

End-to-End Unified Data Analytics Pipeline



Know Your Data



User Guide For Data Engineers

Version: V-1.0

Prepared By: Sharawan Kumar Thapa

Table of Contents

Introduction

Architecture Overview

Step-by-Step Process

- Step 1: Data Input Form
- Step 2: AWS RDS (MySQL)
- Step 3: MySQL Workbench Usage
- Step 4: AWS Glue (ETL)
- Step 5: Amazon S3 – Data Lake
- Step 6: Amazon Athena
- Step 7: ODBC Connection to Athena (Power BI Desktop)
- Step 8: Power BI Service – Publishing & Scheduled Refresh
- Step 9: Power Automate – Notifications

In-Progress Tasks and Upcoming version updates [V 2.0]

Introduction

Purpose: This document guides Data Engineers in managing and scaling an end-to-end cloud-based data pipeline, ensuring smooth data flow from data input to business insights through Power BI dashboards.

Scope: The pipeline covers the entire process from data collection to reporting, including data ingestion, storage, transformation, querying, and automated dashboard refresh.

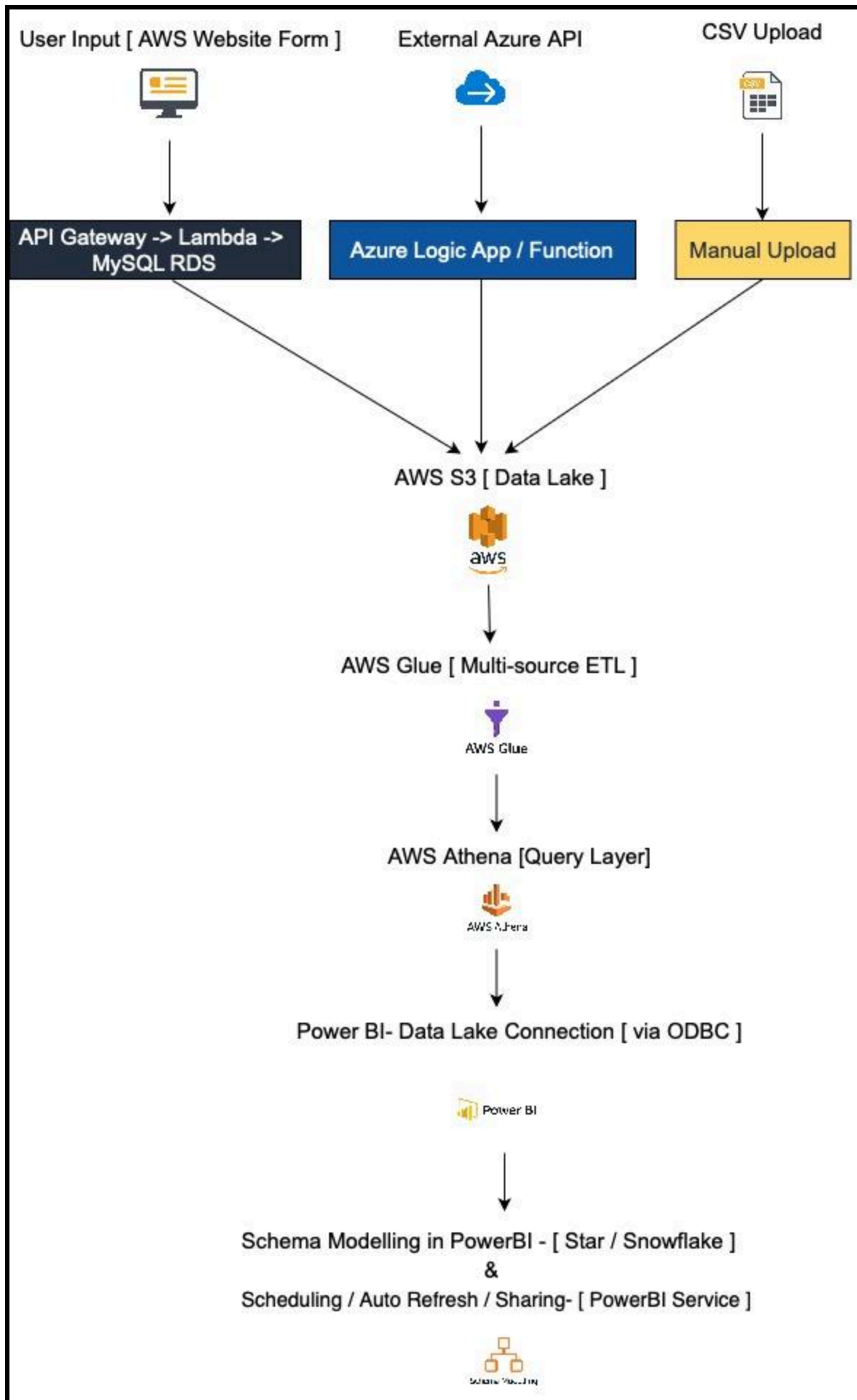
Data Sources: Data are collected from three sources i.e. AWS RDS, Azure services and Manual data ingestion through CSV files.

