# **Data Engineering Project 2:**

# **Transactions and Loan Data for Customers**

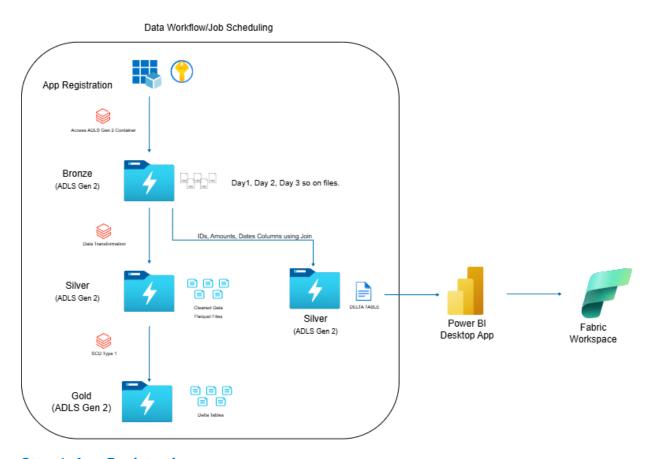
# **Table of Contents**

Objective:	2
Project Architecture Diagram:	2
Step 1: App Registration	2
Step 2: Mount Point and Functions Notebooks	14
Step 3: Upload files into the ADLS Gen 2 Folder	17
Step 4: Bronze to Silver Notebook	19
Step 5: DELTA Table for Data Visualization	20
Step 6: Create a scheduled job for daily run	29
Step 7: Power BI Visualization on DELTA File in Silver Layer	35
Step 8: Publish it to the Fabric Workspace	38
Conclusion:	39
Points to remember:	40

# **Objective:**

This project aims to design and implement a robust data pipeline for processing customer account data. This includes copying data from ADLS GEN2 (Bronze layer) and transforming the data in the Silver layer using <u>Data bricks Notebooks</u> and GOLD Storage into the SCDType 1 Delta Table in ADLS GEN2. The pipeline aims to ensure efficient, accurate, and scalable data processing to support downstream analytics and reporting needs.

# **Project Architecture Diagram:**



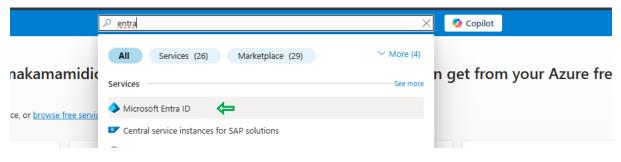
# **Step 1: App Registration**

## Go to Microsoft Entra ID

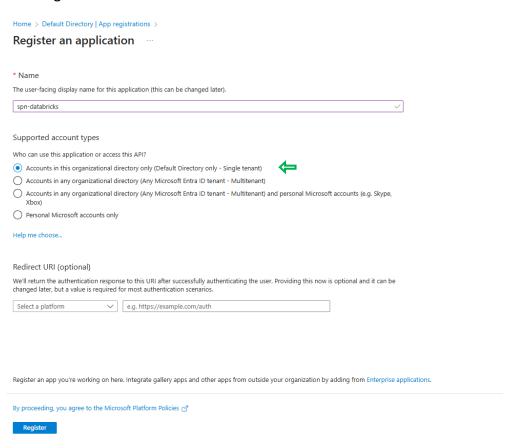
Microsoft Entra ID  $\rightarrow$  Manage  $\rightarrow$  App Registration  $\rightarrow$  New Registration

Name: spn-databricks

Select: Single tenant



# Click Register



Here we need the application(client) ID and Directory(tenant) ID:

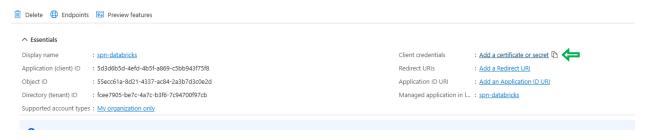
Copy these two IDS into Notepad, as we will need them later

#### Harpalsinh Vaghela

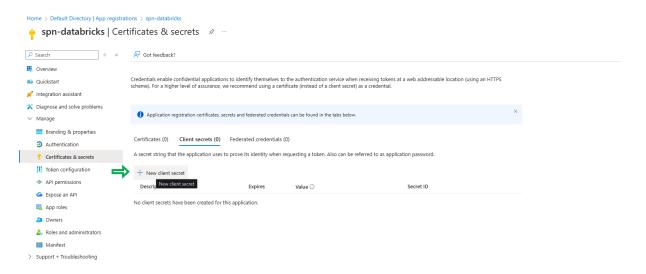
#### Project 2: Transactions and Loan Data for Customers



#### Click on Client Credentials as shown below

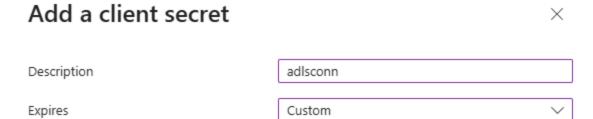


#### Click on New Client Secret



Write a description, Expires, Start and End as shown below

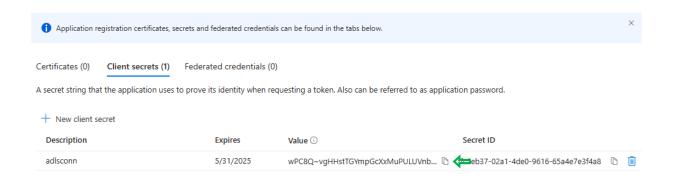
End



Start 04/23/2025

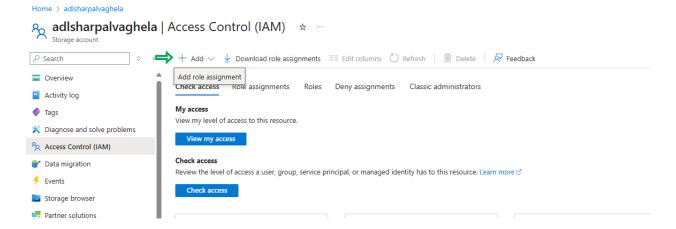
05/31/2025

Copy the **Value** from this newly created secret, paste it into a notepad to use later on.

