# Stream Processing and Analysis

## Uber Data Analytics

Submitted by

**S Yogesh Reddy    - (2022BIFT07AED026)**

*In partial fulfilment for the award of the degree of*

**BACHELOR OF TECHNOLOGY**

In

**INFORMATION TECHNOLOGY**

*Under the Supervision of*

**Dr. Chetan J Shelke**



**DEPARTMENT OF INFORMATION TECHNOLOGY**

**ALLIANCE SCHOOL OF ADVANCED COMPUTING**

**ALLIANCE UNIVERSITY, BENGALURU**

**MARCH – 2025**

**Abstract**

This project showcases the development of an end-to-end data engineering pipeline using Microsoft Azure. The objective is to design a streamlined data flow, from ingestion to visualization, enabling real-time insights into business data. By integrating services like Azure Data Factory, Azure Data Lake, Databricks, and Power BI, the pipeline automates data handling and transformation, delivering valuable business intelligence. The data used in this project is from Ola Cabs, focusing on ride analytics to uncover patterns in demand, pricing, and customer behaviour.

## INTRODUCTION

Modern businesses generate massive volumes of data, making it essential to have a robust data pipeline that can collect, transform, and visualize information. The ability to efficiently process and analyze data is critical for making informed business decisions and staying competitive. This project explores how Azure's cloud ecosystem can handle these tasks, transforming raw data into actionable insights that drive strategic decisions.

The data pipeline is designed to ingest Ola Cabs ride data from GitHub, where ride records, fare amounts, customer ratings, pickup and drop-off locations, and timestamps are stored. Using Azure Data Factory, the data is automatically extracted and loaded into Azure Data Lake Storage Gen2. This raw data is then processed using Azure Databricks, where Spark-based transformations clean, aggregate, and enrich the data to prepare it for analysis.
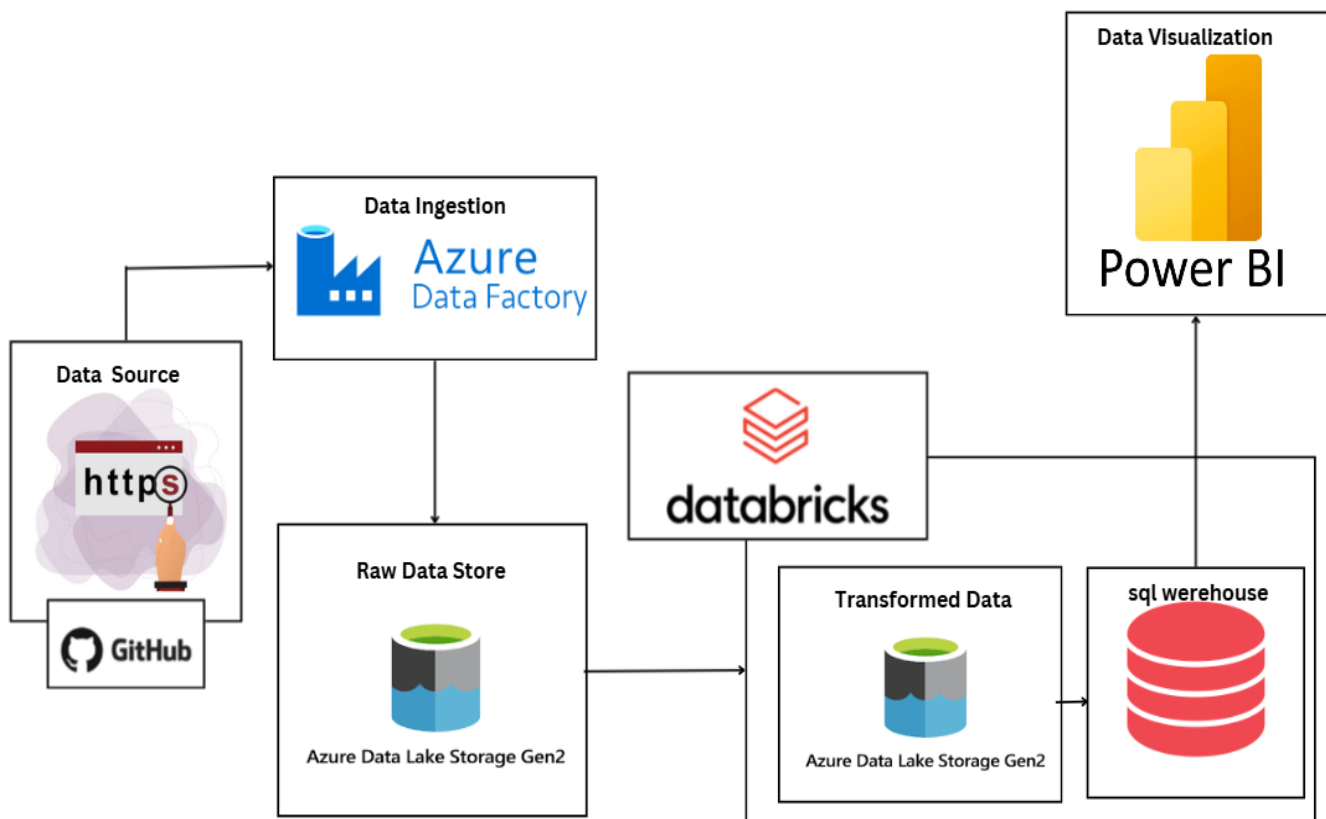
Once transformed, the data is loaded into a SQL warehouse within Databricks, which acts as a central repository for structured analytical data. The SQL warehouse enables rapid querying and complex analytics, making it easy to derive valuable insights. Power BI is connected to this warehouse to create

