob storage to Azure SQL database''.
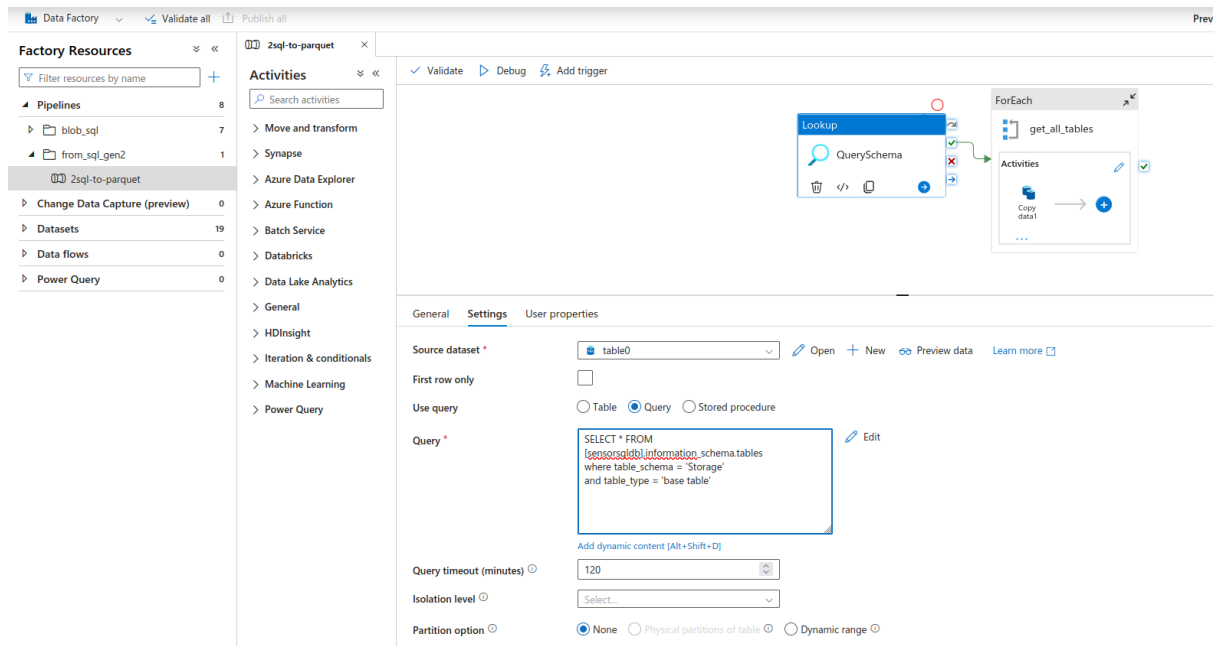
1. **Create pipeline to collect data from database to be processed in the medallion layers**

● The system requires getting data from a SQL database and storing it in a data lake (that was created previously).
● A query is required to retrieve all tables needed from the database schema.

**Pipeline to Copy data from SQL database and save as parquet file into storage Gen2**



● In the pipeline a ForEach activity is established to iterate through each table retrieved from the SQL query.
● Within the ForEach loop, a Copy Data activity is configured to **copy the data** from each table into the bronze layer (**source**) to the data lake (**sink**). The data will be saved as **.parquet**
● Parameters for schema name and table name are dynamically passed to the datasets and activities. However, if the schema does not match you can change the Json file manually. I have added instructions on how to do that in the doc **From_blob_storage_to_Azure_SQL_db.**
● A table is created for the destination in the Azure Data Lake Storage Gen2, with parameters for folder name and file name.
● Folder and file names are dynamically constructed using schema name, table name.

**Souce (SQL database):**



**Sink (Storage Gen2):**



## 2. Databricks data processing
- Next steps involve populating the data into a Databricks notebook for further processing.
- A new Databricks workspace is created within Azure.

