

Azure Databricks end to end Project with Unity Catalog CI/CD

© Shubham Wadekar

Real-Time Data Lakehouse for Traffic and Roads Analytics

Project Overview

This project builds a real-time data pipeline using Azure Data Lake and Delta Lake to process and analyze traffic and roads datasets. Data is manually loaded into a landing zone to simulate incremental ingestion. It then flows into the bronze layer using Spark Structured Streaming with Auto Loader, capturing all raw records.

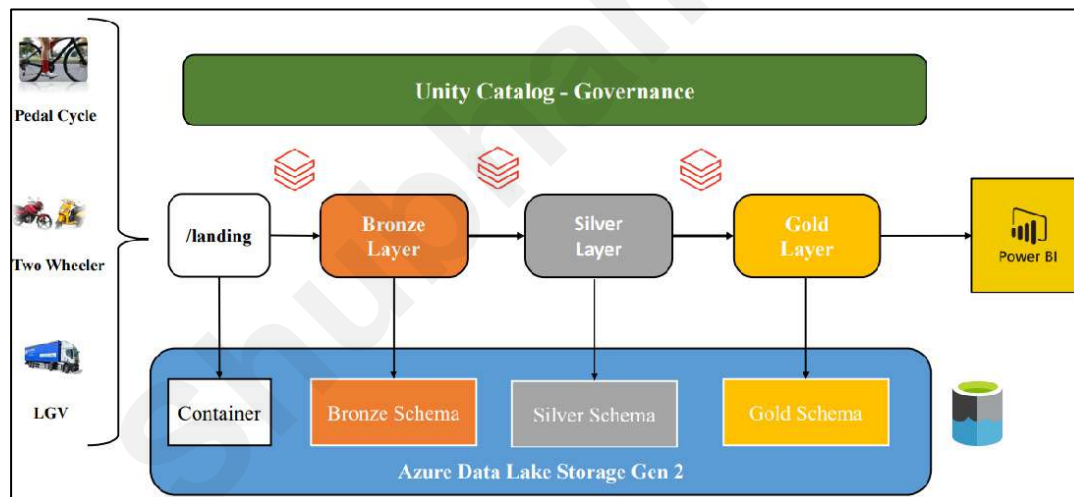
From the bronze layer, only new data is transformed and moved to the silver layer for cleaning and enrichment. The final curated data is stored in the gold layer as optimized tables or views, ready for reporting and data science.

Key technologies include:

- Delta Lake for reliable storage and ACID transactions
- Auto Loader for incremental data ingestion
- Power BI for reporting and visualization
- Azure DevOps for CI/CD
- Access controls to simulate enterprise-grade security

This end-to-end pipeline supports real-time analytics in a scalable, cloud-native environment.

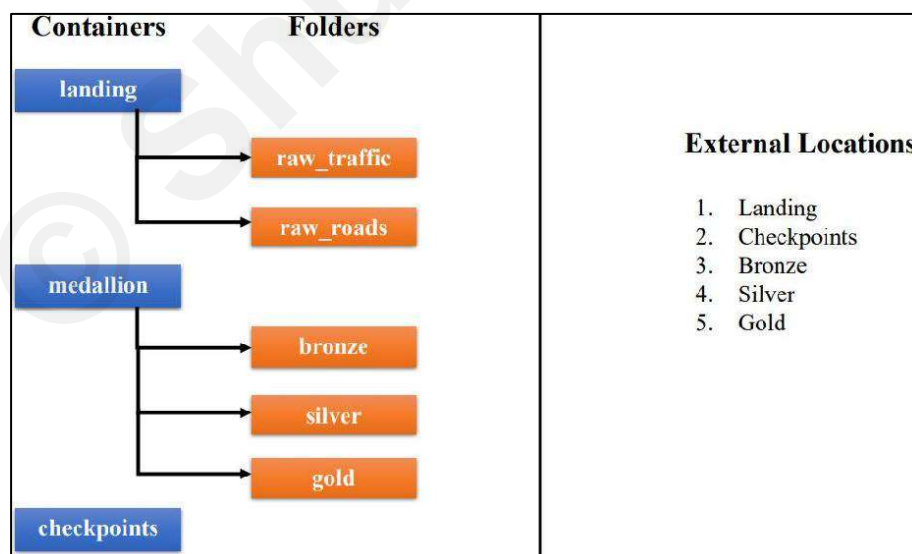
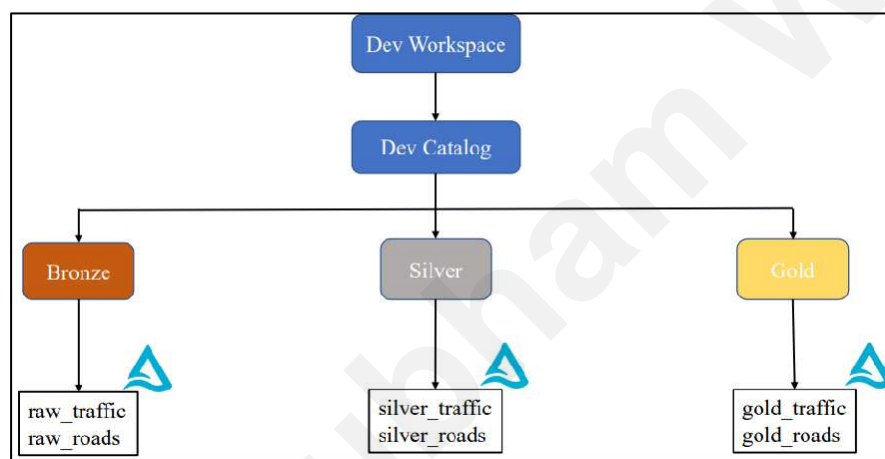
Project Architecture



The project uses a Medallion Architecture (Bronze → Silver → Gold) on Azure Databricks with Delta Lake and Unity Catalog for governance. Data flows through the following stages:

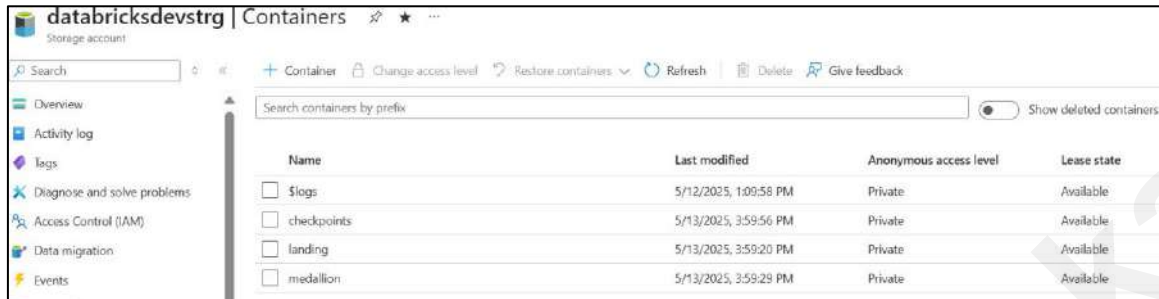
- **LandingZone:** Raw traffic and roads data are manually ingested into Azure Data Lake Storage Gen2 (/landing folder) to simulate streaming input.
- **Bronze Layer:** Ingested incrementally using Auto Loader with Structured Streaming. Raw data is stored as Delta tables (raw_traffic, raw_roads) under the bronze schema.
- **Silver Layer:** Transforms include renaming columns, deriving Electric_Vehicles_Count, Motor_Vehicles_Count, and categorizing road types. Cleaned datasets (silver_traffic, silver_roads) are stored in the silver schema.
- **Gold Layer:** Final business logic is applied (e.g., Vehicle_Intensity), and tables are optimized for reporting (gold_traffic, gold_roads). These are consumed in Power BI for insights.
- **Governance:** Managed via Unity Catalog with a three-tier namespace (catalog.schema.table) and fine-grained access control.
- **CI/CD:** Implemented using Azure DevOps for deploying code and configurations across environments (Dev → UAT → Prod).

Project Setup



This project is organized using Azure Data Lake Storage Gen2, structured into containers and folders representing different stages of the data pipeline.

Containers and Folders

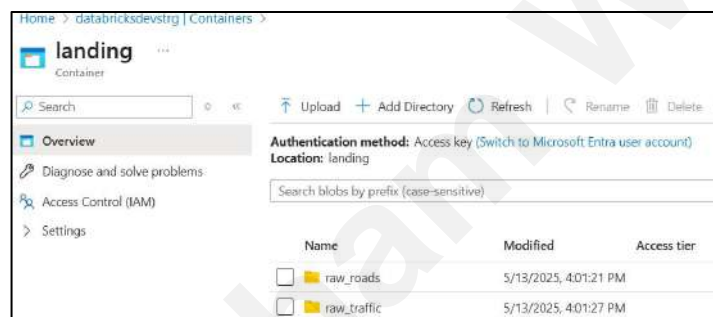


Name	Last modified	Anonymous access level	Lease state
<input type="checkbox"/> \$logs	5/12/2025, 1:09:58 PM	Private	Available
<input type="checkbox"/> checkpoints	5/13/2025, 3:59:56 PM	Private	Available
<input type="checkbox"/> landing	5/13/2025, 3:59:20 PM	Private	Available
<input type="checkbox"/> medallion	5/13/2025, 3:59:29 PM	Private	Available

We have 3 containers in Azure Data Lake Storage:

1. Landing:

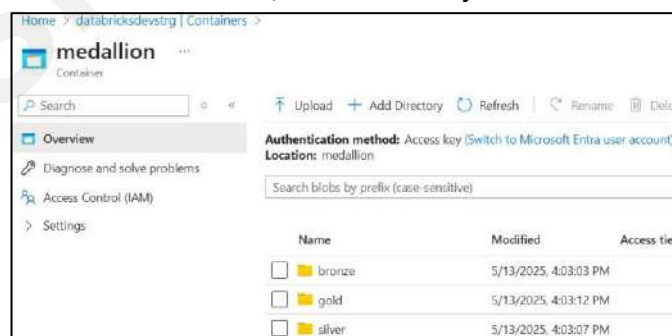
- Holds the raw input data.
- Contains 2 sub folders: raw_traffic (for traffic dataset), raw_roads (for road dataset).



Name	Modified	Access tier
<input type="checkbox"/> raw_roads	5/13/2025, 4:01:21 PM	
<input type="checkbox"/> raw_traffic	5/13/2025, 4:01:27 PM	

2. Medallion:

- Implements the medallion architecture with subfolders.
- Bronze – stores raw ingested data
- Silver – stores cleaned and transformed data
- Gold – stores curated, business-ready data



Name	Modified	Access tier
<input type="checkbox"/> bronze	5/13/2025, 4:03:03 PM	
<input type="checkbox"/> gold	5/13/2025, 4:03:12 PM	
<input type="checkbox"/> silver	5/13/2025, 4:03:07 PM	

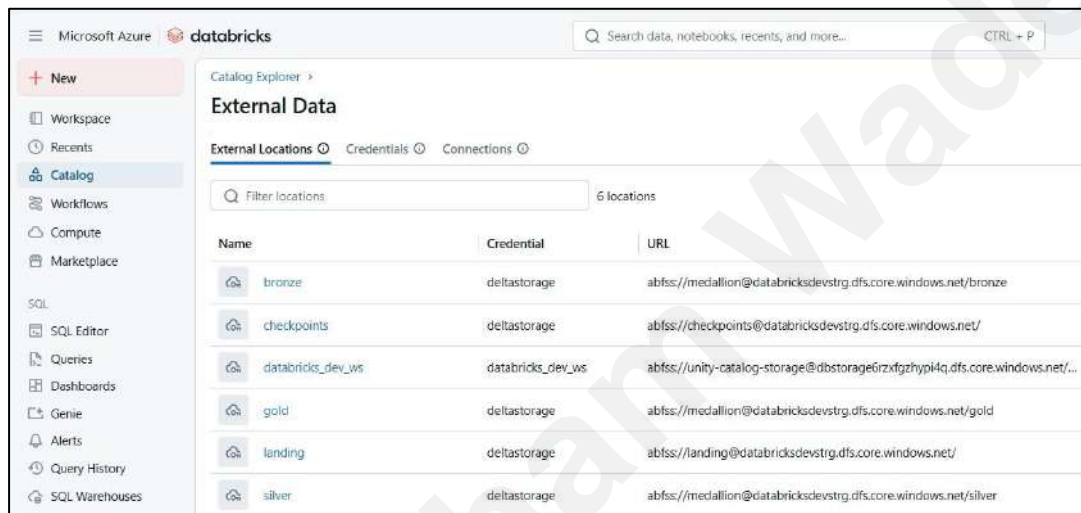
3. Checkpoints:

- a. Stores streaming query checkpoints to support fault tolerance and exactly once processing in Structured Streaming.

External Locations

The following external locations are registered in Unity Catalog to enable secure and governed access to data:

1. Landing
2. Checkpoints
3. Bronze
4. Silver
5. Gold



The screenshot shows the Databricks Catalog Explorer interface. The 'External Data' section is active, displaying a table of registered external locations. The table has three columns: Name, Credential, and URL. There are six locations listed: bronze, checkpoints, databricks_dev_ws, gold, landing, and silver. Each location is associated with a 'deltastorage' credential and a specific URL in the 'abfss' protocol.

Name	Credential	URL
bronze	deltastorage	abfss://medallion@databricksdevstrg.dfs.core.windows.net/bronze
checkpoints	deltastorage	abfss://checkpoints@databricksdevstrg.dfs.core.windows.net/
databricks_dev_ws	databricks_dev_ws	abfss://unity-catalog-storage@dbstorage60rzdfgzhyp4q.dfs.core.windows.net/...
gold	deltastorage	abfss://medallion@databricksdevstrg.dfs.core.windows.net/gold
landing	deltastorage	abfss://landing@databricksdevstrg.dfs.core.windows.net/
silver	deltastorage	abfss://medallion@databricksdevstrg.dfs.core.windows.net/silver

Each location corresponds to a specific path in the data lake and is linked with appropriate storage credentials and access controls using Unity Catalog.

Data Sources

The datasets are from **Kaggle. Raw**

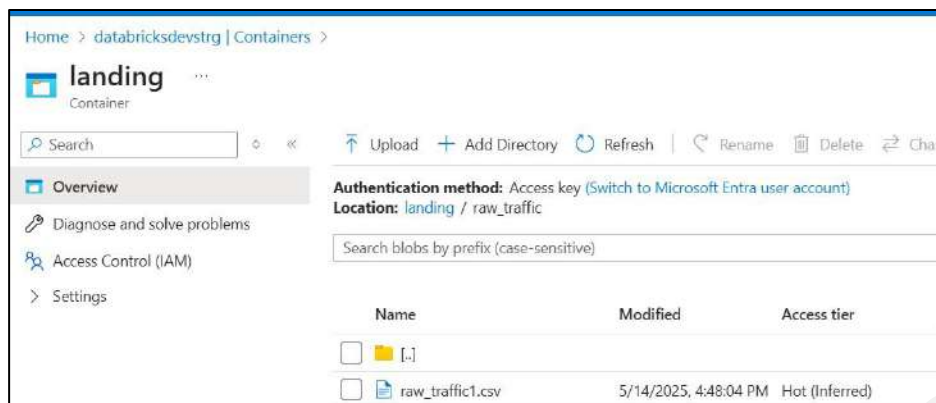
Traffic Dataset

The Raw Traffic Dataset is one of the core inputs for this project. It has actual data collected by trained enumerators to feed data into road traffic estimates. It contains structured information collected from traffic monitoring points across UK roads. This dataset has a raw count for the number of vehicles of each type that flowed past at each point of day, broken by direction and an hour.

