# End-to-End ETL Pipeline Using Microsoft Fabric Lakehouse, Dataflows, and PowerBI

## Introduction

In this project, I designed and implemented a complete, automated ETL (Extract, Transform, Load) pipeline.

The goal was to create a robust system for handling monthly sales data, ensuring data quality, eliminating duplicates, and refreshing Power BI reports without manual intervention.

I used Microsoft Fabric Lakehouse, Dataflows, Power BI Desktop and Service, and Pipelines to achieve this.

This document outlines the challenges faced, solutions implemented, workflow architecture, and detailed steps of the ETL process.
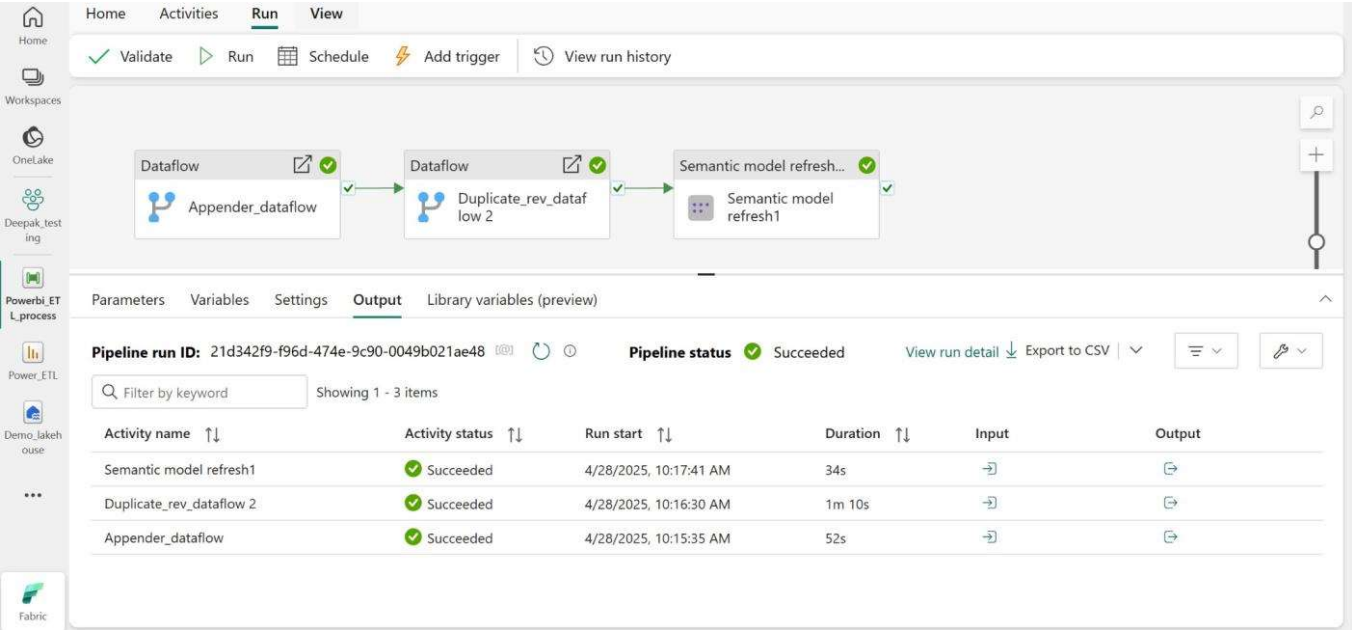
## Challenges Faced

- Data Silos: Monthly sales files were stored separately, making aggregation difficult.

- Manual Data Refresh: Each month, new data needed manual update and validation in Power BI.

- Duplicate Records: Repeated entries from different CSV files led to data inconsistency.

- Scalability Issues: The system had to handle growing sales data over time.

- Automation Need: Reducing human error and building an automated, reliable process was critical.
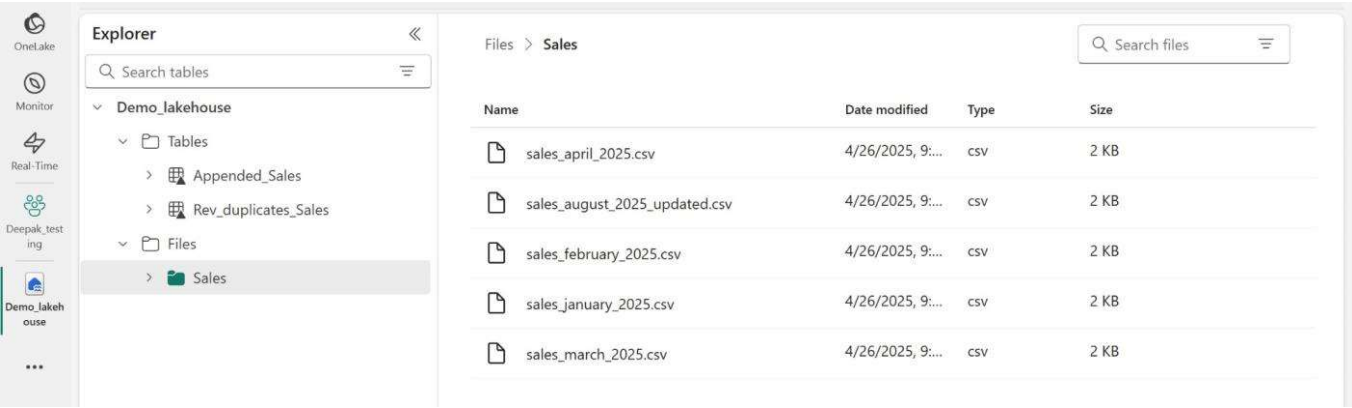
## Solution Overview

The solution was to build a modular, automated, and scalable ETL system:

- Store monthly sales data files in a Lakehouse under a structured folder.

