

# Project Assignment: Data Pipeline for Customer Account Analysis

This document outlines a data pipeline for processing customer account data, including copying data from a backend team's storage account, performing transformations in Databricks, and exposing data for analysis in Azure Synapse Analytics.

## Step 1: Data Ingestion (Backend Storage to Raw(Bronze) Container)

Data ingestion pipeline is a crucial component of modern data architecture, enabling businesses to efficiently manage and utilize their data. It's the process of importing, transferring, loading, and processing data for later use or storage in a database.

### 1.1. Configure Azure Data Factory (ADF) for Data Copy

1. **Sign in to the Azure Portal:**
  - Navigate to the [Azure Portal](#).
2. **Create an Azure Data Factory Instance:**
  - In the Azure Portal, search for "**Data Factory**" and click "**Create**".

- Fill in the necessary details, such as:
  - **Resource Group:** Select an existing group or create a new one.
  - **Name:** Provide a unique name for your Data Factory instance.
  - **Region:** Choose the region where you want to deploy.
- Click "**Review + Create**" and then "**Create**".

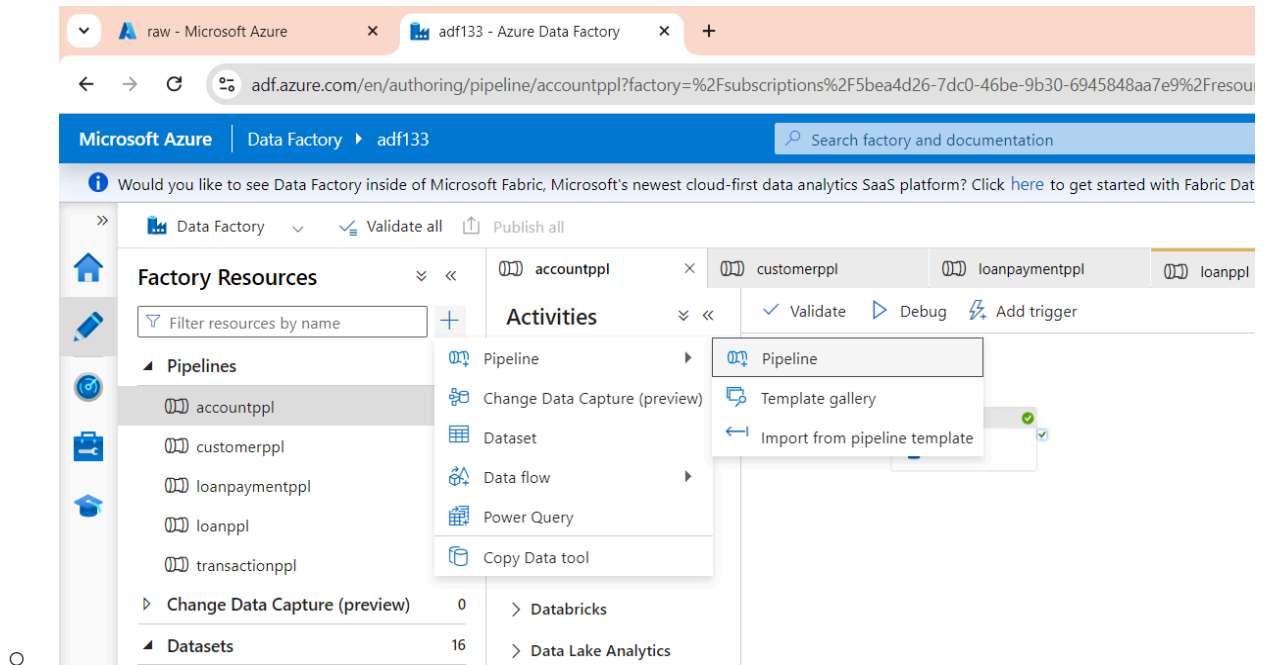
### 3. Set Up Linked Services:

- In the Data Factory, go to **Manage > Linked Services > New**.
- **Create a Linked Service for Backend Storage:**
  - Choose **Azure Blob Storage** as the data store.
  - Enter the storage account details of the backend team's storage account.
  - Use either **Account Key** or **SAS token** for authentication.
- **Create a Linked Service for Your Data Lake Storage:**
  - Repeat the steps above for your own Data Lake Storage account.

