# Analyzing Uber Data Using GCP, Python, and Looker Studio
*Akshit Anand*

---

## 1. Introduction

### 1.1 Objective
The objective of this project was to utilize cloud technology and data analytics to analyze Uber's public datasets. By leveraging Google Cloud Platform (GCP), Python, Mage for ETL, BigQuery, and Looker Studio, the project aimed to uncover patterns and insights in Uber data that could support strategic decision-making.

### 1.2 Data Overview
The Uber dataset consists of millions of Uber pickups across New York City over a specific period. Key features include:

- **Date/Time of Pickup**: Indicates when each ride occurred.
- **Location of Pickup**: Provides coordinates for pickup locations.
- **Ride Details**: Additional attributes relevant to each trip.



### 1.3 Technologies Used
This project used several cloud-based tools and technologies:

- **Google Cloud Storage**: For durable and scalable storage of the large dataset.
- **Python**: For data manipulation, pre-processing, and initial analysis.
- **Mage (ETL Tool)**: To extract, transform, and load data into BigQuery.
- **BigQuery**: For data warehousing, SQL querying, and storage.
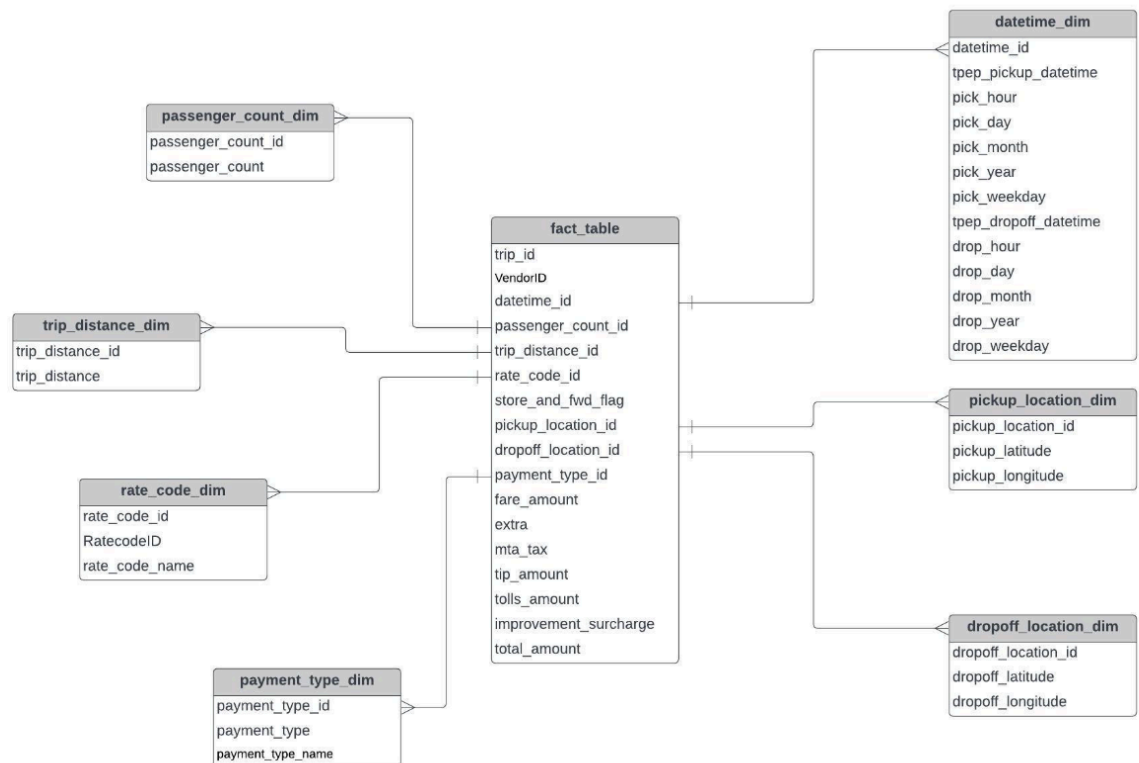- **Looker Studio**: To visualize data insights interactively.



# 2. Data Modeling

### 2.1 Schema Design
A **star schema** was designed to structure the data efficiently, facilitating fast query performance:

- **Fact Table**: Contains core information about each ride, including timestamps, pickup locations, and trip details.
- **Dimension Tables**: Separate tables for time, location, and ride characteristics, which allow for easier filtering and aggregations.
- **Screenshot Suggestion**: Include a **screenshot of the schema diagram** (e.g., from Lucidchart or another modeling tool) here. Place it below this section to visually depict

the data structure used in the project.



## 3. Data Ingestion and Storage

### 3.1 Uploading Data to Google Cloud Storage
The initial dataset (in CSV format) was uploaded to **Google Cloud Storage**. This storage service provides reliability, durability, and scalability, making it suitable for handling large datasets efficiently.