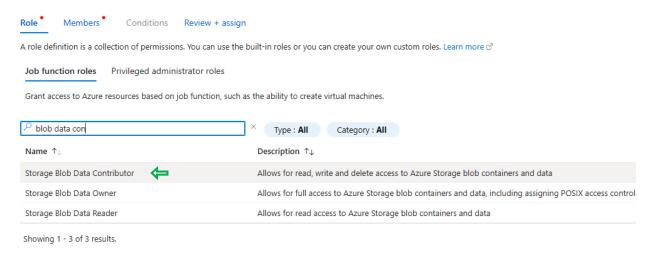


#### Go to the ADLS Gen 2 account:

Access Control (IAM) → Add role assignment



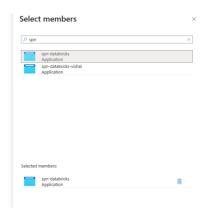
# Select Storage Blob Data Contributor role



# Members Tab: User, group or service principal

## Select member: search spn-databricks





#### Review and Assign

Home > adlsharpalvaghela | Access Control (IAM) >

# Add role assignment



### Azure Key Vault

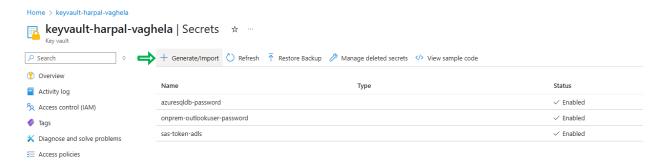
Azure Key Vault is a cloud service provided by Microsoft Azure that helps securely store and manage sensitive information such as keys, secrets, and certificates. It acts as a secure locker where you can store and control access to things like database connection strings, API keys, or passwords, which are crucial for your applications and services but should not be exposed or hard-coded in your code.

# Key Features of Azure Key Vault:

- Security: Azure Key Vault is designed to keep data secure using encryption, ensuring that sensitive information is never exposed.
- Access Control: It allows you to control who and what can access your keys and secrets. You can
  configure permissions specifically tailored to the needs of your organization.
- Audit Trails: It provides logging capabilities, which means you can monitor who accessed what information and when, adding an extra layer of security and compliance.

# Go to Key Vault in the Azure Portal

### Secrets → Generate/Import



# Name: appid

Select the appid from Notepad, and in the Secret value field, paste it there

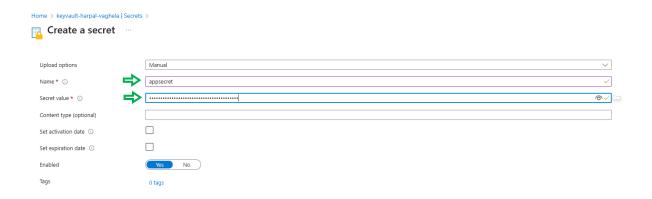


Again, click on Generate/Import → Name: appsecret

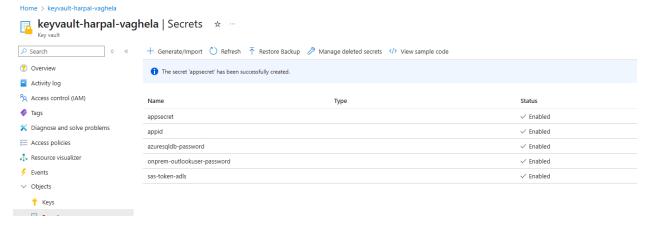
Secret value, paste it here (client secret which we copied earlier while creating spn-databricks)

#### Harpalsinh Vaghela

#### Project 2: Transactions and Loan Data for Customers

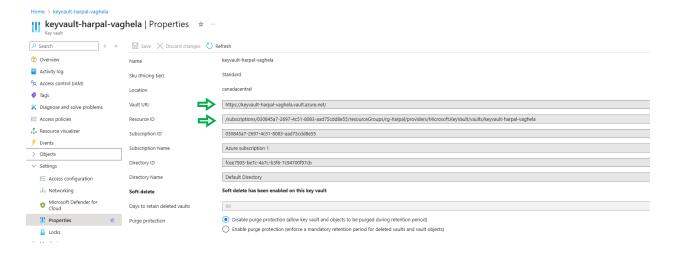


# Keyvault -> Objects -> Secrets page will look like this



# Go to Key Vault-> Properties

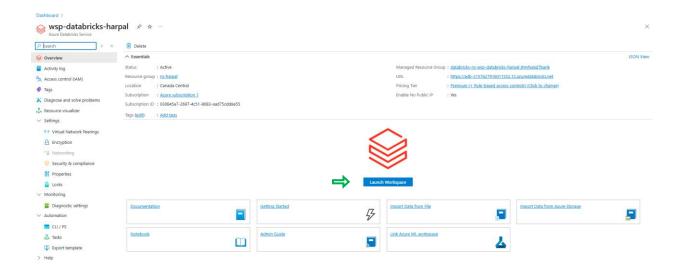
Here we need Vault URI and Resource ID, copy it to Notepad



Go to the Databricks Notebook tab in the browser

#### Harpalsinh Vaghela

#### Project 2: Transactions and Loan Data for Customers



In the browser URL, remove the URL after the .net word, like this

https://adb-2157627916011552.12.azuredatabricks.net/#secrets/createScope

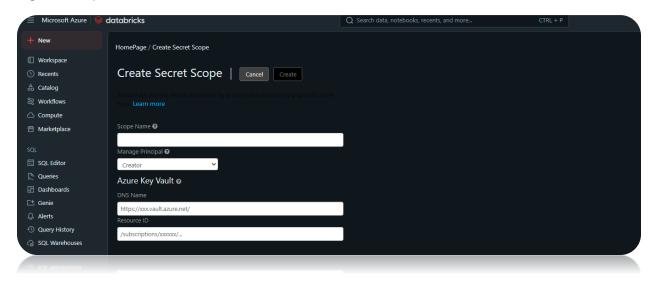
https://adb-2157627916011552.12.azuredatabricks.net/#secrets/createScope

https://adb-2157627916011552.12.azuredatabricks.net/#secrets/createScope - Bing Search

