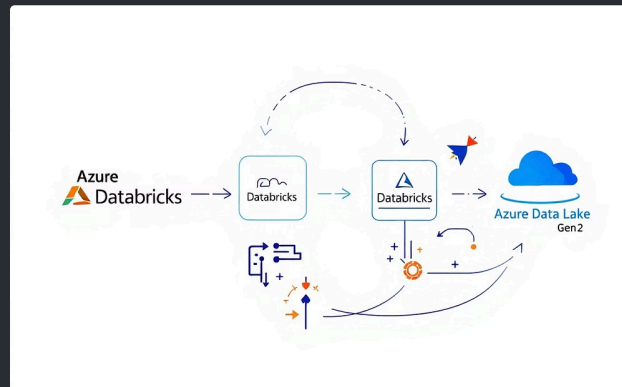# Olympics Data Analytics Pipeline using Azure

This project demonstrates a complete end-to-end **modern data pipeline** using Microsoft Azure.

This comprehensive solution showcases the power of Azure's cloud ecosystem for handling large-scale data analytics, from ingestion through transformation to visualization, creating a robust foundation for Olympic historical data analysis.







### Azure Cloud Foundation

Data acquisition from GitHub repositories and initial staging in Azure Storage Accounts to establish a robust and scalable data lake foundation.

### Streamlined Data Flow

Efficient data transformation using Azure Databricks for processing and refining raw data, ensuring its readiness for analytical workloads, and storing in Azure Data Lake Gen2.

### Insightful Analytics & Visualization

Advanced querying capabilities with Azure Synapse Analytics for complex data analysis, and interactive visualization through Power BI dashboards to extract actionable insights.

# Project Overview

### Data Source

Olympics dataset from Kaggle, featuring 120 years of Olympic history with athlete and results data

### Storage Layer

Azure Blob Storage and Azure Data Lake Storage Gen2 for raw and processed data

### Processing Engine

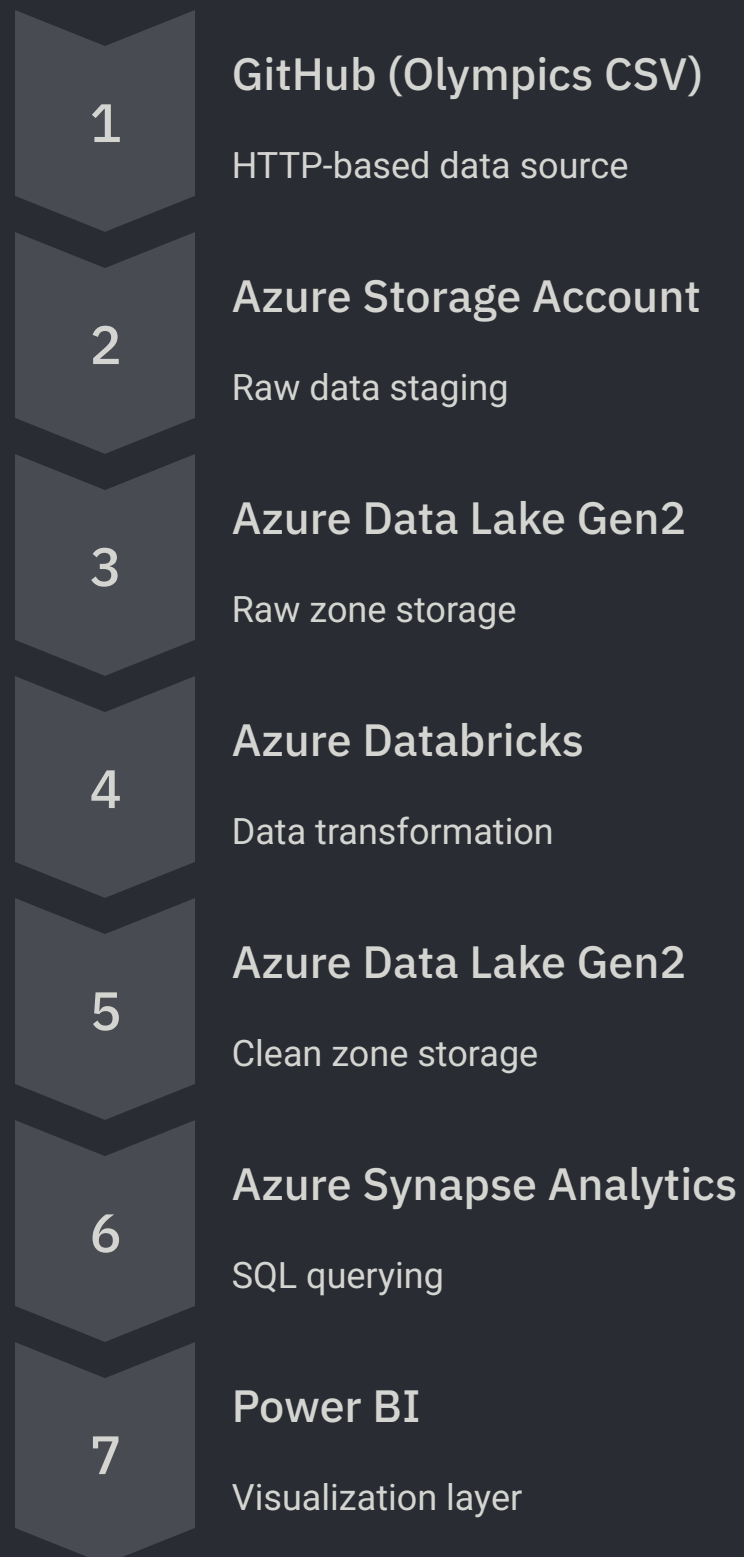Azure Databricks for data transformation, cleaning, and advanced analytics