dynamic, interactive dashboards, allowing stakeholders to visualize trends like peak booking hours, revenue distribution across cities, and customer satisfaction metrics.

This end-to-end pipeline not only automates data ingestion and transformation but also ensures that reports are always up to date, giving business leaders access to the most current insights. By leveraging Azure's scalability and integration capabilities, the solution is well-equipped to handle increasing data volumes as the business grows, providing a future-proof foundation for data-driven decision-making.

## METHODOLOGY

**The Flowchart Diagram:**

**Step 1: Data Source**

**GitHub link for Data: https://github.com/Yogeshreddy07/OlaCabs-Data-Analytics-Azure.git**

- The data originates from Ola Cabs, accessed via GitHub over HTTPS.
- The dataset includes information about ride bookings, trip durations, pickup locations, drop-off points, fare amounts, and customer ratings.
- GitHub serves as a central repository for version-controlled data, making it a reliable source for real-time data ingestion.

**Step 2: Resource Group:**(Create the **-** Azure Data Factory, Storge Account, Databricks**)**



**Step 3: Data Ingestion:**

- Azure Data Factory (ADF) orchestrates the data ingestion process.
- ADF pipelines are set up to extract Ola Cabs ride data and load it into Azure Data Lake Storage Gen2.

- Triggers and schedules are configured to automate the ingestion process, ensuring data freshness.

**Pipeline Connection for Data Ingestion** (Data Lake Storage Gen2)**:**



## Step 4: Raw Data Storage:

- Ingested data is stored in Azure Data Lake Storage Gen2.
- The lakehouse architecture is used, with raw data landing in the "bronze" layer.
- The data lake provides a scalable and cost-efficient storage solution, handling large volumes of ride and customer data.

**Data loaded into Storage account- Containers:**



## Step 5: Data Transformation:

- Azure Databricks, powered by Apache Spark, processes the raw data.
- Notebooks are created in Databricks to clean, enrich, and aggregate the ride data.
- The transformation process includes handling missing values, calculating ride distances, aggregating revenue by location, and analyzing peak hours.
- The data is organized into bronze (raw), silver (cleaned), and gold (aggregated) layers, following the medallion architecture.