And write this: #secrets/createScope

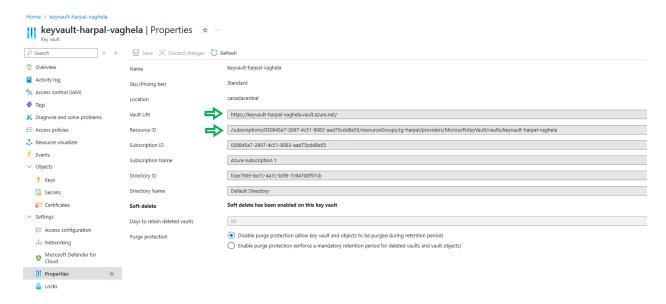
Press Enter on the keyboard

Page will be opened as shown below:

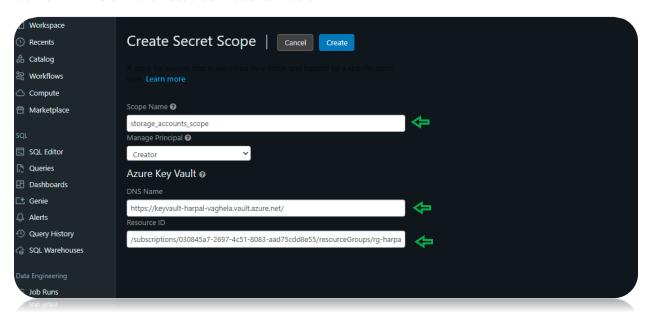


Write the Scope Name, also, we need the DNS Name and Resource ID here

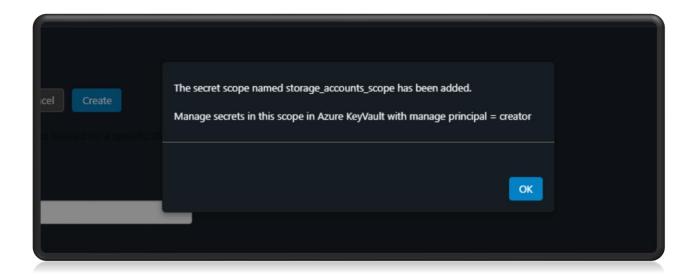
So we have to go to Key vault -> Properties tab and copy it to the scope window



Vault URI in DNS URL and Resource ID as shown below:



Click on Create



# Go to Databricks Workspace:

# Go to Compute -> Create Compute



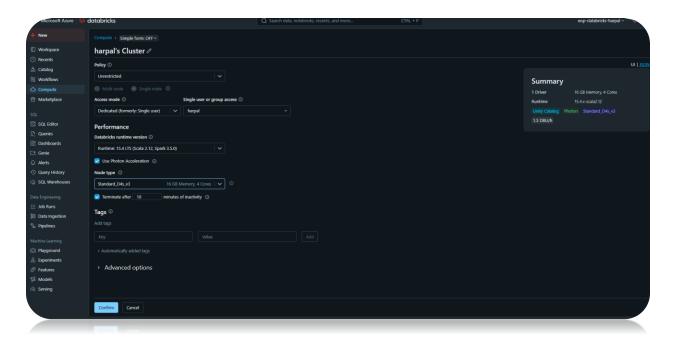
Use this configuration for creating a cluster

Policy: Unrestricted, Single Mode Access mode: Dedicated (Single User) Performance:

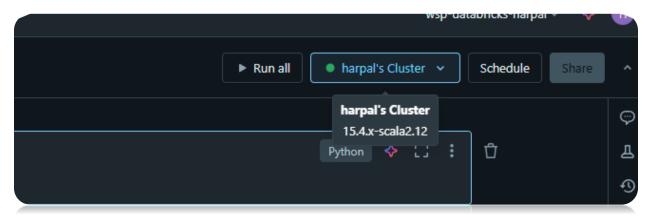
- Databricks runtime version:
   Runtime 15.4 LTS (Scala 2.12, Spark 3.5.0)
- Node Type: Standard\_D4s\_v3
- o Terminate after 10 minutes of inactivity



# Click on Create Compute



Create a new notebook and attach that cluster



Follow this URL to get the code for the mounting process

Code to mount ADLS Gen 2 in Databricks

https://learn.microsoft.com/en-us/azure/databricks/dbfs/mounts

Go to Databricks Notebook, attach and start the cluster

Check for all mounted locations using the command below

dbutils.fs.mounts()

Check all secret scopes

### dbutils.secrets.listScopes()

```
dbutils.secrets.listScopes()

[SecretScope(name='storage_accounts_scope')]
```

Execute this to check all secrets from the key vault

# dbutils.secrets.list("storage\_accounts\_scope")

```
dbutils.secrets.list("storage_accounts_scope")

[SecretMetadata(key='appid'),
    SecretMetadata(key='appsecret'),
    SecretMetadata(key='azuresqldb-password'),
    SecretMetadata(key='onprem-outlookuser-password'),
    SecretMetadata(key='sas-token-adls')]
```

Make a change in that mount code and modify as below

We will create a total of 3 notebook files:

Let's go with the first one, in which we already have the mount code.

# **Step 2: Mount Point and Functions Notebooks**

- Gets today's date in "YYYY-MM-DD" format using datetime.today().strftime.
- Splits the date into year, month, and day components.
- Builds a dynamic bronze layer path based on today's date.
- Sets static paths for silver and gold layers.

