# Retail Store Reporting Using Azure Data Engineering

## Project Overview

This project is designed to generate **Revenue and Stock Wastage Reports** for an **FMCG Retail Company**. The **Revenue Reports** provide insights on a **monthly basis**, categorized by **product and city**. Additionally, the reporting includes **stock wastage analysis**, calculating the **financial loss** incurred due to different suppliers, which helps in determining the **supplier rating**.

The data pipeline is built using **Azure Data Factory (ADF)**, **Azure Data Lake Storage (ADLS)**, and **Azure Databricks Unity Catalog** for efficient data processing and management.

## Requirement

The following reports need to be developed,

- Monthly Revenue for past one year.
- Top 5 City Revenue distribution.
- Revenue by different product category.
- Loss due to Stocky Wastage by different Supplier – The stock wastage is calculated by calculating the difference between the current day opening stock count and previous day closing stock count.
- Wastage over different days of the weekdays.
- List of people eligible for loyalty program – A loyalty score is calculated identifying the customer with most number of high value transactions over the period of one year

## Tech Stack

- **Azure Data Factory (ADF)** – Data ingestion and orchestration.

- **Azure Data Lake Storage (ADLS)** – Raw, Bronze and Databricks unity catalog metastore.

- **Azure Databricks** – Data transformation (Silver Layer) and Aggregated data for reporting in Unity Catalog (Gold Layer).

- **Power BI** – Visualization and reporting.

## Source Datasets Used:

## Dimension Datasets:

1. **Customer:**

| Customer | Name | Age | Gender | City | SignupDate | LastUpdated |
|---|---|---|---|---|---|---|
| 1 | Ashley Hall | 34 | Other | Port Victoria | 05-06-2022 | 04-11-2023 |
| 2 | Adrian Blankenship | 43 | Male | South Brent | 01-03-2022 | 08-06-2024 |
| 3 | Elizabeth Stanley | 63 | Female | Tamaraton | 22-11-2023 | 15-08-2024 |
| 4 | Shannon Jackson | 73 | Male | South Markbury | 28-12-2024 | 21-01-2025 |

## 2. Product

| ProductID | Name | Category | Price | SupplierID | LastUpdated |
|---|---|---|---|---|---|
| 1 | Practice Cultural | Home | 233.53 | 98 | 15-03-2024 |
| 2 | Blood Leg | Clothing | 515.13 | 43 | 10-09-2023 |
| 3 | American Phone | Clothing | 22.42 | 84 | 27-06-2023 |

## 3. Store

| StoreID | Name | City | Region | LastUpdated |
|---|---|---|---|---|
| 1 | Store Black Inc | Smithshire | North | 24-07-2023 |
| 2 | Store Nicholson Ltd | Lake Kathyfort | South | 14-12-2022 |
| 3 | Store Reese Ltd | Elaineton | West | 09-08-2022 |
| 4 | Store Collier PLC | North Sarah | West | 14-12-2023 |

## 4. Supplier

| SupplierID | Name | ContactInfo | Rating | LastUpdated |
|---|---|---|---|---|
| 1 | May Ltd | +1-909-912-6461x15282 | 1.2 | 01-01-2025 |
| 2 | Walls, Bush and Davis | (712)791-4202x302 | 2.1 | 15-12-2023 |
| 3 | Mann LLC | 9972320647 | 2.6 | 10-03-2024 |

## 5. Time

| Date | Year | Quarter | Month | Weekday | DayOfWeek | HolidayFlag |
|---|---|---|---|---|---|---|
| 01-01-2023 | 2023 | 1 | 1 | Sunday | 6 | TRUE |
| 02-01-2023 | 2023 | 1 | 1 | Monday | 0 | FALSE |
| 03-01-2023 | 2023 | 1 | 1 | Tuesday | 1 | FALSE |
| 04-01-2023 | 2023 | 1 | 1 | Wednesday | 2 | FALSE |
| 05-01-2023 | 2023 | 1 | 1 | Thursday | 3 | FALSE |