[Data Pipeline from Azure to AWS Data Lake is in Progress]

Technologies involved:

AWS hosted website Input Form → RDS → MySQL Workbench → Glue → S3 Data Lake (AWS RDS + CSV + Azure Service) → Athena → ODBC → Power BI → Power BI Service Schedule → Power Automate

Architecture Overview



[Figure: Data pipeline Architecture Overview](#)

Step-by-Step Process

Step 1: Data Input Form

Purpose: Collects user-entered data through a web form hosted on AWS [https://mydata.knowyour*****.com/] using HTML Form / Web App / API Gateway / Lambda / Route 53

Process:

Created a web form with fields like **Invoice ID**, **Branch**, **Date**, etc.

01. Connect the form backend to RDS using APIs (Lambda + API Gateway).
02. Validate data before writing to RDS.

← → ↻ https://mydata.knowyour.com

Invoice Form

Invoice ID:
INV78140

Branch:
A

City:
Tasmania

Customer type:
Member

Gender:
Female

Product line:
Food and beverages

Unit price:
14.06

Quantity:
4

Tax 5%:
2.81

Total:
59.05

Date:
09/02/2024

Time:
11:52 pm

Payment:
Cash

COGS:
56.24

Gross margin %:
4.76

Gross income:
2.81

Rating:
6.6

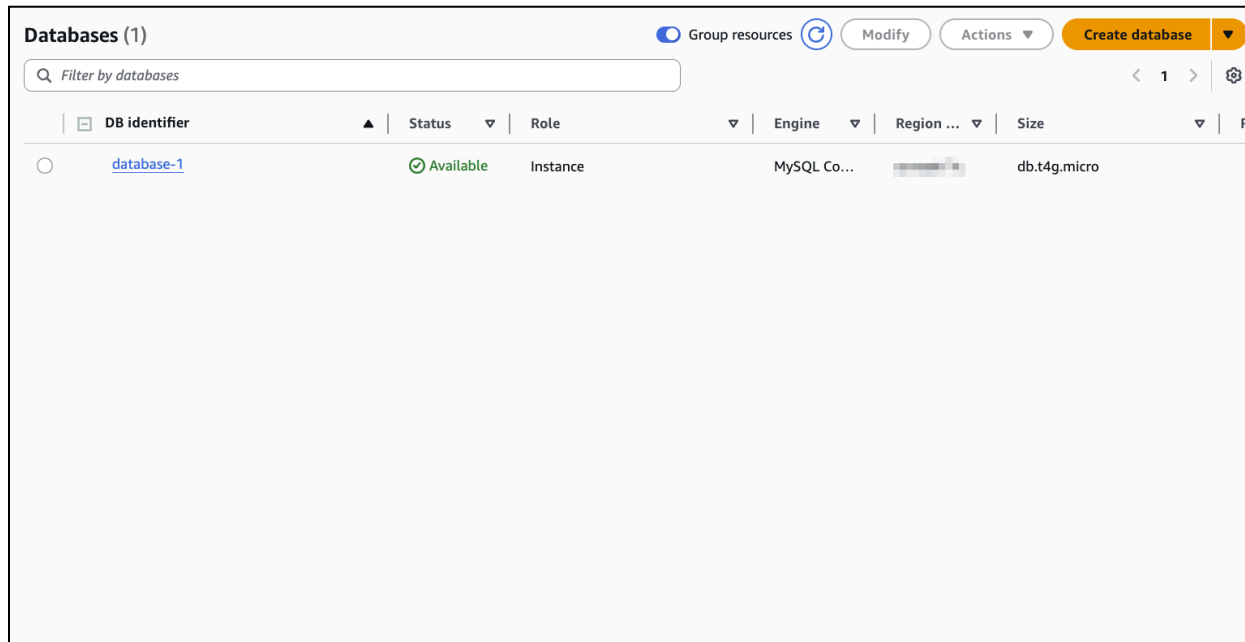
Submit

Step 2: AWS RDS (MySQL)

Purpose: Relational database to store raw transactional data.

Process:

01. Set up an RDS instance with MySQL.
02. Configure security groups (allow form/web app to write, Glue to read).
03. Use MySQL Workbench to Query data, Validate entries, Perform backups or schema changes if needed.

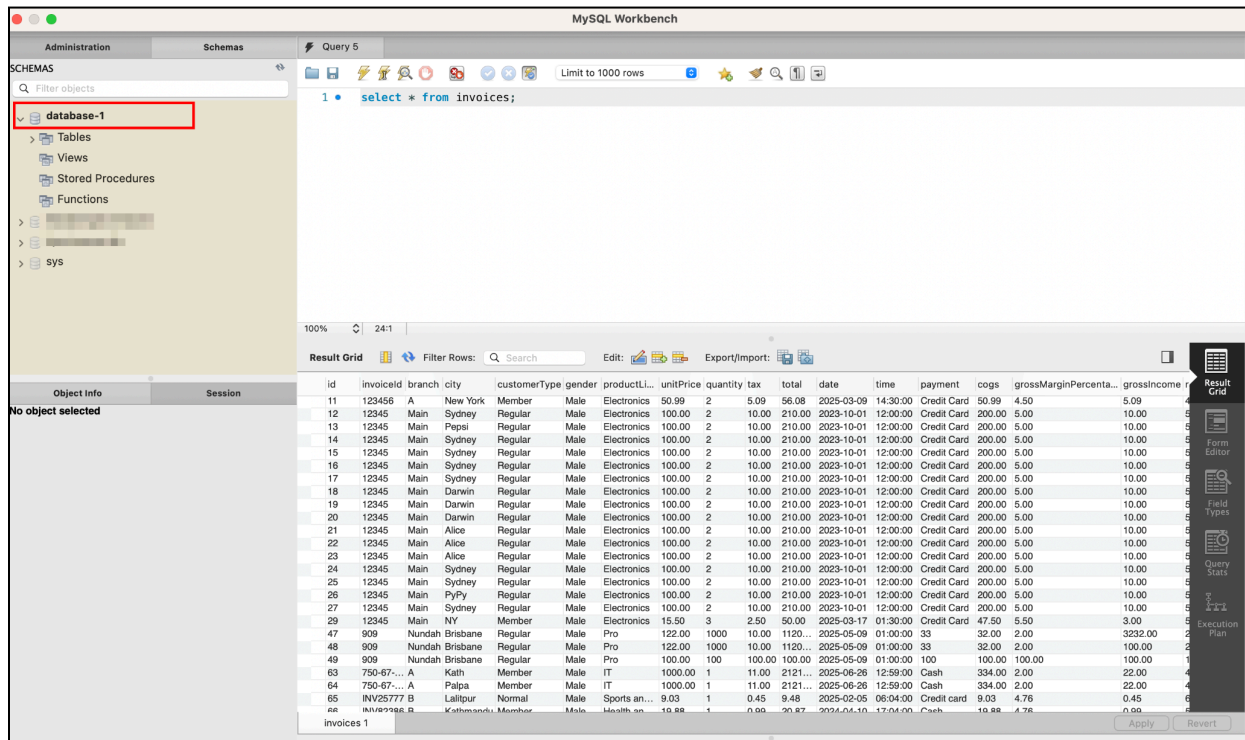


Step 3: MySQL Workbench Usage

Purpose: Connect securely to AWS RDS

Process:

01. Manual queries
02. Troubleshooting
03. Schema updates

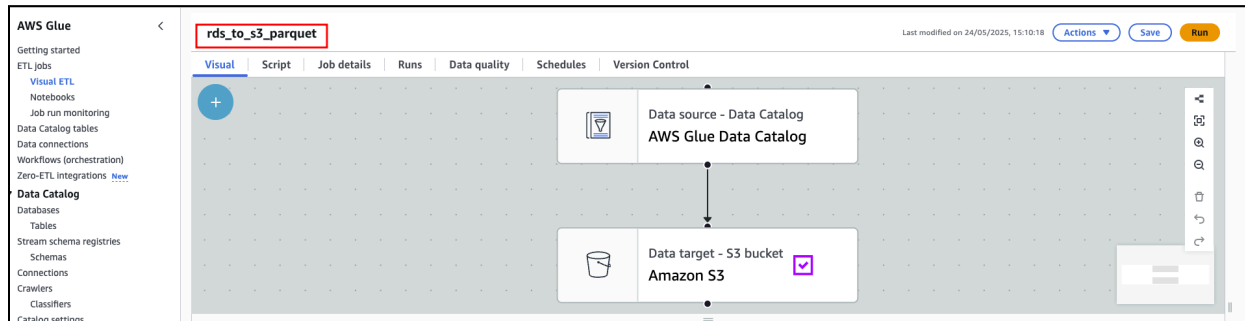


Step 4: AWS Glue (ETL)

Purpose: Extract from MySQL → Transform → Load to S3 in Parquet format. Performs ETL (Extract, Transform, Load) processes — moves, cleans, and transforms data into a query-ready format in S3.

Process:

01. Set up a Glue Crawler to connect to RDS (JDBC Connection).
02. Create a Glue Job to export data from RDS to S3:
 - a. Output format: Parquet
 - b. Folder structure: e.g.,
`s3://data-lake/sales_data/year=2025/month=05/`



03. Schedule Glue jobs (hourly, daily, etc.)

The screenshot shows the 'Schedules' tab for the 'rds_to_s3_parquet' job. The 'mysql to athena update trigger' is listed with a schedule of 'At 55 minutes past the hour' and a cron expression of '55 0/1 * * * *'. The status is 'Activated'.

Description	Schedule	Cron Expression	Status
-	At 55 minutes past the hour	55 0/1 * * * *	Activated

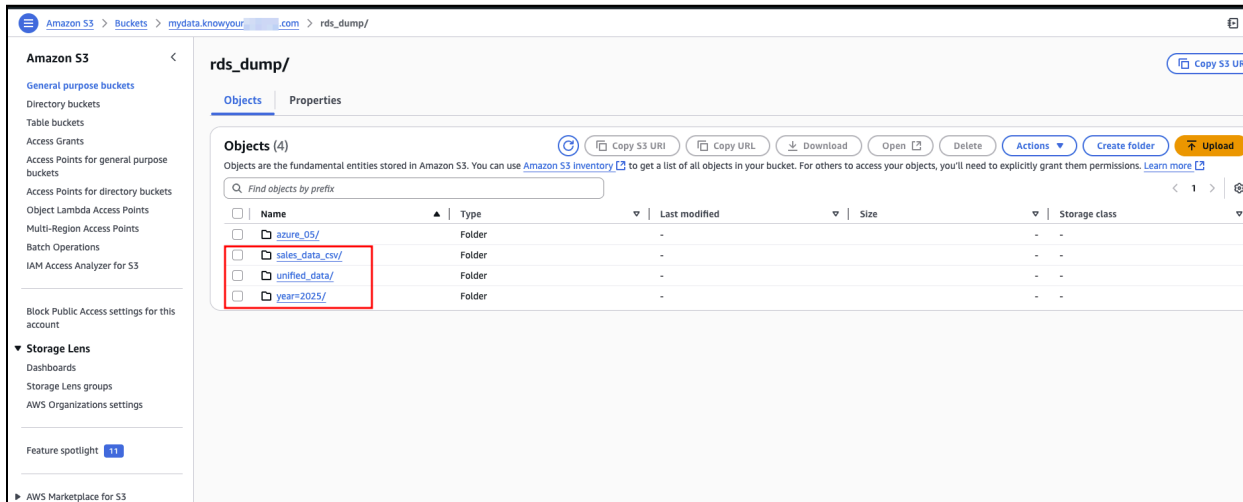
Step 5: Amazon S3 – Data Lake

Purpose: Central repository for transformed data in Parquet. Central storage for raw, structured, and semi-structured data from multiple sources [AWS RDS, Static CSV upload, Azure]. Create Glue crawlers that scan S3 folders to create Athena tables.