- Imports data types like StringType, DoubleType, DateType, and IntegerType to define DataFrame schemas.
- Imports functions like col, when, lit, concat\_ws, concat, current\_timestamp, and crc32 for data transformations.
- Imports DeltaTable to work with Delta Lasske tables for updates, merges, and upserts

```
from pyspark.sql.types import StructType, StructField, StringType, DoubleType, DateType, IntegerType
from pyspark.sql.functions import col, when, lit, concat_ws, concat, col, current_timestamp, crc32
from delta.tables import DeltaTable
```

- Defines a function read\_silver\_layer\_files that takes file\_name as input.
- Reads a Parquet file from the silver layer path using the given file name.
- Returns the loaded DataFrame for further processings

- Replaces .csv in the file name to create a base name for output.
- Writes the DataFrame as a single .parquet file by coalescing into one partition and saving to a temp path.
- Finds the actual .parquet part file generated in the temp folder.
- Moves and renames the part file to the final path in the silver layer.
- Deletes the temporary folder after moving the file to clean up.

- Defines a function get\_schema that returns a specific schema based on the file name.
- Each schema is created using StructType with a list of StructField for column names and data types.
- Handles five known files: accounts.csv, customers.csv, loans.csv, loan\_payments.csv, and transactions.csv.
- If an unknown file name is passed, the function raises an error to prevent incorrect schema assignment.

```
| Table | Description | Descri
```

# Backup File logic:

- path\_exists checks if a given path exists without raising an error.
- The main function loops through .csv files in the bronze path if it exists.
- Each file is read into a DataFrame and written to a temporary backup folder as a single .csv.
- Renames the part-xxxxx.csv file to the original file name.
- Moves the renamed file to the final backup folder and deletes the temp folder.

Prints the backup status for each file; skips if the bronze path doesn't exist.

```
elif file_name == 'loan_payments.csv':

return StructField("payment_id", IntegerType(), True),

StructField("payment_date", DateType(), True),

StructField("payment_date", DateType(), True),

StructField("payment_mount", DoubleType(), True)

structField("payment_mount", DoubleType(), True)

lif file_name == 'transactions.csv':

return StructType([

StructField("transaction_id", IntegerType(), True),

StructField("transaction_date", DateType(), True),

StructField("transaction_type", StringType(), True)

structField("transaction_type", StringType(), True)

raise ValueError(f"Unknown file_name: {file_name}")
```

- Defines a function list\_csv\_files\_in\_bronze to get all .csv files from tshe given bronze path.
- Uses dbutils.fs.ls to list all files in the directory.
- Filters only files that end with .csv and returns a list of matching file names.

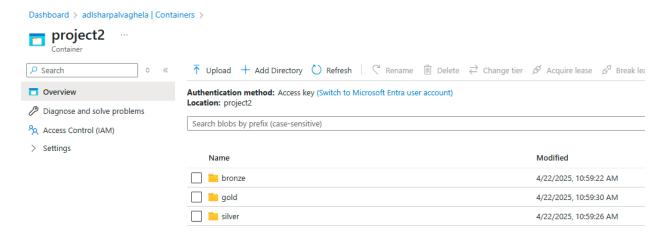
- filter\_nulls: Filters out rows where any of the specified ID columns are null.
- replace\_nulls\_with\_defaults: Replaces nulls in specified columns with default values using a dictionary.

- remove\_duplicates: Removes duplicate rows based on the given ID columns.s
- Each function returns a cleaned DataFrame for the next transformation step.

# Step 3: Upload files into the ADLS Gen 2 Folder

ADLS Gen 2 folders:

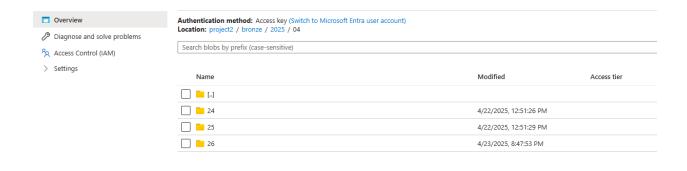
Create Bronze, Silver and Gold folders in this storage account



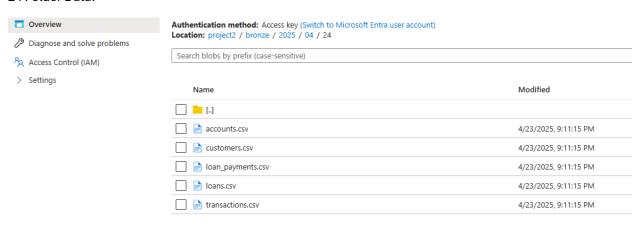
In the Bronze Folder in ADLS Gen 2, we have created a day-wise folder to upload the raw data files

## Harpalsinh Vaghela

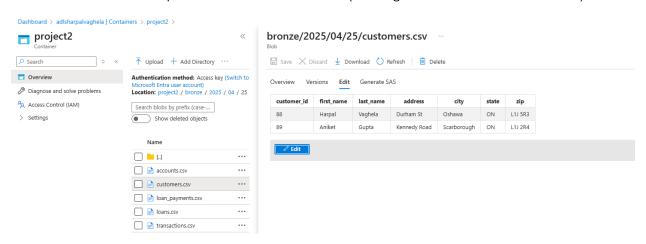
## Project 2: Transactions and Loan Data for Customers



#### 24 Folder Data:



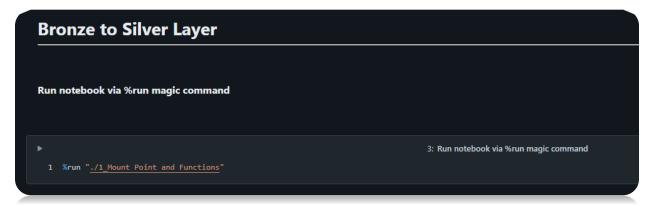
25 Folder Data with some updated values and new records (showing data for the customers CSV file):



# **Step 4: Bronze to Silver Notebook**

Let's write code into the bronze to silver notebook now (this is our second notebook)

- %run is a Databricks notebook command used to include and run another notebook's code.
- It imports all variables, functions, and logic defined in that notebook into the current one.



