# BOOTCAMP PROJECT 4

Project Title: Incremental Data Loading and Automated Notifications using Microsoft Fabric

## Problem Statement:

In modern data ecosystems, organizations need to efficiently ingest, transform, and load data from various sources into centralized platforms for analytics, while also ensuring timely monitoring and notification upon successful data refreshes. This project addresses the challenge of incrementally loading data from on-premises sources to Microsoft Fabric Lakehouse, processing it through a structured transformation pipeline, and triggering automated notifications upon successful execution.
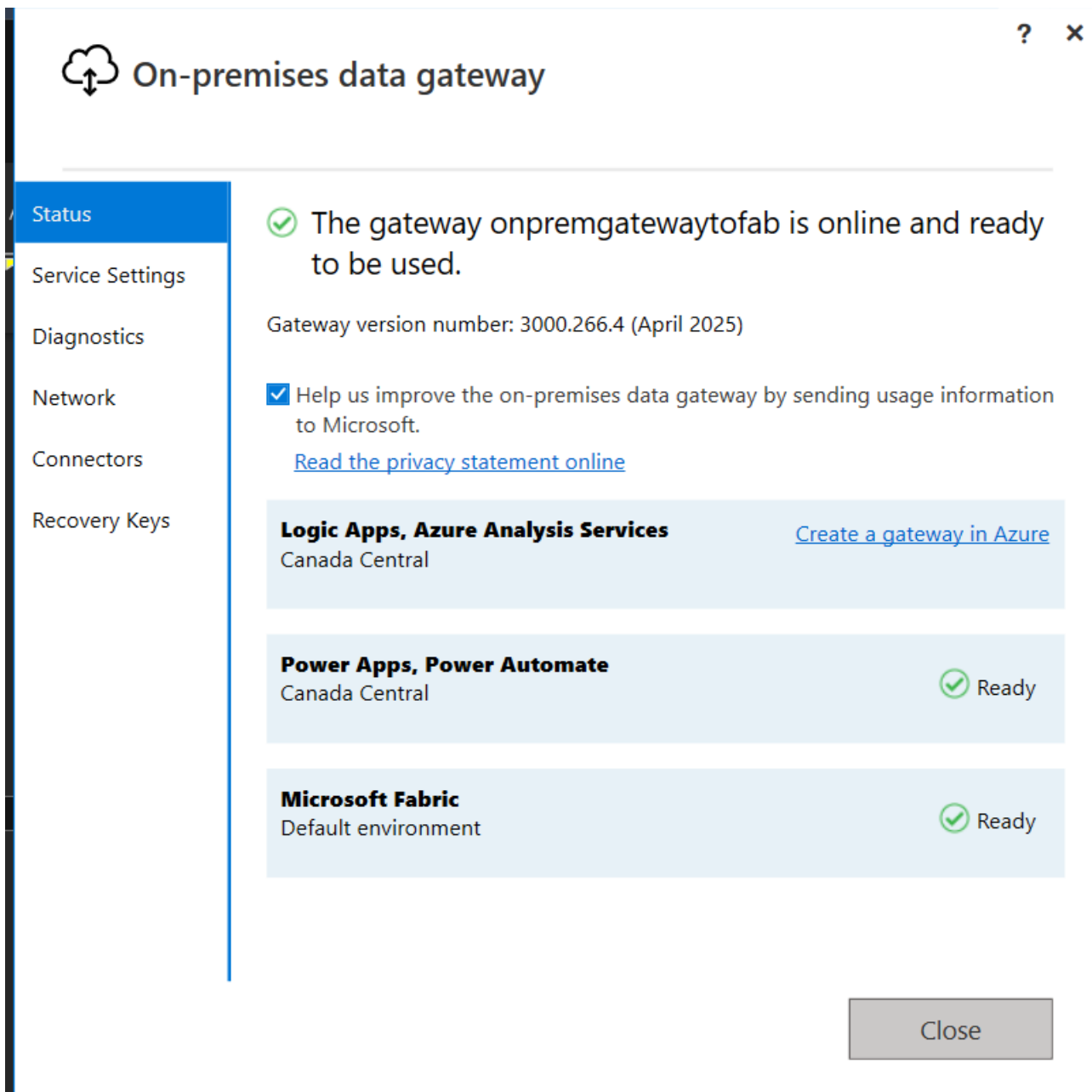
## Project Architecture Diagram:

## Tools and Technology used :

• Microsoft Fabric

• On-Premises Data Gateway

• Fabric Lakehouse and Warehouse

• Fabric Dataflow Gen 1

• Fabric Notebook

• Email Notification Task (in-built)