- A secret key is set up in Databricks using Azure Key Vault to securely store access key, enabling secure access to Azure resources from Databricks notebooks.
- Next, in Databricks create a **'Compute'** for data processing, then mount the silver and gold storage accounts and confirm their accessibility.
- Create the Compute and specify the configuration needed [for this example a single node, 12.2 LTS (includes Apache Spark 3.3.2, Scala 2.12) and Standard_DS3_v2 were selected for processing].
- Create a notebook and mount to the storage to execute scripts.

```python
bronze_mount_point = '/mnt/bronze'
silver_mount_point = '/mnt/silver'
gold_mount_point = '/mnt/gold'

# Check if bronze mount point is already mounted
if not any(mount.mountPoint == bronze_mount_point for mount in dbutils.fs.mounts()):
    dbutils.fs.mount(
        source='wasbs://bronze@envsensorstorageacct.blob.core.windows.net',
        mount_point=bronze_mount_point,
        extra_configs={
            'fs.azure.account.key.envsensorstorageacct.blob.core.windows.net': dbutils.secrets.get('databricksScope', 'storageAccountKey')
        }
    )

# Check if silver mount point is already mounted
if not any(mount.mountPoint == silver_mount_point for mount in dbutils.fs.mounts()):
    dbutils.fs.mount(
        source='wasbs://silver@envsensorstorageacct.blob.core.windows.net',
        mount_point=silver_mount_point,
        extra_configs={
            'fs.azure.account.key.envsensorstorageacct.blob.core.windows.net': dbutils.secrets.get('databricksScope', 'storageAccountKey')
        }
    )

# Check if gold mount point is already mounted
if not any(mount.mountPoint == gold_mount_point for mount in dbutils.fs.mounts()):
    dbutils.fs.mount(
        source='wasbs://gold@envsensorstorageacct.blob.core.windows.net',
        mount_point=gold_mount_point,
        extra_configs={
            'fs.azure.account.key.envsensorstorageacct.blob.core.windows.net': dbutils.secrets.get('databricksScope', 'storageAccountKey')
        }
    )
```
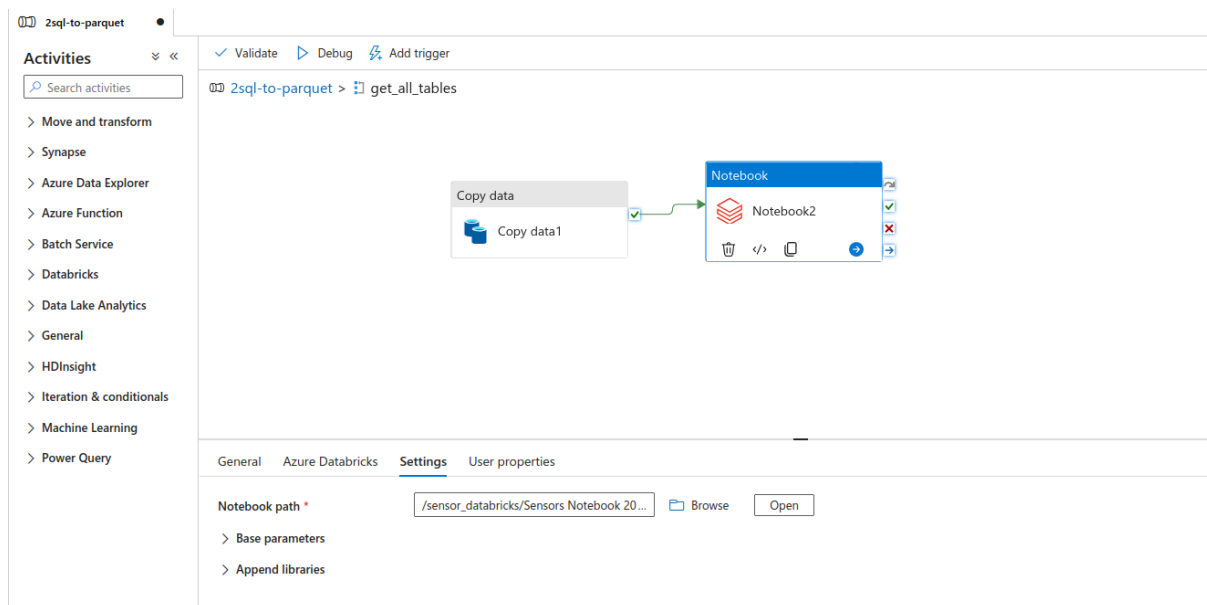
Establish connections between Data Factory and Databricks, the notebook can be added inside the pipeline.

