

# Transportation & Logistics Data Processing Project – Documentation

---

## Project Summary

This project focuses on analyzing data from the transportation and logistics domain to generate meaningful insights related to deliveries, routes, vehicles, and drivers. The overall objective is to:

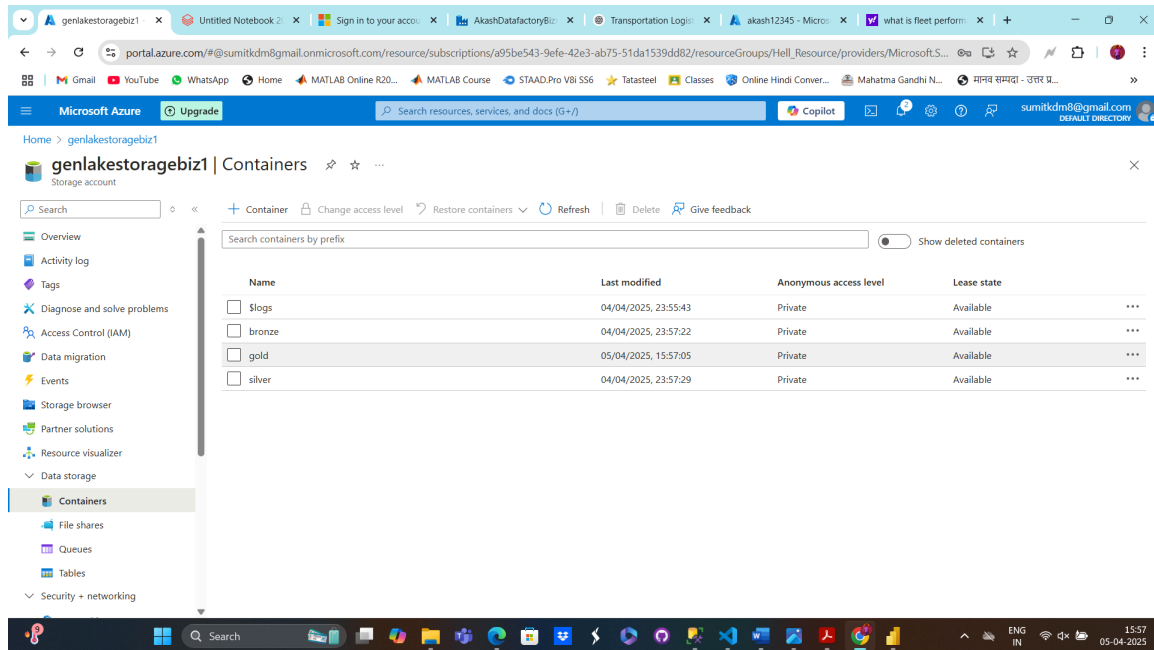
- Improve route efficiency
- Monitor fleet and driver performance
- Track and reduce fuel consumption
- Visualize operations through dashboards

To achieve this, we use a structured data pipeline built using PySpark for data processing, MySQL for storing cleaned and aggregated data, and Power BI for visual reporting. The pipeline follows the Medallion Architecture model consisting of Bronze (raw data), Silver (cleaned and enriched), and Gold (aggregated for reporting) layers.

## Tools & Technologies Used

Tool/Technology	Purpose/Usage
PySpark	Used for reading, cleaning, transforming, and enriching large CSV datasets.
Jupyter Notebook	To write and run PySpark scripts interactively.
MySQL	Stores Silver (clean) and Gold (aggregated) layer data for further reporting.
Power BI	Connects to MySQL Gold layer to build dashboards and visualizations.
CSV Files	Raw input data files containing delivery, route, vehicle, and driver information.

## Project Architecture – Medallion Approach



Name	Last modified	Anonymous access level	Lease state
<input type="checkbox"/> slogs	04/04/2025, 23:55:43	Private	Available
<input type="checkbox"/> bronze	04/04/2025, 23:57:22	Private	Available
<input checked="" type="checkbox"/> gold	05/04/2025, 15:57:05	Private	Available
<input type="checkbox"/> silver	04/04/2025, 23:57:29	Private	Available

### 1. Bronze Layer – Raw Data Ingestion

The Bronze layer is the foundation where raw data is ingested. CSV files from multiple sources are loaded and converted to a uniform Parquet format. No major transformation is applied at this layer. Metadata fields such as ingestion date and source filename are added for auditing.