### Analytics & Visualization

Azure Synapse Analytics for querying with Power BI dashboards for insights

**Data Hosting:** Dataset uploaded to GitHub [AkshathD2298] for HTTP-based access, enabling seamless integration with Azure services and automated data pipeline workflows.

# Architecture & Tech Stack

## 📁 Architecture Flow

**1** GitHub (Olympics CSV)

HTTP-based data source

**2** Azure Storage Account

Raw data staging

**3** Azure Data Lake Gen2

Raw zone storage

**4** Azure Databricks

Data transformation

**5** Azure Data Lake Gen2

Clean zone storage

**6** Azure Synapse Analytics

SQL querying

**7** Power BI

Visualization layer

## Technology Components

**GitHub (HTTP URL):** Hosting Olympics CSV file for HTTP-based ingestion

**Azure Blob Storage:** Initial staging of raw data

**Azure Data Lake Gen2:** Raw + Processed data storage

**Azure Databricks:** Data transformation and cleansing

**Azure Synapse Analytics:** Analytics SQL querying on clean data

**Azure Data Factory:** Orchestrating data movement and pipeline automation

**Power BI:** Dashboard and data visualization

**Python / PySpark:** Data wrangling and transformation within Databricks

# Data Pipeline Flow

01

## Data Upload

The original Olympics CSV file is hosted on GitHub for HTTP access

02

## Ingestion

Azure Data Factory fetches the CSV and loads it into an Azure Storage account

03

## Raw Zone

Data is moved into the raw zone of Azure Data Lake Gen2

04

## Transformation

Azure Databricks reads the raw data, performs cleaning (nulls, duplicates, column formatting), and writes to the clean zone

05

## Querying

Azure Synapse Analytics connects to the clean data zone for querying
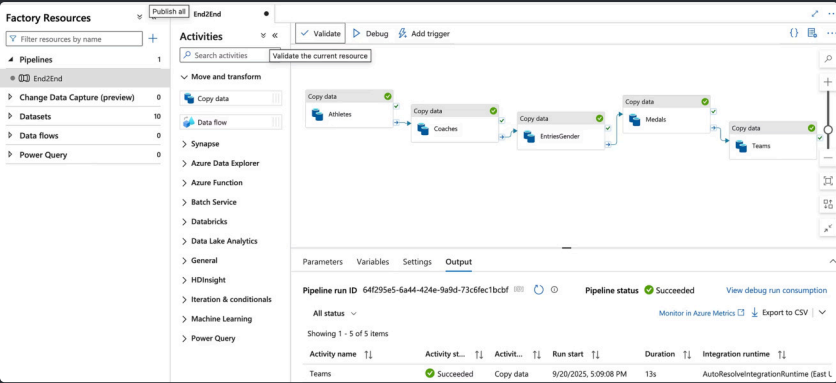
06