# Transactional Datasets

## 1. Sales Transactions

| TransactionID | Date | CustomerID | StoreID | ProductID | Quantity | TotalAmount | PaymentMethod |
|---|---|---|---|---|---|---|---|
| ea471e57-4a4e-4605-bae3-a588bbf12cef | 01-01-2025 | 4131 | 90 | 380 | 4 | 334.28 | Debit Card |
| 2031c1d7-f54f-4550-96bd-7ddcdc278d74 | 01-01-2025 | 2725 | 41 | 138 | 1 | 470.72 | Credit Card |
| 6c75c49d-fd42-40cb-a7ac-048f93dfa030 | 01-01-2025 | 2329 | 43 | 366 | 3 | 116.38 | Cash |
| e34712dd-5452-40ab-9397-eff60ce16556 | 01-01-2025 | 3962 | 66 | 365 | 9 | 389.99 | Online Payment |

## 2. Stock Movement

| StockID | Date | StoreID | ProductID | OpeningStock | ReceivedStock | SoldStock | ClosingStock |
|---|---|---|---|---|---|---|---|
| 62eeb7d6-906c-45ab-a475-0754439b619f | 01-01-2025 | 50 | 313 | 194 | 48 | 44 | 198 |
| 9a3cabf8-4abc-41ca-a4ec-a444ab78f548 | 01-01-2025 | 68 | 450 | 188 | 60 | 20 | 228 |
| 5760ba71-1735-4de7-aef7-653e43e61d0f | 01-01-2025 | 11 | 376 | 338 | 62 | 29 | 371 |
| 0ad0103d-7fdb-4514-94ce-87aa6b63075c | 01-01-2025 | 64 | 384 | 223 | 48 | 28 | 243 |
| 035431ac-b003-40c6-8cd5-f9fb42163344 | 01-01-2025 | 61 | 154 | 447 | 62 | 18 | 491 |

# Data Model



# Architecture



# Detailed Implementation:

**Azure Data Lake Storage:**

The below containers are created for the respective usage,

- Source – to store the raw csv source files.

NOTE: This action is implemented to mimic the data ingestion from source to ADLS

- Bronze – to store the parquet file ingested by ADF.

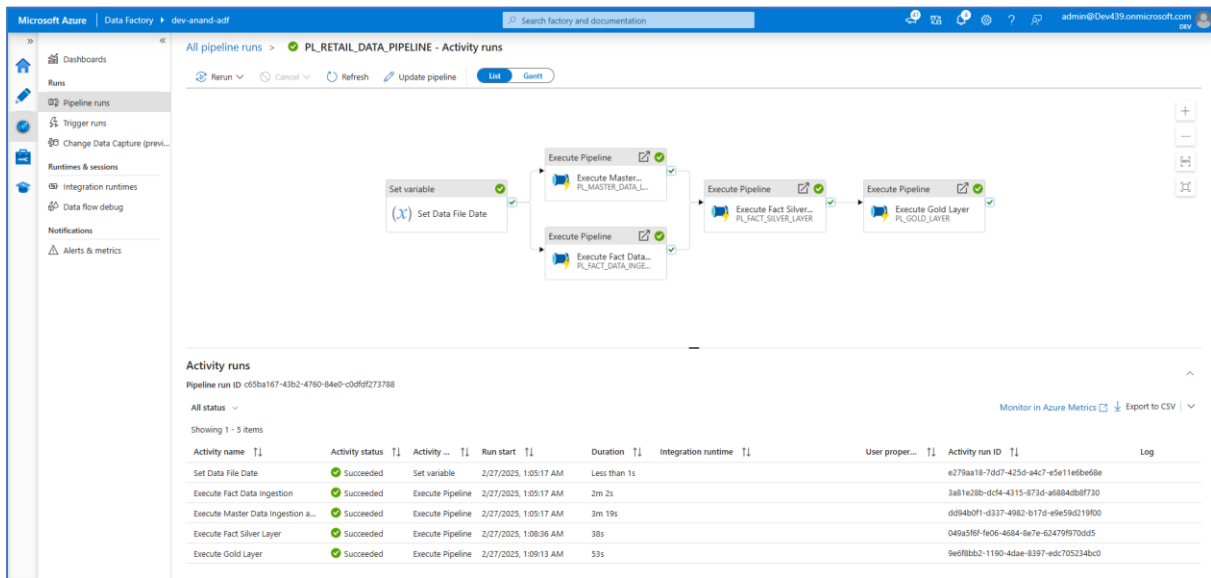The transaction data files are stored in YYYY/MM/DD date folder format.