Process:

Create Folder Structure in data lake:

```
s3://data-lake/  
├── sales_data_csv/  
│   ├── unified_data/  
│   │   └── year=2025/
```

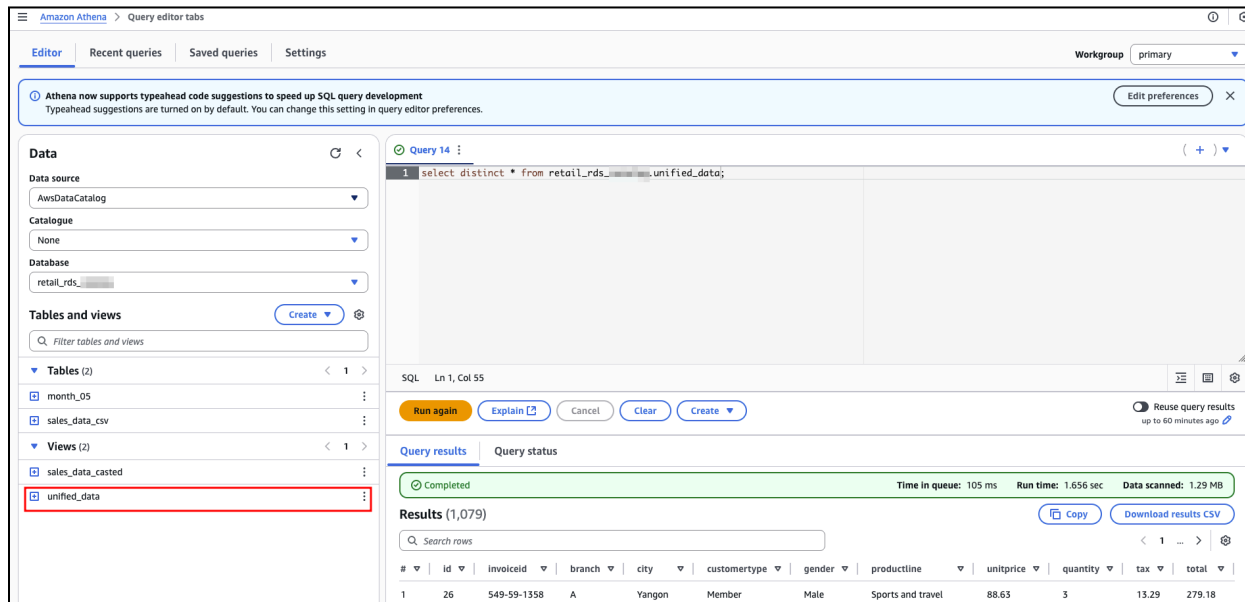


Step 6: Amazon Athena

Purpose: SQL engine to query data directly from S3; acts as the data access layer for reporting tools.

Process:

01. Write SQL queries that read from the S3 Data Lake, and can include filters, exclusions, or even create new “cleaned” tables or views for downstream use.
02. Create views like `unified_data` that merge multiple sources (`month_05` + `sales_data_csv`).