## 1.2. Create a Data Factory Pipeline

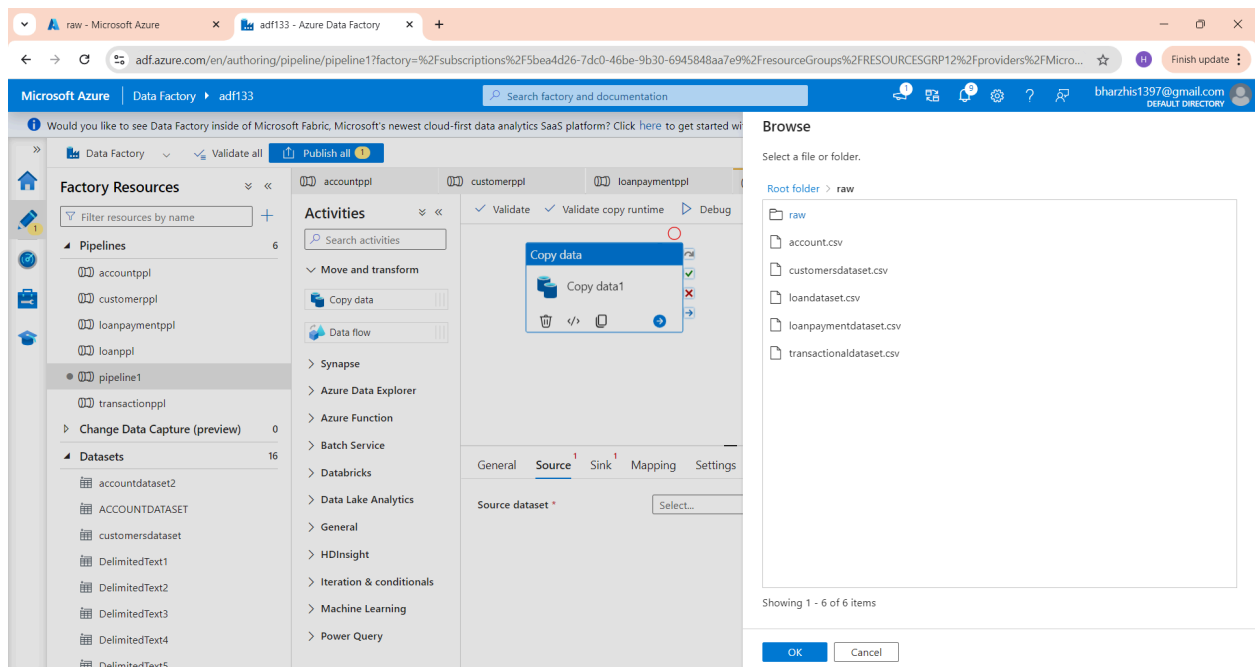
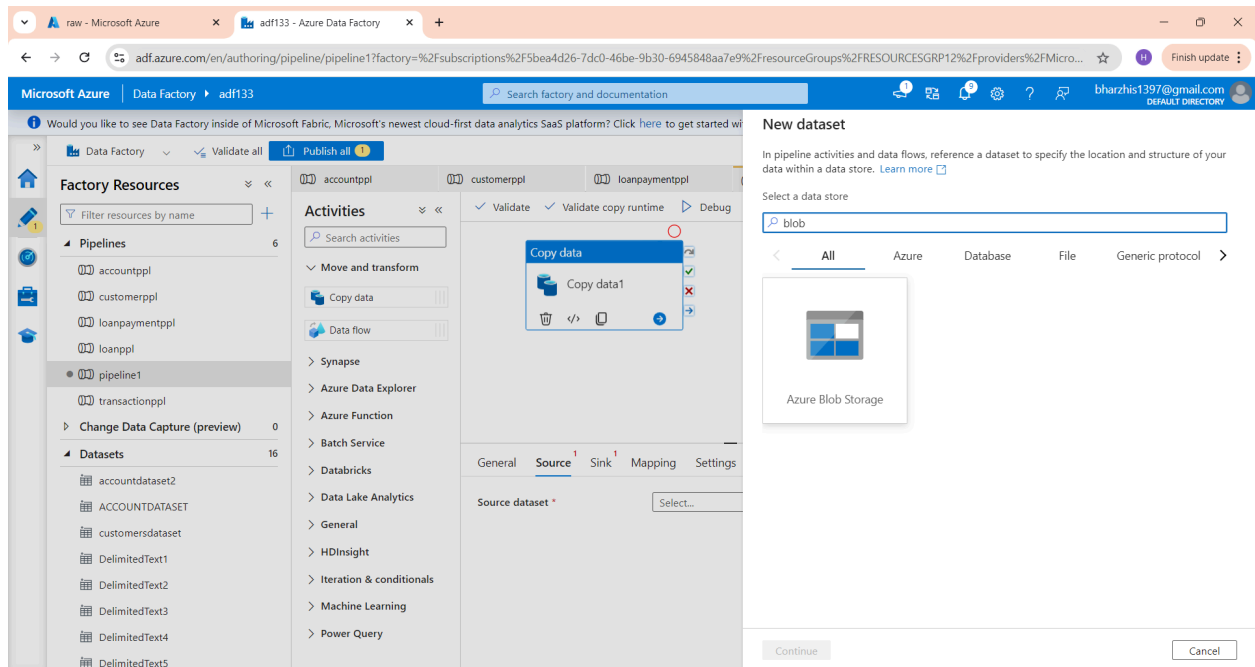
### 1. Navigate to the Author Tab:

- Click "**Author**" > "+" > "**Pipeline**".



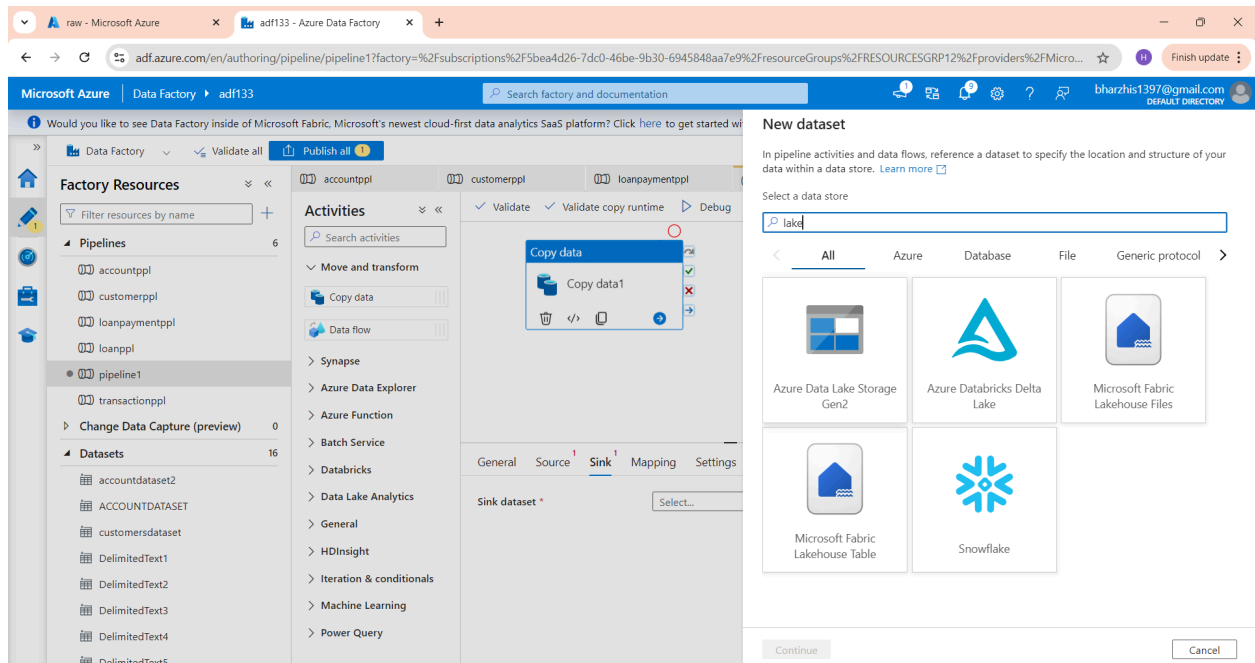
## 2. Add a Copy Data Activity:

- Drag and drop the **Copy Data** activity onto the canvas.
- **Configure the Source:**
  - Select **Source** in the Copy Data activity settings.
  - Click **New** to add a dataset pointing to the backend storage account.
  - Choose **DelimitedText** for CSV files.
  - Specify the path for each file:
    - `accounts.csv`
    - `customers.csv`
    - `loan_payments.csv`
    - `loans.csv`
    - `transactions.csv`