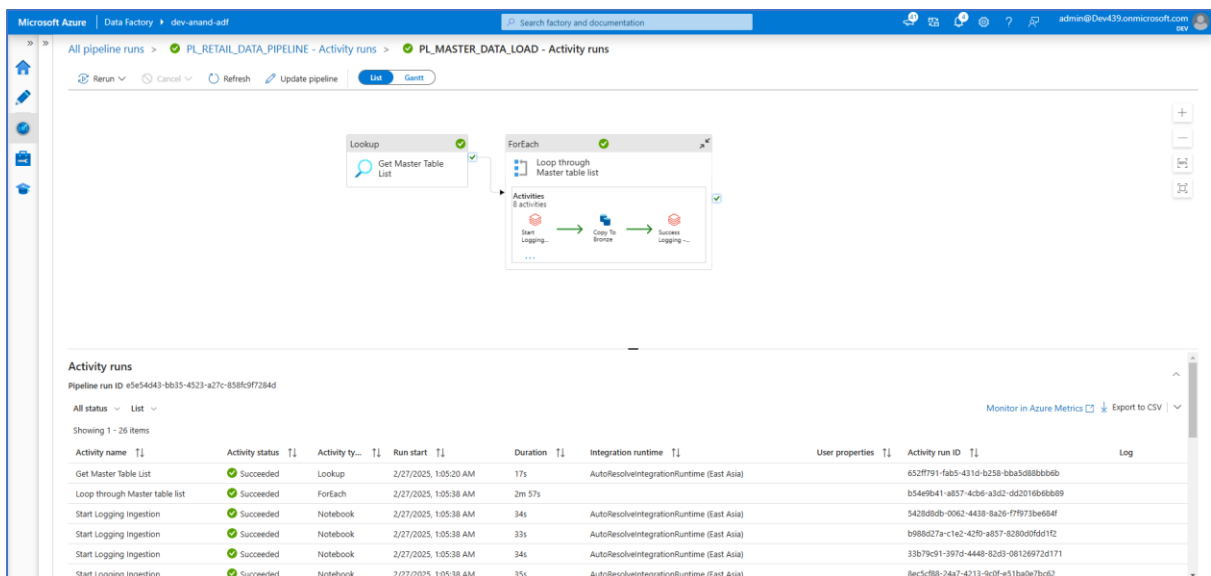- **devanand-metastore** – Metastore for the Databricks to create the unity catalog.



**Azure Data Factory:**

The ADF contains a main pipeline which starts with setting the load date (current date -1 day). It then executes Dimension data and Fact Data Ingestion and Curation pipelines.

PL_RETAIL_DATA_PIPELINE (main pipeline):



PL_MASTER_DATA_LOAD: Ingests all the dimension data on Incremental basis.

Time dimension is populated programmatically.

1. The Data from the source is first ingested in the bronze layer.
2. The Silver Layer notebook will read the data, performs cleanups and load the data to the silver layer tables on SCD2 basis.

The above two steps are run in a loop with all the entities passed dynamically to the for each activity using the configuration table.



The pipeline also performs the necessary logging. The Silver layer notebook is also developed in the form of framework to process each entity dynamically which will be discussed below in the Databricks sections.

PL_FACT_DATA_INGESTION: The Fact table load is also made dynamic. The fact tables are configured in the above configuration table. The load is incremental as the data at the source are stored in YYYY/MM/DD folder format



Incremental Logic:

PL_FACT_SILVER_LAYER – Then the fact tables are processed and loaded to the silver tables



PL_GOLD_LAYER: The final pipeline will aggregate the data as per report requirement and load the data on full load basis to the gold layer tables.



**Databricks:**

**Utils:**

- Pipeline Logging Start - This notebook is called before the copy/notebook activity to start the logging.
- Pipeline Logging End – This notebook is called once the copy/notebook activity is completed to update the status of the activity.

The log data is stored in the below table,



**Silver Layer:**

- Dimension Tables – This notebook dynamically accepts the entity and the respective entity's key column and mutable columns and loads the data into the silver layer tables in SCD2 mode.