Expected files:

- delivery\_data.csv
- vehicle\_data.csv
- route\_data.csv
- driver\_data.csv

### 2. Silver Layer – Cleaned & Enriched Data

In this layer, the raw data is cleaned by removing records with nulls, invalid formats, or duplicates. The tables are then enriched by joining related datasets such as vehicle information, route data, and driver details

geniakestoragebiz1 x Untitled Notebook 2 x Sign in to your acco... x AkashDatafactoryBiz... x Transportation Logis... x akash12345 - Micro... x what is fleet perform... x + -

adb-542866418164150.10.azuredatabricks.net/editor/notebooks/3696802500255949?o=542866418164150#command/7015074710219969

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P AkashBizmetric

New Workspace Recents Catalog Workflows Compute Marketplace SQL SQL Editor Queries Dashboards Genie Alerts Query History SQL Warehouses Data Engineering Job Runs Data Ingestion Pipelines Machine Learning

Untitled Notebook 2025-04-04 23:59:07 Python Last edit was 4 hours ago Run all sumit kadam's Cluster Schedule Share

Catalog Type to search... For you All My organization akashbizmetric system Delta Shares Received samples Legacy hive\_metastore

```
df = spark.read.parquet("/mnt/adls_storage/bronze/delivery_data.parquet")
df = df.dropna()
display(df)
df.write.jdbc(
    url="jdbc:sqlserver://akash12345.database.windows.net:1433;databaseName=gold",
    table="silver_db.delivery_data",
    mode="overwrite",
    properties={
        "user": "akash12345",
        "password": "Naddraj1998@",
        "driver": "com.microsoft.sqlserver.jdbc.SQLServerDriver"
    })
```

(3) Spark Jobs

pyspark.sql.dataframe.DataFrame = [delivery\_id: string, vehicle\_id: string ... 6 more fields]

	delivery_id	vehicle_id	route_id	driver_id	delivery_date	delivery_time	de
38	189	329	32/	29	09-11-2024	4.52	91/91
39	191	436	1279	180	27-07-2024	12	185.56
40	193	448	1737	24	29-07-2024	12	500
41	206	472	760	183	10-07-2024	4.2	500
42	212	29	1930	139	19-04-2024	12	500
43	213	314	879	287	16-12-2024	1.38	162.84
44	217	108	1326	255	11-11-2024	12	131.07

15:57 05-04-2025

geniakestoragebiz1 x Untitled Notebook 2 x Sign in to your acco... x AkashDatafactoryBiz... x Transportation Logis... x akash12345 - Micro... x what is fleet perform... x + -

adb-542866418164150.10.azuredatabricks.net/editor/notebooks/3696802500255949?o=542866418164150#command/7015074710219969

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P AkashBizmetric

New Workspace Recents Catalog Workflows Compute Marketplace SQL SQL Editor Queries Dashboards Genie Alerts Query History SQL Warehouses Data Engineering Job Runs Data Ingestion Pipelines Machine Learning

Untitled Notebook 2025-04-04 23:59:07 Python Last edit was 4 hours ago Run all sumit kadam's Cluster Schedule Share

Catalog Type to search... For you All My organization akashbizmetric system Delta Shares Received samples Legacy hive\_metastore

50	244	83	533	64	24-07-2024	12	500
51	252	315	668	196	15-02-2025	12	197.14
52	267	373	675	49	26-01-2025	12	95.66

204 rows | 4.04s runtime Refreshed 4 hours ago

```
df1 = spark.read.parquet("/mnt/adls_storage/bronze/driver_data.parquet")
df1.show()
df1.write.jdbc(
    url="jdbc:sqlserver://akash12345.database.windows.net:1433;databaseName=gold",
    table="silver_db.vechile_data",
    mode="overwrite",
    properties={
        "user": "akash12345",
        "password": "Naddraj1998@",
        "driver": "com.microsoft.sqlserver.jdbc.SQLServerDriver"
    })
```

(3) Spark Jobs

pyspark.sql.dataframe.DataFrame = [driver\_id: string, driver\_name: string ... 2 more fields]

15:58 05-04-2025

geniakstoragebiz1 x Untitled Notebook 2 x Sign in to your acco... x AkashDatafactoryBiz x Transportation Logis x akash12345 - Micros x what is fleet perform x + -

adb-542866418164150.10.azuredatabricks.net/editor/notebooks/3696802500255949?o=542866418164150#command/7015074710219969

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P AkashBizmetric

Untitled Notebook 2025-04-04 23:59:07 Python

File Edit View Run Help Last edit was 4 hours ago

Run all sumit kadam's Cluster Schedule Share

Catalog

Type to search...