• Calling the Backup function to take all files backup to a different folder in csv files

- Retrieves a list of all .csv files from the bronze layer and prints them.
- For each file, reads it with an assigned schema and header using Spark.
- Based on the file name, defines primary ID columns and default values for null replacements.
- Applies three transformations: filters out rows with null IDs, replaces other nulls with defaults, and removes duplicates.
- Writes the cleaned data as a single .parquet file into the silver layer using the helper function.

```
If file_name_list = list_cve_file=_in_bronze(pronze_base_path)

print(frile inuse_list : (file_name_list))

# storp through such file

# clicent_dfs = []

# for soften in file_name_list:

# soft soften for the file

# soft soften for the soft for the file

# soft soften for the soft for the file

# soft soften for the soft for the file

# soft soften for the soft for the soft for the file

# soft soft for the soft for the soft for the file

# soft soft for the soft for the soft for the file

# soft soft for the soft for the soft for the file

# soft soft for the soft for the soft for the file

# soft soft for the soft for the soft for the file

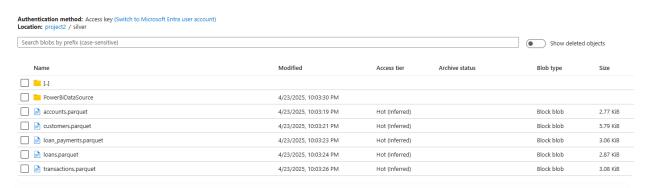
# soft soft for the soft for the soft for the file

# soft soft for the soft for the soft for the file

# soft soft for the soft for the soft for the file

# soft soft file_name = "soft soft for the so
```

### Silver Folder data in ADLS Gen 2



# **Step 5: DELTA Table for Data Visualization**

- Reads five cleaned .parquet files from the silver layer and loads them into individual DataFrames.
- Joins accounts with customers using customer\_id.
- Joins the result with transactions using account\_id.

- Further joins with loans on customer\_id, then with loan\_payments on loan\_id.
- Combines all related data into a single DataFrame df\_combined for visualization or reporting.

- Creates a new DataFrame df\_selected\_columns with only the required fields.
- Renames the selected columns using alias to more readable names.

- Writes df\_selected\_columns to the silver layer in Delta format for efficient querying and updates.
- Uses mode("overwrite") to replace any existing data at the path.
- Saves the output to the folder PowerBiDataSource inside the silver layer path.

```
12: Store it in DELTA format in Silver Layer

1 df_selected_columns.write.format("delta").mode("overwrite").save(f"{silver_base_path}PowerBiDataSource")
```

- Runs the notebook "./3\_Silver to Gold Layer SCD Type 1" inside the current notebook.
- Perform SCD Type 1 logic and write cleaned records to the gold layer.



Let's go with a 3rd notebook now

# Step 5: SCD Type 1 Logic

- Executes a SQL command using %sql magic in Databricks.
- Creates a new database named bankdb under the hive\_metastore if it doesn't already exist.



- Creates a Delta table named transactions in the bankdb database if it doesn't exist.
- Defines schema with fields like transaction\_id, account\_id, and audit columns such as CreatedDate and HashKey.
- Uses the DELTA format for efficient updates and versioning.
- Specifies the physical storage location in the gold layer: /mnt/project2/gold/transactions.

- Creates a Delta table named loan\_payments in the bankdb database if it doesn't already exist.
- Defines the table schema with columns like payment\_id, loan\_id, payment\_amount, and audit fields.
- Uses DELTA format to enable ACID transactions and efficient data updates.
- Stores the table data at the specified gold layer path: /mnt/project2/gold/loan\_payments.

- Creates a Delta table named loans in the bankdb database if it doesn't already exist.
- Defines schema with loan-related fields and audit columns like CreatedDate, UpdatedBy, and HashKey.
- Uses DELTA format for optimized performance and support for upserts.
- Stores the data in the gold layer path: /mnt/project2/gold/loans.

- Creates a Delta table named customers in the bankdb database if it doesn't already exist.
- Defines schema with customer details and audit fields for tracking changes.
- Uses DELTA format to enable efficient updates and data consistency.
- Saves the table data in the gold layer at /mnt/project2/gold/customers.

- Creates a Delta table named accounts in the bankdb database if it doesn't already exist.
- Defines columns for account info along with audit fields like CreatedBy and HashKey.
- Uses DELTA format to support versioning, ACID compliance, and upserts.
- Stores the table at the gold layer path: /mnt/project2/gold/accounts.

## Check for new Data and load it into the tables

- Reads the silver data and adds audit columns (CreatedDate, UpdatedDate, CreatedBy, UpdatedBy) and a HashKey for change detection using crc32.
- Performs an anti-join with the gold table to identify new or changed records (based on join\_key and HashKey).
- Prepares update\_set and insert\_values dictionaries for the merge operation, applying updates or inserts accordingly.
- Runs a Delta Lake merge using when Matched Update and when Not Matched Insert to apply SCD Type 1 logic (overwrite on change, insert if new).

```
| Jef scd_types_logic(| silver_path: str, | silver_path: | silver_path: str, | silver_
```

# Calling the function: scd\_type1\_logic

- Lists all .csv files from the bronze path and extracts their base names (e.g., accounts, customers).
- For each known file type, sets the join\_key and relevant columns for update and insert operations.
- Builds full paths for silver (input Parquet file) and gold (Delta table output) layers dynamically.
- Calls scd\_type1\_logic function for each dataset to apply SCD Type 1 logic and update the gold layer.
- Skips any unknown file types and logs a message for awareness.