## Creation of Databricks:



## Databricks Spark integration for loading Transformed data to silver:

**SQL Warehouse:**

- Transformed data is loaded into a SQL warehouse within Databricks.
- The SQL warehouse enables fast, efficient querying, supporting complex analytical workloads.
- It acts as the central repository for analytical data, ready to be consumed by business intelligence tools.
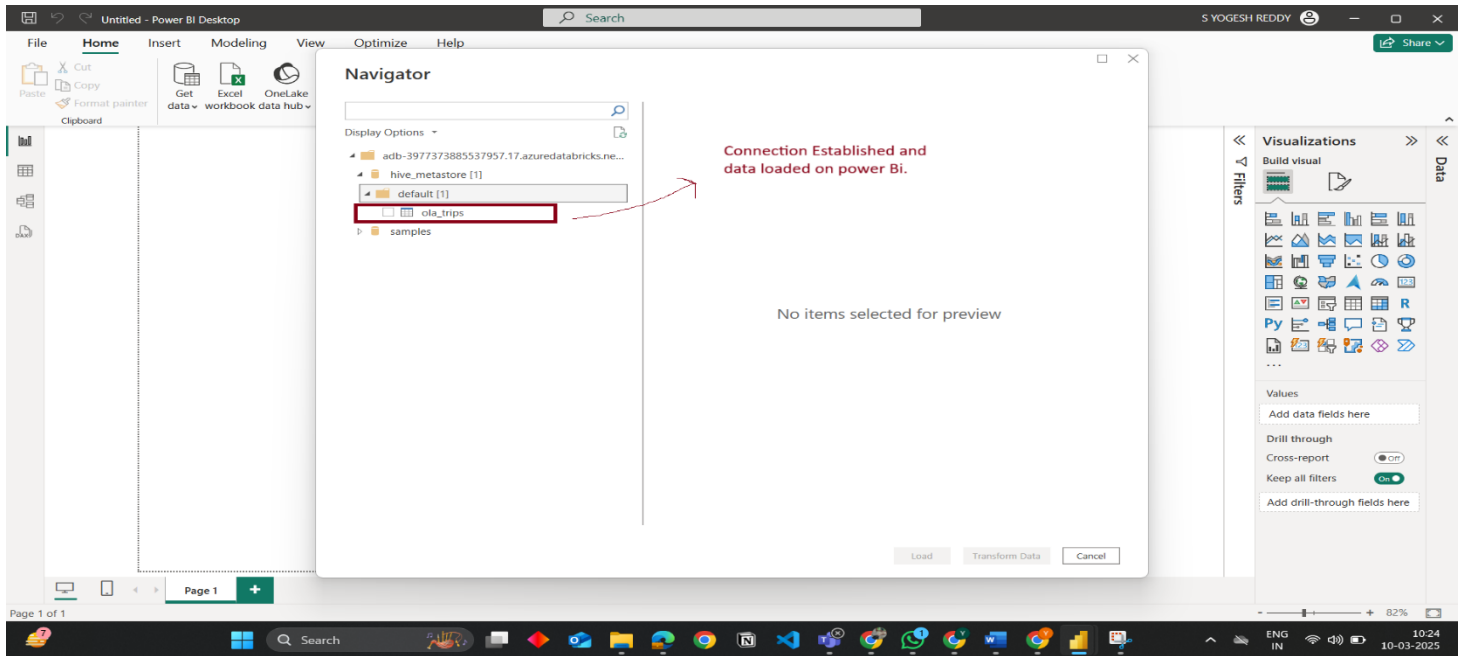
**SQL Warehouse Server Integration**