For you All

My organization

akashbizmetric

system

Delta Shares Received

samples

Legacy

hive\_metastore

df1: pyspark.sql.dataframe.DataFrame = [driver\_id: string, driver\_name: string ... 2 more fields]

df2 = spark.read.parquet("/mnt/adls\_storage/bronze/route\_data.parquet")  
df2.show()  
df2.write.jdbc(  
url="jdbc:sqlserver://akash12345.database.windows.net:1433;databaseName=gold",  
table="silver\_db.route\_data",  
mode="overwrite",  
properties={  
"user": "akash12345",  
"password": "Haddraj1998@",  
"driver": "com.microsoft.sqlserver.jdbc.SQLServerDriver"  
})

df2: pyspark.sql.dataframe.DataFrame = [route\_id: string, start\_location: string ... 2 more fields]

df3 = spark.read.parquet("/mnt/adls\_storage/bronze/vehicle\_data.parquet")  
df3.show()  
df3.write.jdbc(  
url="jdbc:sqlserver://akash12345.database.windows.net:1433;databaseName=gold",  
table="silver\_db.vehicle\_data")

15:58 05-04-2025

geniakstoragebiz1 x Untitled Notebook 2 x Sign in to your acco... x AkashDatafactoryBiz x Transportation Logis x akash12345 - Micros x what is fleet perform x + -

adb-542866418164150.10.azuredatabricks.net/editor/notebooks/3696802500255949?o=542866418164150#command/7015074710219969

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P AkashBizmetric

Untitled Notebook 2025-04-04 23:59:07 Python

File Edit View Run Help Last edit was 4 hours ago

Run all sumit kadam's Cluster Schedule Share

Catalog

Type to search...

For you All

My organization

akashbizmetric

system

Delta Shares Received

samples

Legacy

hive\_metastore

df3: pyspark.sql.dataframe.DataFrame = [vehicle\_id: string, vehicle\_type: string ... 2 more fields]

df\_silver=df.join(df1, "driver\_id", "inner") \  
.join(df2, "route\_id", "inner") \  
.join(df3, "vehicle\_id", "inner") \  
.withColumn("fuel\_consumed", col("distance\_covered") / col("fuel\_efficiency"))

display(df\_silver)

df\_silver: pyspark.sql.dataframe.DataFrame = [vehicle\_id: string, route\_id: string ... 16 more fields]

	vehicle_id	route_id	driver_id	delivery_id	delivery_date	delivery_time	fuel_efficiency
10	360	1297	43	78	14-10-2024	2:59	187.11
11	192	960	173	84	21-02-2025	12	107.4
12	311	1578	193	93	12-11-2024	12	500
13	435	127	216	107	27-03-2024	5:01	500
14	218	818	254	123	01-08-2024	12	500
15	30	1237	298	126	05-10-2024	12	77.09
16	100	757	225	130	02-03-2025	12	500

15:58 05-04-2025

geniakestoragebiz1 x Untitled Notebook 2 x Sign in to your acco... x AkashDatafactoryBio... x Transportation Logis... x akash12345 - Micro... x what is fleet perform... x + -

adb-542866418164150.10.azuredatabricks.net/editor/notebooks/3696802500255949?o=542866418164150#command/7015074710219969

Gmail YouTube WhatsApp Home MATLAB Online R20... MATLAB Course STAAD.Pro V8i 556 Tatabsteel Classes Online Hindi Conver... Mahatma Gandhi N... मानव सम्पदा - उत्तर प्र...

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P AkashBizmetric

New Workspace Recents Catalog Workflows Compute Marketplace SQL SQL Editor Queries Dashboards Genie Alerts Query History SQL Warehouses Data Engineering Job Runs Data Ingestion Pipelines Machine Learning

Untitled Notebook 2025-04-04 23:59:07 Python File Edit View Run Help Last edit was 4 hours ago Run all sumit kadam's Cluster Schedule Share

Catalog Type to search... For you All My organization akashbizmetric system Delta Shares Received samples Legacy hive\_metastore

40	253	724	215	261	25-02-2025	12	500
41	260	618	74	262	08-04-2024	12	500
42	89	900	88	264	29-03-2024	4.96	183.06
43							

140 rows | 0.80s runtime Refreshed 4 hours ago

```
from pyspark.sql.functions import *

df_silver= df_silver.withColumn("route_name",concat_ws(" to ", col("start_location"), col("end_location")))
df_silver=df_silver.withColumn("processed_date",current_date())

display(df_silver)
```