## Visualization

Power BI imports data from Synapse to create visual analytics and dashboards

# Project Implementation: Visual Walkthrough

This section provides a visual walkthrough of the key components and configurations within the Azure portal, Databricks notebooks, Synapse Analytics, and Power BI dashboards, demonstrating the end-to-end data pipeline in action.
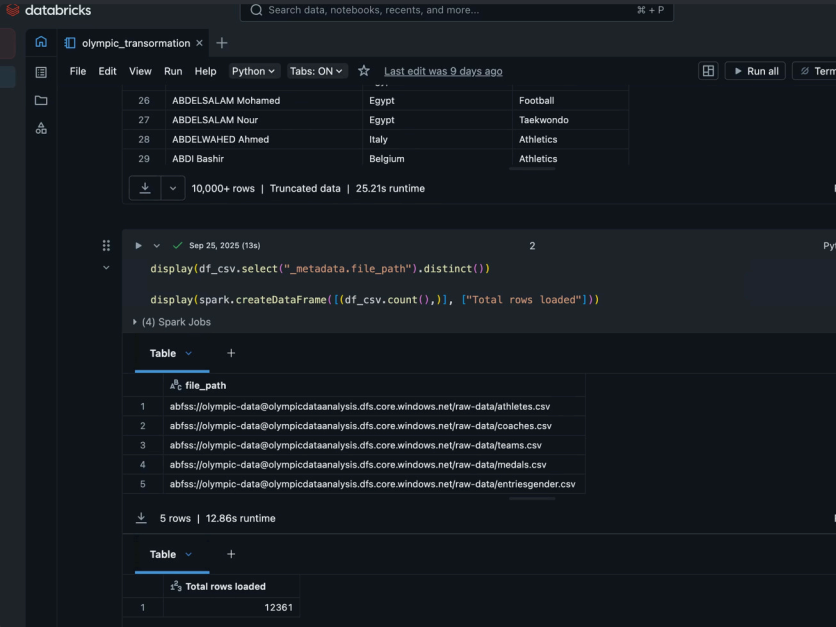
## Azure Data Factory: Ingestion Pipeline

Observe the Azure Data Factory pipeline, meticulously configured to ingest the Olympics CSV data directly from GitHub. This pipeline skillfully orchestrates data movement, ensuring efficient loading of raw data into Azure Storage for subsequent processing.
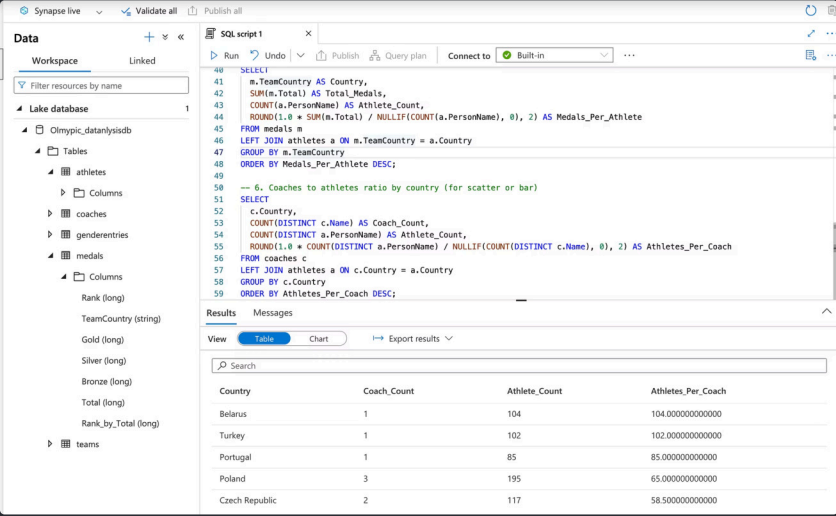


## Azure Databricks: Data Transformation Notebook

Explore this Azure Databricks notebook, where PySpark code is applied for robust data cleaning, transformation, and enrichment. Steps encompass handling null values, removing duplicates, and standardizing column formats, all prior to writing the refined data to the clean zone in Azure Data Lake Gen2.



## Azure Synapse Analytics: SQL Querying

Witness Azure Synapse Analytics in action, querying the meticulously cleaned data. This environment facilitates interactive analysis of the processed Olympics data using SQL, providing a powerful interface for data exploration and preparing it for dynamic visualizations.



## Raw vs. Transformed Data

These visuals illustrate the critical data transformation process, showcasing the initial raw Olympics dataset alongside the cleaned and transformed results. Observe the improvements in data quality, consistency, and structure after Databricks processing, ready for advanced analysis and reporting.