It has pedal cycles, 2-wheeler vehicles, buses and coaches, LGV (Large Good Vehicles) and HGV (Heavy goods vehicles), and electric vehicles. So, we need to find out at a given time within an hour, how many vehicles are recorded in the raw traffic count dataset. We will analyze the type of vehicle travelling at a given point along with roads.

Source and Storage

- The dataset is manually placed in the /landing/raw_traffic folder of Azure Data Lake



Storage Gen2.

- It simulates real-time ingestion and is later processed incrementally using Spark Structured Streaming with Auto Loader.
- Once ingested, it is first stored in the bronze layer as the table raw_traffic.

Purpose in the pipeline

This dataset serves multiple purposes:

- Provides the raw measurements of traffic flow, needed to calculate derived metrics.
- Enables tracking of hourly, daily, and yearly trends in vehicle movement.
- Acts as a base for data enrichment and transformation in the silver layer, eventually powering analytical dashboards and reports on the gold layer.

Data Dictionary

Data Dictionary has all column names. It defines what information that each column has.

Data Dictionary	
1. Record ID	Vehicle flow point
2. Count point id	
3. Direction of travel	
4. Year	
5. Count date	Travel info of vehicle
6. hour	
7. Region id	
8. Region name	
9. Local authority name	Count of types of vehicle
10. Road name	
11. Road Category ID	
12. Start junction road name	
13. End junction road name	
14. Latitude	
15. Longitude	
16. Link length km	
17. Pedal cycles	
18. Two wheeled motor vehicles	
19. Cars and taxis	
20. Buses and coaches	
21. LGV Type	
22. HGV Type	
23. EV Car	
24. EV Bike	

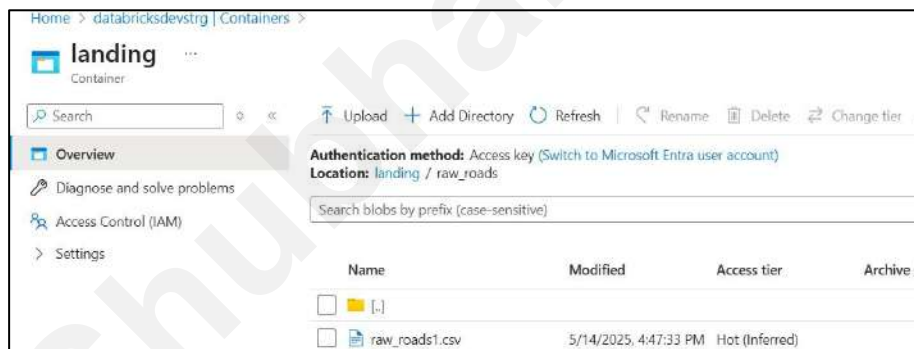
Data Dictionary		
1. Record ID	=	Uniquely identifies a record
2. Count point id	=	A unique reference for the road link
3. Direction of travel	=	Direction of travel
4. Year	=	Year it happened
5. Count date	=	The date when the actual count took place
6. hour	=	Hour 7 represents from 7am to 8am, and 17 tells from 5pm to 6pm
7. Region id	=	Website region identifier
8. Region name	=	The name of the Region that travel took place
9. Local authority name	=	Local authority that region
10. Road name	=	This is the road name (for instance M25 or A3).
11. Road Category ID	=	Uniquely identifies road ID
12. Start junction road name	=	The road name of the start junction of the link
13. End junction road name	=	The road name of the end junction of the link
14. Latitude	=	Latitude of the Location
15. Longitude	=	Longitude of the Location
16. Link length km	=	Total length of the network road link
17. Pedal cycles	=	Counts for pedal cycles
18. Two wheeled motor vehicles	=	Counts of Two wheeled motor vehicles
19. Cars and taxis	=	Counts of Cars and taxis
20. Buses and coaches	=	Counts of Buses and coaches
21. LGV Type	=	Counts of LGV Type
22. HGV Type	=	Counts of HGV Type
23. EV Car	=	Counts of EV Car
24. EV Bike	=	Counts of EV Bike

Raw Roads Dataset

The raw roads dataset defines the road category. It provides essential metadata about the road network across different regions. It includes classifications, measurements, and summaries of road segments, which are later used for enriching traffic data and performing spatial analysis.

Source and storage

- The dataset is manually ingested into the /landing/raw_roads folder within Azure Data Lake Storage Gen2.



- It simulates external data ingestion and is incrementally loaded into the bronze layer using Spark Structured Streaming with Auto Loader.
- Stored as a Delta table named raw_roads in the bronze schema.

Purpose in the pipeline

The Raw Roads Dataset:

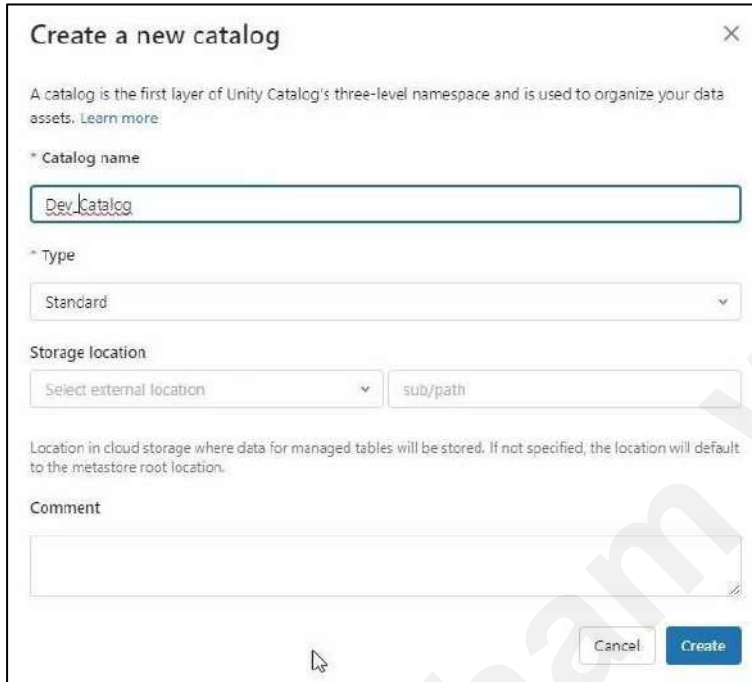
- Provides road classification and geographical context for traffic analysis.
- Helps join with traffic datasets using common region or road category IDs.

- Supports the derivation of road type attributes used in visualization and aggregation in the gold layer.

The common link or common column from both the datasets is Road Category.

Creating Databricks Dev Catalog

We created the dev workspace for Azure Databricks (**databricks-dev-ws**).



Create a new catalog

A catalog is the first layer of Unity Catalog's three-level namespace and is used to organize your data assets. [Learn more](#).

* Catalog name
Dev Catalog

* Type
Standard

Storage location
Select external location sub/path

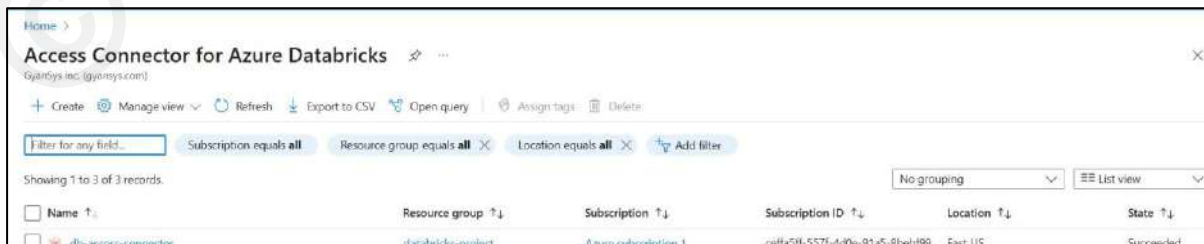
Location in cloud storage where data for managed tables will be stored. If not specified, the location will default to the metastore root location.

Comment

Cancel Create

Access Connector for Databricks

In this project, the Access Connector for Azure Databricks is used to securely access Azure Data Lake Storage Gen2, where all data layers (landing, bronze, silver, gold, checkpoints) are stored. It enables Databricks to read **raw traffic** and **raw roads** datasets and write Delta tables without using storage keys, by leveraging managed identity authentication. This ensures secure, role-based access control and is required for integrating with Unity Catalog's external locations.



Home > Access Connector for Azure Databricks

Gyrfors Inc. (@gyrfors.com)

+ Create Manage view Refresh Export to CSV Open query Assign tags Delete

Filter for any field... Subscription equals all Resource group equals all Location equals all Add filter

Showing 1 to 3 of 3 records. No grouping List view

Name	Resource group	Subscription	Subscription ID	Location	State
db-access-connector	databricks-project	Azure subscription 1	ceffa3ff-557f-4d0e-91a5-6beb199...	East US	Succeeded

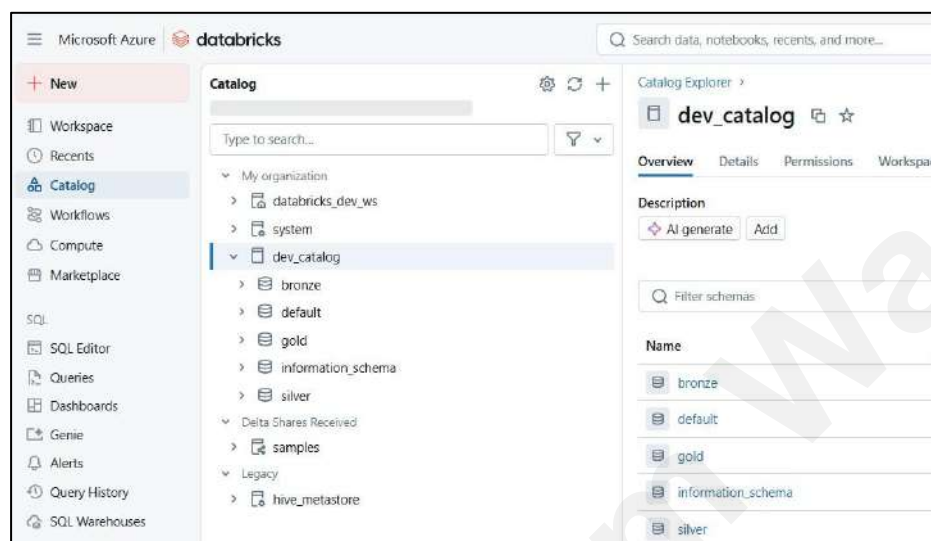
Metastore Creation

A metastore is a top-level container for data in Unity Catalog. Within a metastore, Unity Catalog provides a 3-level namespace for organizing data (catalogs, schemas, tables/views).

If we do not assign the workspace to a metastore we will not be able to create a catalog or schema. After assigning the workspace to the metastore we need to enable the Unity Catalog.

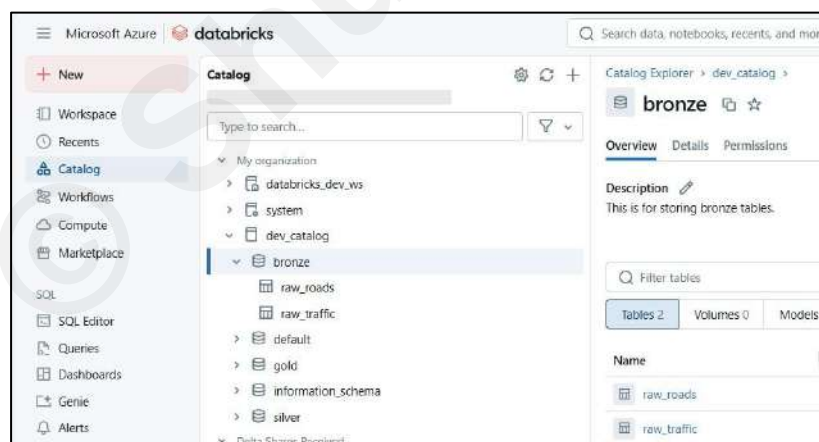
Creating all the Schemas in Dev Catalog

We have created 3 schemas in the dev catalog (bronze, silver, gold).



Ingestion to Bronze Layer

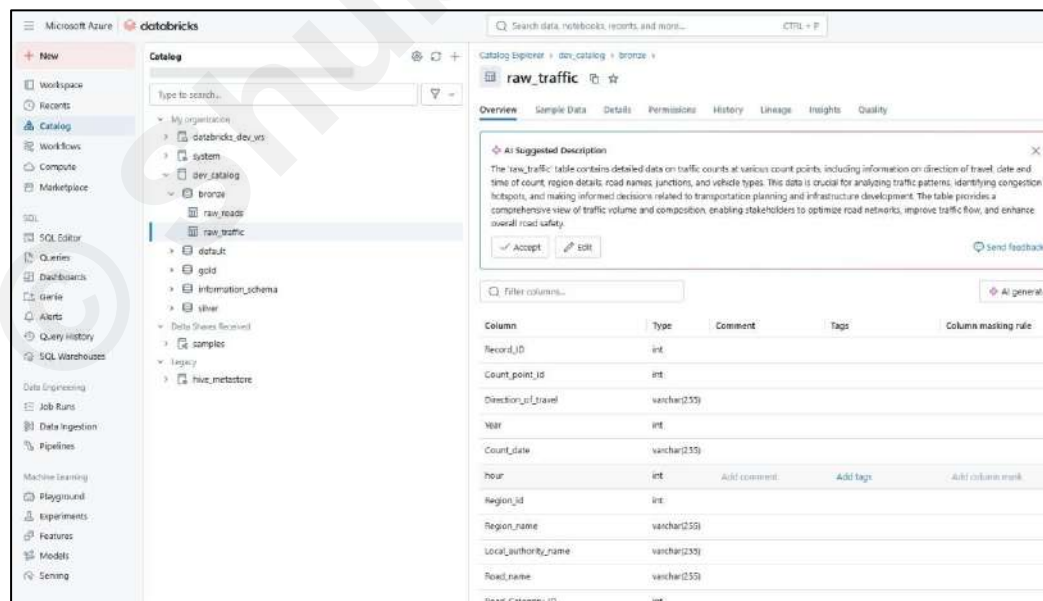
The bronze layer is the first structured zone in the medallion architecture where raw data is ingested from the landing zone and stored as Delta tables. It serves as the source of the truth, capturing unprocessed data exactly as received.