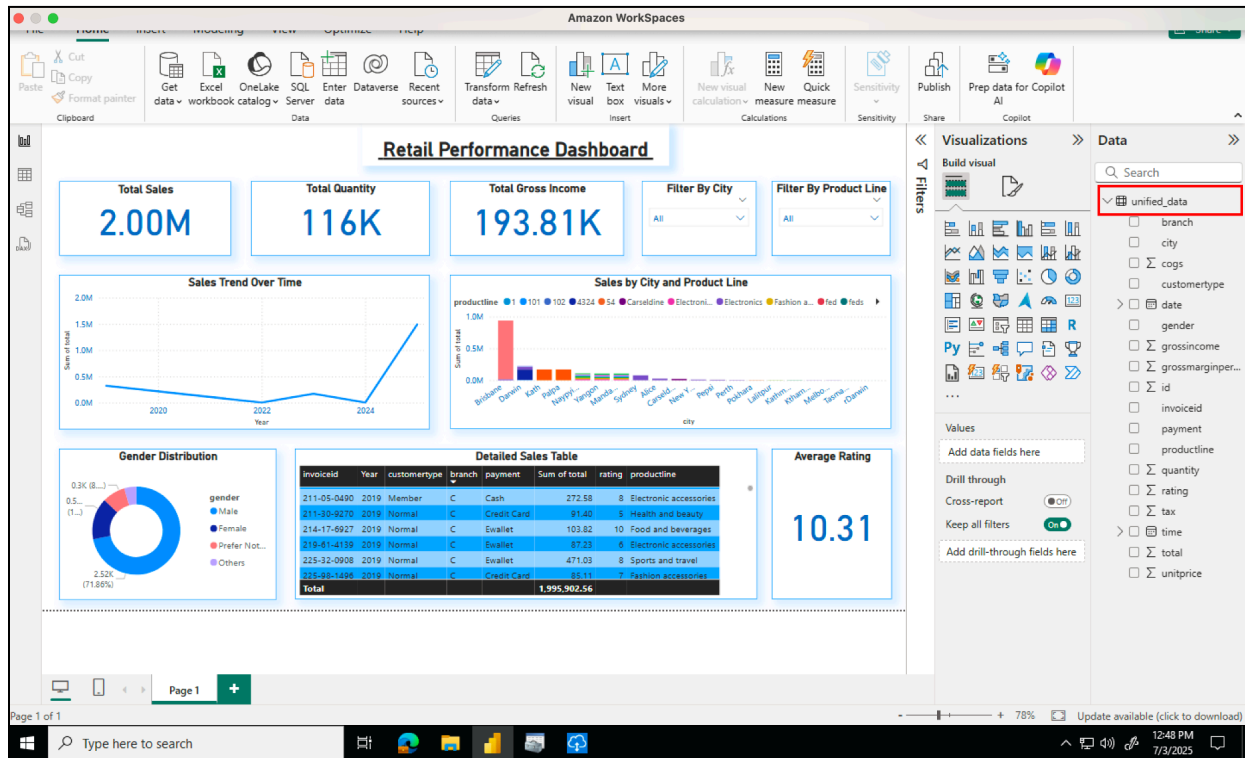


Step 7: ODBC Connection to Athena (Power BI Desktop)

Purpose: Connects to Athena to fetch query results in PowerBI Desktop, enables report building and data visualization.

Process:

01. In Power BI: **Get Data** → **ODBC** → **Athena DSN** → **Select tables/views like unified_data**.
02. Build reports, charts, KPIs.