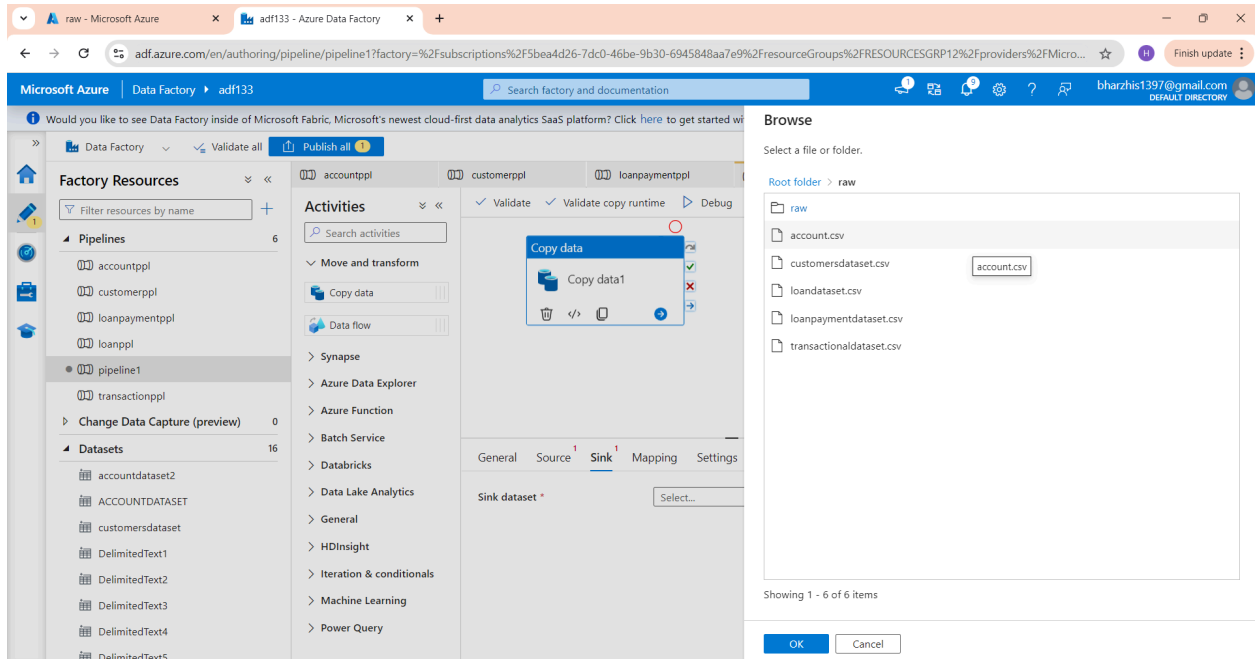


- **Configure the Sink:**

- Select **Sink** in the Copy Data activity settings.
- Click **New** to add a dataset pointing to your Data Lake Storage Raw (Bronze) container.



- Specify the destination paths (e.g., `raw/accounts.csv`, `raw/customers.csv`).