- Create two Dataflows:

    - Dataflow 1 (Appender): Append all CSV files into a single dataset.

    - Dataflow 2 (Duplicate Removal): Remove duplicate records from the appended data.

- Load cleaned data into Power BI Desktop for report creation.

- Publish the report to Power BI Service.

- Create a Fabric Pipeline to automate the refresh of: Dataflow 1 -> Dataflow 2 -> Power BI Semantic Model.
- Thus, every new file added to Lakehouse is automatically picked up and reflected in the Power BI   dashboard.
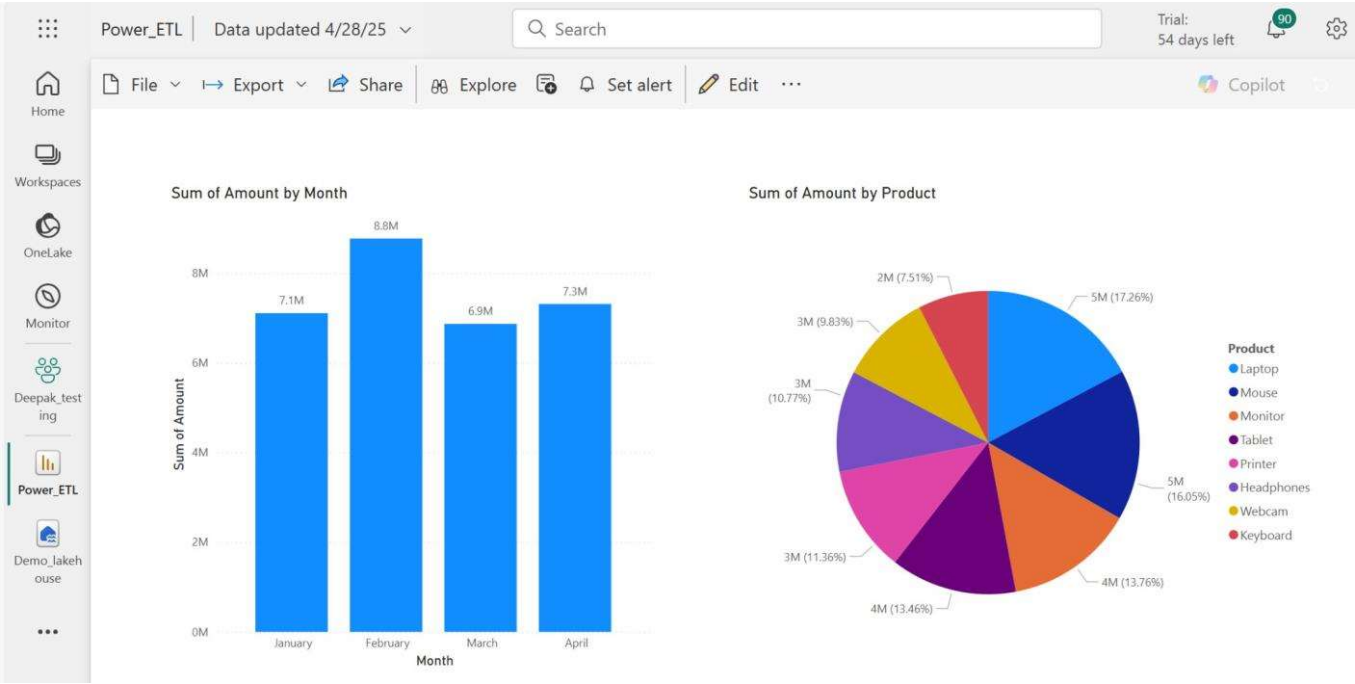
# Architecture Diagram



# Data Storage Setup (Lakehouse Structure)

# ETL Process Deep Dive



# Power BI Integration

- Connected Power BI Desktop directly to Dataflow 2 (cleaned data).

- Built interactive visualizations and KPIs on top of the clean dataset.

- Published the Power BI Report to Power BI Service workspace for online access and sharing.

# Pipeline Setup for Automation

# Conclusion

Through the implementation of this end-to-end ETL pipeline, I successfully:

- Automated the ingestion, transformation, and loading of monthly sales data.

- Eliminated duplicate records ensuring data reliability.

- Achieved seamless, automatic refresh of Power BI reports.

- Built a highly scalable and modular architecture.

- Improved efficiency, minimized human intervention, and enhanced overall data quality.

This project showcases strong skills in Microsoft Fabric, Data Engineering, ETL Design, and Power BI development