- **Screenshot Suggestion**: Add a **screenshot showing the successful upload of the data to Google Cloud Storage**. This screenshot should display the dataset in GCP, confirming its availability for further processing.



A screenshot that the data is successfully uploaded to GCP storage.
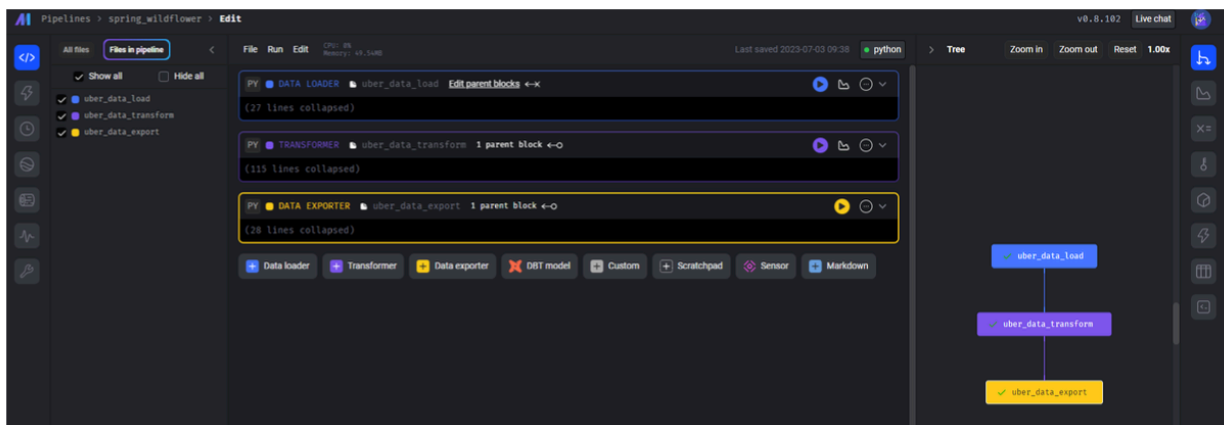
## 4. ETL Pipeline

**4.1 Deploying Mage on Google Compute Engine**

Mage, an ETL tool, was deployed on **Google Compute Engine** to handle data extraction, transformation, and loading processes. This setup allowed us to efficiently manage large data volumes and automate ETL tasks.

**4.2 ETL Process Using Mage**

Mage was configured to:

1. **Extract** the data from Google Cloud Storage.
2. **Transform** the data according to the pre-designed schema, including data cleansing steps.
3. **Load** the transformed data into BigQuery for further analysis.
● **Screenshot Suggestion**: Include a **screenshot of Mage deployed on Compute Engine** or a **screenshot of an ETL job in progress** to show the data pipeline setup.



# 5. Data Cleansing in BigQuery

**5.1 Data Import to BigQuery**

Once the ETL process was complete, the data was imported into BigQuery. This allowed us to utilize SQL-based querying to manage and analyze the data quickly and effectively.

**5.2 Data Cleansing Operations**

We performed several cleansing operations in BigQuery to ensure data quality:

● **Handling Missing Values**: Rows with missing data were removed or imputed to ensure complete records.
● **Removing Duplicates**: Duplicated rows were identified and removed to prevent double-counting.

- **Standardizing Formats**: Consistent formatting was applied to ensure data uniformity across all entries.



# 6. Data Visualization with Looker Studio

### 6.1 Visualization Process
Using **Looker Studio**, we created a series of interactive dashboards to visualize key insights and trends. Looker Studio's drag-and-drop interface enabled easy customization, and we designed the following visualizations:

- **Heatmaps** of pickup locations to identify high-demand areas.
- **Time-based trend graphs** showing peak hours for Uber pickups.

### 6.2 Insights from Visualizations
The visualizations provided a comprehensive view of Uber data patterns, including:

- **Peak Hours**: High demand during specific hours of the day.
- **Hotspot Locations**: High-demand areas were visible in the heatmaps, revealing popular pickup locations within New York City.
- **Screenshot Suggestion**: Add **screenshots of key Looker Studio dashboards** such as heatmaps and time-series graphs here. Place them below this section to illustrate insights visually.

## 7. Conclusion

### 7.1 Summary of the Project
This project demonstrated the effectiveness of cloud-based data analytics tools in handling large datasets. By combining GCP, BigQuery, and Looker Studio, we successfully processed, cleansed, and visualized Uber data, gaining valuable insights.

### 7.2 Key Insights
Our analysis uncovered essential insights, such as:

- **Peak Ride Times**: Certain hours showed higher ride frequencies, useful for optimizing driver allocation.
- **High-Demand Areas**: Identified through location-based heatmaps, helping pinpoint areas where Uber services are frequently required.

### 7.3 Future Scope and Enhancements
Potential extensions of this project could include:

- **Incorporating Additional Data**: Using data sources such as weather or local events to uncover more nuanced patterns.
- **Machine Learning for Predictive Analysis**: Applying machine learning algorithms to predict demand in specific locations or times for strategic planning.