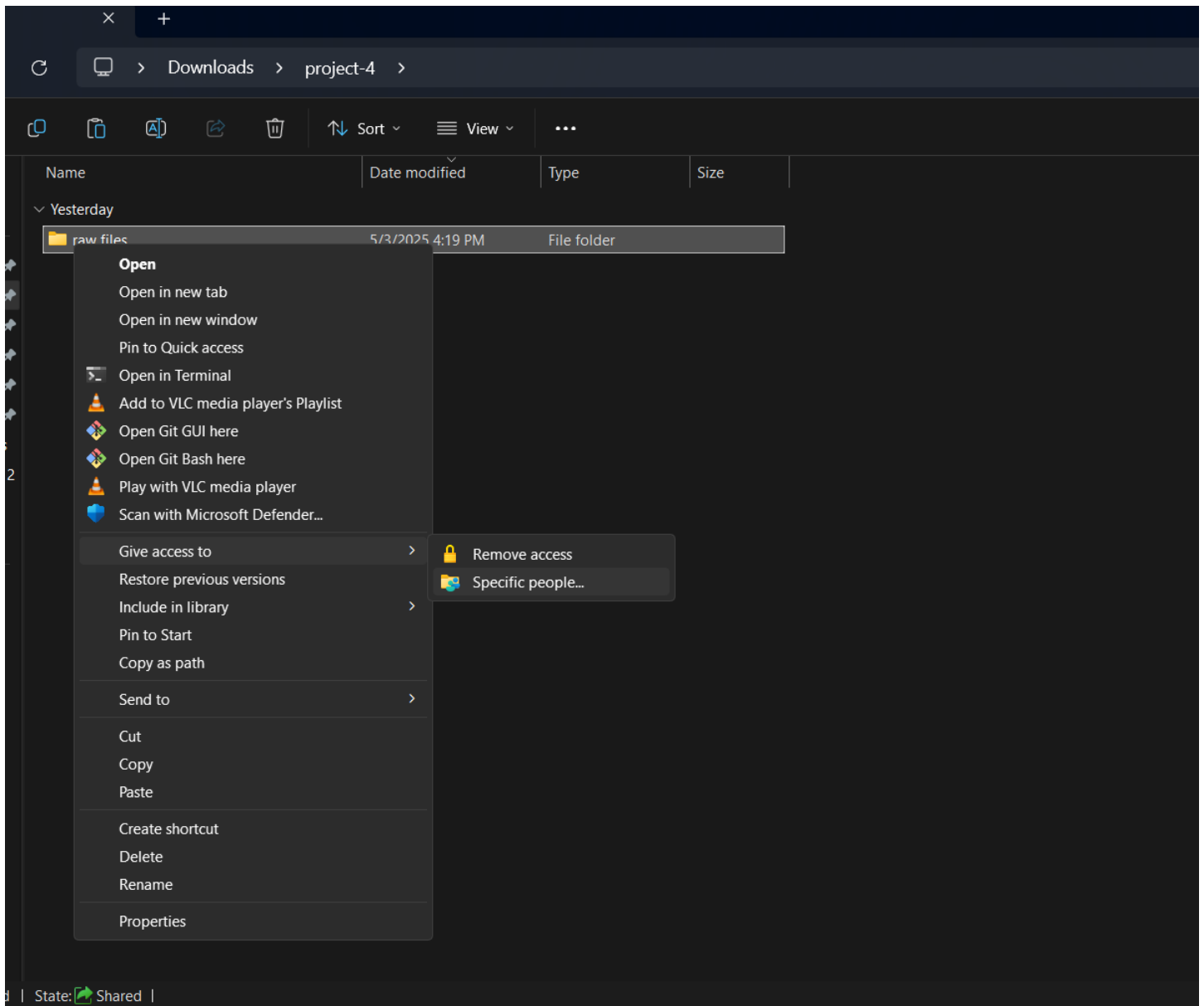• Draw.io / Visio for architecture diagram