Displaying the output of the accounts table:

Displaying the output of the customers table:

```
14
1 df_gold_customers = spark.read.format("delta").option("header", True).load(f"{gold_base_path}customers")
2 print("-----customers.csv file Data-----")
3 display(df_gold_customers)
```

Displaying the output of the loans table:

```
15
1 df_gold_loans = spark.read.format("delta").option("header", True).load(f"{gold_base_path}loans")
2 print("-----loans.csv file Data-----")
3 display(df_gold_loans)
```

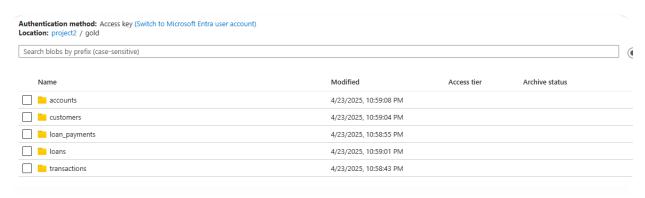
Displaying the output of the loan\_payments table:

```
16
1 df_gold_loan_payments = spark.read.format("delta").option("header", True).load(f"{gold_base_path}loan_payments")
2 print("-----loan_payments.csv file Data-----")
3 display(df_gold_loan_payments)
```

Displaying the output of the transactions table:

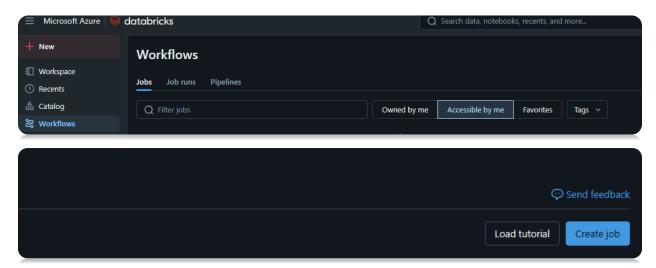
```
17
1 df_gold_transactions = spark.read.format("delta").option("header", True).load(f"{gold_base_path}transactions")
2 print("-----transactions.csv file Data-----")
3 display(df_gold_transactions)
```

## Gold Folder in ADLS Gen 2:



# Step 6: Create a scheduled job for daily run

Go to the Workflow tab -> Jobs -> Create Job

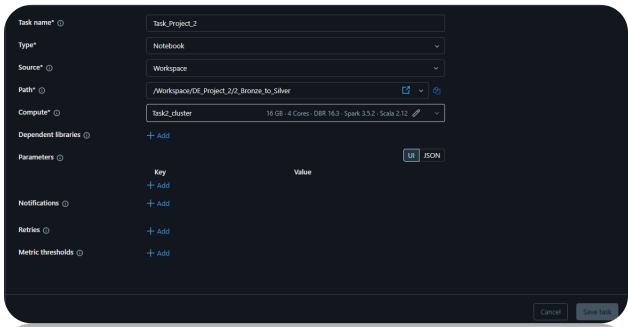


## Create a new task:

## Mention this information

- Task Name: Name of this new task
- Type: Notebook
- Source Workspace (Select the main notebook which we want to run on a schedule trigger)
- Path: Path of that notebook
- Computer: Select a job compute cluster/Create a new one based on the requirement

# Click on Save Task



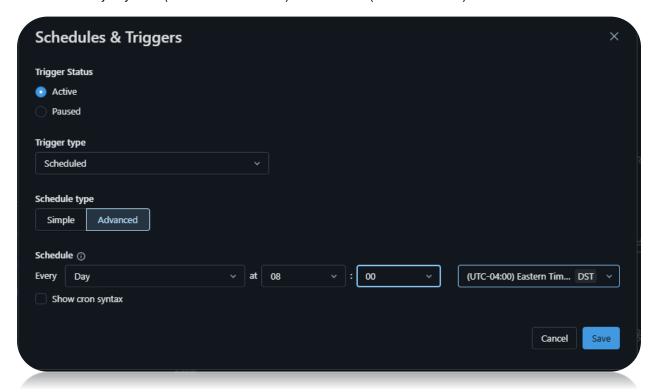
On the right-side panel, click on add schedule trigger and go to advanced options

Here we have selected

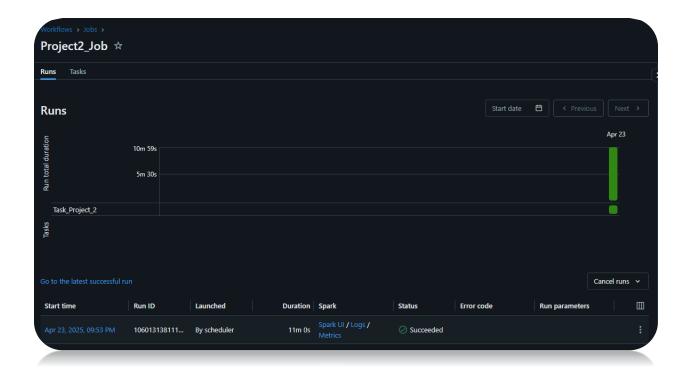
Trigger type: Schedule

Schedule type: Advanced

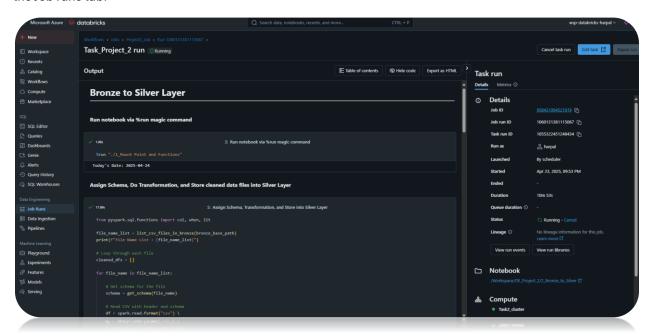
Schedule: Every day 8 AM (UTC -4:00 Time zone) Eastern Time (US and Canada)