(4) Spark Jobs

df\_silver: pyspark.sql.dataframe.DataFrame = [vehicle\_id: string, route\_id: string ... 18 more fields]

	vehicle_id	route_id	driver_id	delivery_id	delivery_date	delivery_time	de
20	371	69	41	156	14-05-2024	12	500
21	8	1050	171	166	23-09-2024	1.48	84.71
22	247	1742	30	173	18-03-2025	4.32	500

15:38 05-04-2025

geniakstoragebiz1 x Untitled Notebook 2 x Sign in to your acco... x AkashDatafactoryBiz x Transportation Logis x akash12345 - Micro x what is fleet perform x + -

adb-542866418164150.10.azuredatabricks.net/editor/notebooks/3696802500255949?o=542866418164150#command/7015074710219969

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P AkashBizmetric

Untitled Notebook 2025-04-04 23:59:07 Python

File Edit View Run Help Last edit was 4 hours ago

Run all submit kadam's Cluster Schedule Share

Catalog

Type to search...

For you All

- My organization
  - akashbizmetric
  - system
- Delta Shares Received
  - samples
- Legacy
  - hive\_metastore

12:24 PM (1g) 10

```
df_silver=df_silver.select("delivery_id","vehicle_type","driver_name","route_name","delivery_time",
"distance_covered","delivery_status","fuel_consumed","processed_date")
display(df_silver)
```

(4) Spark Jobs

df\_silver: pyspark.sql.dataframe.DataFrame = [delivery\_id: string, vehicle\_type: string ... 7 more fields]

	delivery_id	vehicle_type	driver_name	route_name	delivery_time
1	5	Truck	Joseph Wilson	New Rogerton to North Jennifer	12
2	6	Van	Elizabeth Stout	North Susan to Bellport	12
3	9	Van	Mary Watson	Norriston to East Karenville	12
4	10	Truck	Cheryl Davis	Garrettsstad to Shellaport	1.39
5	16	Truck	Alan Steele	Lake Denise	2.3
6	34	Van	Paul Foster	Pagechester	12
7	52	Car	Mr. Patrick Adams III	North John to Robinbury	4.72
8	60	Bus	Spencer Randolph	Lake Michele	12
9	77	Truck	Sara Pham	Port Laurenton to Francistown	2.54
10	78	Bus	Jeffrey Lindsey	Elizabethfort to Angelaland	2.59
11	84	Car	Roy Price	Walkerton to Jeffreyport	12
12	93	Van	Mrs. Lisa Clark	South Tylerchester to Luidand	12
13	107	Truck		East Melanie to Garretthaven	5.01

15:38 05-04-2025

geniakstoragebiz1 x Untitled Notebook 2 x Sign in to your acco... x AkashDatafactoryBiz x Transportation Logis x akash12345 - Micro x what is fleet perform x + -

adb-542866418164150.10.azuredatabricks.net/editor/notebooks/3696802500255949?o=542866418164150#command/7015074710219969

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P AkashBizmetric

Untitled Notebook 2025-04-04 23:59:07 Python

File Edit View Run Help Last edit was 4 hours ago

Run all submit kadam's Cluster Schedule Share

Catalog

Type to search...

For you All

- My organization
  - akashbizmetric
  - system
- Delta Shares Received
  - samples
- Legacy
  - hive\_metastore

12:24 PM (12g) 12

```
df_silver.write.mode("overwrite").parquet("/mnt/adls_storage/silver/silver_data.parquet")
```

(4) Spark Jobs

True

12:25 PM (4g) 13

```
df_silver.write.mode("overwrite").parquet("/mnt/adls_storage/silver/silver_data.parquet")
```

(4) Spark Jobs

12:25 PM (3g) 14

```
df_silver.write.jdbc(
url="jdbc:sqlserver://akash12345.database.windows.net:1433;databaseName=gold",
table="silver_db.Silver_table",
mode="overwrite",
properties={
"user": "akash12345",
"password": "Naddraj1998@",
"driver": "com.microsoft.sqlserver.jdbc.SQLServerDriver"
})
```

(4) Spark Jobs

15:38 05-04-2025

Calculated fields such as fuel consumed are added to improve analysis. Data is stored in both Parquet format and a MySQL table called `silver\_db.delivery\_data\_silver`.