# PROJECT FLOW

1. DATA INGESTION FROM LOCAL MACHINE TO FABRIC LAKEHOUSE USING ON-PREMISES GATEWAY

- In this phase, data gets transferred on local windows OS machine to fabric Lakehouse using On-premises gateway. Here are the steps I followed:
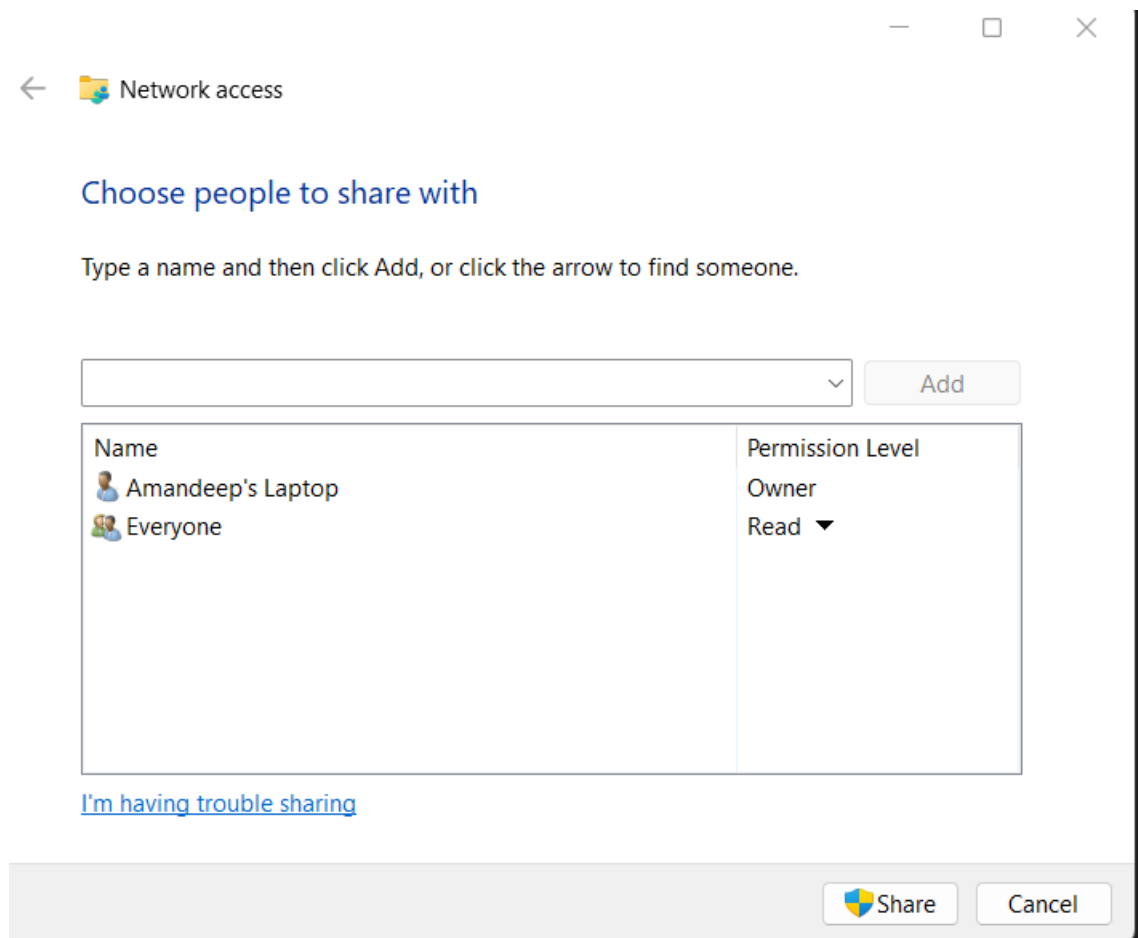


Install On-prem gateway and login with same credentials used in fabric. Setup up recouvery key in case of emergency.

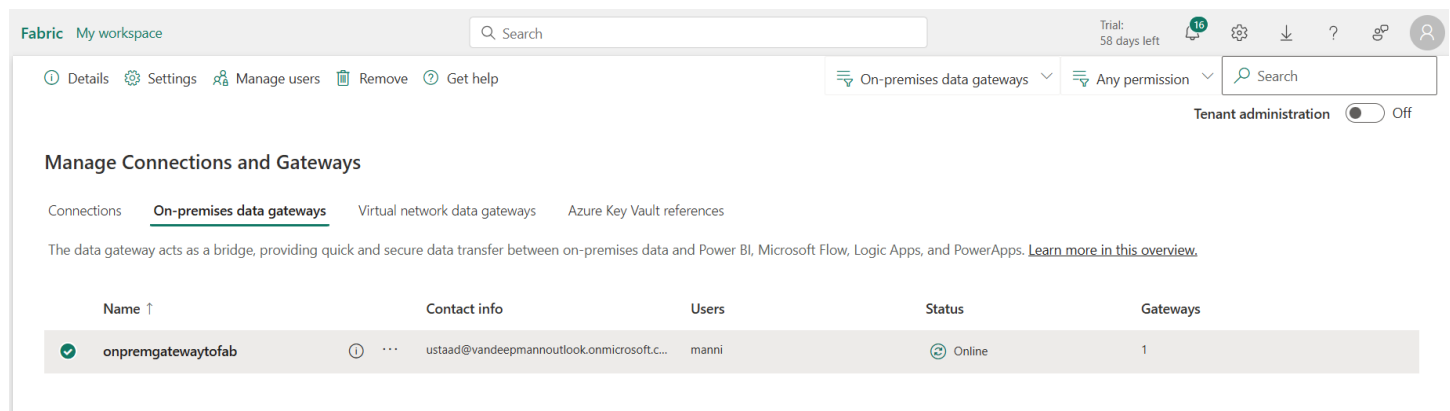Go to the dir, which you want to share and give read only access to everyone

Now check it gateway status is online in the fabric.