**Data Visualization:**

- Power BI connects to the SQL warehouse to build interactive dashboards.
- Visualizations include ride volume trends, revenue by city, average fare amounts, customer ratings, and trip heatmaps.
- The dashboard provides filtering capabilities, allowing users to explore data by date, location, fare range, and customer feedback.
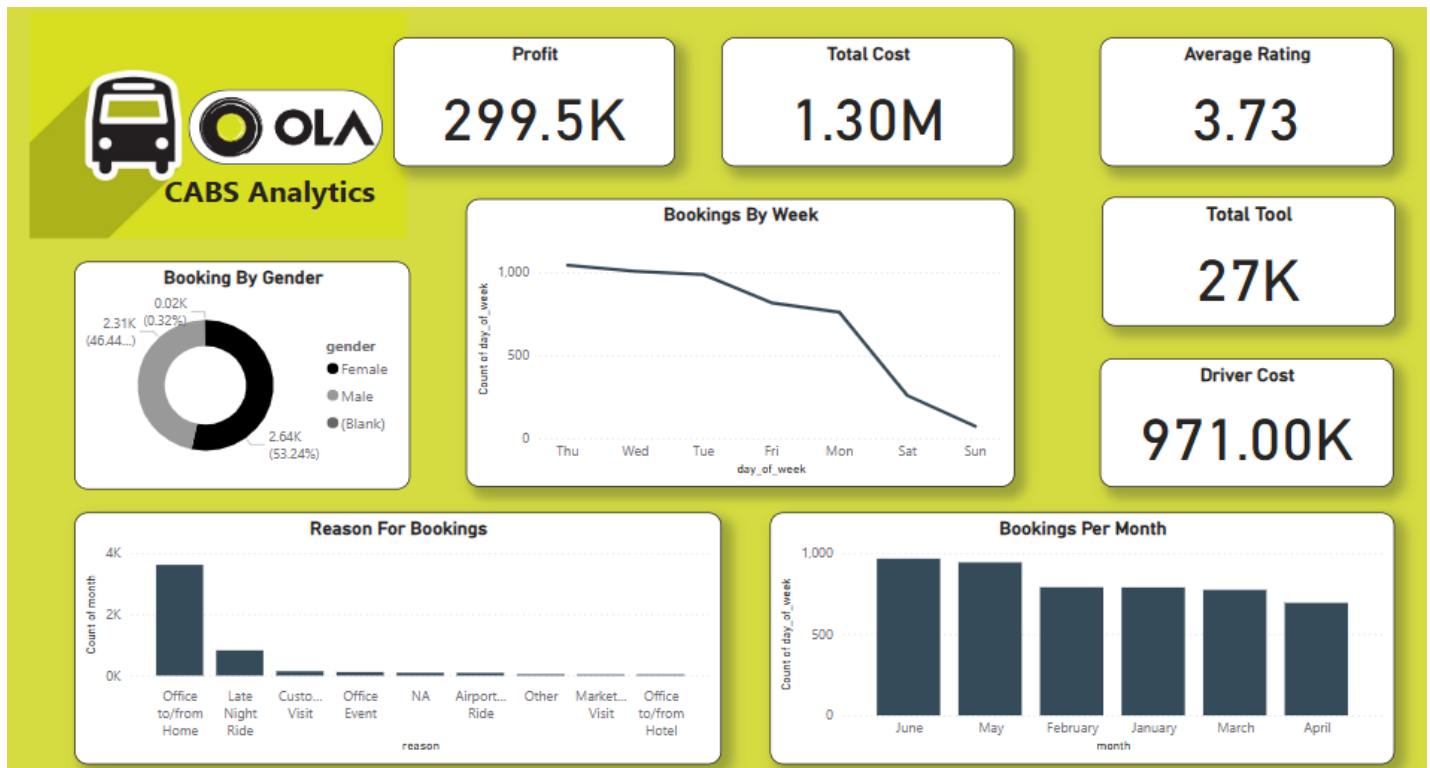
**Data integrated to Power BI** :



# RESULTS

The pipeline delivers a fully automated solution, enabling real-time data analysis for Ola Cabs. Stakeholders can access up-to-date dashboards that provide a clear view of operational metrics. The SQL warehouse ensures rapid query execution, while Power BI offers dynamic, easy-to-understand visualizations. The solution is scalable, allowing future growth in data volume and complexity.

**Data Loaded in to Power BI:**

**The Power BI Output :**



## Conclusion

This project highlights the capabilities of Azure's integrated services in building a scalable, efficient, and fully automated data pipeline. By leveraging tools like Azure Data Factory, Databricks, and Power BI, the solution transforms raw Ola Cabs data into valuable business insights. It serves as an excellent learning experience, reinforcing core data engineering concepts and best practices in cloud architecture.