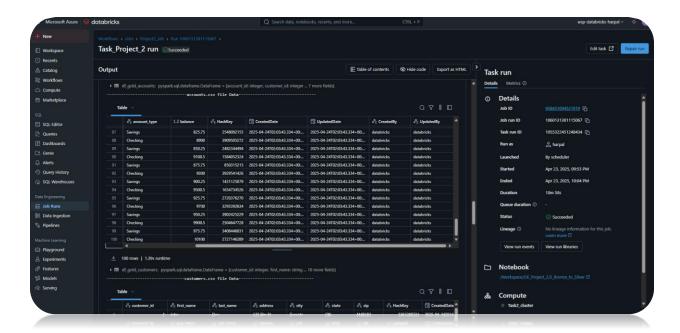


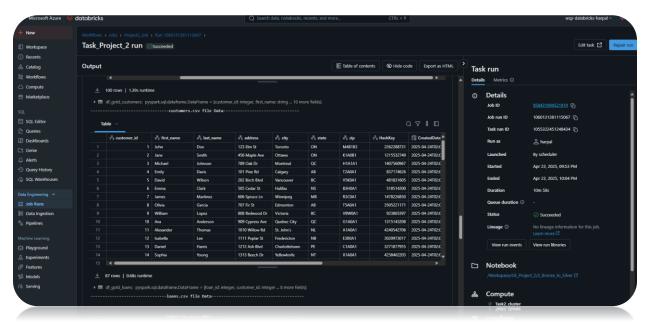
Once the job runs by the scheduler, it will be shown in the Runs tab in the workflow section:



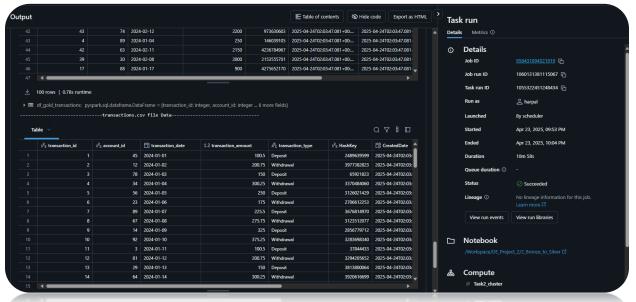
When the job runs on schedule time we can check the process side by side, by clicking on the start time in the Job runs tab:





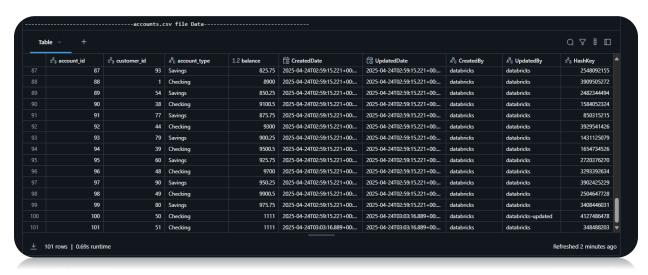


At the end Job cluster will be terminated by itself.



## Day 2 Schedule Job run output:

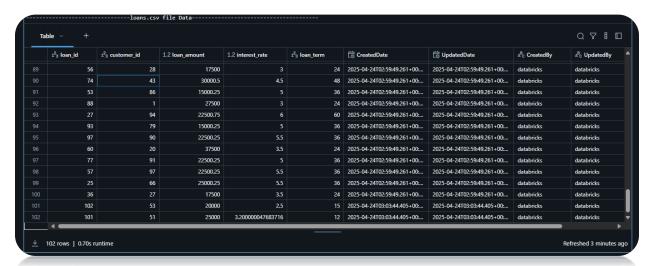
Accounts:



## Customers:



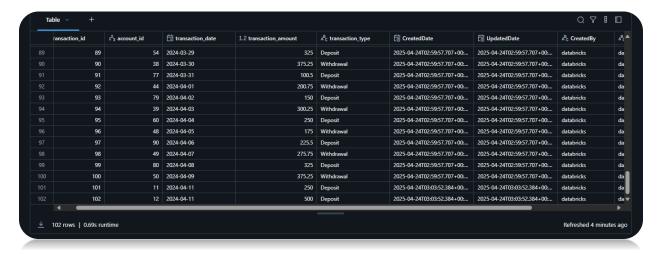
#### Loans:



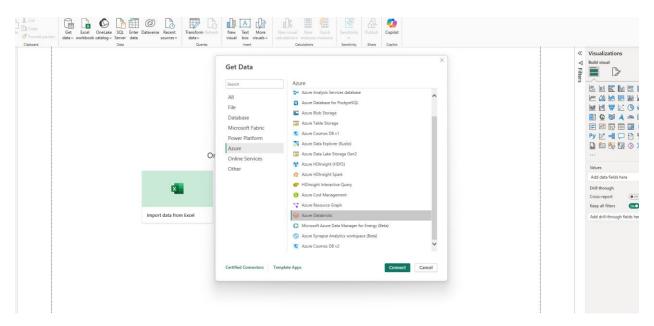
# Loan Payments:



## Transactions:

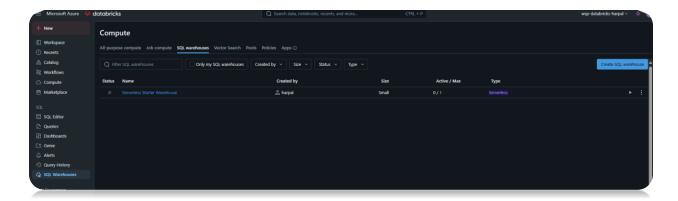


Step 7: Power BI Visualization on DELTA File in Silver Layer

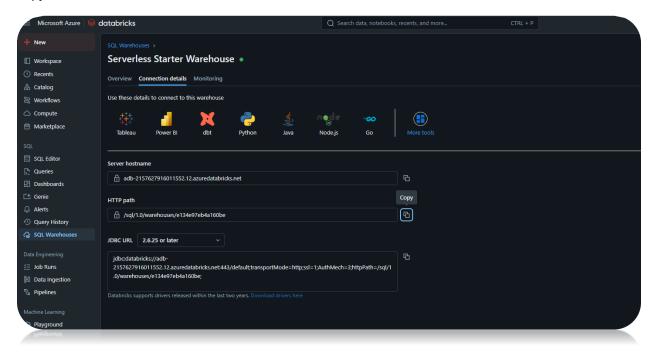


Click on Connect

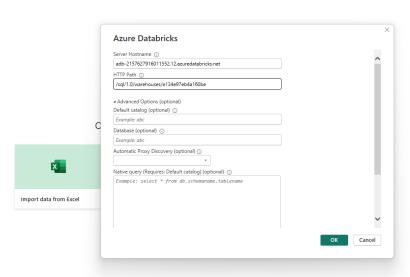
Now, to Databricks SQL Warehouse, start the warehouse