Ingestion Process

- Source:
 - Data is manually placed in:
 - /landing/raw_traffic (for traffic data)
 - /landing/raw_roads (for road data)
- Streaming Ingestion with Auto Loader:
 - Spark reads incoming files using Auto Loader. It improves pipeline efficiency by only processing newly arrived data.
 - Used to incrementally ingest raw traffic and road CSV files from the /landing folder.
 - We have created 2 autoloaders. One for raw_roads and the other for raw_traffic.
 - Enables real-time data ingestion into the bronze layer using Structured Streaming.
 - Reads data using .format("cloudFiles") with cloudFiles.format = "csv".
 - Stores schema information using cloudFiles.schemaLocation for schema inference.
 - Tracks progress using a checkpoint directory to ensure fault tolerance and supports automatic detection of new files, eliminating the need for manual triggers.
 - Loads data into Delta tables: bronze.raw_traffic and bronze.raw_roads.
- Schema: Bronze
- Tables created:
 - raw_traffic
 - raw_roads

Ingestion of raw_traffic table

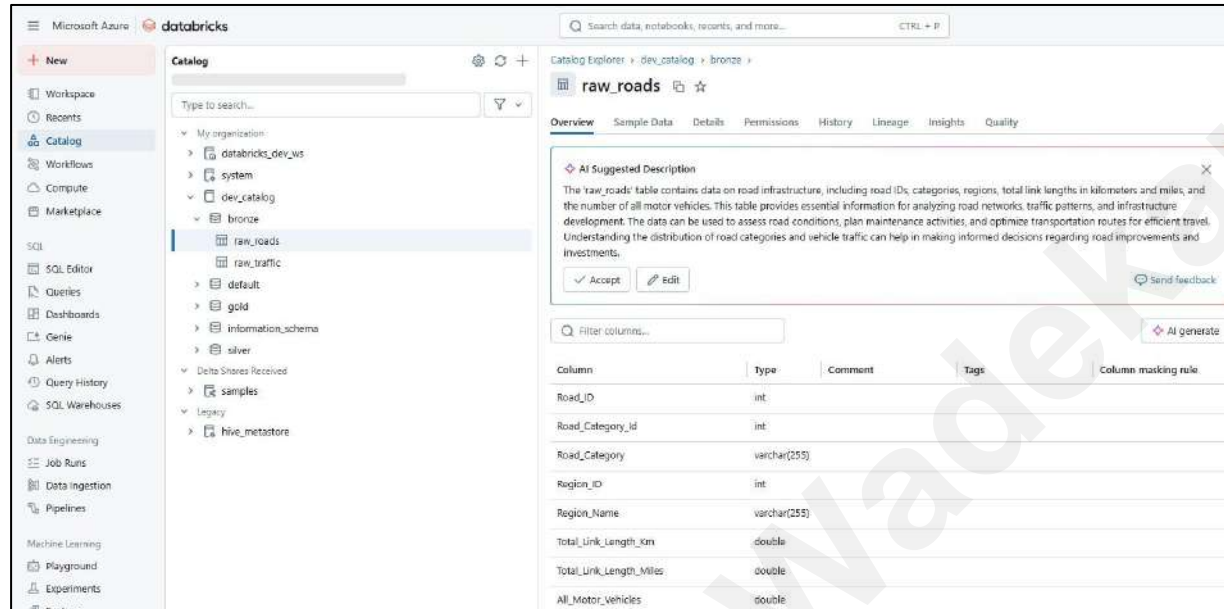


The screenshot displays the Databricks Catalog Explorer interface. On the left, a sidebar shows the navigation menu with options like Workspace, Recent, Catalog, Workflows, Compute, and Marketplace. The main panel shows the 'raw_traffic' table selected under the 'dev_catalog' and 'bronze' schema. The table's structure is detailed in the 'Overview' tab, showing columns such as Record_ID, Count_point_ID, Direction_of_travel, Year, Count_date, hour, Region_ID, Region_name, Local_authority_name, Road_name, and Road_Category_ID. The table is described as containing detailed data on traffic counts at various count points, including information on direction of travel, date and time of count, region details, road names, junctions, and vehicle types. The interface also includes a search bar, a filter column dropdown, and a table of column details with types and comments.

Column	Type	Comment	Tags	Column masking rule
Record_ID	int			
Count_point_ID	int			
Direction_of_travel	varchar(255)			
Year	int			
Count_date	varchar(255)			
hour	int			
Region_ID	int			
Region_name	varchar(255)			
Local_authority_name	varchar(255)			
Road_name	varchar(255)			
Road_Category_ID	int			

We have ingested the raw_traffic table into the bronze schema within the dev_catalog in the Databricks workspace (dataricks_dev_ws).

Ingestion of raw_roads table



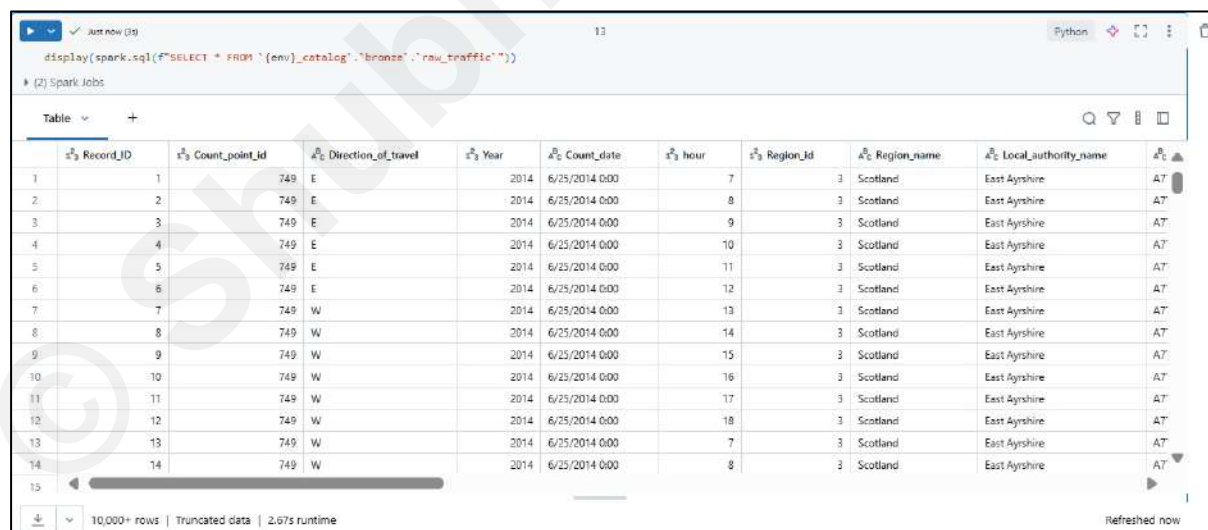
The screenshot shows the Databricks Catalog Explorer interface. On the left, the 'Catalog' sidebar lists the workspace structure, including 'dev_catalog' and its 'bronze' schema. The 'raw_roads' table is highlighted. The main panel displays the 'Overview' tab for the 'raw_roads' table, which includes an 'AI Suggested Description' and a table of columns with their data types.

Column	Type	Comment	Tags	Column masking rule
Road_ID	int			
Road_Category_Id	int			
Road_Category	varchar(255)			
Region_ID	int			
Region_Name	varchar(255)			
Total_Link_Length_Km	double			
Total_Link_Length_Miles	double			
All_Motor_Vehicles	double			

We have ingested the raw_roads table into the bronze schema within the dev_catalog in the Databricks workspace (dataricks_dev_ws).

Loading data to Bronze Tables

Data written for raw_traffic:



The screenshot shows a Databricks SQL query result. The query is: `display(spark.sql("SELECT * FROM 'dev_catalog'.bronze.raw_traffic"))`. The result is a table with 14 rows and 11 columns. The columns are: Record_ID, Count_point_Id, Direction_of_travel, Year, Count_date, hour, Region_Id, Region_name, Local_authority_name, and two unlabeled columns. The data shows traffic counts for various locations in Scotland, including East Ayrshire.

Record_ID	Count_point_Id	Direction_of_travel	Year	Count_date	hour	Region_Id	Region_name	Local_authority_name		
1	1	749 E	2014	6/25/2014 0:00	7	3	Scotland	East Ayrshire	A7	
2	2	749 E	2014	6/25/2014 0:00	8	3	Scotland	East Ayrshire	A7	
3	3	749 E	2014	6/25/2014 0:00	9	3	Scotland	East Ayrshire	A7	
4	4	749 E	2014	6/25/2014 0:00	10	3	Scotland	East Ayrshire	A7	
5	5	749 E	2014	6/25/2014 0:00	11	3	Scotland	East Ayrshire	A7	
6	6	749 E	2014	6/25/2014 0:00	12	3	Scotland	East Ayrshire	A7	
7	7	749 W	2014	6/25/2014 0:00	13	3	Scotland	East Ayrshire	A7	
8	8	749 W	2014	6/25/2014 0:00	14	3	Scotland	East Ayrshire	A7	
9	9	749 W	2014	6/25/2014 0:00	15	3	Scotland	East Ayrshire	A7	
10	10	749 W	2014	6/25/2014 0:00	16	3	Scotland	East Ayrshire	A7	
11	11	749 W	2014	6/25/2014 0:00	17	3	Scotland	East Ayrshire	A7	
12	12	749 W	2014	6/25/2014 0:00	18	3	Scotland	East Ayrshire	A7	
13	13	749 W	2014	6/25/2014 0:00	7	3	Scotland	East Ayrshire	A7	
14	14	749 W	2014	6/25/2014 0:00	8	3	Scotland	East Ayrshire	A7	

After defining the schema and reading the raw_traffic CSV file from the landing zone in using **Auto Loader**, the data was successfully written to the bronze layer in Delta format.

Data written for raw_roads:

Table

	1.1 Road_ID	1.2 Road_Category_Id	1.3 Road_Category	1.4 Region_ID	1.5 Region_Name	1.6 Total_Link_Length_Km	1.7 Total_Link_Length_Miles	1.8 All_Motor_Vehicles
1	1	1	TM	1	South West	301.339	187.24	346594
2	2	3	TA	1	South West	993.586	617.39	348471
3	3	4	PA	1	South West	3874.924	2407.77	779400
4	4	5	M	1	South West	43581.7	27080.41	911200
5	5	1	TM	2	East Midlands	178.609	110.98	273666
6	6	3	TA	2	East Midlands	1219.231	757.6	493670
7	7	4	PA	2	East Midlands	2571.21	1597.68	553572
8	8	5	M	2	East Midlands	26712.7	16598.5	708336
9	9	1	TM	3	Scotland	335.263	208.32	248546
10	10	3	TA	3	Scotland	2820.302	1752.45	505292
11	11	4	PA	3	Scotland	7405.483	4601.55	651443
12	12	5	M	3	Scotland	47971.6	29808.17	765643
13	13	1	TM	4	Wales	130.291	80.96	130658
14	14	3	TA	4	Wales	1509.635	938.04	312394

76 rows | 2.33s runtime

After defining the schema and reading the raw_roads CSV file from the landing zone using **Auto Loader**, the data was written to the bronze layer in Delta format.

Proving Autoloader handles incremental loading

Autoloader tracks based on the checkpoint location, and it is going to write this data incrementally by reading the only newly added rows.

Table

	1.1 count(1)
1	18546

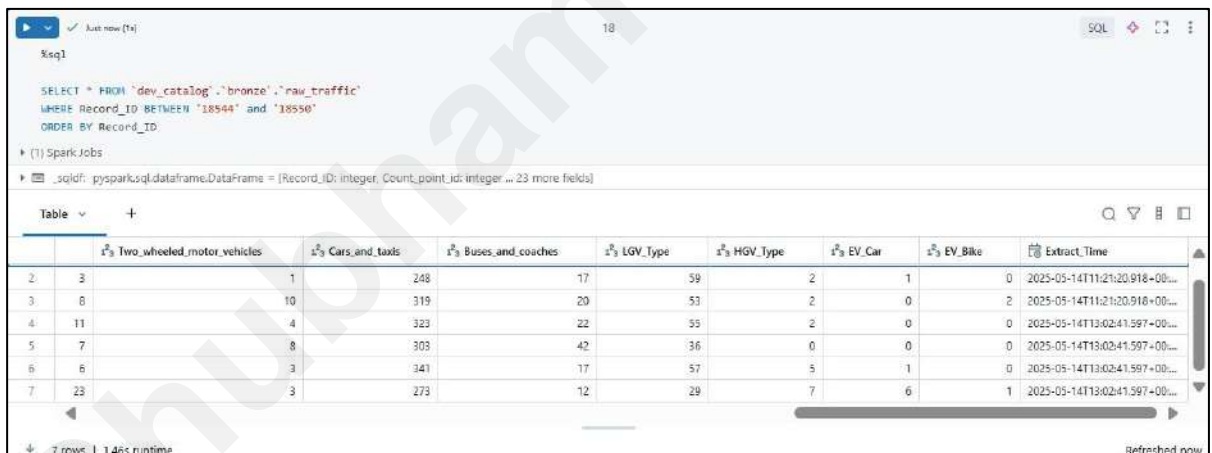
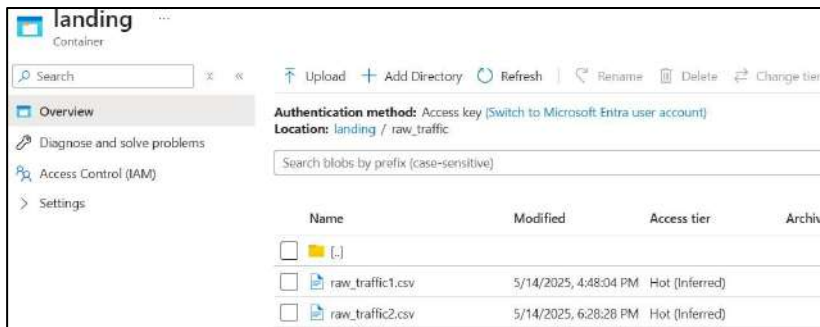
- Checking the row count of the raw_traffic data
- The same we can say that the last record_id column value would also be the same

Table

	1.1 max(Record_ID)
1	18546

number as the above row count.

- Now if we upload another file, those many new records will now be having new extract time value after 18546 records. There will be difference in the extract time value.



The count changed from 18546 to 37092. We can check the timestamp between the last few record_id till some of the newly added record_id.

- So, this proves that this is going to take that data based on a micro batch. For each micro batch, it is going to process all records which are available and is going to write the last record information somewhere in the checkpoint. When we upload another data, it goes to the checkpoint, and it is going to see where exactly the previous load was done and is going to compare that and then do the next load.

This proves that autoloader is capable to do the incremental loading.

- To run the notebook to process the newly added data, we need to monitor if there is any new file available. We can have it in a scheduled manner as well like twice a day or thrice a day. For now, we have just manually run the notebook to check the results.

Transforming Data into Silver Layer

Raw Traffic and Raw Roads data from the bronze layer is cleaned and enriched here.

- Schema: Silver
- Tables: silver_traffic, silver_roads

Transformations on raw_traffic silver_traffic:

- Renamed columns for easier querying and readability (e.g., Count point id → Count_point_id).
- Removed duplicates.
- Created Electric_Vehicles_Count = EV_Car + EV_Bike. It combines electric vehicle types to get total EV presence at a location.
- Created Motor_Vehicles_Count = Two_wheeled_motor_vehicles + Cars_and_taxis + Buses_and_coaches + LGV_Type + HGV_Type + Electric_Vehicles_Count. It calculates the total number of motorized vehicles for a given record.
- Derived Vehicle_Intensity = Motor_Vehicles_Count / Link_length_km to measure traffic density.
- Added timestamp columns like Extract_Time (from bronze) to track ingestion time.