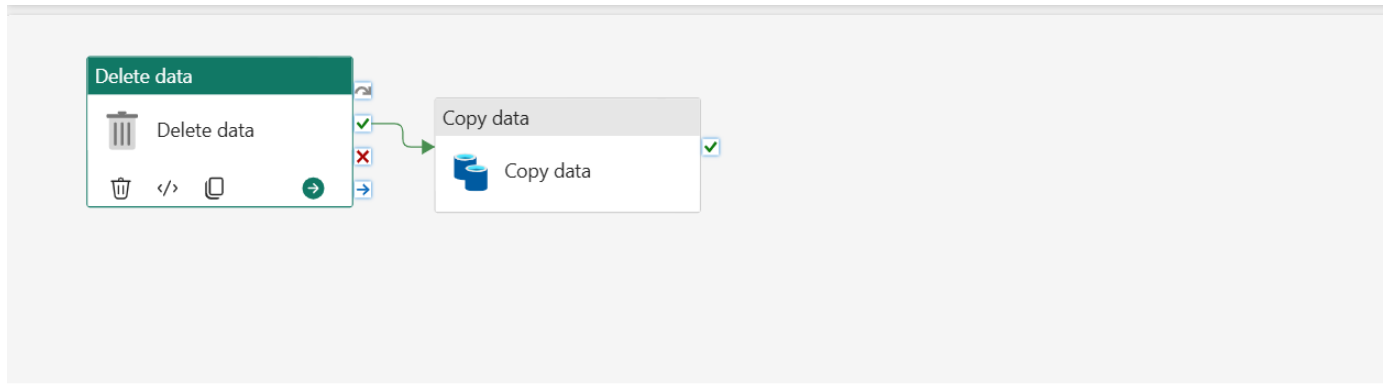Go to fabric>settings>manage connections and gateways>> gateways and check if it is online.



## Create new Data Pipeline for data ingestion

- Create delete activity and copy activity in the pipeline.
- Delete activity is to delete the previous csv files already present in lakehouse directory.

- In the copy activity, in source section, I have selected "folder" data source.
- In the destination section, I have used lakehouse as sink.

General    **Source**    Logging settings

Connection *                    📄 dev_lakehouse          ⌄        🔄 Refresh      ⬀ Open

File path type                  ⦿ File path    ○ Wildcard file path    ○ List of files ⓘ

File path          Files  /    project-4/raw-files          /    File name           📁 Browse  | ⌄

Recursively ⓘ                  ☑

> Advanced



General    **Source**    Destination    Mapping    Settings

Connection *                    📄 local_connection        ⌄      🔄 Refresh    ⚡ Test connection    ✎ Edit