**Dimension dynamic notebook:**



- Sales Transactions – The sales transaction fact table is processed in this notebook. The highest transaction for each day is identified here to calculate the customer eligibility score for loyalty programs
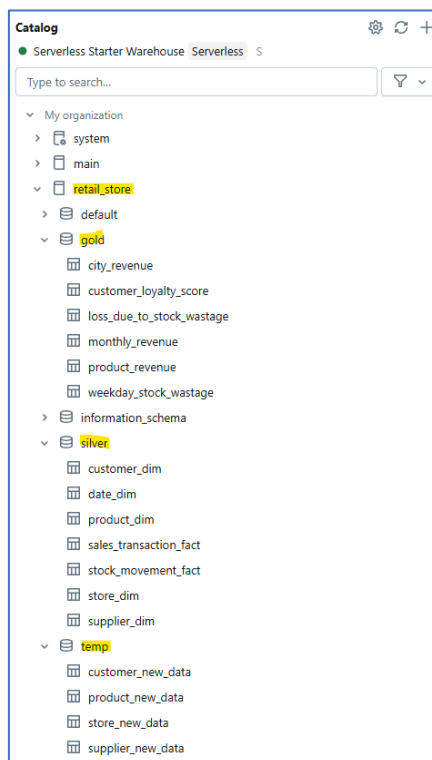
- Stock Movement – This Stock Movement fact table is processed to calculate the wastage on each day for each product.
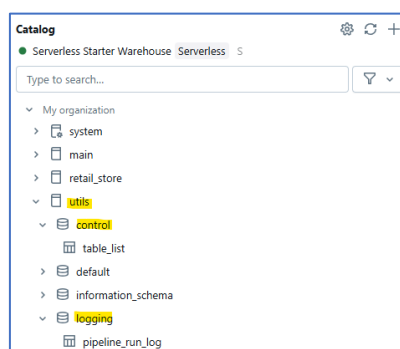
**Gold Layer:**

- Revenue Analysis – This notebook calculates the Monthly revenue, Revenue by product category and revenue by cities
- Wastage Analysis – This notebook calculates the wastage on a daily basis - the difference between the current day opening stock count and previous day closing stock count
- Customer Loyalty program eligibility – This notebook calculates the score for each customers and picks the top 10 for loyalty programs.

**Unity Catalog:**

- The Retail Store Catalog has the silver schema, gold schema and temp schema (required for curation process)
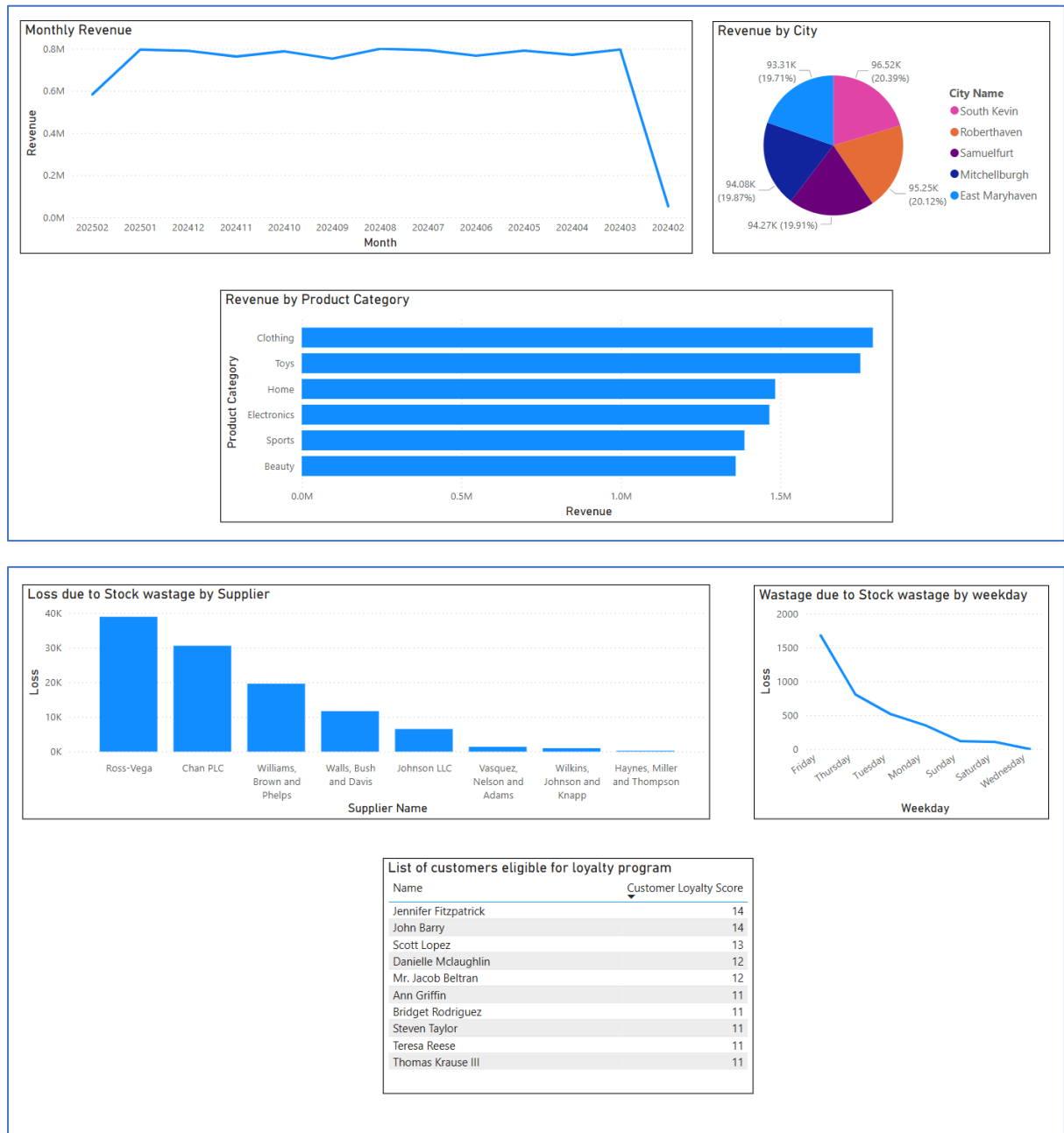