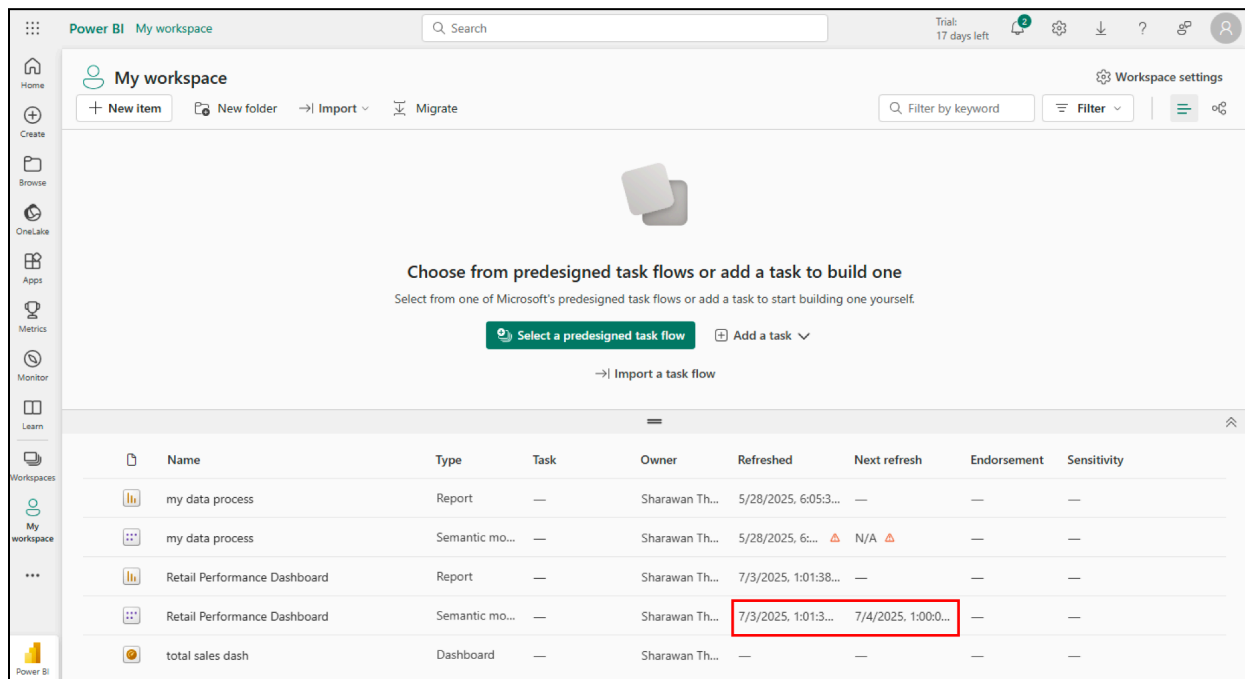


Step 8: Power BI Service – Publishing & Scheduled Refresh

Purpose: To publish reports to Power BI Service for collaboration, and automate data updates through scheduled refresh

Process:

01. Publish .pbix report to Power BI Service workspace
02. Configure Athena ODBC credentials in Power BI Service (via Gateway)
03. Set up **Scheduled Refresh**: Frequency (e.g., hourly, daily), Alerts for refresh failures