File path type                  ⦿ File path    ○ File filter    ○ Wildcard file path    ○ List of files ⓘ
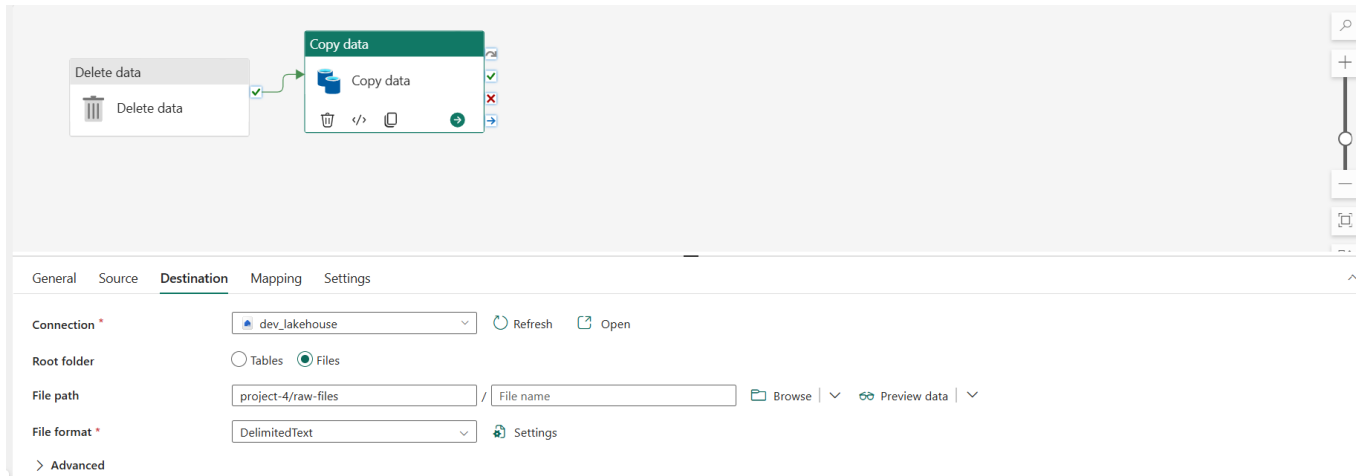
File path          raw files                          /    File name          📁 Browse  | ⌄    👓 Preview data
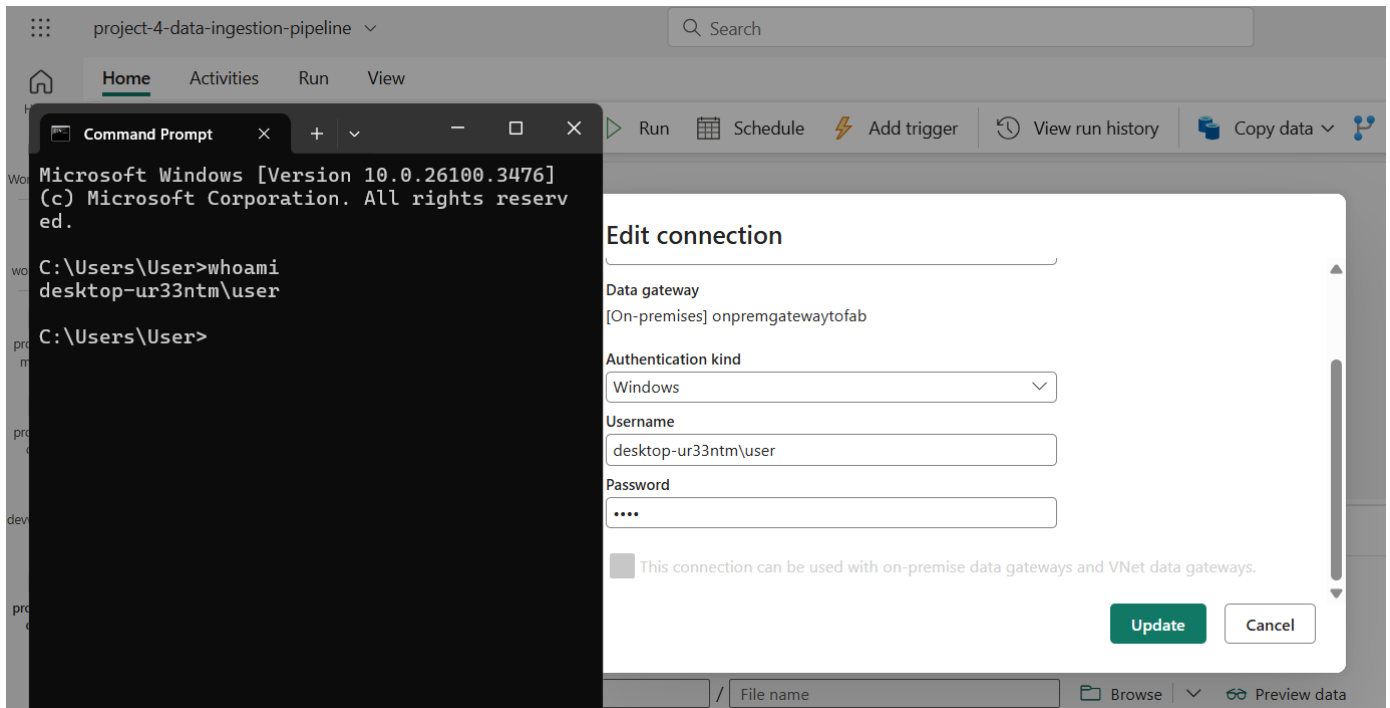
Recursively ⓘ                  ☑

File format *                   DelimitedText             ⌄      ⚙ Settings

> Advanced

- In source section, I've used same gateway to connect with local PC.



Over here, we have to use windows credentials and Administrator's account name.

- Select the list of files options to copy all csv files from  a dir

Here is the completion status of data ingestion pipeline:



The data get copied from LOCAL Machine to Fabric Lakehouse:

## 2. DATA CLEANING AND DATA TRANSFORMATION USING DATAFLOW GEN1 AND STORING IN FABRIC DATA WAREHOUSE

- Now, data got ingested into fabric lakehouse in raw files ( csv files).
- It's time to clean the data and transforming it according to requirements.
- DATAFLOW gen1 will be used to clean and transform the data and all csv files data will be saved into Tables in fabric- data warehouse.