**Steps:**
- A pipeline named "get-all-tables" is created, consisting of activities to get table names and loop through each table for further processing.
- A Lookup activity retrieves the table names using the dataset created earlier, and a ForEach activity is added to iterate through each table.
- Within the ForEach loop, a Copy Data activity is added for each table to copy data from the SQL database into the Bronze layer in Azure Data Lake Storage Gen2.
- Configure the settings such as source and sink tables to Copy Data and ensure the data is copied accurately.
- The pipeline is structured to adhere to the Medallion architecture's Bronze-Silver-Gold layers, with data flowing from the Bronze layer (SQL database) to the Silver layer (Azure Data Lake Storage Gen2) for further processing.

Once the pipeline execution is successful, verify if the data is correctly copied into the bronze layer and if the databases and tables are created as expected in the Azure Databricks environment.



**3. Processing the data in dbt-Databricks**
- Transitioning to the **dbt** (Data Build Tool) for data transformation, a new project is created.
- The project name is set as "de-proj-az" with Python 3.9.
- Dependencies including Databricks and dbt are installed with "pip install" commands.
- Configuration of Databricks CLI is done to establish a connection with Databricks on the cloud.
- Testing the Databricks connection and file system listing ensures successful connection.

```
● (toxin) river@river:~/de-proj-az/dbt_de$ databricks secrets list-scope
  s
  Scope              Backend          KeyVault URL
  -----------------  ---------------  -----------------------------------
  ---
  databricks-vscode  AZURE_KEYVAULT   https://de-project-az.vault.azure.n
  et/
  databricksscope    AZURE_KEYVAULT   https://de-project-az.vault.azure.n
  et/
  dbt-databrics      AZURE_KEYVAULT   https://de-project-az.vault.azure.n
  et/
● (toxin) river@river:~/de-proj-az/dbt_de$ databricks fs ls
  Volume
  Volumes
  databricks-datasets
  databricks-results
  mnt
  user
  volume
  volumes
○ (toxin) river@river:~/de-proj-az/dbt_de$ █
```

🗗 bash
▷ Python dbt.

