# Data Architecture Design for Global Superstore Dataset

DAMG7370 Final Project - Group 6

Team Members:
Nikhila Devi Maddala 002878220
Sai Puja Kamatam 002818267
Likhitha Velagapudi 002727757

# Roles and Responsibilities

Nikhila Devi Maddala - Azure till silver stage

Sai Pujitha Kamatam - Azure till gold stage, connecting to Power BI

Likhitha Velagapudi - Power BI

# Objective

**Project** **Objective:**
**To** build a scalable, end-to-end data pipeline on Microsoft Azure to ingest, transform, and visualize the Global Superstore dataset, delivering actionable retail insights through interactive dashboards.

**Project Summary:**

**Dataset**: Global Superstore (sales, customers, products, returns) from Kaggle.

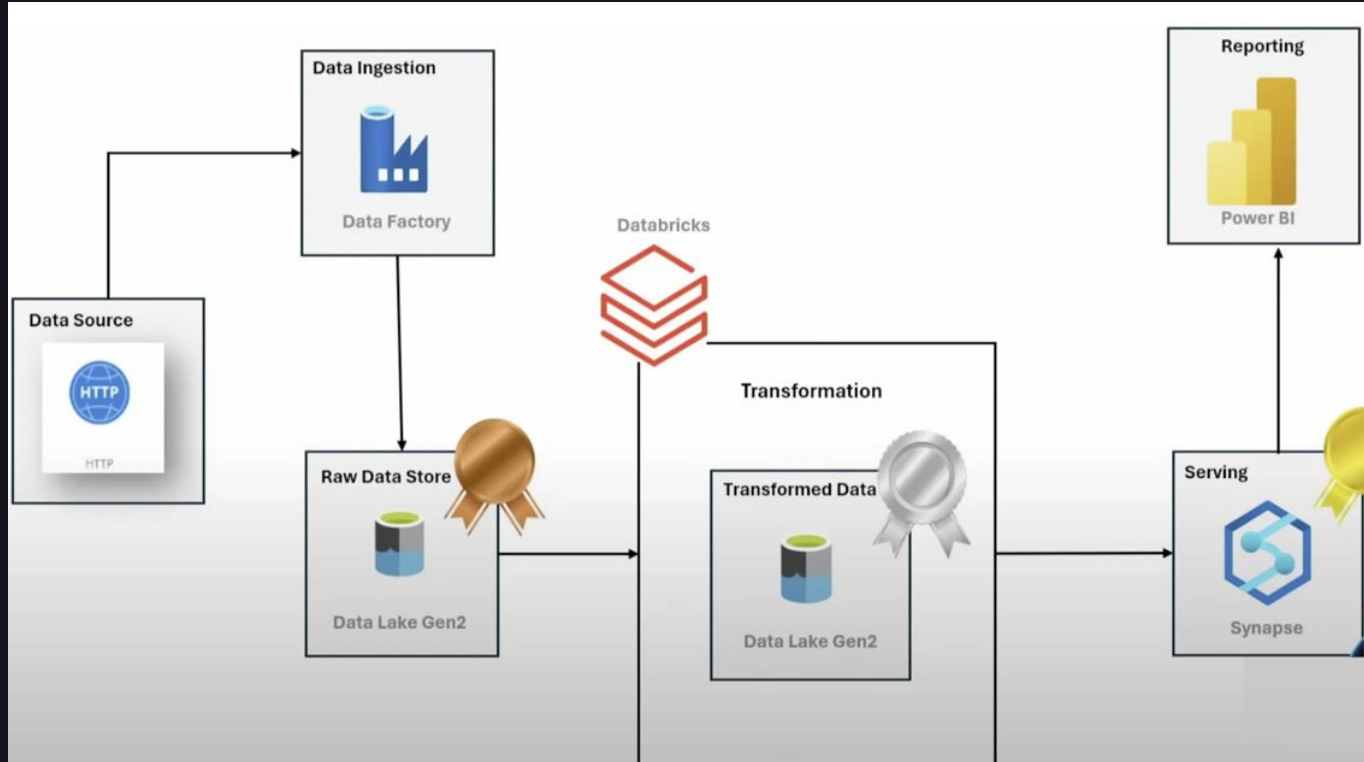**Architecture**: Leverages Azure Blob Storage, Data Factory, Databricks, Synapse Analytics, and Power BI.