Here the data pipleine is created where dataflow is invoked, which helps in data cleaning and data transformation.
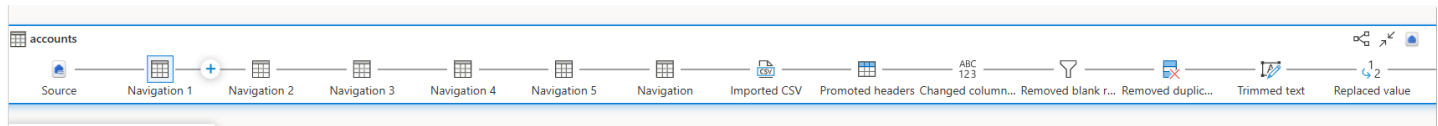
In the dataflow, I'm doing following tasks:

- Handle Missing Values
- Removing duplicate records
- Filtering the records
- Transforms date data type
- Trimming string column values

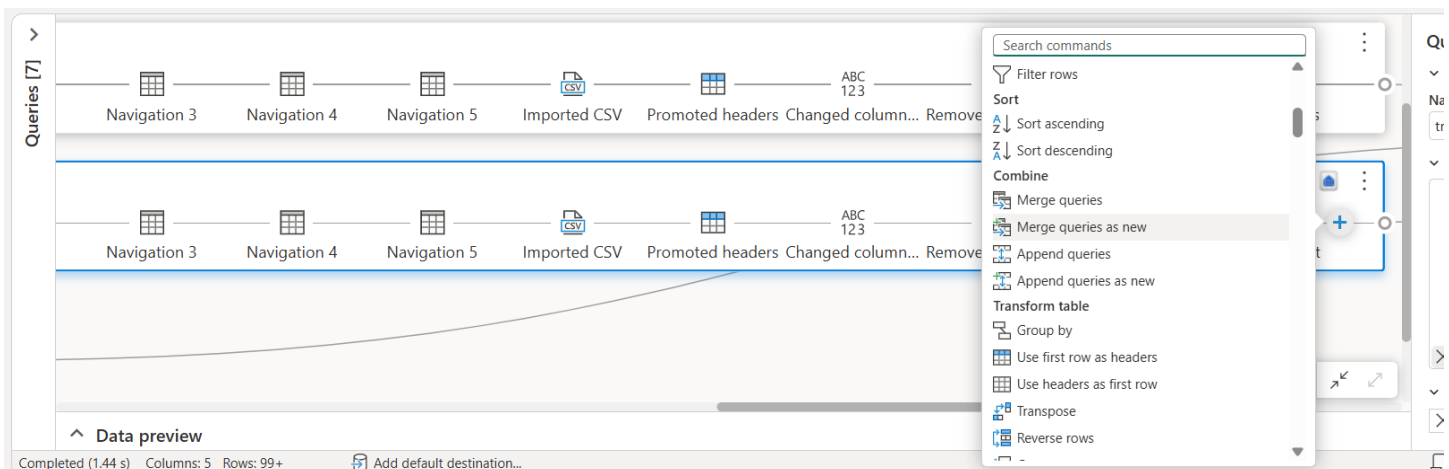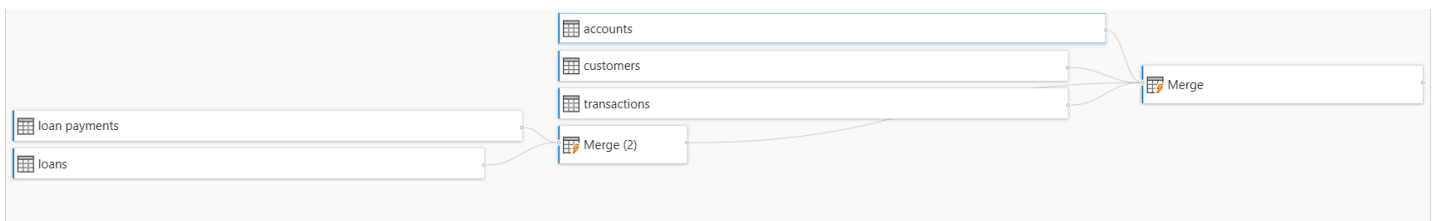For each table. So there will 5 data sources and 5 sinks.

Project-4  by Amandeep Singh

In addition to that, I'm creating new table, in which data from 5 tables get combined into single table named : "cleaned_common_table"

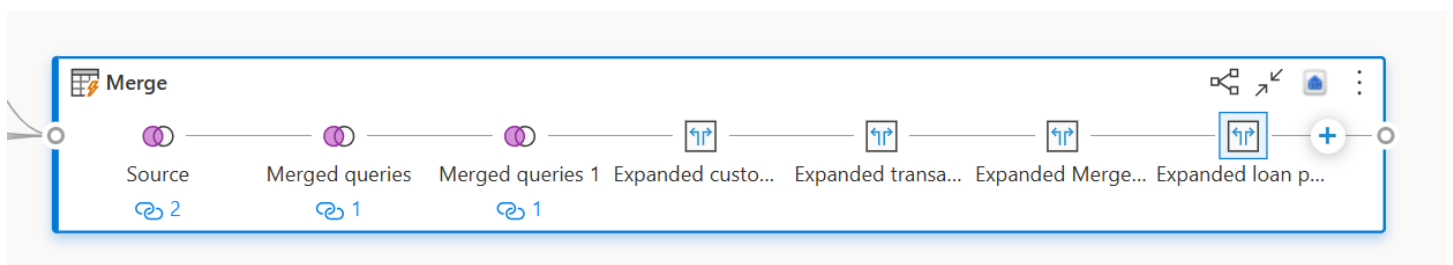Here is the query view of 1 csv file :