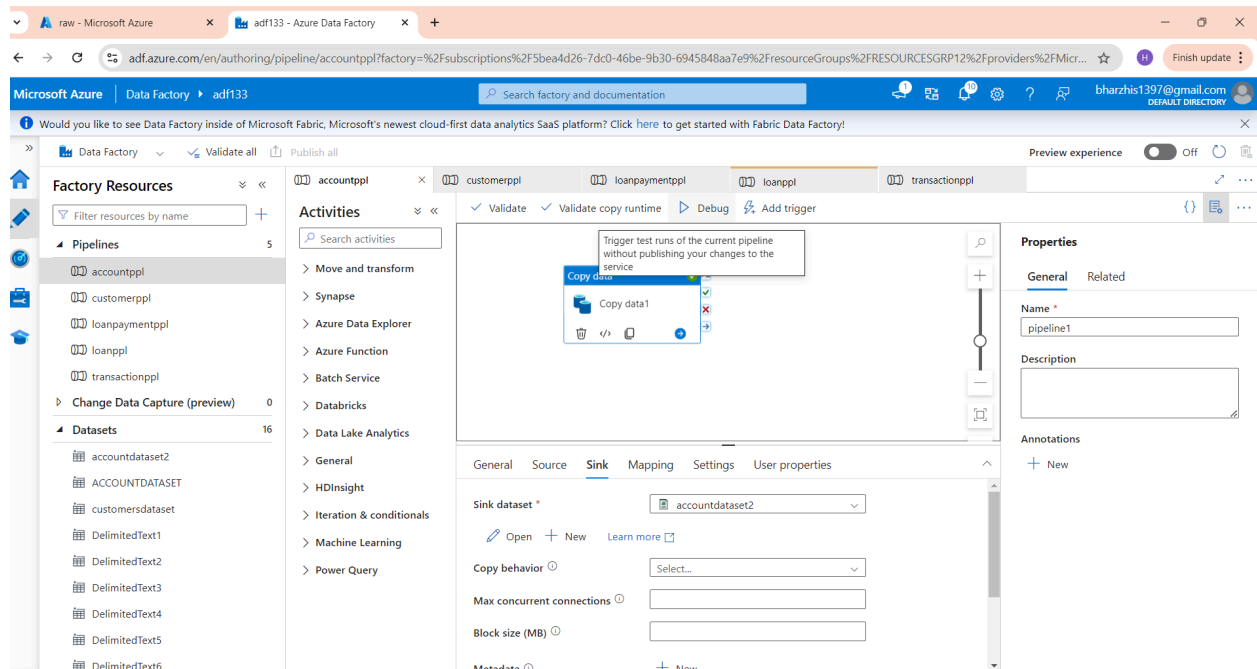


### 3. Set Up Parameters (Optional for Dynamic Configurations):

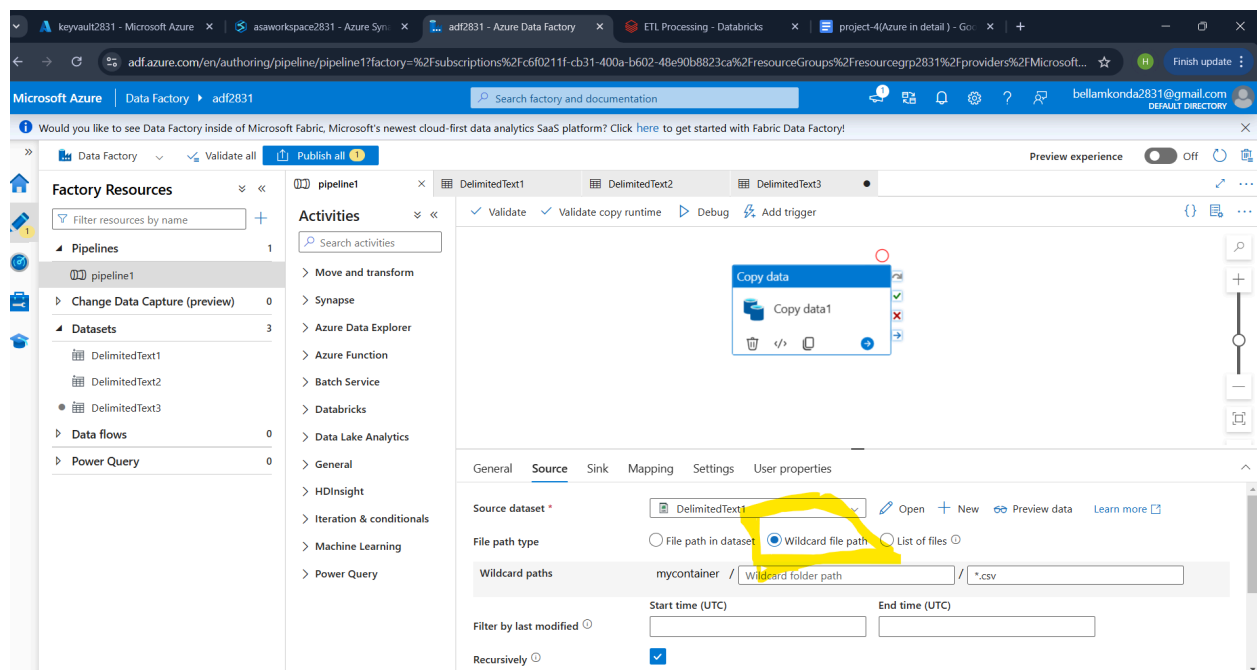
- Use **parameters** to define file paths, making your pipeline flexible and easily configurable.

### 4. Debug and Publish:

- Click **Debug** to test the pipeline.
- If successful, click **Publish All** to save your pipeline.



Note: If you would like to push N number of files into the pipeline, you can achieve it by the option called "wildcard file path" in the source tab as shown in the below screenshot.