Example Calculation:

$$\text{fuel\_consumed} = \text{distance\_covered} / \text{fuel\_efficiency}$$

[Space for Silver Layer diagram/image]

### **3. Gold Layer – Aggregated Data for Reporting**

This layer prepares data for final reporting. Key metrics are calculated by aggregating the Silver layer data. Examples include:

- Total number of deliveries per route, driver, and vehicle
- Average delivery time
- Total fuel consumption and efficiency
- Driver ratings and performance over time

Stored in MySQL table: `gold\_db.transportation\_gold`

geniakestoragebiz1 x Untitled Notebook 2 x Sign in to your acco... x AkashDatafactoryBiz... x Transportation Logi... x gold (akash12345/g... x what is fleet perform... x + -

portal.azure.com/#@sumitkdm8@gmail.onmicrosoft.com/resource/subscriptions/a95be543-9efe-42e3-ab75-51da1539dd82/resourceGroups/Hell\_Resource/providers/Microsoft.S...

Microsoft Azure Upgrade Search resources, services, and docs (G+V) Copilot sumitkdm8@gmail.com DEFAULT DIRECTORY

Home > gold (akash12345/gold) | Query editor (preview) ☆ ...

SQL database

Search Login + New Query Open query Feedback Getting started

Overview Activity log Tags Diagnose and solve problems Query editor (preview) Mirror database in Fabric (preview) Resource visualizer Settings Data management Integrations Power Platform Security Intelligent performance Monitoring Automation Help

gold.driver\_summary gold.fleet\_avg\_distance gold.fleet\_avg\_distance1 gold.fleet\_delivery\_summary gold.fleet\_performance gold.fleet\_performance\_summar... gold.fleet\_performance1 gold.fleet\_performance2 gold.route\_optimization\_summa... gold.total\_deliveries\_per\_driver gold.total\_deliveries\_per\_route gold.total\_deliveries\_per\_vehicle... gold.transportation\_gold silver\_db.delivery\_data silver\_db.driver\_data silver\_db.route\_data silver\_db.Silver\_table silver\_db.vechile\_data Views Stored Procedures

Query 1 x Query 2 x

Run Cancel query Save query Export data as Show all

Results Messages

Search to filter items...

delivery_id	vehicle_type	driver_name	route_name	delivery_tim
5	Truck	Joseph Wilson	New Rogerton to North Jennifer	12
6	Van	Elizabeth Stout	North Susan to Bellport	12
9	Van	Mary Watson	Norriston to East Kareville	12
10		Cheryl Davis	Garrettstad to Sheilaport	1.39
16	Truck	Alan Steele	Lake Denise	2.3
34	Van	Paul Foster	Pagechester	12
52	Car	Mr. Patrick Adams III	North John to Robinbury	4.72
60	Bus	Spencer Randolph	Lake Michele	12

Query succeeded | 1s

geniakestoragebiz1 x Untitled Notebook 2 x Sign in to your acco... x AkashDatafactoryBiz... x Transportation Logi... x gold (akash12345/g... x what is fleet perform... x + -

portal.azure.com/#@sumitkdm8@gmail.onmicrosoft.com/resource/subscriptions/a95be543-9efe-42e3-ab75-51da1539dd82/resourceGroups/Hell\_Resource/providers/Microsoft.S...

Microsoft Azure Upgrade Search resources, services, and docs (G+V) Copilot sumitkdm8@gmail.com DEFAULT DIRECTORY

Home > gold (akash12345/gold) | Query editor (preview) ☆ ...

SQL database

Search Login + New Query Open query Feedback Getting started

Overview Activity log Tags Diagnose and solve problems Query editor (preview) Mirror database in Fabric (preview) Resource visualizer Settings Data management Integrations Power Platform Security Intelligent performance Monitoring Automation Help

gold.driver\_summary gold.fleet\_avg\_distance gold.fleet\_avg\_distance1 gold.fleet\_delivery\_summary gold.fleet\_performance gold.fleet\_performance\_summar... gold.fleet\_performance1 gold.fleet\_performance2 gold.route\_optimization\_summa... gold.total\_deliveries\_per\_driver gold.total\_deliveries\_per\_route gold.total\_deliveries\_per\_vehicle... gold.transportation\_gold silver\_db.delivery\_data silver\_db.driver\_data silver\_db.route\_data silver\_db.Silver\_table silver\_db.vechile\_data Views Stored Procedures

Query 1 x Query 2 x

Run Cancel query Save query Export data as Show all

Results Messages

Search to filter items...