Here is the diagram view of those dataflow 1 :





I'm using merge queries as new and merge queries activity to join two data stream into new stream.
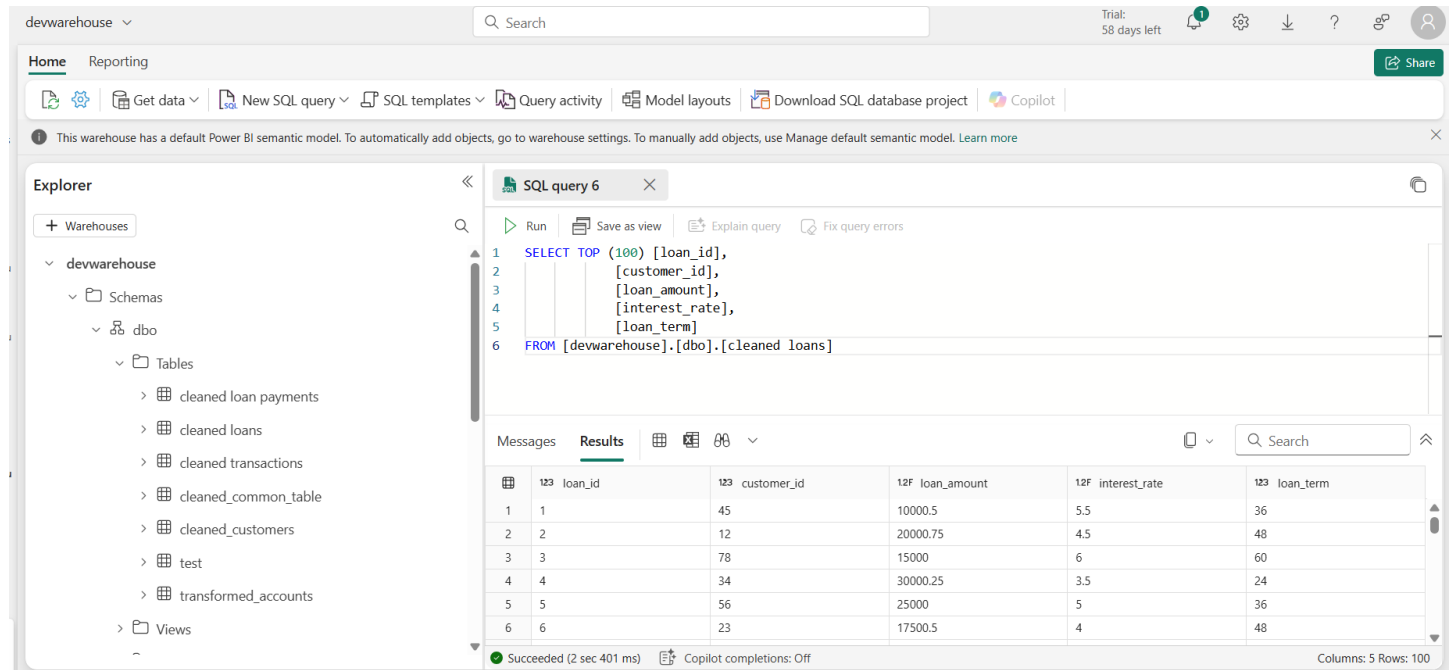
In last, I've created destination table for that common stream too.

Project-4  by Amandeep Singh

I'm using choose column activity and table expansion activity to explode the data into common stream.

Here the resulted data from dataflow 1 get stored as tables in devwarehouse:



## 3. TRANSFORMING THE DATA INTO SCD TYPE 1 TABLES:

- For this, we are using jupyter notebooks in fabric.
- To access the tables from fabric warehouse, I'm using shortcut of lakehouse, which I've created in the lakehouse.

- 
- This is the notebook, where I was accessing warehouse's tables using shortcut.