Then there is utils catalog which has logging schema and control schema (for configuration table).
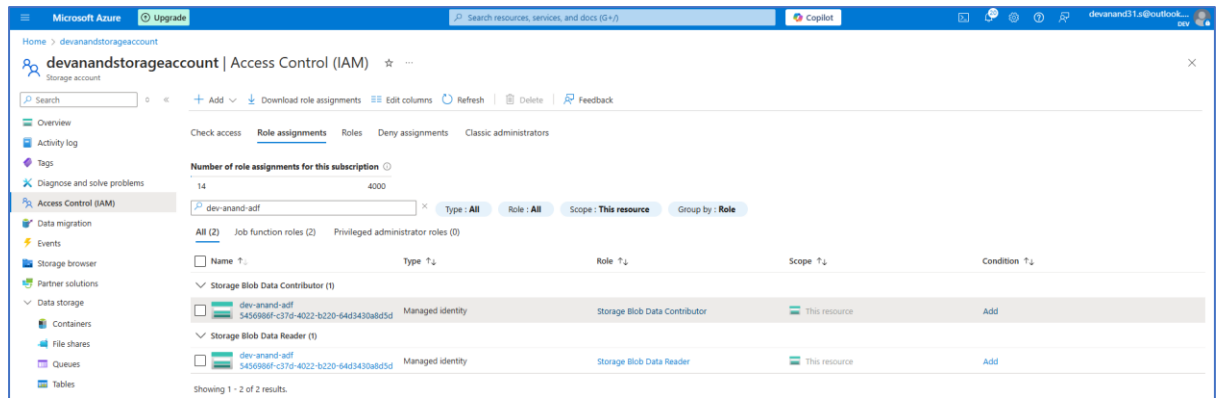
**Power BI:**

Finally, the Power BI is connected to the gold schema to fetch all the table data for the below reports.
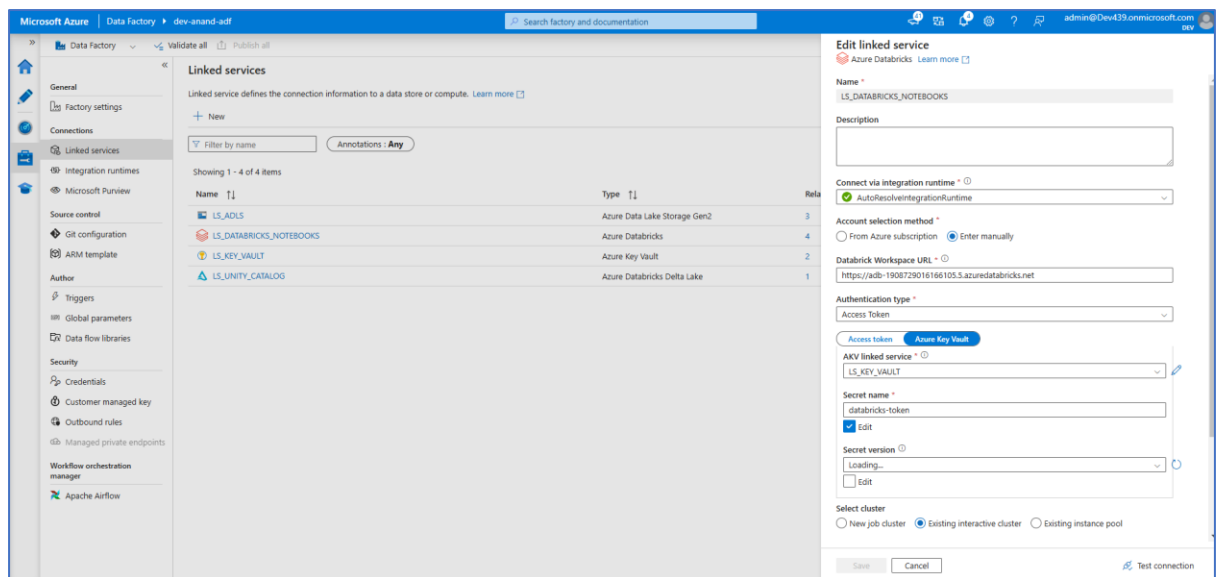