route_name	total_deliveries	avg_delivery_time	avg_fuel_consumed	total_distanc
Arroyofort	1	12	62.34413965087282	500
Brandonport to Shannonview	1	12	43.36513443191674	500
Buckborough to Randystad	1	12	12.760406091370559	125.69
Burnsmouth to East William	1	12		175.98
Carolineland to Kevinfurt	1	4.37		63.9
Charlesshire to Port Ashley	1	3.55	48.590864917395535	500
Christopherburgh to West Alicia	1	1.93	70.72135785007072	500
Crystalhaven to Michellestad	1	12	9.812039312039312	79.87

Query succeeded | 0s



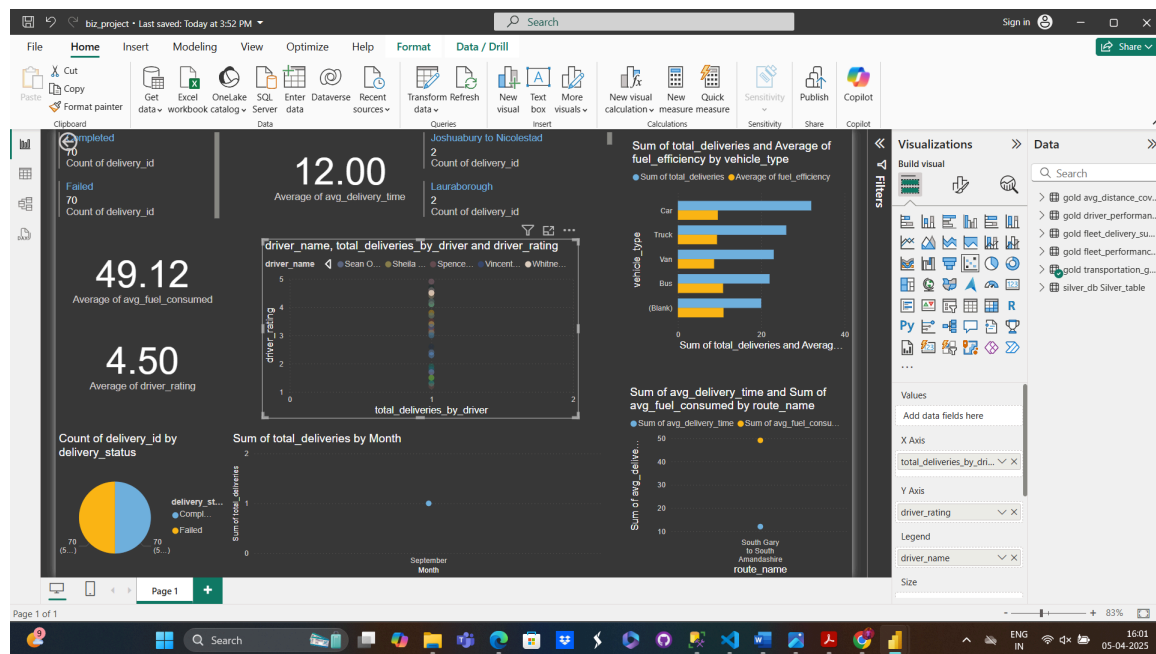
## Power BI Dashboard – Visual Insights

Power BI is used to build interactive visual dashboards based on the Gold layer data stored in MySQL. It includes key performance indicators (KPIs), trend lines, and comparative analysis for better decision-making.

Connected Data Source: MySQL – gold\_db.transportation\_gold

Key Visuals and KPIs:

1. Line Chart – Average delivery time & fuel usage per route
2. Bar Chart – Total deliveries per vehicle and their efficiency
3. Scatter Plot – Driver performance vs ratings
4. Pie Chart – Completed vs Failed deliveries
5. Line Graph – Delivery trends by week/month



## Final Deliverables

1. MySQL Databases:
  - Silver layer table: silver\_db.delivery\_data\_silver
  - Gold layer table: gold\_db.transportation\_gold
2. Power BI Dashboard:
  - Includes visual KPIs, filters, and export options.
3. Documentation:

- Detailed Word or PDF file explaining each step in the pipeline.

#### 4. Git Repository:

- Contains PySpark code, MySQL scripts, and Power BI .pbix file.

[https://github.com/Akash05111998/Project\\_Bizmetric.git](https://github.com/Akash05111998/Project_Bizmetric.git)

#### **Automation & Best Practices**

- Pipeline is fully automated from raw data ingestion to dashboard refresh.
- Consistent naming conventions across files and tables.
- Modular code with reusable functions.
- Data audit logs maintained at each transformation stage.
- Easy to scale and integrate with more data sources.