▶

Just now (1s)

23

display(spark.sql(f"SELECT * FROM {env}._ceteolog.`silver`.`silver_traffic`"))

▶ (1) Spark Jobs

Table

▼

+

🔍

🔍

🔍

🔍

	🔍 LGV_Type	🔍 HGV_Type	🔍 EV_Car	🔍 EV_Bike	🕒 Extract_Time	🔍 Electric_Vehicles_Count	🔍 Motor_Vehicles_Count	🕒 Transformed_Time
1	9	209	38	12	92025-05-14T11:21:20.918+00...	21	1362	2025-05-14T12:17:32.859+00...
2	8	634	109	20	132025-05-14T11:21:20.918+00...	33	3083	2025-05-14T12:17:32.859+00...
3	15	55	3	0	22025-05-14T11:21:20.918+00...	2	623	2025-05-14T12:17:32.859+00...
4	18	256	58	10	62025-05-14T11:21:20.918+00...	16	2159	2025-05-14T12:17:32.859+00...
5	3	99	19	1	42025-05-14T11:21:20.918+00...	5	553	2025-05-14T12:17:32.859+00...
6	17	80	6	1	02025-05-14T11:21:20.918+00...	1	526	2025-05-14T12:17:32.859+00...
7	5	278	35	11	202025-05-14T11:21:20.918+00...	31	3119	2025-05-14T12:17:32.859+00...
8	20	270	53	10	52025-05-14T11:21:20.918+00...	15	2229	2025-05-14T12:17:32.859+00...
9	8	350	63	17	252025-05-14T11:21:20.918+00...	42	1745	2025-05-14T12:17:32.859+00...
10	11	589	140	33	402025-05-14T11:21:20.918+00...	73	3222	2025-05-14T12:17:32.859+00...
11	5	128	25	2	72025-05-14T11:21:20.918+00...	9	1061	2025-05-14T12:17:32.859+00...
12	7	94	20	0	02025-05-14T11:21:20.918+00...	0	2575	2025-05-14T12:17:32.859+00...
13	1	78	18	5	92025-05-14T11:21:20.918+00...	14	502	2025-05-14T12:17:32.859+00...
14	6	87	0	0	02025-05-14T11:21:20.918+00...	0	841	2025-05-14T12:17:32.859+00...
15								

▼

9,705+ rows | Truncated data | 1.24s runtime

Refreshed now

Transformations on raw_roads silver_roads:

- Renamed fields (e.g., Road category → Road_category).
- Created Road_Category_Name using mappings. It converts road category codes into human-readable names to improve dashboard clarity.
 - TA → Class A Trunk Road
 - TM → Class A Trunk Motorway
 - PA → Class A Principal Road
 - PM → Class A Principal Motorway
 - M → Class B Road
- Derived Road_Type groups road categories into broader classifications for filtering and aggregation.
 - If Road_Category_Name contains "Class A" → "Major"

- If Road_Category_Name contains "Class B" → "Minor"

```

%sql
SELECT * FROM `dev_catalog`.`silver`.`silver_roads`
  
```

(2) Spark Jobs

_sqlidf: pyspark.sql.dataframe.DataFrame = [Road_ID: integer, Road_Category_Id: integer ... 8 more fields]

Category	Region_ID	Region_Name	1.2 Total Link Length_Km	1.2 Total Link Length_Miles	1.2 All_Motor_Vehicles	Road_Category_Name	Road_Type
1	7	East of England	245.9	152.8	4120104199	Class A Trunk Motor	Major
2	10	West Midlands	27956.1	17371.12	9927697665	Class B road	Minor
3	7	East of England	2719.344	1689.72	7493977513	Class A Principal road	Major
4	8	Yorkshire and the Humber	2559.104	1590.15	6365559059	Class A Principal road	Major
5	10	West Midlands	351.963	218.7	5280343918	Class A Trunk Motor	Major
6	9	South East	40851.1	25383.7	15597621429	Class B road	Minor
7	2	East Midlands	1215.788	755.46	5056721450	Class A Trunk Road	Major
8	1	South West	3882.154	2412.26	8028746829	Class A Principal road	Major
9	6	London	66.386	41.25	1311689053	Class A Trunk Motor	Major
10	5	North West	558.639	347.25	7370068617	Class A Trunk Motor	Major
11	2	East Midlands	178.609	110.98	2736661520	Class A Trunk Motor	Major
12	2	East Midlands	1219.231	757.6	4936702686	Class A Trunk Road	Major
13	10	West Midlands	6.895	4.28	1464122624	Class A Principal Motorway	Major
14	3	Scotland	2813.698	1748.35	5192586960	Class A Trunk Road	Major

Transformations on Incrementally Loaded Data

Here we can see that only the newly added records were taken to process the data. It is because the bronze tables can have thousands of records every time, the incremental loading will be taken place from the landing zone to bronze, where incremental data will be appended to the bronze table and in point of time somewhere the records may be a million records. And if we are trying to do this silver layer transformation by creating a new column and applying the data that should not be applied on the entire data set each time, this should be applied only to the changed data, which means the rows which are newly added.

```

%sql
SELECT COUNT(*) FROM `dev_catalog`.`silver`.`silver_traffic`
  
```

(4) Spark Jobs

_sqlidf: pyspark.sql.dataframe.DataFrame = [count(1): long]

count(1)
37092

- Checking the count of the current records. It is 37092 after adding the 2nd traffic file.
- When we add the 3rd traffic file, the count changes to 55638. Now when we query the silver transformed data for traffic data, we can see the changes in transformed time for the new records only. The previous loaded data was the previous time when it was transformed. **So, this proves that the data was transformed incrementally and did not transform the old data.** This is possible because we are using the spark structure streaming and there is a delta lake for this table.

1 minute ago (2s)

19

SQL

```
SELECT * FROM 'dev_catalog','gold','gold_traffic'
```

(2) Spark Jobs

_sqldf: pyspark.sql.DataFrame = [Record_ID: integer, Count_point_id: integer ... 26 more fields]

		EV_Bike	Extract_Time	Electric_Vehicles_Count	Motor_Vehicles_Count	Transformed_Time	Vehicle_Intensity	Load_Time
1	12	9	2025-05-14T11:21:20.918+00...	21	1382	2025-05-14T12:17:32.859+00...	1151.6666666666667	2025-05-15T07:28:18.927+00...
2	20	13	2025-05-14T11:21:20.918+00...	33	3083	2025-05-14T12:17:32.859+00...	416.6216216216216	2025-05-15T07:28:18.927+00...
3	0	2	2025-05-14T11:21:20.918+00...	2	623	2025-05-14T12:17:32.859+00...	327.89473664210526	2025-05-15T07:28:18.927+00...
4	10	6	2025-05-14T11:21:20.918+00...	16	2159	2025-05-14T12:17:32.859+00...	674.6875	2025-05-15T07:28:18.927+00...
5	1	4	2025-05-14T11:21:20.918+00...	5	553	2025-05-14T12:17:32.859+00...	71.81818181818181	2025-05-15T07:28:18.927+00...
6	1	0	2025-05-14T11:21:20.918+00...	1	526	2025-05-14T12:17:32.859+00...	2630	2025-05-15T07:28:18.927+00...
7	11	20	2025-05-14T11:21:20.918+00...	31	3119	2025-05-14T12:17:32.859+00...	1356.0069565217392	2025-05-15T07:28:18.927+00...
8	10	5	2025-05-14T11:21:20.918+00...	15	2229	2025-05-14T12:17:32.859+00...	602.4324324324324	2025-05-15T07:28:18.927+00...
9	17	25	2025-05-14T11:21:20.918+00...	42	1745	2025-05-14T12:17:32.859+00...	1090.625	2025-05-15T07:28:18.927+00...
10	33	40	2025-05-14T11:21:20.918+00...	73	3222	2025-05-14T12:17:32.859+00...	847.8947366421053	2025-05-15T07:28:18.927+00...
11	2	7	2025-05-14T11:21:20.918+00...	9	1061	2025-05-14T12:17:32.859+00...	221.04166666666669	2025-05-15T07:28:18.927+00...
12	0	0	2025-05-14T11:21:20.918+00...	0	2575	2025-05-14T12:17:32.859+00...	1514.7058823529412	2025-05-15T07:28:18.927+00...
13	5	9	2025-05-14T11:21:20.918+00...	14	502	2025-05-14T12:17:32.859+00...	193.07692307692307	2025-05-15T07:28:18.927+00...
14	0	0	2025-05-14T11:21:20.918+00...	0	841	2025-05-14T12:17:32.859+00...	290	2025-05-15T07:28:18.927+00...
15								

2 minutes ago (1s)

18

SQL

```
SELECT * FROM 'dev_catalog','gold','gold_roads'
```

(2) Spark Jobs

_sqldf: pyspark.sql.DataFrame = [Road_ID: integer, Road_Category_ID: integer ... 9 more fields]

	Region_Name	1.2 Total_Link_Length_Km	1.2 Total_Link_Length_Miles	1.2 All_Motor_Vehicles	Road_Category_Name	Road_Type	Load_Time
1	7 East of England	245.9	152.8	4120104199	Class A Trunk Motor	Major	2025-05-15T07:28:34.918+00...
2	10 West Midlands	27956.1	17371.12	9927697665	Class B road	Minor	2025-05-15T07:28:34.918+00...
3	7 East of England	2719.344	1689.72	7493977513	Class A Principal road	Major	2025-05-15T07:28:34.918+00...
4	8 Yorkshire and the Humber	2559.104	1590.15	6365559059	Class A Principal road	Major	2025-05-15T07:28:34.918+00...
5	10 West Midlands	351.963	218.7	5280343918	Class A Trunk Motor	Major	2025-05-15T07:28:34.918+00...
6	9 South East	40851.1	25383.7	15597621429	Class B road	Minor	2025-05-15T07:28:34.918+00...
7	2 East Midlands	1215.788	755.46	5056721458	Class A Trunk Road	Major	2025-05-15T07:28:34.918+00...
8	1 South West	3882.154	2412.26	8028746829	Class A Principal road	Major	2025-05-15T07:28:34.918+00...
9	6 London	66.386	41.25	1311689053	Class A Trunk Motor	Major	2025-05-15T07:28:34.918+00...
10	5 North West	558.839	347.25	7370068517	Class A Trunk Motor	Major	2025-05-15T07:28:34.918+00...
11	2 East Midlands	178.609	110.98	2736661520	Class A Trunk Motor	Major	2025-05-15T07:28:34.918+00...
12	2 East Midlands	1219.231	757.6	4936702686	Class A Trunk Road	Major	2025-05-15T07:28:34.918+00...
13	10 West Midlands	6.895	4.28	146412262.4	Class A Principal Motorway	Major	2025-05-15T07:28:34.918+00...
14	3 Scotland	3813.698	1748.35	5192586980	Class A Trunk Road	Major	2025-05-15T07:28:34.918+00...

Jobs to orchestrate the flow

- There are certain notebooks which need to run daily to get the data, and some notebooks need not run daily. Here also we need to load data to bronze table, silver layer transformations and the gold transformations. So, all these notebooks need to be executed one after another. Based on the cadence when the data arrives, we need to run these notebooks in a flow.
- So, to run these notebooks in a flow, we need to give a job. A job will orchestrate all these notebooks, and it will run all these notebooks in a flow.
- Created a playground notebook that has the count of all the records. Counts in gold layer for both datasets:

- gold_traffic: 55638

Just now (1s)

-- Count of Gold Traffic Rows

```
SELECT COUNT(*) FROM `dev_catalog`.`gold`.gold_traffic
```

(3) Spark Jobs

Table ▾ +

	count(1)
1	55638

- gold_roads: 76

1 minute ago (1s)

-- Count of Gold Road Rows

```
SELECT COUNT(*) FROM `dev_catalog`.`gold`.gold_roads
```

(3) Spark Jobs

Table ▾ +

	count(1)
1	76

- Created a job named ETL Flow having various tasks:
 - TaskName: Load_to_Bronze
Cluster Used: Job Cluster (Once completed, it terminates)
 - TaskName: Silver_Traffic
 - TaskName: Silver_Roads
 - Task_Name: Gold



- We added new csv files in the landing zone for raw_traffic and raw_roads. Now we just need to run this ETL Flow workflow and check the counts. The new records will be added.



The screenshot shows the 'Runs' tab for the 'ETL Flow' workflow. A table lists the runs, with the most recent one being successful.

Start time	Run ID	Launched	Duration	Status	Error code	Run parameters
May 15, 2025, 02:53 PM	193446035213140	Manually	2m 23s	Succeeded		

The screenshot shows the results of a SQL query in Databricks. The query counts the number of rows in the 'gold_traffic' table.

```
-- Count of Gold Traffic Rows
SELECT COUNT(*) FROM `dev_catalog`.`gold`.gold_traffic
```

count(1)
74184

The screenshot shows the results of a SQL query in Databricks. The query counts the number of rows in the 'gold_roads' table.

```
-- Count of Gold Road Rows
SELECT COUNT(*) FROM `dev_catalog`.`gold`.gold_roads
```

count(1)
152

- The counts changed after running the ETL Flow.

- Added Triggers:
 - Raw_Traffic will be updated twice a day, not sure for the time so we use **file arrival** trigger here. (tracking for new file path – landing external location/raw_traffic/)

Schedules & Triggers

Trigger Status

☒ Active

☐ Paused

Trigger type

File arrival

This job does not have any failure notifications. Consider adding email or webhook notifications to alert if trigger evaluation fails.

File arrival triggers monitor cloud storage paths of up to 10,000 files for new files. These paths are either Unity Catalog volumes or external locations managed through Unity Catalog.

Storage location

abfss://landing@databricksdevstrg.dfs.core.windows.net/raw_traffic/

Advanced

Success Cancel Save

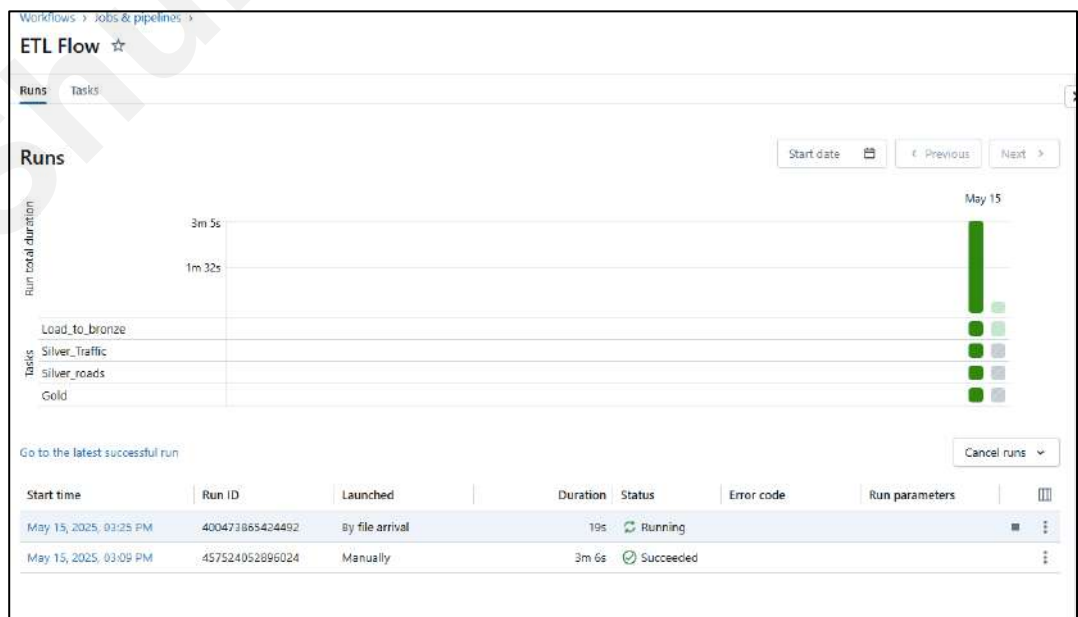
Schedules & Triggers

File arrival

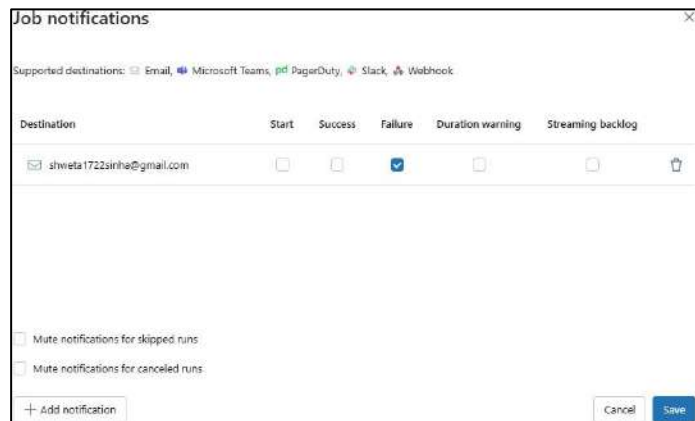
abfss://landing@databricksdevstrg.dfs.core.windows.net/raw_traffic/ (Every minute)

Edit trigger Pause Delete

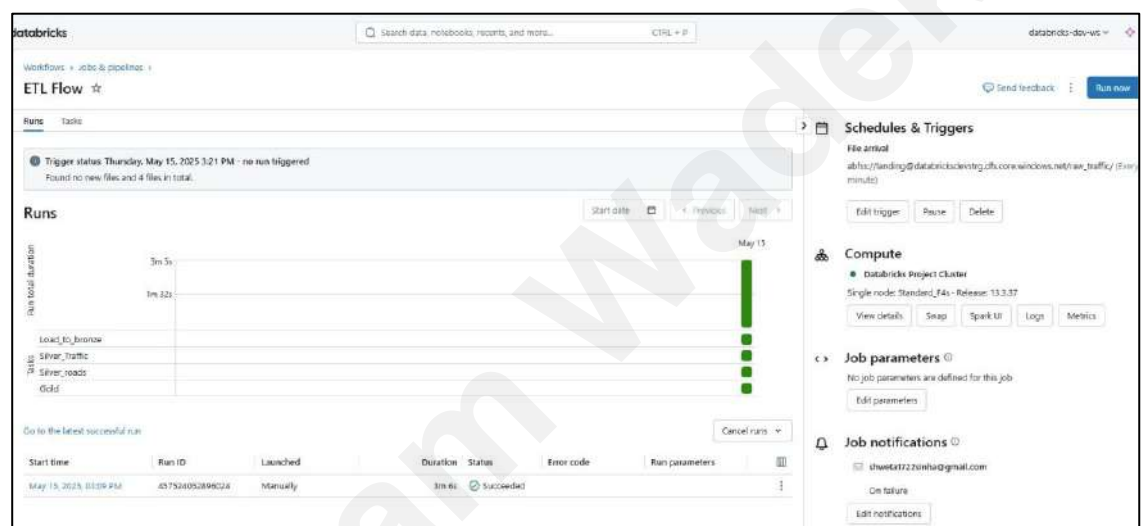
After every one minute it checks for the file as we have the file arrival trigger available. When we upload a new csv file for raw_traffic in the landing zone it will automatically run the ETL Flow job. In the Launched section we can see the job started running by the file arrival.