## Step 2: Databricks Activity (Incremental/Delta Processing)

### 2.1. Set Up Databricks

**Incremental and delta processing in Databricks** allows for the processing of data in a way that is more efficient and cost-effective than repeated batch jobs.

#### 1. Create an Azure Databricks Workspace:

- In the Azure Portal, search for **"Azure Databricks"** and click **"Create"**.
- Provide the required details and click **"Review + Create"**, then **"Create"**.
- Once created, launch the **Databricks workspace**.

#### 2. Create a Databricks Cluster:

- In the Databricks workspace, go to **Clusters > Create Cluster**.
- Configure the cluster settings (e.g., name, node types) and create

### 2.2. Create a Databricks Notebook for Incremental Processing

#### 1. Set Up a Notebook:

- In Databricks, click **"Create" > "Notebook"**.
- Name the notebook (e.g., **Incremental\_Processing**).
- Choose **Language** as PySpark.

#### 2. Read Data from Raw (Bronze) Container:



Incremental and processing Python

Free trial ends in 10 days. Upgrade to Premium in Azure Portal

4 days ago (2)

```
storage_account_name = "strgacc2831"
storage_account_key = "nK93TB+8Q47M4lV45af/QN9W01RMS1NJYTe3ITBFE/Qt5i7jp2jAtmC/DsNenMt5aNgGjDBxR+ASTA8o5SA=="

# Configure Spark to use the storage account key
spark.conf.set(f"fs.azure.account.key.strgacc2831.dfs.core.windows.net","nK93TB+8Q47M4lV45af/QN9W01RMS1NJYTe3ITBFE/Qt5i7jp2jAtmC/DsNenMt5aNgGjDBxR+ASTA8o5SA==")
```

Add a cell title

4 days ago (4)

```
# Assuming both files are in CSV format
accounts_df = spark.read.format("csv").option("header", "true").load("abfss://bronze@strgacc2831.dfs.core.windows.net/accounts.csv")
customers_df = spark.read.format("csv").option("header", "true").load("abfss://bronze@strgacc2831.dfs.core.windows.net/customers.csv")
loan_payments_df = spark.read.format("csv").option("header", "true").load("abfss://bronze@strgacc2831.dfs.core.windows.net/loan_payments.csv")
loans_df = spark.read.format("csv").option("header", "true").load("abfss://bronze@strgacc2831.dfs.core.windows.net/loans.csv")
transactions_df = spark.read.format("csv").option("header", "true").load("abfss://bronze@strgacc2831.dfs.core.windows.net/transactions.csv")

accounts_df.show(5)
customers_df.show(5)
loan_payments_df.show(5)
loans_df.show(5)
transactions_df.show(5)
```

(10) Spark Jobs

Incremental and processing Python

only showing top 5 rows

4 days ago (4)

```
# Select only the customer_id column from the customer DataFrame
valid_customers = customers_df.select("customer_id").distinct()

# Join accounts with valid_customers on customer_id
cleaned_accounts_df = accounts_df.join(valid_customers, on="customer_id", how="inner")

cleaned_accounts_df.show(20)
```

(3) Spark Jobs

valid_customers: pyspark.sql.dataframe.DataFrame = [customer_id: string]
cleaned_accounts_df: pyspark.sql.dataframe.DataFrame = [customer_id: string, account_id: string ... 2 more fields]

customer_id	account_id	transaction_type	amount
151	381	Checking	3900.50
541	891	Savings	850.25
111	241	Checking	2600.00
291	131	Savings	1300.25
691	651	Savings	550.25
421	541	Checking	5500.50
731	711	Savings	625.75
871	731	Savings	650.25
641	141	Checking	3200.50
31	111	Savings	1100.75
301	801	Checking	8100.00
341	41	Checking	3000.25
591	691	Savings	600.25
81	681	Checking	1000.00

keyvault2831 - Microsoft Azure | asaworkspace2831 - Azure Synapse | adf2831 - Azure Data Factory | Incremental and processing - D | project-4(Azure in detail) - Go |

adb-3761879947263093.13.azuredatabricks.net/editor/notebooks/2589490834336633?o=3761879947263093#command/2589490834336634

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | brickspace2831 | Finish update

New | Workspace | Recents | Catalog | Workflows | Compute | SQL | SQL Editor | Queries | Dashboards | Alerts | Query History | SQL Warehouses | Data Engineering | Job Runs | Data Ingestion | Delta Live Tables | Machine Learning | Playground | Experiments | Features