**Process**:

- **Ingestion**: Raw data to Bronze layer via Data Factory.
- **Transformation**: Advanced processing in Databricks to Silver layer.
- **Analytics**: Gold layer views in Synapse for optimized querying.
- **Visualization**: Interactive Power BI dashboards for sales and customer insights.

**Outcome**: Efficient, scalable pipeline enabling real-time retail analytics.

# Tools Used in the Global Superstore Data Pipeline

- **Azure Blob Storage / Data Lake Storage Gen2**
  - Stored raw (Bronze), transformed (Silver), and analytics-ready (Gold) data.
  - Enabled hierarchical namespace for efficient file management.
- **Azure Data Factory**
  - Orchestrated data ingestion from Kaggle to Bronze layer.
  - Used static and dynamic pipelines for scalable data copying.
- **Azure Databricks**
  - Performed advanced data transformations (e.g., calculating Shipping Delay Days).
  - Processed data into Parquet files for the Silver layer.
- **Azure Synapse Analytics**
  - Created Gold layer views (e.g., sales, customers) for optimized querying.
  - Supported SQL-based analytics and external tables.
- **Power BI**
  - Built interactive dashboards for sales and customer insights.
  - Connected to Synapse via DirectQuery for real-time visualization.

# Data Architecture

# Data Source – Global Superstore Dataset

**Content**:

- **Source**: Kaggle Dataset
  - URL: https://www.kaggle.com/datasets/apoorvaappz/global-super-store-dataset
- **Format**: Multiple CSV files
  - Files: Orders, Returns, People
- **Size**: 12.09 MB
- **Upload Method**:
  - Downloaded locally from Kaggle then uploaded to github ..
  - Uploaded to Azure Blob Storage (Bronze container) via Azure Data Factory.
- **Description**: Contains sales, customer, product, and returns data for a global retail business, ideal for analyzing sales trends and operational insights.

## Data Ingestion Layer

- **Tool**: Azure Data Factory
  - Instance: neu-aw-project (Version 2, East US region).
- **Source**: Kaggle Dataset (Global Superstore)
  - Files: Orders, Returns, People (CSV, 12.09 MB).
  - Downloaded locally for upload.
- **Destination**: Azure Blob Storage (Bronze Container)
  - Storage Account: neudamgdatalake.
  - Stored raw, untouched data for downstream processing.
- **Process**:
  - Created a static pipeline to copy CSV files to Bronze container.
  - Configured dynamic pipelines for scalability (iterative file copying).
- **Outcome**: Reliable, scalable ingestion of raw retail data into the Bronze layer.

# Data Cleaning and Transformation- Azure Databricks



- **Tool**: Azure Databricks
  - Workspace configured in resource group DAMG7370.
  - Processed data using Spark DataFrames.
- **Input**: Bronze Container (Azure Blob Storage)
  - CSV files: Orders, Returns, People (12.09 MB).
- **Transformations**:
  - **Orders**: Added Order Month, Order Year from Order Date; calculated Shipping Delay Days (Ship Date - Order Date).
  - **Returns**: Standardized Return Reason for consistency.
  - **People**: Unified Region names (e.g., "US-East" to "East US").
- **Output**: Silver Container
  - Saved transformed DataFrames as Parquet files for optimized analytics.
- **Outcome**: Clean, structured data ready for analytics in the Gold layer

# Data Cleaning and Transformation- Azure Databricks

- Raw CSV data loaded into PySpark DataFrames from bronze container.
- Silver transformations: added Order Month, Order Year, Shipping Delay Days; standardized Market, region; validated Order ID.
- Transformed data saved as optimized Parquet in the silver container.
- Gold layer views created: calendar, customers, products, returns, sales, subcat, territories.
- Ensures schema consistency and query performance for analysis.

# Analytics and Query Layer - Azure Synapse

Tools used: Azure Synapse Serverless SQL

Data Source: External data source to silver container

Views Created:

- Gold.orders
- Gold.sales

# Analytics and Query Layer - Azure Synapse

- Utilized Managed Identity for secure data access
- SQL queries can be executed directly on silver or gold layer data without needing to copy or move it.
- Azure Synapse supports high-performance operations like joins, aggregations, and filtering on large datasets.

# Power BI Dashboard

Visualizations created:

THANK YOU