Added notification email on failure.

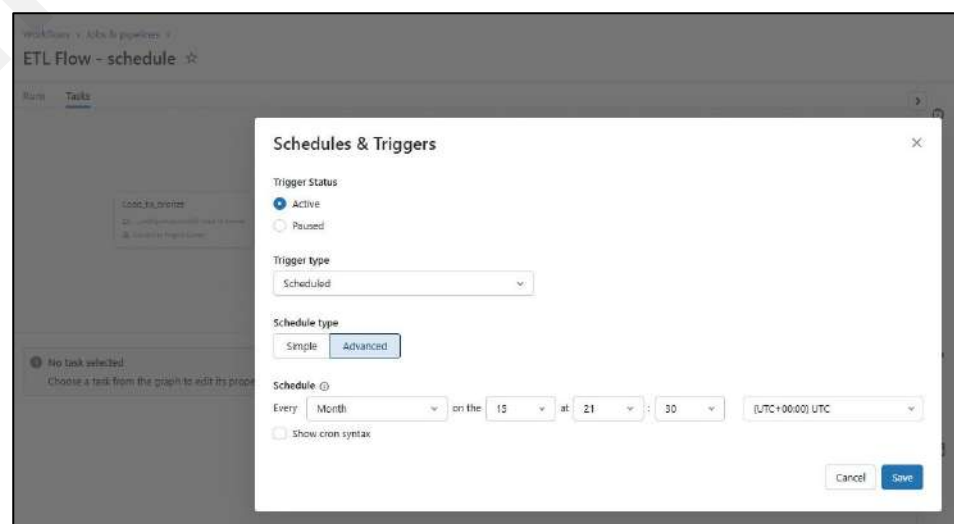


Job notifications configuration window. Supported destinations: Email, Microsoft Teams, PagerDuty, Slack, Webhook. Destination: shweta1722sirhe@gmail.com. Columns: Start, Success, Failure, Duration warning, Streaming backlog. Checkboxes: Mute notifications for skipped runs, Mute notifications for canceled runs. Buttons: + Add notification, Cancel, Save.

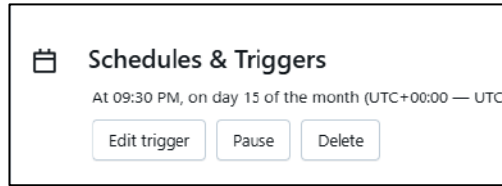


Databricks ETL Flow workflow overview. Trigger status: Thursday, May 15, 2025 3:21 PM - no run triggered. Found no new files and 4 files in total. Runs: Load_Cs_browse, Silver_traffic, Silver_roads, Gold. Schedules & Triggers: File arrival, abhish/landing@databricksdevtrg.cbf.com/windows.net/new_traffic/ (5 min). Compute: Databricks Project Cluster, Single node: Standard_f4s - Release: 13.3.37. Job parameters: No job parameters are defined for this job. Job notifications: shweta1722sirhe@gmail.com, On failure.

- Raw_Roads will be updated monthly basis, because this is kind of fixed (example – road length, type of road). We cloned the ETL Flow workflow and edited the trigger to scheduled type for every month at a particular time, because we cannot add 2 triggers within the same workflow.



ETL Flow - schedule configuration window. Schedules & Triggers. Trigger Status: Active. Trigger type: Scheduled. Schedule type: Simple. Schedule: Every Month on the 15 at 21:30 (UTC+00:00) UTC. Buttons: Cancel, Save.

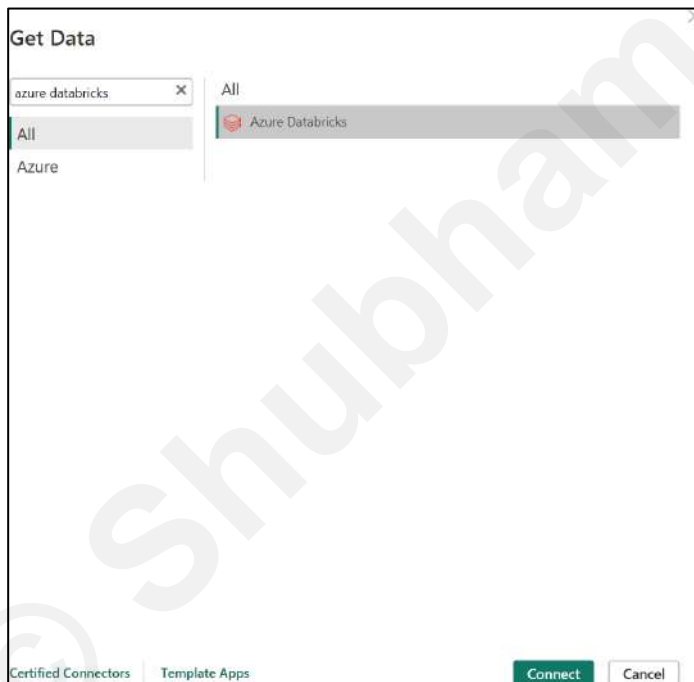


Reporting Data to PowerBI

PowerBI is used as the reporting and visualization tool to consume and present the Gold Layer data stored in Azure Data Lake via Azure Databricks. It provides interactive dashboards and data-driven insights from traffic and roads data. This will help to support decision-makers in analyzing traffic patterns, road utilization, and vehicle trends across different regions and timeframes.

PowerBI connection with Azure Databricks

To get the gold data, select the 'Get Data' in PowerBI and search for Azure Databricks. PowerBI connects to the Gold Layer Delta tables (gold_traffic, gold_roads).



Using the Databricks compute details we connect PowerBI to Azure Databricks as shown below.

Azure Databricks

Server Hostname ⓘ
 adb-3978535877632182.2.azuredatabricks.net

HTTP Path ⓘ
 sql/protocolv1/o/3978535877632182/0513-105057-5otdpysm

Advanced Options (optional)

Default catalog (optional) ⓘ
 dev_catalog

Database (optional) ⓘ
 Example: abc

Automatic Proxy Discovery (optional) ⓘ
 [Dropdown]

Native query (Requires: Default catalog) (optional) ⓘ
 Example: select * from db.schemaname.tablename

OK Cancel

Navigator

Display Options ⓘ

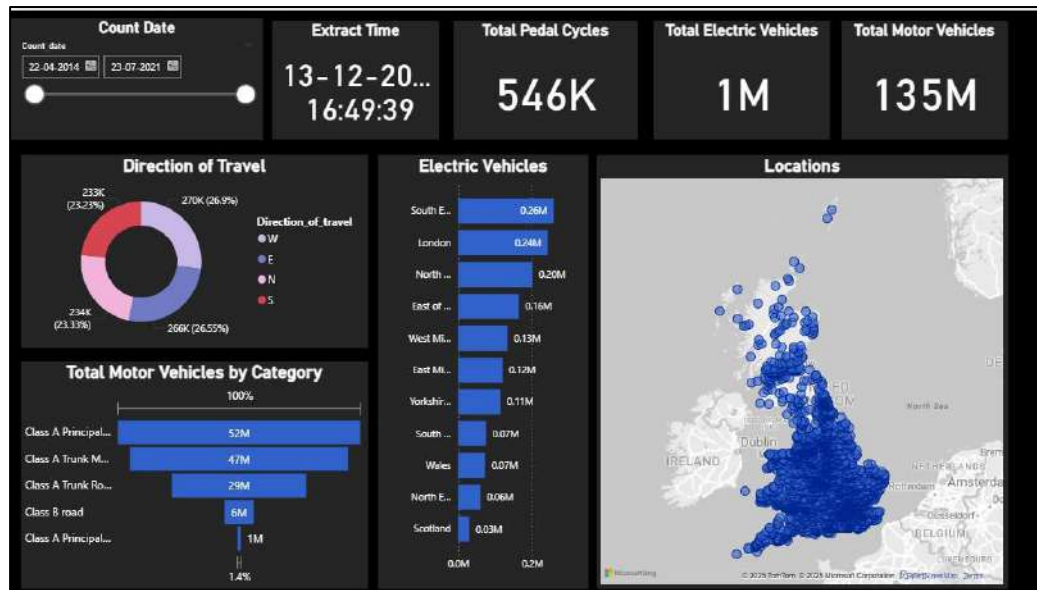
- adb-3978535877632182.2.azuredatabricks.net...
 - databricks_dev_ws
 - dev_catalog [4]
 - bronze
 - default
 - gold [2]
 - ☒ gold_roads
 - ☒ gold_traffic
 - silver
 - hive_metastore
 - samples

gold_traffic

Record_ID	Count_point_id	Direction_of_travel	Year	Count_date
255	540	S	2014	5/9/2014 0:00
297	6005	S	2016	9/7/2016 0:00
805	581	W	2015	9/23/2015 0:00
1232	16594	E	2014	7/1/2014 0:00
1331	7091	S	2014	6/5/2014 0:00
1376	64	N	2014	9/18/2014 0:00
1383	559	E	2014	9/24/2014 0:00
1576	16468	S	2014	5/8/2014 0:00
1674	18222	N	2014	6/25/2014 0:00
2100	6007	E	2016	7/5/2016 0:00
2309	26284	N	2014	5/23/2014 0:00
2488	6928	N	2014	10/16/2014 0:00
2659	18577	N	2014	6/13/2014 0:00
2721	6167	N	2016	5/18/2016 0:00
3096	6005	N	2014	7/10/2014 0:00
3535	16532	W	2014	4/25/2014 0:00

The data in the preview has been truncated due to size limits.

Load Transform Data Cancel

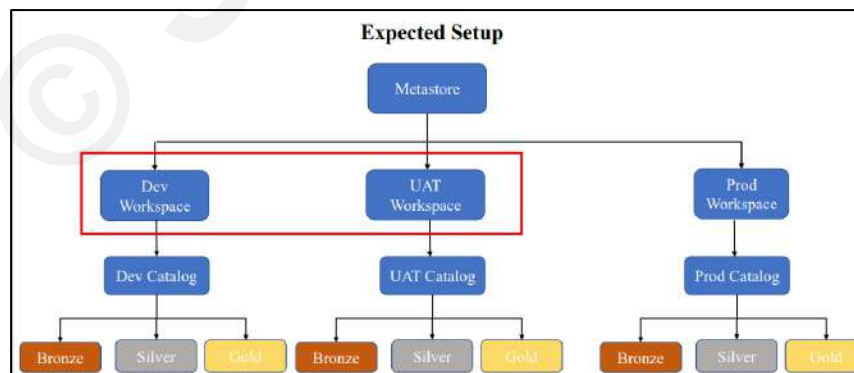


This PowerBI dashboard presents key insights from the gold layer of the project using curated traffic and road data.

Key points of the dashboard:

- Date filters: Users can filter data by count date and view when the data was last ingested.
- Extract Time: Shows when the dashboard is refreshed.
- KPI Tiles: Show total counts of:
 - Pedal Cycles
 - Electric Vehicles
 - Motor Vehicles
- Direction of Travel (Donut Chart): Shows vehicle distribution by travel direction (N, S, E, W).
- Electric Vehicles by Region (Bar Chart): Highlights regional EV usage, with Southeast and London leading.
- Motor Vehicles by Road Category: Displays total vehicle count across road types (e.g., Class A Principal, Class B).
- Location Map: Geospatial view of all traffic count points across the UK.

Continuous Integration & Continuous Deployment (CI/CD)



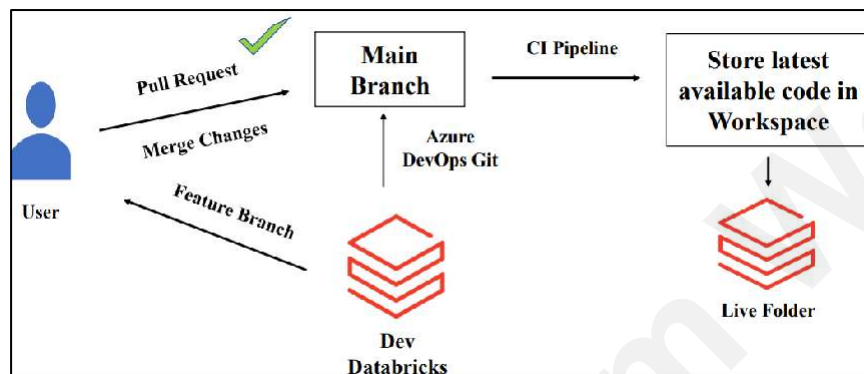
In this project, we have used Azure DevOps git as our repository. We have created the UAT workspace for testing and we have implemented the CI/CD. This project uses Azure DevOps to implement a robust Continuous Integration and Continuous Deployment (CI/CD) pipeline for managing notebooks, workflows, and configurations in Databricks. To ensure that code updates are version-controlled and tested before deployment.

Since this is a sample project, we have created the UAT workspace and we have implemented the CI/CD, where we have all the data from the dev workspace to the UAT workspace.

We have created the catalog and dynamically created all the schemas to represent the medallion architecture.

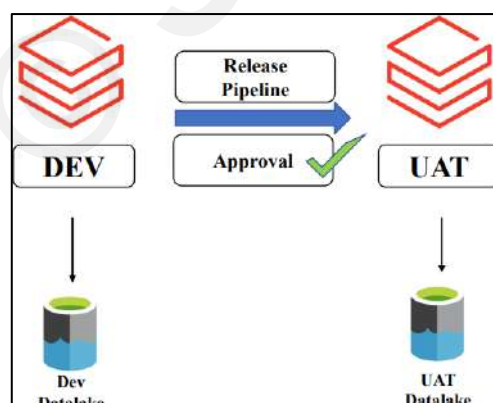
Data from dev Data to UAT

Continuous Integration



- Git repository used: Azure DevOps git
- Main branch holds all the changes done to the project. Whatever we work on notebooks, it is stored in a centralized place that is called the main branch.
- Continuous Integration lets multiple developers work and all the changes are merged with the main branch repository.
- Pull requests need to be approved by a technical lead or technical head in the team, who has access to Azure DevOps.