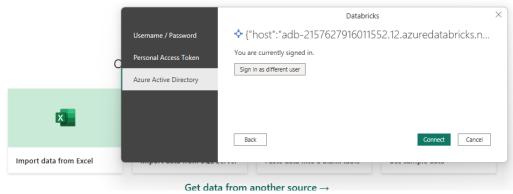
# Copy Hostname and HTTP Path



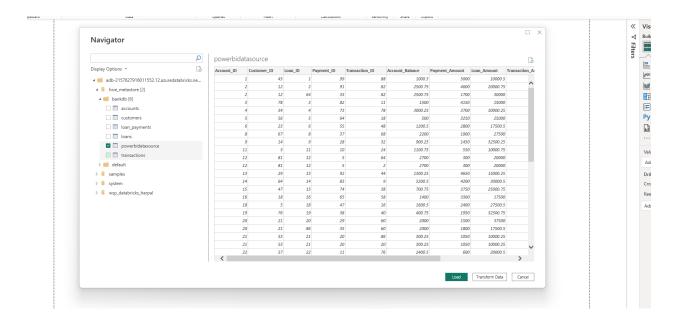
Paste that Hostname and HTTP Path in Power BI window as shown below:



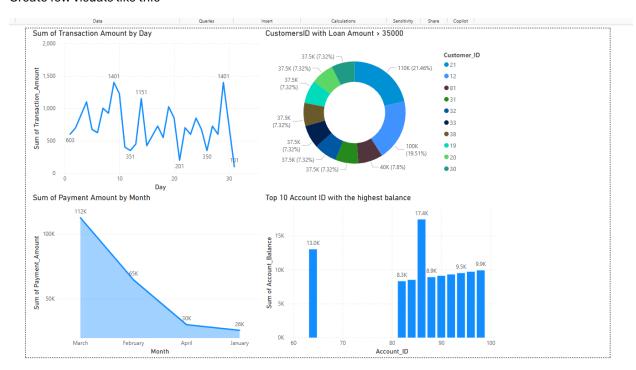
# Click on Connect



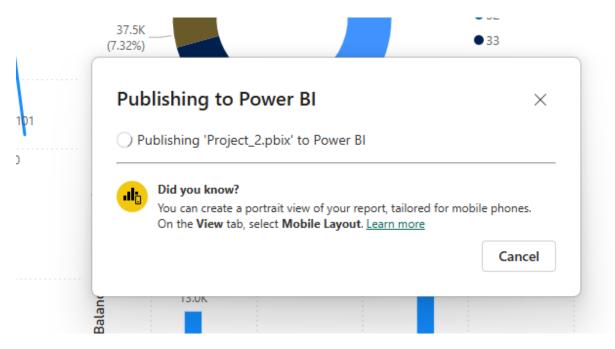
Select the table



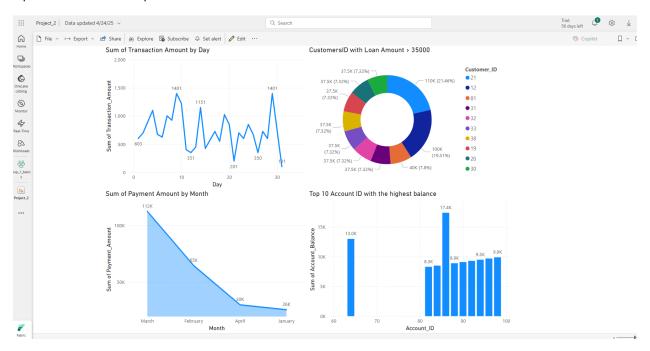
## Create few visuals like this



Step 8: Publish it to the Fabric Workspace



# Open the Fabric Workspace and check the visuals there:



# **Conclusion:**

This project successfully designed and implemented a robust data pipeline in Databricks to process customer account, loan, transaction, and payment data. Using ADLS Gen2, Databricks notebooks, and Delta/Parquet formats, we established a layered architecture (Bronze  $\rightarrow$  Silver  $\rightarrow$  Gold) for scalable transformation, SCD handling, and analytics. Final outputs were integrated with Power BI, enabling powerful visual reporting and business insights.

## Points to remember:

- Parquet files can be overwritten using .mode("overwrite"), but they don't support merge or versioning like
   Delta.
- Use Delta format when you need SCD Type 1, incremental processing, or Power BI connectivity, it supports ACID transactions and time travel.
- Use .saveAsTable("db.table") for managed Delta tables, but never provide a file path with it.
- When overwriting **CSV** or **Parquet** files, use .coalesce(1) to avoid multiple part-\*.csv files if you want a single output.
- To rename output from Spark (e.g., to accounts.csv), write to a temp folder, then rename part-\*.csv, and remove the temp folder.
- Implementing SCD Type 1 logic requires MERGE INTO with Delta format, Parquet files can't handle this.
- Use **functions** like transform\_dataframe() and load\_to\_silver() to modularize and reuse your transformation logic across different files.
- Always create date-based backup folders like /client\_backup\_files/YYYY/MM/DD/ to keep historical snapshots organized.
- In Power BI, we can **connect to Delta tables via Databricks SQL Warehouses**, and use **refresh** to update visuals after new data loads.
- When reading .csv files, always use .schema() instead of inferSchema=True to avoid type mismatches and ensure consistent transformations.
- Implement reusable utility functions for: schema loading, SCD Type 1 logic, safe reads, backup routines, and path building. This saves a lot of code duplication.
- For Power BI, use sum, average, and count aggregations in visuals (not raw tables), and apply filters like "Top N" or "greater than" directly in the Filters pane for performance and clarity.