### Incremental and processing

Python | File | Edit | View | Run | Help | Last edit was 4 days ago | Run all | HARSHITHA BELLAM's ... | Schedule | Share

```
cleaned_accounts_df.write.format("csv").mode("overwrite").option("header", "true").save("abfss://silver@strgacc2831.dfs.core.windows.net/cleaned_accounts.csv")
accounts_df.show(10)
```

▶ (4) Spark Jobs

account_id	customer_id	account_type	balance
1	45	Savings	1000.50
2	12	Checking	2500.75
3	78	Savings	1500.00
4	34	Checking	3000.25
5	56	Savings	500.00
6	23	Checking	1200.50
7	89	Savings	800.75
8	67	Checking	2200.00
9	14	Savings	900.25
10	92	Checking	1800.50

only showing top 10 rows

keyvault2831 - Microsoft Azure | asaworkspace2831 - Azure Synapse | adf2831 - Azure Data Factory | Incremental and processing - D | project-4(Azure in detail) - Go |

adb-3761879947263093.13.azuredatabricks.net/editor/notebooks/2589490834336633?o=3761879947263093#command/2589490834336634

Microsoft Azure | databricks | Search data, notebooks, recents, and more... | CTRL + P | brickspace2831 | Finish update

New | Workspace | Recents | Catalog | Workflows | Compute | SQL | SQL Editor | Queries | Dashboards | Alerts | Query History | SQL Warehouses | Data Engineering | Job Runs | Data Ingestion | Delta Live Tables | Machine Learning | Playground | Experiments | Features

### Incremental and processing

Python | File | Edit | View | Run | Help | Last edit was 4 days ago | Run all | HARSHITHA BELLAM's ... | Schedule | Share

```
# Remove rows with null values in important columns
accounts_df = accounts_df.dropna(subset=["account_id", "balance"])

# Remove duplicate rows
accounts_df = accounts_df.dropDuplicates(["account_id"])
accounts_df.show(10)
```

▶ (2) Spark Jobs

accounts\_df: pyspark.sql.dataframe.DataFrame = [account\_id: string, customer\_id: string ... 2 more fields]

account_id	customer_id	account_type	balance
7	89	Savings	800.75
51	72	Savings	375.75
15	47	Savings	700.75
54	42	Checking	5500.50
11	3	Savings	1100.75
29	58	Savings	75.25
69	59	Savings	600.25
42	36	Checking	4300.50
73	87	Savings	650.25
87	93	Savings	825.75

only showing top 10 rows

keyvault2831 - Microsoft Azure | asaworkspace2831 - Azure Synapse | adb2831 - Azure Data Factory | Incremental and processing - D | project-4(Azure in detail) - Go |

adb-3761879947263093.13.azuredatabricks.net/editor/notebooks/2589490834336633?o=3761879947263093#command/2589490834336634

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P bricspace2831

New Incremental and processing Python Last edit was 4 days ago Run all HARSHITHA BELLAM's ... Schedule Share

Workspace Recents Catalog Workflows Compute SQL SQL Editor Queries Dashboards Alerts Query History SQL Warehouses Data Engineering Job Runs Data Ingestion Delta Live Tables Machine Learning Playground Experiments Features

```
# Filter out accounts with negative balances if not required
accounts_df = accounts_df.filter(accounts_df.balance >= 0)
accounts_df.show(10)

(2) Spark Jobs

accounts_df: pyspark.sql.dataframe.DataFrame = [account_id: string, customer_id: string ... 2 more fields]

+-----+-----+-----+-----+
|account_id|customer_id|account_type|balance|
+-----+-----+-----+-----+
| 7| 89| Savings| 800.75|
| 51| 72| Savings| 375.75|
| 15| 47| Savings| 700.75|
| 54| 42| Checking|5500.50|
| 11| 3| Savings|1100.75|
| 29| 58| Savings| 75.25|
| 69| 59| Savings| 600.25|
| 42| 36| Checking|4300.50|
| 73| 87| Savings| 650.25|
| 87| 93| Savings| 825.75|
+-----+-----+-----+-----+
only showing top 10 rows
```

```
# Save cleaned data as Delta format for efficient storage
```

keyvault2831 - Microsoft Azure | asaworkspace2831 - Azure Synapse | adb2831 - Azure Data Factory | Incremental and processing - D | project-4(Azure in detail) - Go |

adb-3761879947263093.13.azuredatabricks.net/editor/notebooks/2589490834336633?o=3761879947263093#command/2589490834336634