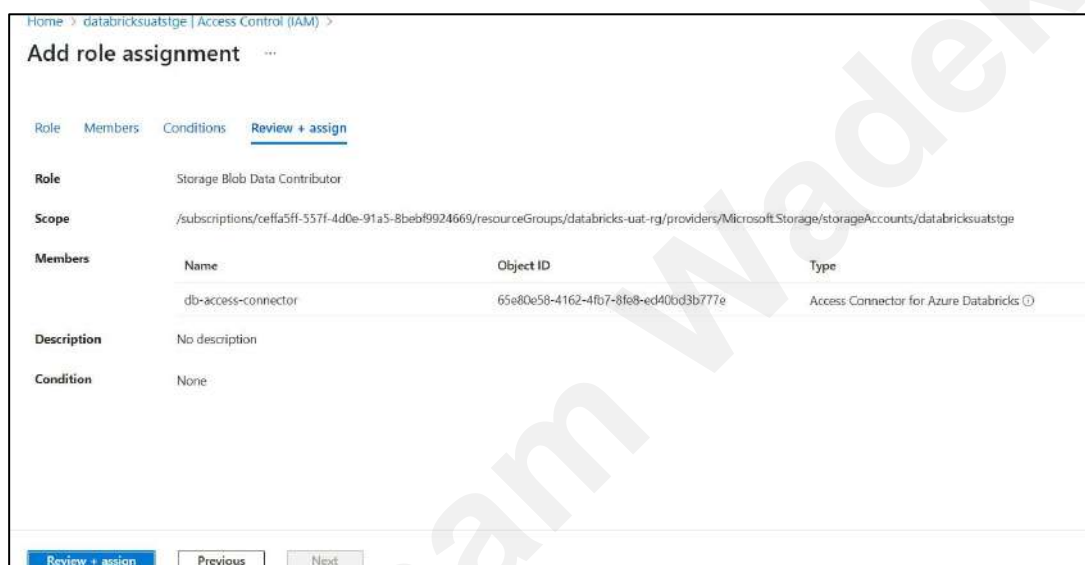
Continuous Deployment



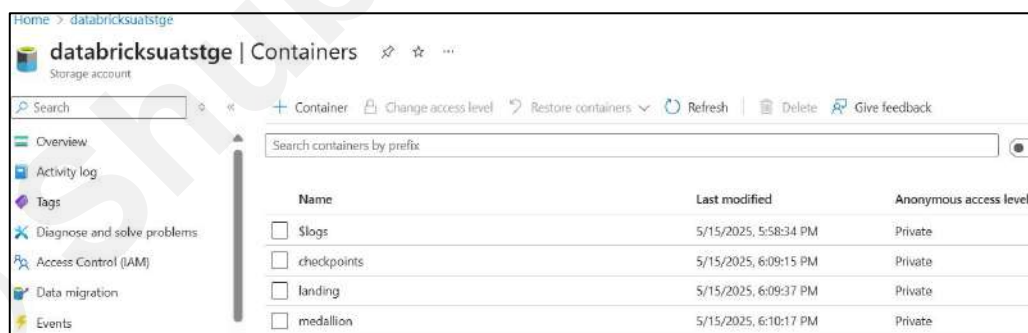
Release Pipeline: It gets all the changes in a live folder to UAT workspace. This will go on after an approval system has been done.

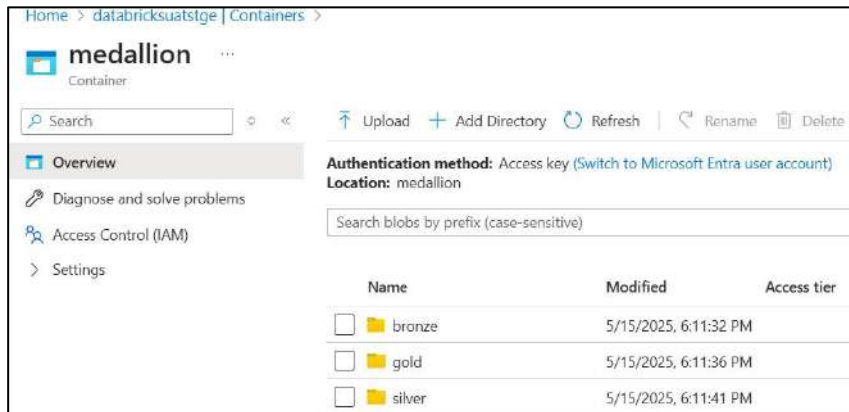
UAT Resources

- Resource Group: databricks-uat-rg
- Databricks workspace: databricks-uat-wsp
- Storage account: databricksuatstge
- Provided role assignment as storage blob data contributor to db-access-connector (Access connector for Azure Databricks) and gave the managed identity. Now UAT workspace will be a part of Unity Catalog.

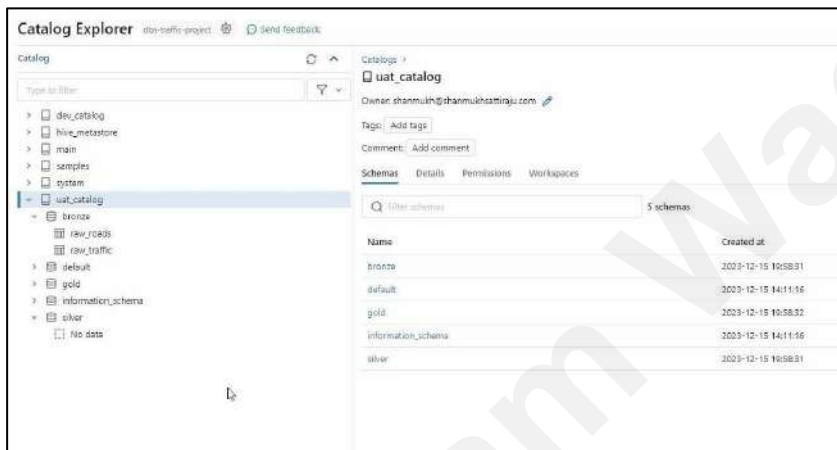


- Created the same containers as we created for the dev workspace.





- Created uat_catalog with all the schemas.



- Created another credential for creating external locations.

Create a new credential

A storage credential represents an authentication and authorization mechanism for accessing data stored on your cloud tenant. [Learn more](#)

Credential Type*

Azure Managed Identity

Credential name*

UAT-Access

Access connector ID [Learn more](#)

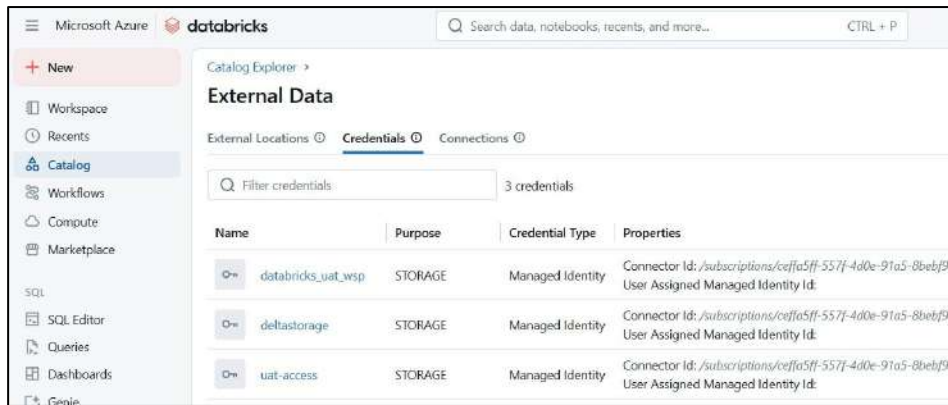
/subscriptions/ceffa5ff-557f-4d0e-91a5-8bebf9924669/resourceGroups/databricks-project/prov

User assigned managed identity ID (optional)

Comment

Cancel

Create



- Created 5 external locations for UAT also
 - bronze-uat
 - silver-uat
 - gold-uat
 - checkpoints-uat
 - landing-uat

External Locations		
Name	Credential	URL
bronze	deltastorage	abfss://medallion@databricksdevstg.dfs.core.windows.net/bronze
bronze-uat	uat-access	abfss://medallion@databricksuatstg.dfs.core.windows.net/bronze
checkpoints	deltastorage	abfss://checkpoints@databricksdevstg.dfs.core.windows.net/
checkpoints-uat	uat-access	abfss://checkpoints@databricksuatstg.dfs.core.windows.net/
databricks_uat_wsp	databricks_uat_wsp	abfss://unity-catalog-storage@dbstorage7jgw6cscatqum.dfs.core.wind...
gold	deltastorage	abfss://medallion@databricksdevstg.dfs.core.windows.net/gold
gold-uat	uat-access	abfss://medallion@databricksuatstg.dfs.core.windows.net/gold
landing	deltastorage	abfss://landing@databricksdevstg.dfs.core.windows.net/
landing-uat	uat-access	abfss://landing@databricksuatstg.dfs.core.windows.net/
silver	deltastorage	abfss://medallion@databricksdevstg.dfs.core.windows.net/silver
silver-uat	uat-access	abfss://medallion@databricksuatstg.dfs.core.windows.net/silver

Integrating Azure DevOps with Azure Databricks

Continuous Integration:

Created a new project in Azure DevOps named as databricks-traffic.

Create new project

Project name *

databricks-traffic

Description

Visibility

Public

Anyone on the internet can view the project. Certain features like TFVC are not supported.

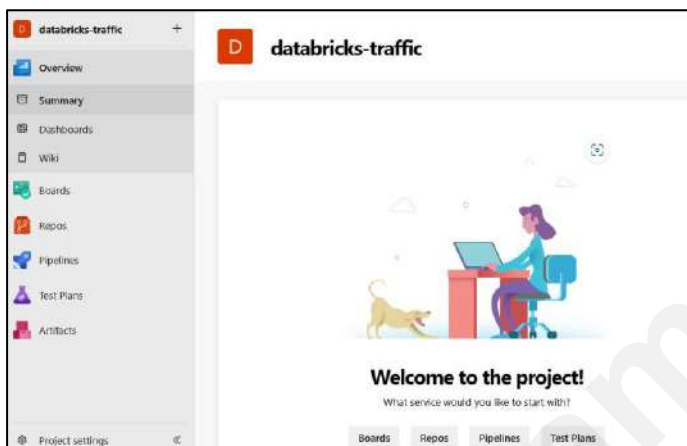
Private

Only people you give access to will be able to view this project.

Public projects are disabled for your organization. You can turn on public visibility with [organization policies](#).

Advanced

Cancel Create



- Created a repository in Azure DevOps.

Create a repository

Repository type

Git

Repository name *

dbproject

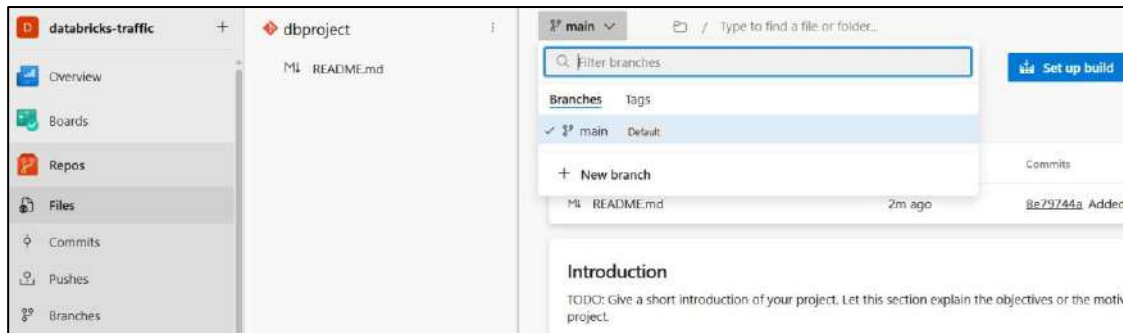
☒ Add a README

Add a .gitignore: None

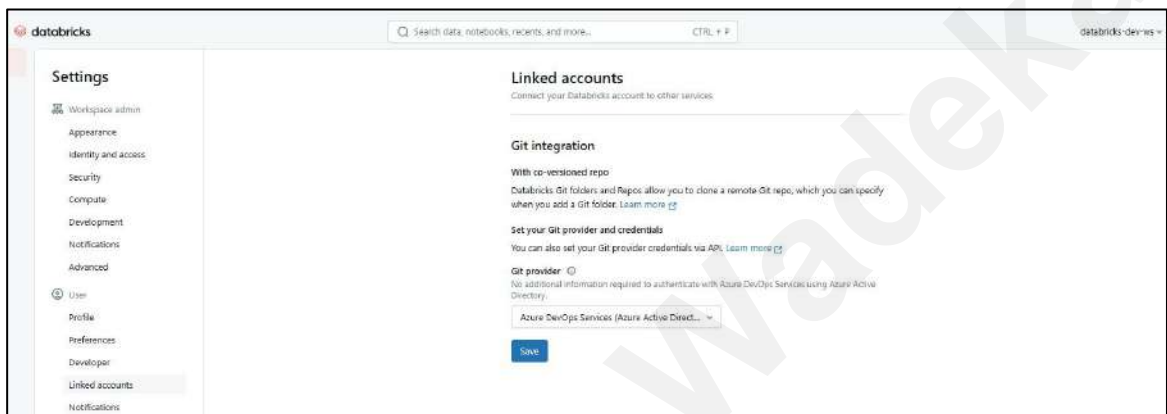
Your repository will be initialized with a `main` branch.

Cancel Create

We can see the main branch in our repository (dbproject) as shown below. This is the place where every code will be copied. This is used as the central repository.



- Link your Azure databricks dev workspace to Azure DevOps Services (Azure Active

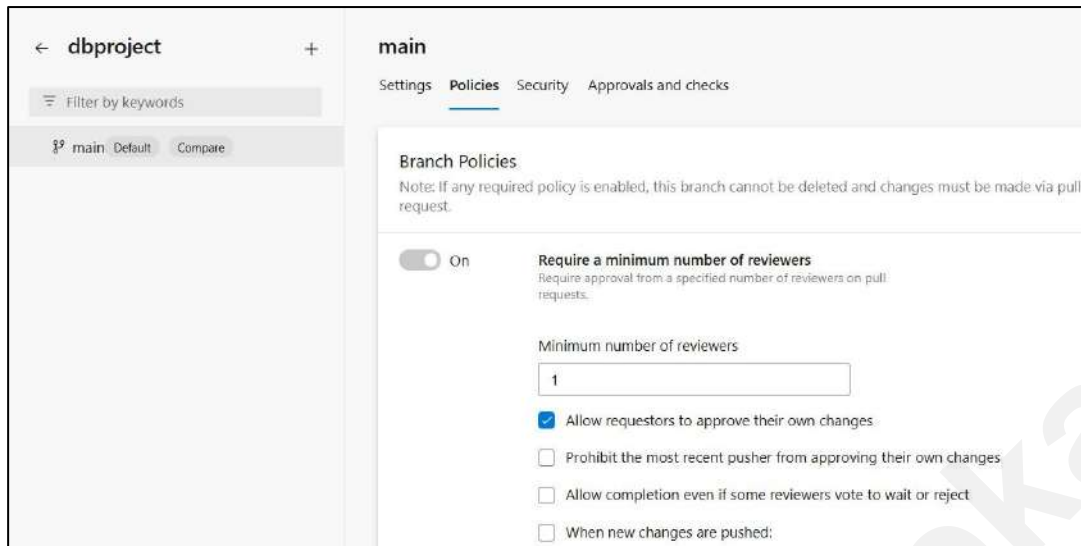


Directory) in User Settings.