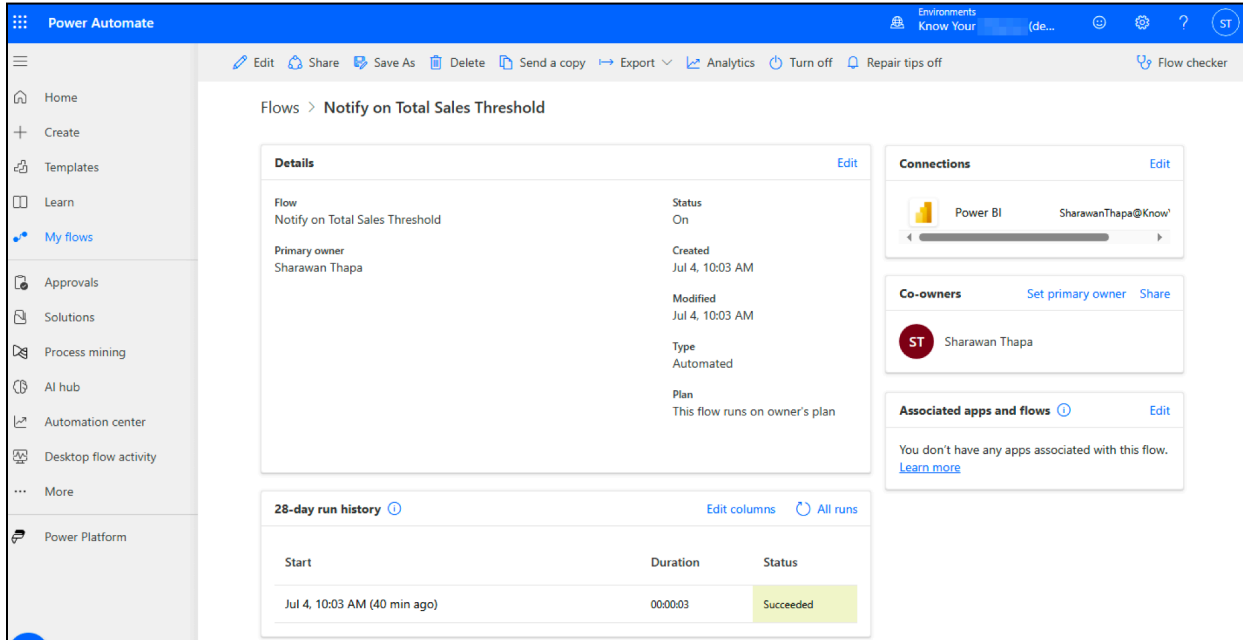


Step 9: Power Automate – Notifications

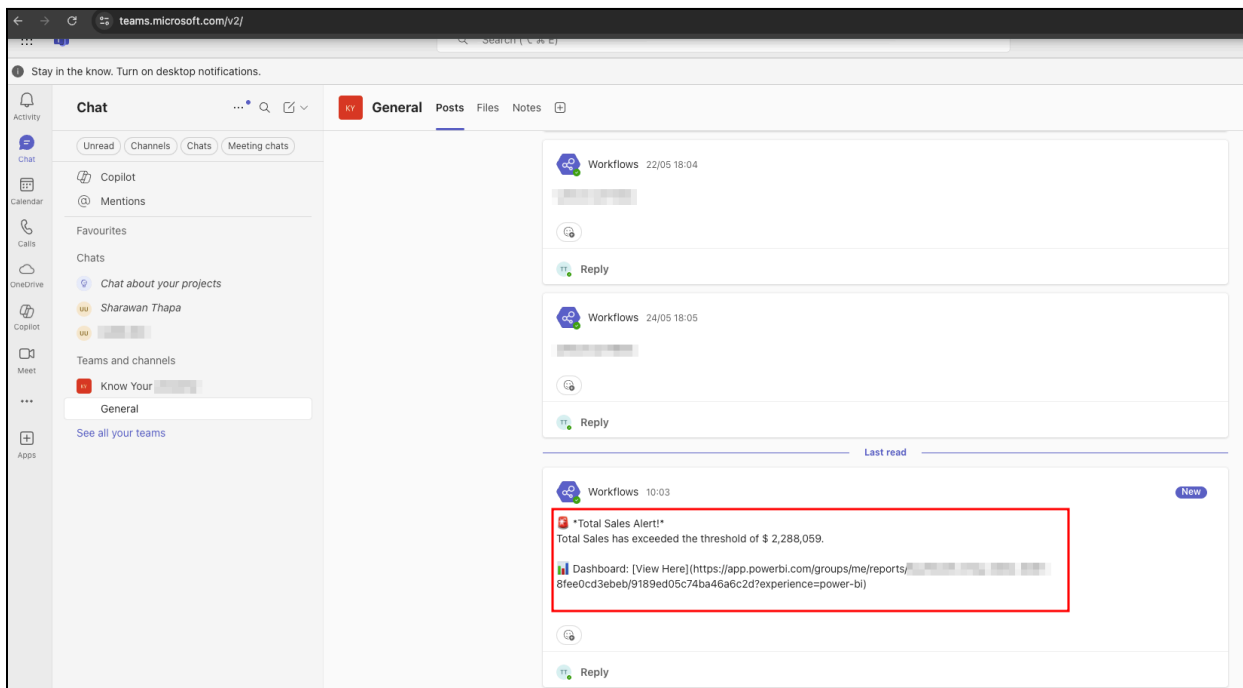
Purpose: Trigger email or Teams notification when Power BI refresh fails when the defined Power BI KPI threshold is exceeded, ensuring timely awareness and action.

Process:

01. Build a flow in Power Automate using the When a Power BI data-driven alert is triggered connector.
02. Configure the flow to post a message to a specified Teams channel or send an email.



03. Include relevant details in the notification, such as current metric value and a link to the Power BI dashboard.



In-Progress Tasks and Upcoming version updates [V 2.0]

01. Integrate external Azure Data Pipeline with Data Lake and ML on S3-based data
02. Implement Predictive Modelling
03. Star/ Snowflake Schema Modelling
04. Real-time streaming pipelines (e.g., Kinesis/Event Hubs → Data Lake)
05. CI/CD pipelines for data workflows (version control & automated tests)