# Azure Data Engineering Project
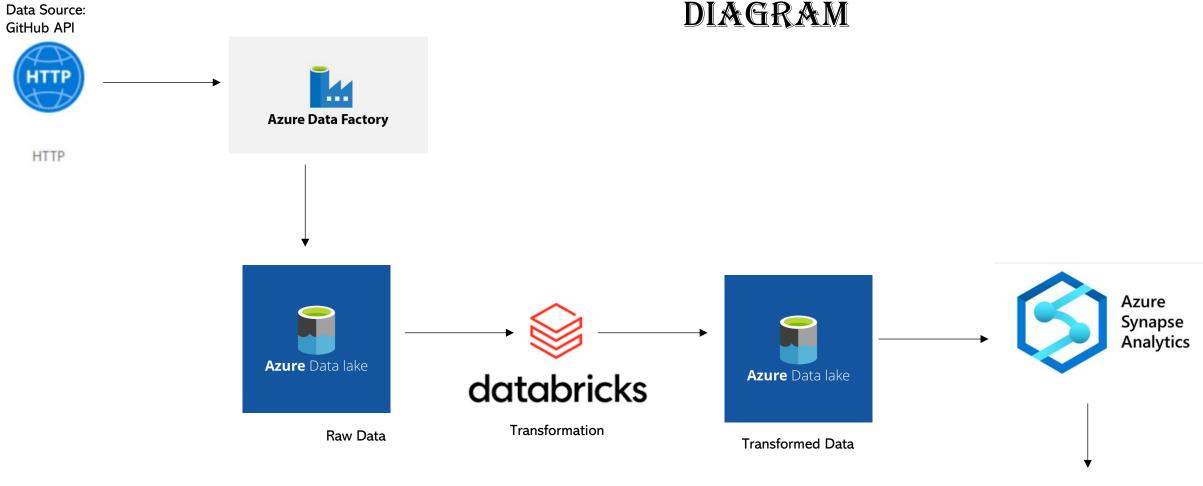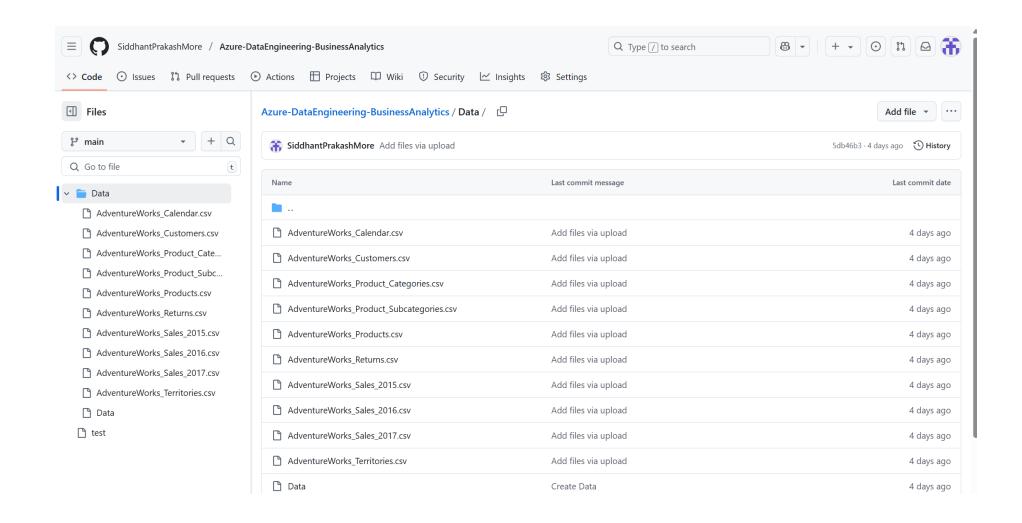
~ Siddhant Prakash More

# Architecture Diagram
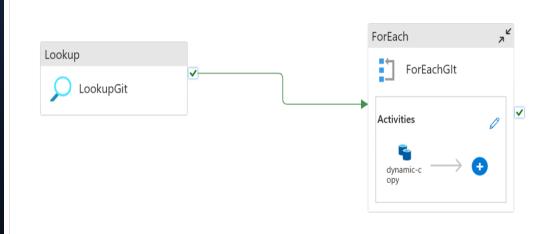
Data Source:
GitHub API
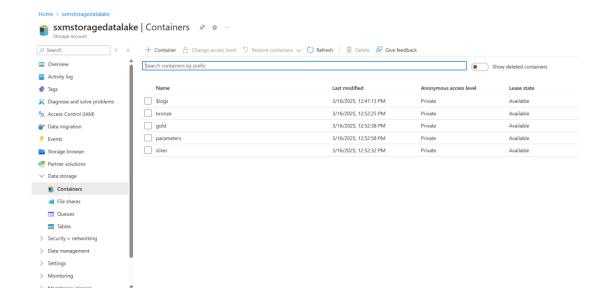
HTTP

**Azure Data Factory**

Azure Data lake

Raw Data

databricks

Transformation

Azure Data lake

Transformed Data

Azure Synapse Analytics

Power BI

# DATA SOURCE

# DATA INGESTION:

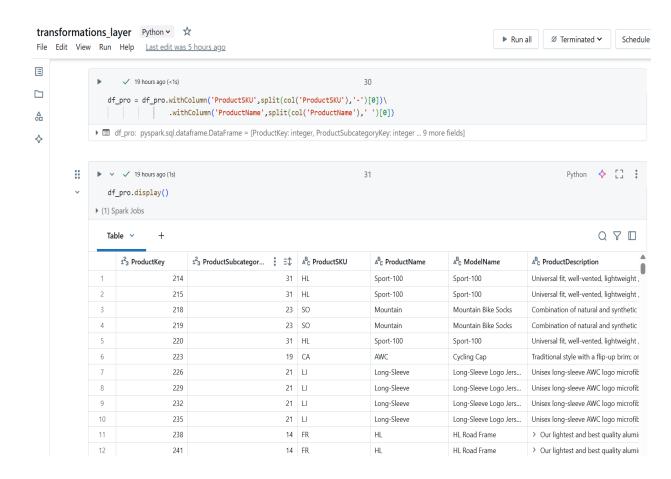Using **Azure Data Factory (ADF)** to fetch data from the **GitHub API** via HTTP

# STORAGE: AZURE DATA LAKE

- The storage architecture follows a **medallion design** using **Azure Data Lake**, ensuring efficient data management across three layers. The **bronze layer** stores raw, unprocessed data directly ingested from the source. The **silver layer** contains cleaned and transformed data, making it structured and ready for analysis. The **gold layer** holds the final, optimized dataset, designed for efficient querying and reporting in **Azure Synapse Analytics** and **Power BI**.

# DATA TRANSFORMATION: DATABRICKS

- Used **Databricks with PySpark** for data transformation, processing raw data from the **bronze layer** and refining it into the **silver layer**. The transformation included data cleaning, structuring, and enrichment to ensure high-quality, usable data. This process optimized the data for further analysis and integration with downstream systems.

# Analytics & BI:

- **Azure Synapse Analytics** – Used for running analytical queries on transformed data and integrating with BI tools.

**Power BI** – Utilized for creating dashboards and visual reports based on the processed data.