## Security:

### ADF <-> ADLS:

The Connection for the ADF to read and write files from and to ADLS is performed using managed identity with below access roles.
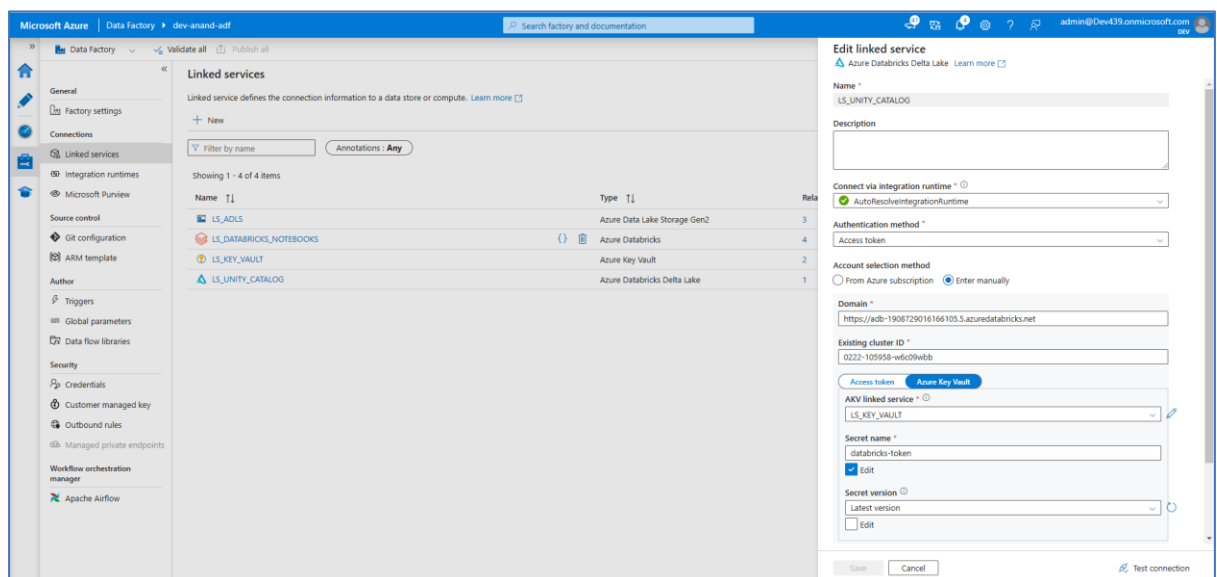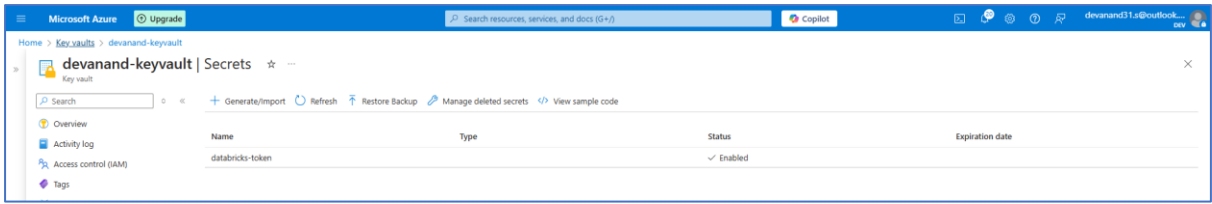
**ADF <-> Databricks:**

The Connection for the ADF to Databricks to execute the notebook is created using the Access Token stored in the Key vault.