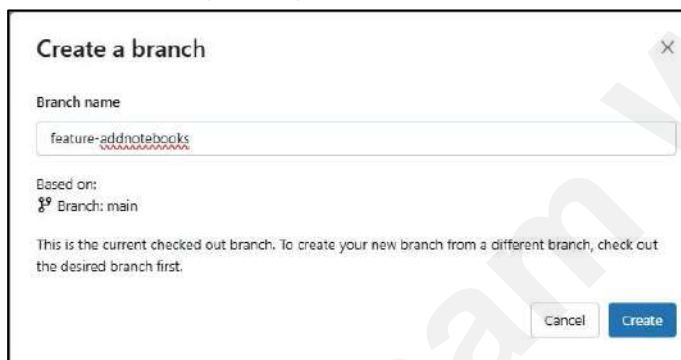
- Created our own repository in Azure Databricks so that we can integrate Azure Databricks with Azure DevOps. We cloned our new project created in Azure DevOps and copied the http link and pasted while creating a Repo to get all its details.



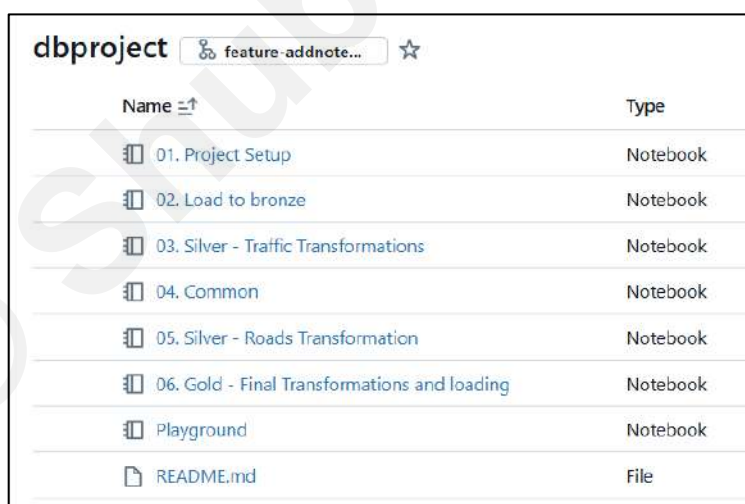
- In Azure DevOps we have set the minimum number of reviewers to be 1. So, when we have a pull request, it needs to be approved by any technical person in this project which I am doing. Mostly in every project we have more than 1 so we would be needing more people to approve then. Since I am the only one working here so I will approve my own pull request here.



- We can see we have a main branch and currently we can see it is empty. We need to have all our codes in the main branch, and for that we have created a feature branch in which we will create a pull request.



- After creating our feature branch, we moved all the codes from the user's workspace to

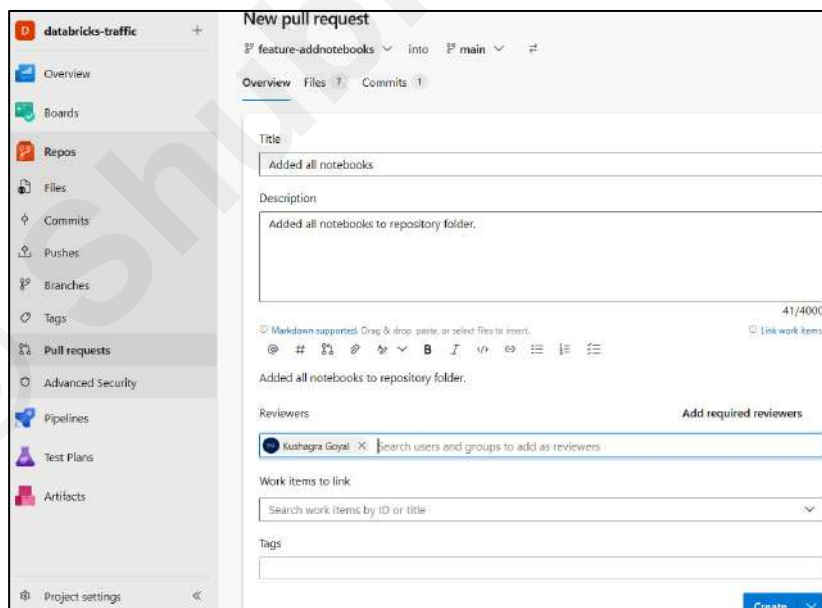
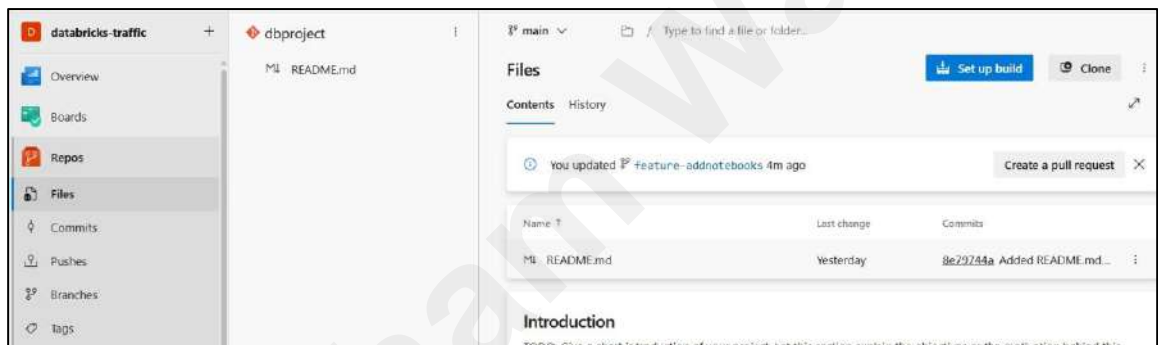


our own repository in our feature branch (feature-addnotebooks) for this project.

- We need to save all these changes to our feature branch. So, we will commit and push the changes to our feature branch as shown below.



- Since we did the commit and push to our feature branch, in Azure DevOps we will get notified in our main branch that a new branch has been created that is having some commit. Based on that commit we will be creating the pull request.



- Now we need to create the pull request.

Complete pull request

×

Merge type

Merge (no fast forward)

Post-completion options

☒ Complete associated work items after merging

☒ Delete feature-branch after merging

☐ Customize merge commit message

- Approved the pull request and completed it.
- Now we can see that all our codes are available in the main branch in Azure DevOps and in

databricks-traffic

Overview

Boards

Repos

Files

Commits

Pushes

Branches

Tags

Pull requests

Advanced Security

Pipelines

Test Plans

dbproject

PY 01. Project Setup.py
PY 02. Load to bronze.py
PY 03. Silver - Traffic Transformati
PY 04. Common.py
PY 05. Silver - Roads Transformat
PY 06. Gold - Final Transformation
Playground.sql
README.md

main

Type to find a file or folder...

Files

Contents History

Name ↑	Last change	Commits
PY 01. Project Setup.py	10m ago	06785da3 Added all notebook
PY 02. Load to bronze.py	10m ago	06785da3 Added all notebook
PY 03. Silver - Traffic Transformations.py	10m ago	06785da3 Added all notebook
PY 04. Common.py	10m ago	06785da3 Added all notebook
PY 05. Silver - Roads Transformation.py	10m ago	06785da3 Added all notebook
PY 06. Gold - Final Transformations and loading.py	10m ago	06785da3 Added all notebook
Playground.sql	10m ago	06785da3 Added all notebook
README.md	Yesterday	8e7974a Added README.m

Azure Databricks as well.

dbproject main ☆	
Name ↕	Type
01. Project Setup	Notebook
02. Load to bronze	Notebook
03. Silver - Traffic Transformations	Notebook
04. Common	Notebook
05. Silver - Roads Transformation	Notebook
06. Gold - Final Transformations and loading	Notebook
Playground	Notebook
README.md	File

Create a branch

Branch name

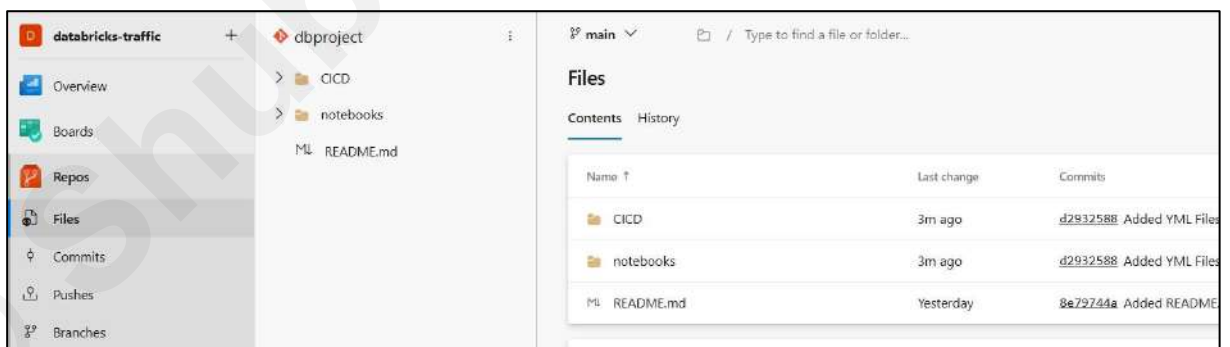
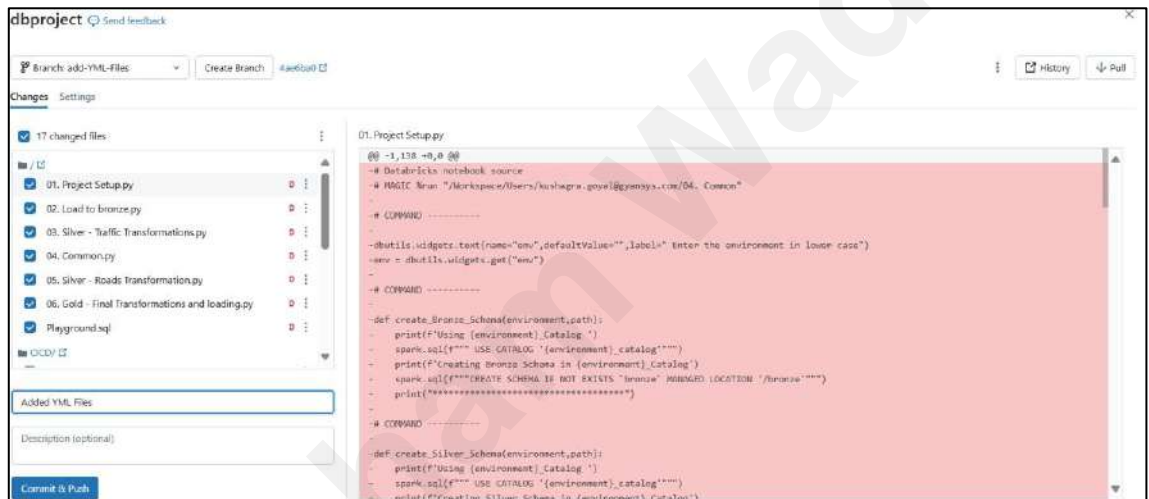
Based on:
 Branch: main

This is the current checked out branch. To create your new branch from a different branch, check out the desired branch first.

Cancel

Create

- Created a feature branch to add some YML files.
- Made separate folders for CI/CD and all notebooks. Committed and pushed the changes to the feature branch and created and completed a pull request in Azure DevOps.



- The CI/CD YML File checks that when there is a change in the main branch, it needs to trigger the CI pipeline, and it needs to deploy notebooks in the live folder.

- Created a library (dev-cicd-grp) and added all the required variables in the library.

The screenshot shows the 'dev-cicd-grp' library configuration page in Databricks. The left sidebar contains navigation links: Overview, Boards, Repos, Pipelines, Environments, Library (selected), Test Plans, and Artifacts. The main content area is titled 'Library > dev-cicd-grp*' and includes tabs for 'Variable group', 'Save', 'Clone', 'Security', 'Pipeline permissions', 'Approvals and checks', and 'Help'. Under the 'Properties' section, the 'Variable group name' is 'dev-cicd-grp' and the 'Description' field is empty. A toggle for 'Link secrets from an Azure key vault as variables' is turned off. The 'Variables' section contains a table with the following data:

Name	Value
dev-env-name	dev-cicd-name
dev-resource-group-name	databricks-project
dev-service-connection-name	dev-service-name

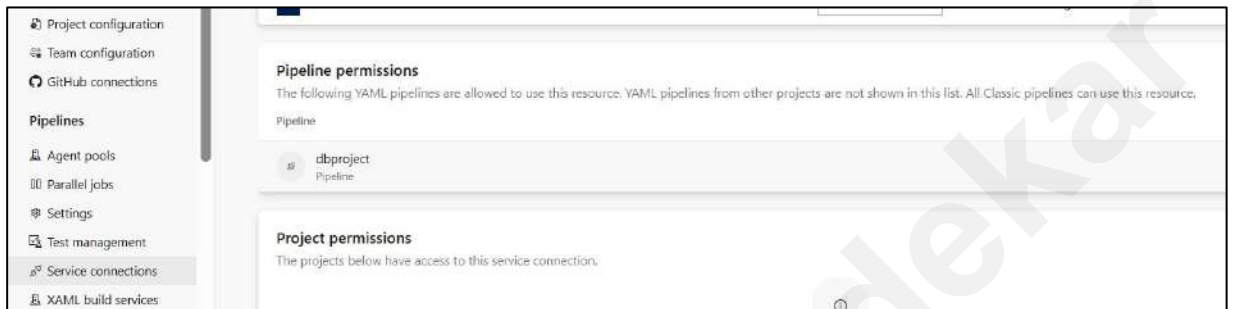
The screenshot shows the 'Configure your pipeline' page in Databricks. The left sidebar is the same as the previous screenshot. The main content area has tabs for 'Connect', 'Select', 'Configure' (selected), and 'Review'. Under 'Configure your pipeline', there are four options: 'Python package', 'Python to Linux Web App on Azure', 'Starter pipeline', and 'Existing Azure Pipelines YAML file'. The 'Existing Azure Pipelines YAML file' option is selected. On the right, a 'Select an existing YAML file' dialog is open, showing 'Branch' as 'main' and 'Path' as '/CICD/cicd-main.yml'. A 'dbproject' icon is visible at the bottom right of the dialog.

- Created a pipeline and saved it.

The screenshot shows the 'dev-cicd-grp' library configuration page, similar to the first screenshot. A 'Pipeline permissions' dialog box is open in the foreground. The dialog has a title bar 'dev-cicd-grp' and a close button. It shows a search bar with 'dbproject' entered. Below the search bar, it says 'No permitted pipelines' and 'This resource cannot be used in a YAML pipeline until at least one pipeline has permission. All Classic pipelines can use this resource. Learn more'. The background shows the same library configuration page as before.

- Gave the permission to the library that we created with the new pipeline.

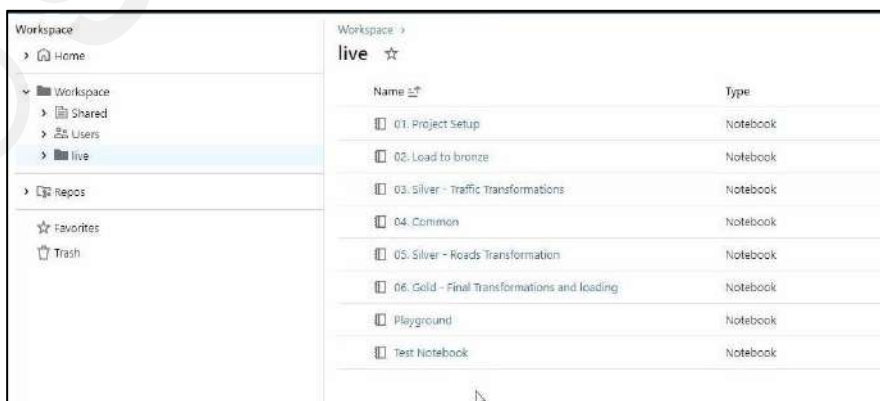
- Added pipeline permission to the environment and service connection in their security option.