- dbt initialization is performed with "dbt init" command, setting project name, database, Databricks host, access token, and default schema.
- Created SQL files for snapshots of various tables: data schema [ ]
- Added YAML configuration files for recognizing source tables and schemas from Azure Databricks.
- Executed the dbt snapshot command after resolving dependencies, generating snapshots for all tables.
- The dbt system keeps track of historical changes in data, including the last valid date before changes.
- Data snapshots are stored in the silver section of the storage account, providing a means to replay data if necessary.
- Transformation of data occurs in the bronze and silver sections, while the final output will be stored in the gold layer .
- SQL scripts are used for transformations, selecting data from snapshots and performing necessary operations.
- The dbt documentation displays the structure of the transformations, including source tables, code, and dependencies, aiding in understanding the data transformation pipeline.

```
      . . . e  which  do  not  apply  to  any  resources.
There  are  1  unused  configuration  paths:
- models.dbt_de.example
17:29:53  Found  3  snapshots,  1  model,  3  sources,  0  exposures,  0  metric
s,  538  macros,  0  groups,  0  semantic  models
17:29:53
17:29:55  Concurrency:  1  threads  (target='dev')
17:29:55
17:29:55  1  of  1  START  sql  table  model  sensor.merged_data  ...........
................. [RUN]
17:29:59  1  of  1  OK  created  sql  table  model  sensor.merged_data  .......
................. [OK  in  4.81s]
17:30:00
17:30:00  Finished  running  1  table  model  in  0  hours  0  minutes  and  6.94
 seconds  (6.94s).
17:30:00
17:30:00  Completed  successfully
17:30:00
17:30:00  Done.  PASS=1  WARN=0  ERROR=0  SKIP=0  TOTAL=1
(toxin)  river@river:~/de-proj-az/dbt_de$
```
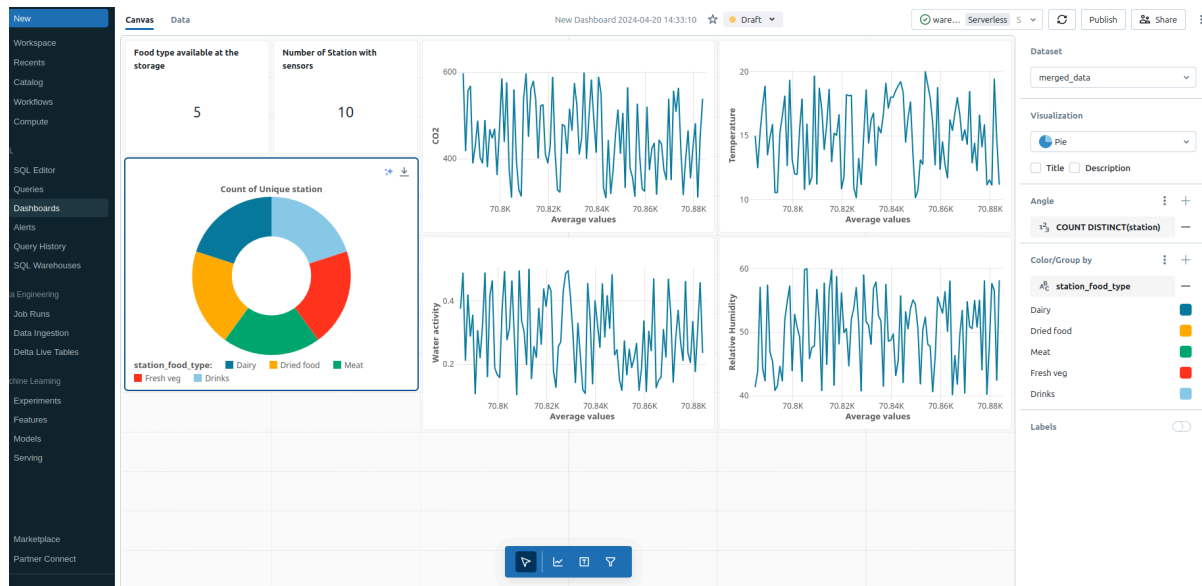
- The structure of the data on the database level mirrors what is seen in the dbt project, with tables like "merged_data" and "snapshot" will be added to Azure database containers (bronze, silver and gold layers) as well as into Databricks storage for further use.

**4. Visualisation dashboard**

Dashboards can be created in Databricks or in connection with Looker



# Summary

This data engineering project was successfully completed as part of the ZoomCamp course. I developed this project on Azure implementing the Medallion architecture. The infrastructure was created with Terraform, and I leveraged Azure Data Factory, Azure Databricks, and Azure Data Lake Storage Gen2 to transform the data efficiently.

The main objective of this project was to streamline data processing and analysis in a cloud environment, enabling efficient decision-making processes. By using Azure services and implementing the Medallion architecture, scalability and flexibility can be ensured. The use of synthetic datasets aided in risk assessment, while pipelines facilitated smooth data transfer.

I use dbt-Databricks for data transformation and analysis, integrating SQL files to capture table snapshots and gain valuable insights. The transformed data was then stored in the appropriate layer (bronze, silver, or gold) based on the processing stage.
Furthermore, a dashboard was created using Databricks and Looker for data visualisation.

Overall, this project exemplified effective data engineering practices within the Azure environment, showcasing the power of cloud-based solutions for complex data projects.