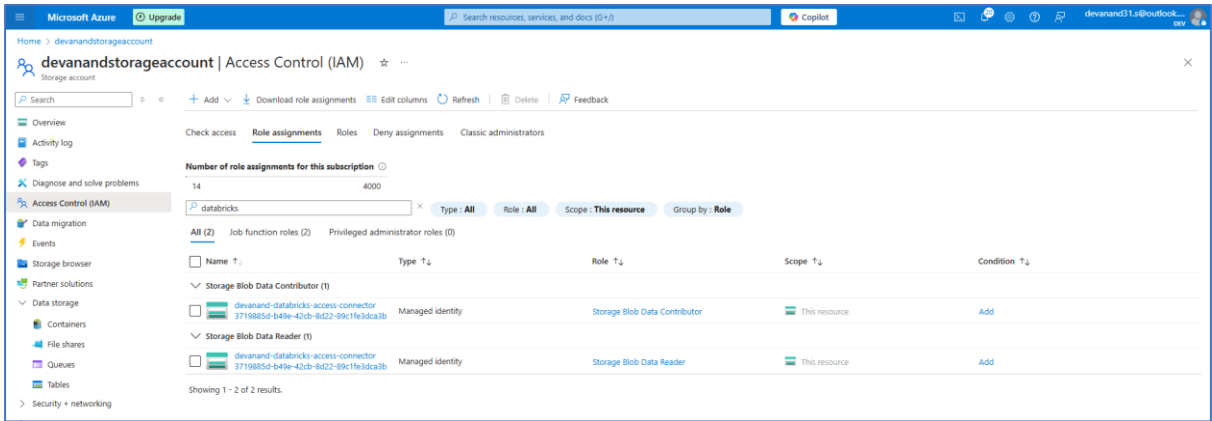
The connection to the unity catalog,
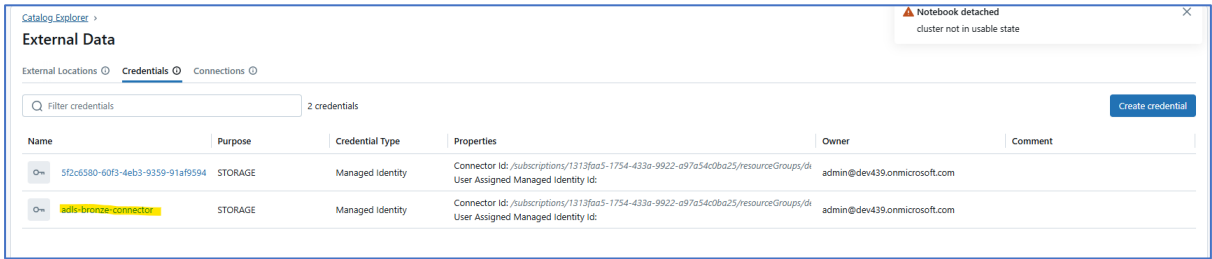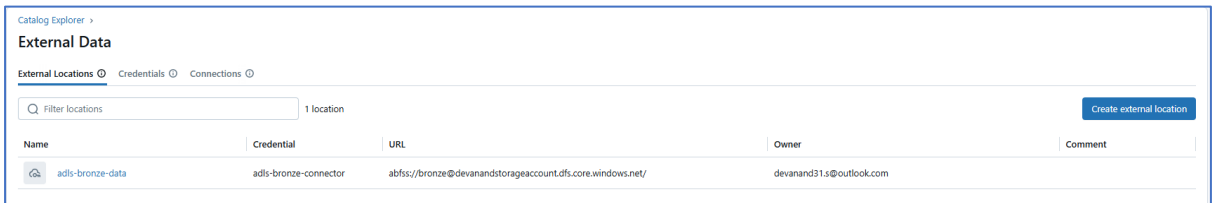
Key vault secrets,



## Databricks <-> ADLS:

The External Storage is created in the Databricks using Databricks Access connector which in-turn is authorized by managed identity.



Databricks credential created in the catalog,



External location access the ADLS via databricks access connector,



## Power BI <-> Databricks unity catalog:

The PowerBI access the unity catalog gold schema using the Access Tokens.