- To test the pipeline that we created in Azure DevOps we created a test notebook in our feature branch and merged it with the main branch. We can see that the pipeline will start running on its own as soon as there is any change in the main branch.

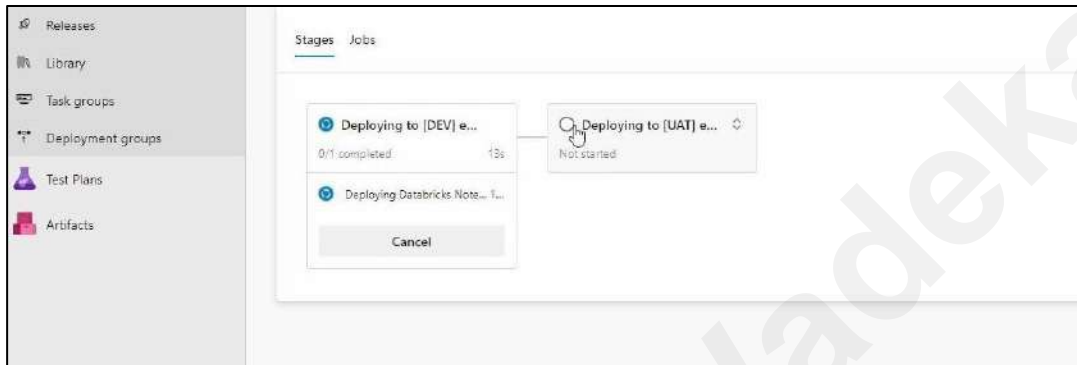


- After running the pipeline, we can see that there is a live folder created in Azure Databricks which has all the notebooks. So, everytime whoever pushes their changes to main branch all the codes will be copied to the live folder.

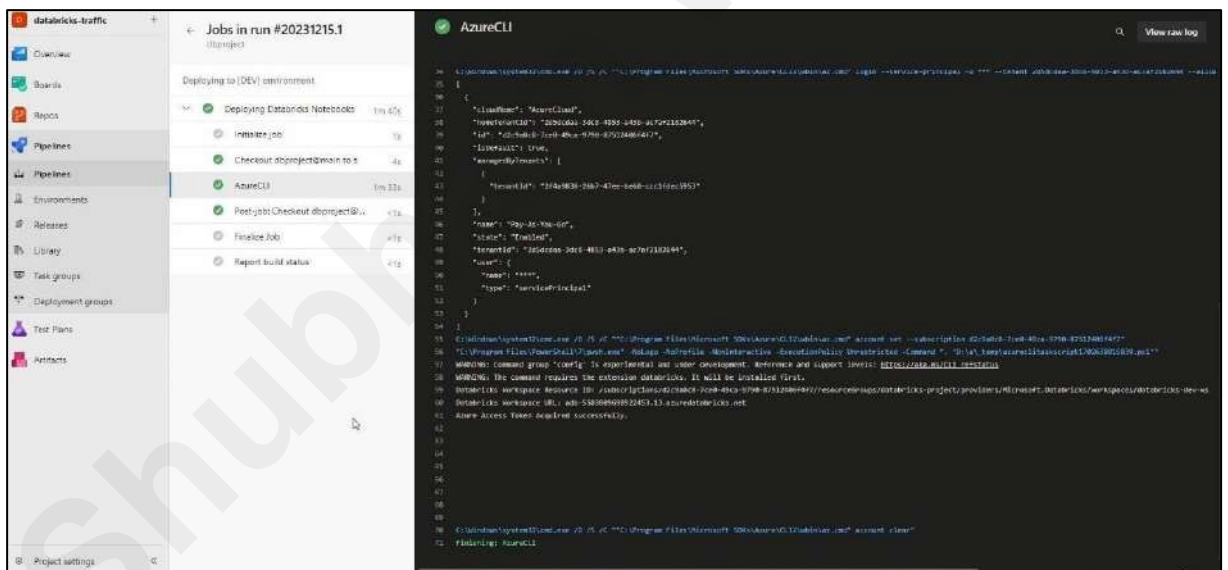


Continuous Deployment:

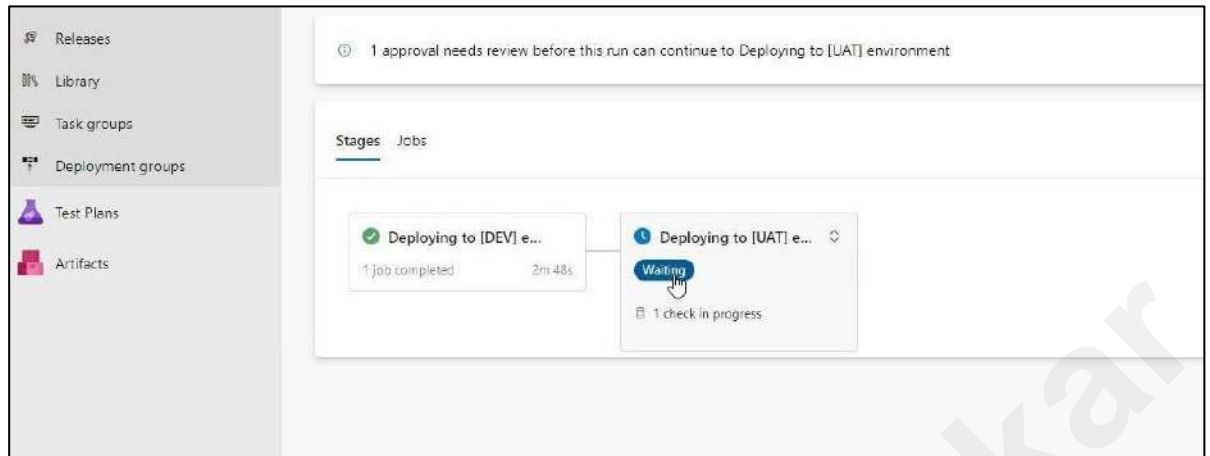
- Previously we added a group where we need to get the credentials of all dev workspaces. Since we want to deploy this to UAT we need to create a group that would hold all the variables of UAT environment. Then create variables for the UAT environment just like we did for dev.
- Tested with a change we can see that the main CI pipeline started running and it got



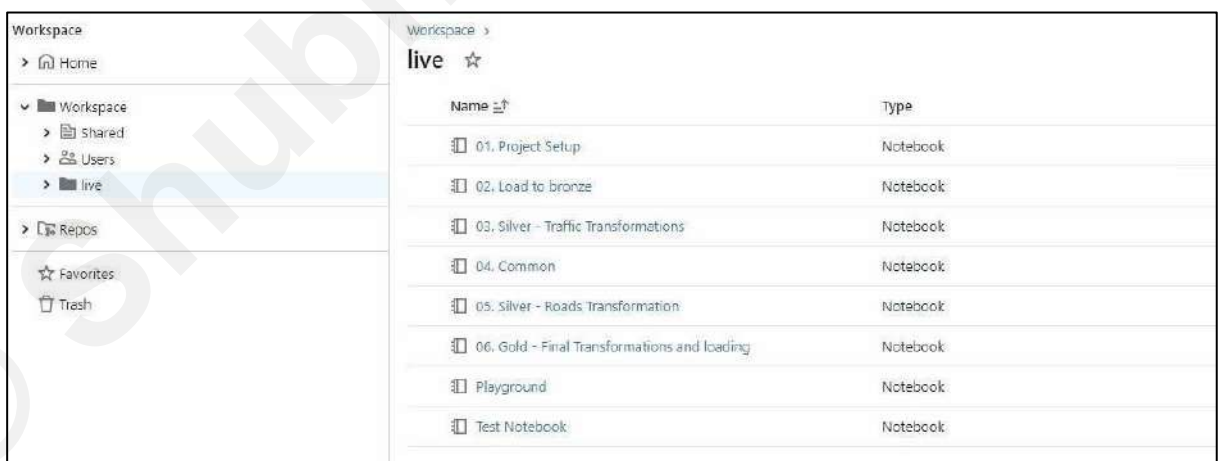
deployed to the dev environment.



- After deploying it to dev, it is waiting for the approval so that it can be deployed to the UAT environment. We can go and confirm this and give the approval of the required user.



- After deploying it to UAT we can see the live folder in Azure Databricks with all the codes in our UAT workspace also.



- After running all the notebooks, we can get all the data in UAT environment.

Road_ID	Road_Category_ID	Road_Category	Region_ID	Region_Name	Total_Link_Length_Km	Total_Link_Length_Miles	All_Motor_Vehicles	Road_Category_Name	Road_Class
99	1	TM	1	South West	300.899	186.99	3779917881	Class A Trunk Motor	Major
79	1	TM	8	Yorkshire and the Humber	335.806	209.32	3904082279	Class A Trunk Motor	Major
80	2	PM	8	Yorkshire and the Humber	5.012	3.49	88178179.43	Class A Principal Motorway	Major
124	1	TM	7	East of England	248.014	154.11	4320804105	Class A Trunk Motor	Major
150	3	PR	1	South West	3879.488	2410.58	8414392699	Class A Principal road	Major
75	1	TM	7	East of England	245.9	152.8	4120104199	Class A Trunk Motor	Major
44	5	M	10	West Midlands	2798.1	17371.12	9937607665	Class B road	Minor
132	5	M	8	Yorkshire and the Humber	2581.1	17166.74	7956272669	Class B road	Minor
198	1	TM	10	West Midlands	358.449	222.73	8721315484	Class A Trunk Motor	Major
28	3	PR	7	East of England	2719.344	1689.72	7493677513	Class A Principal road	Major
33	4	PR	8	Yorkshire and the Humber	2539.934	1589.15	6365559035	Class A Principal road	Major
40	1	TM	10	West Midlands	351.863	218.7	5200343918	Class A Trunk Motor	Major
133	1	TM	9	South East	405.147	251.82	10402000348	Class A Trunk Motor	Major
39	5	M	9	South East	40051.1	250383.7	15087821429	Class B road	Minor
55	2	TA	2	East Midlands	1215.388	754.46	5056721438	Class A Trunk Road	Major

Delta Live Table

Delta Live Tables (DLT) is used in this project to orchestrate, automate, and manage the data pipeline from raw data ingestion to the creation of gold-level analytics tables. It automates the creation of bronze, silver, and gold tables with built-in data quality checks and lineage tracking. It reduces the complexity of manual job scheduling and notebook chaining.

- Created DLT in default schema.
- Used autoloader to handle incremental loading.
- Mode of pipeline: Triggered

Workflows > Jobs & pipelines >

Create pipeline [Send feedback](#)

General

* Pipeline name
DLT execution

☐ Serverless

Product edition
Advanced

[Help me choose](#)

Pipeline mode
☒ Triggered
 ☐ Continuous

Source code Paths

Destination

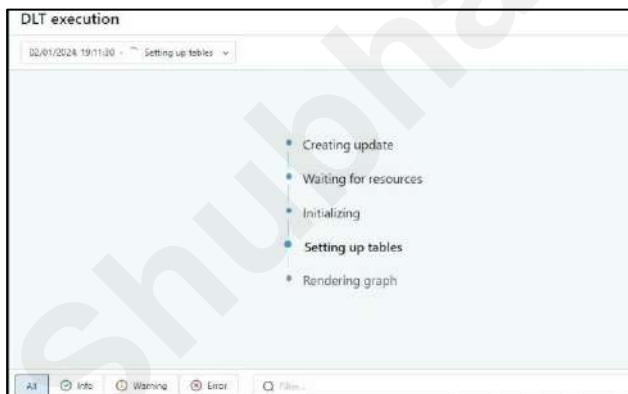
Storage options
☐ Hive Metastore ☒ Unity Catalog [Preview](#)
Default catalog
* Default schema

Compute

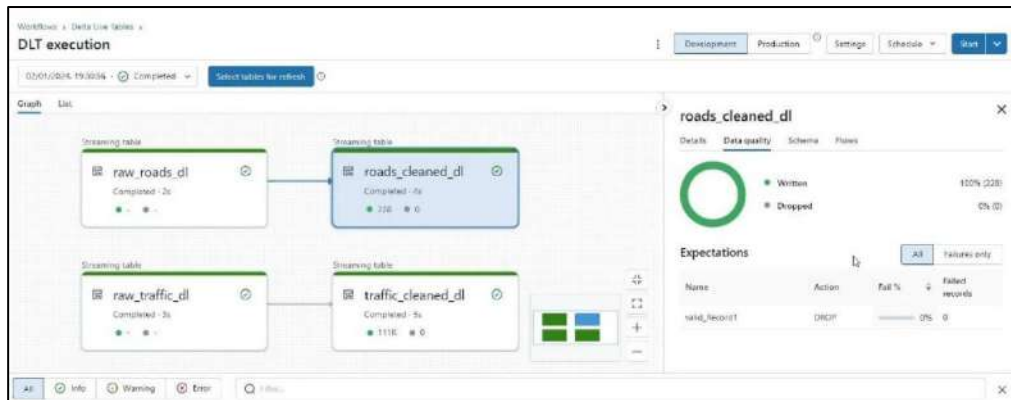
Cluster policy
Cluster mode
[Learn about Enhanced Autoscaling](#)
* Workers

We can see 2 modes in DLT pipeline:

- **Development:** It retains the cluster for 2 hours. It does not make any retry.
- **Production:** It invokes the cluster, and it stops the cluster once the execution is completed. It does retry as well.
- Through this DLT pipeline we need only 1 table and we will join both the tables after passing data quality checks.
- These are certain steps taken by Delta Live Table.



- Data quality metric is also visible in DLT. We can see the name of the constraint as valid, and it shows the percentage of failure.



- In the final gold table, we need only a few selected columns, so we select them and join both the tables in DLT.



The screenshot displays the 'Delta live' configuration page. Under the 'Tasks' tab, a 'Delta' task is shown with the pipeline 'DLT execution'. Below the task card, there is a form to configure the task. The 'Task name' is 'Delta', the 'Type' is 'Delta Live Tables pipeline', and the 'Pipeline' is 'DLT execution'. There are checkboxes for 'Trigger a full refresh on the Delta Live Tables pipeline', 'Notifications', 'Retries', and 'Duration threshold', each with an '+ Add' button.

- We can create a job for DLT as well as shown below.

Conclusion

This project successfully demonstrates the design and implementation of a modern data lakehouse architecture using Azure Databricks, Delta Lake, and Azure Data Lake Storage Gen2 for managing and analyzing traffic and road datasets.

By following the medallion architecture (Bronze → Silver → Gold), we ensured a structured and scalable data pipeline that supports both batch and real-time data ingestion using Spark Structured Streaming with Auto Loader.

Key achievements include:

- Secure and governed access using Access Connector and Unity Catalog.
- Cleaned and transformed data in the Silver Layer with derived metrics like `Electric_Vehicles_Count` and `Vehicle_Intensity`.
- Creation of business-ready Gold Layer tables optimized for analytics and reporting.
- Integration with Power BI to generate interactive dashboards that provide insights into traffic patterns, road usage, and electric vehicle trends.
- Implementation of CI/CD pipelines using Azure DevOps, enabling version control and automated deployments across environments.

This end-to-end pipeline not only enables real-time analytics but also serves as a reusable framework for building similar data-driven solutions in the transportation or smart city domain.