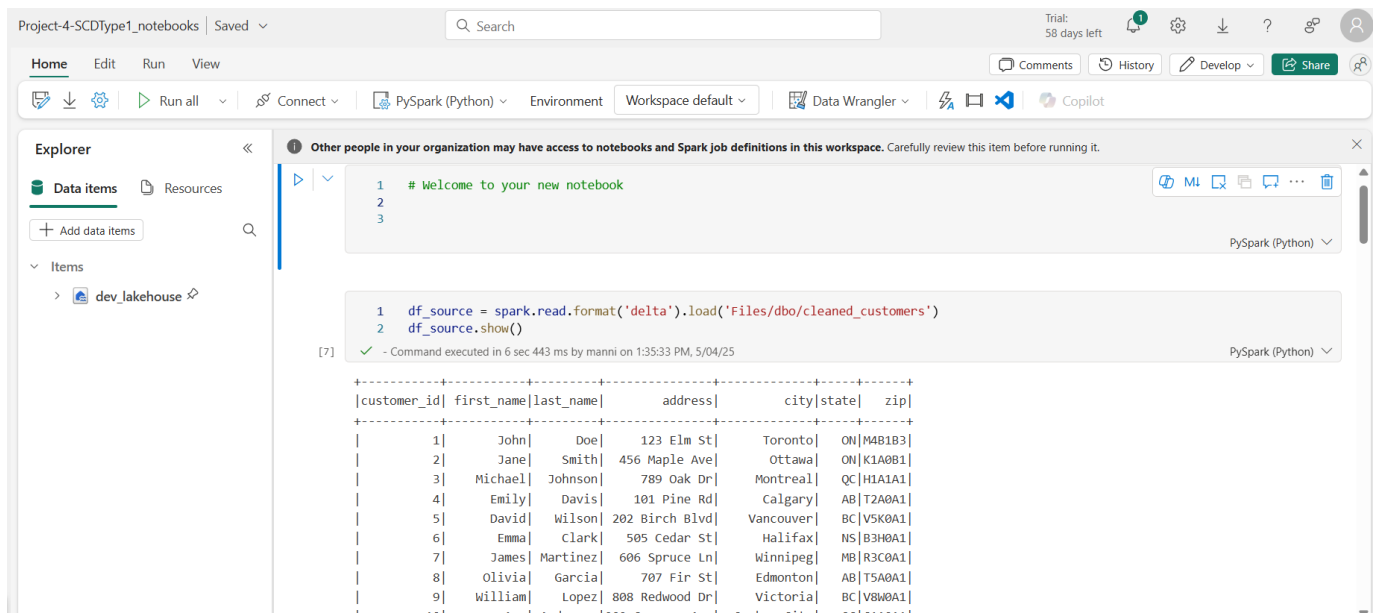Then, I've created SCD tables for each transformed table.

```
 1   create_query = """create table scd_customers
 2       (
 3           customer_id int,
 4           first_name varchar(100),
 5           last_name varchar(100),
 6           address varchar(100),
 7           city varchar(25),
 8           state varchar(50),
 9           zip varchar(50),
10           createdby varchar(50),
11           updatedby varchar(50),
12           created_date timestamp,
13           updated_date timestamp,
14           hash_key BIGINT
15       )
16       using delta
17       location 'Tables/dbo/scd/scd_customers'"""
18   spark.sql(create_query)
```

[12]  ✓  - Command executed in 6 sec 303 ms by manni on 1:44:51 PM, 5/04/25                    PySpark (Python) ∨

DataFrame[]

Following is the code to transform dataframe object into SCD Type 1 and storing into target table:

```
from delta.tables import DeltaTable

target_path = "Tables/dbo/scd/scd_customers"

delta_target = DeltaTable.forPath(spark, target_path)

 from pyspark.sql.functions import *

df_src1= df_source.withColumn("hash_key",crc32(concat(df_source.columns))) display(df_src1)
df_src1=df_src1.alias("src").join(delta_target.toDF().alias("tgt"),((col("src.customer_id")==col
("tgt.customer_id"))&(col("src.hash_key")==col("tgt.hash_key"))),"anti").select(col("src.")) df_src1.show()
from pyspark.sql.functions import col
delta_target.alias("tgt").merge(df_src1.alias("src"),"tgt.customer_id =
src.customer_id").whenMatchedUpdate(set={"tgt.customer_id":"src.customer_id","tgt.first_name":"src.
first_name","tgt.last_name":"src.last_name","tgt.address":"src.address","tgt.city":"src.city","tgt.state":"
src.state","tgt.zip":"src.zip","tgt.hash_key":"src.hash_key","tgt.updated_date":current_timestamp(),"tgt.
updatedby":lit("databricks_Updated") }).whenNotMatchedInsert(values={"tgt.customer_id":"src.custom
er_id","tgt.first_name":"src.first_name","tgt.last_name":"src.last_name","tgt.address":"src.address","tgt
.city":"src.city","tgt.state":"src.state","tgt.zip":"src.zip","tgt.hash_key":"src.hash_key","tgt.created_date"
:current_timestamp(),"tgt.createdby":lit("databricks"),"tgt.updated_date":current_timestamp(),"tgt.upda
tedby":lit("databricks")}).execute()
display(spark.read.format("delta").option("header","true").load(target_path))
```

## 4. SCHEDULING THE PROJECT BY MASTER PIPELINE AND ADDING ETL PROCESS COMPLETION NOTIFICATION:

- After completion of all phases, It's time to schedule the whole process by invoking pipelines and notebook using master pipeline.
- We are also adding e-mail notification that the pipeline has run successfully.