Microsoft Azure databricks Search data, notebooks, recents, and more... CTRL + P bricspace2831

New Incremental and processing Python Last edit was 4 days ago Run all HARSHITHA BELLAM's ... Schedule Share

Workspace Recents Catalog Workflows Compute SQL SQL Editor Queries Dashboards Alerts Query History SQL Warehouses Data Engineering Job Runs Data Ingestion Delta Live Tables Machine Learning Playground Experiments Features

```
+-----+-----+-----+-----+
| 7| 89| Savings| 800.75|
| 51| 72| Savings| 375.75|
| 15| 47| Savings| 700.75|
| 54| 42| Checking|5500.50|
| 11| 3| Savings|1100.75|
| 29| 58| Savings| 75.25|
| 69| 59| Savings| 600.25|
| 42| 36| Checking|4300.50|
| 73| 87| Savings| 650.25|
| 87| 93| Savings| 825.75|
+-----+-----+-----+-----+
only showing top 10 rows
```

```
# Save cleaned data as Delta format for efficient storage
accounts_df.write.format("delta").mode("overwrite").save("/mnt/datalake/silver/accounts_cleaned")
customers_df.write.format("delta").mode("overwrite").save("/mnt/datalake/silver/customers_cleaned")

(8) Spark Jobs
```

[Shift+Enter] to run and move to next cell

## Step 3: Databricks Activity (ETL Processing)

### 3.1. Create Another Databricks Notebook for ETL

#### 1. Set Up a New Notebook:

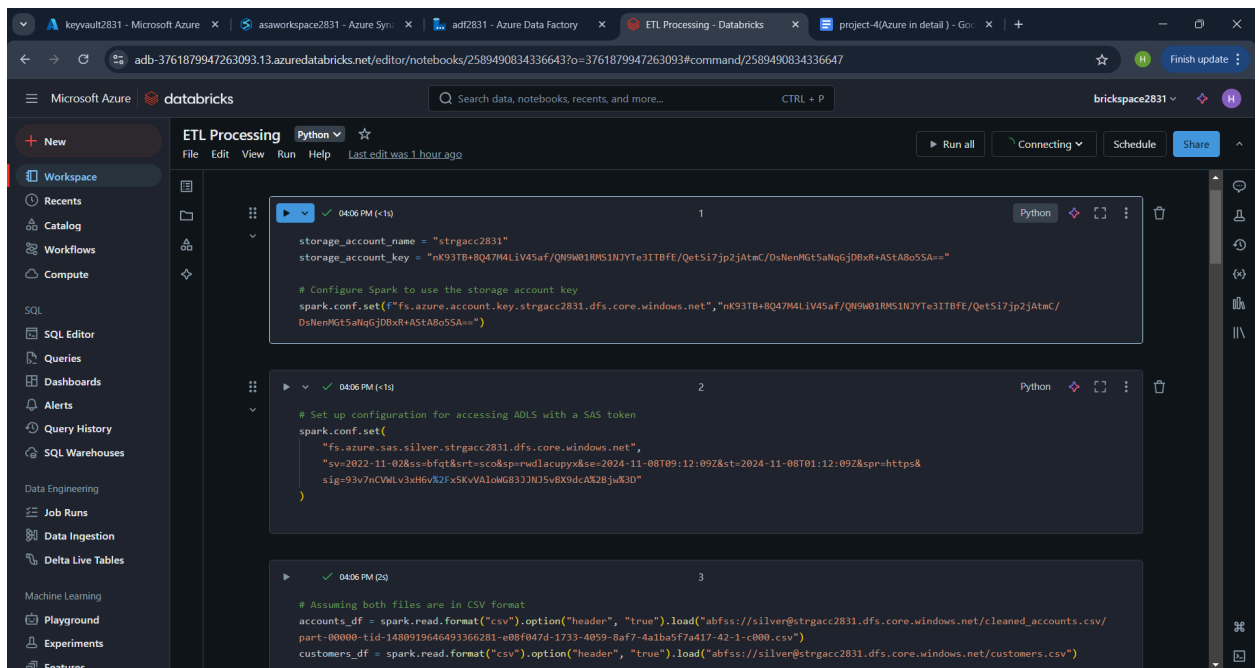
- Create another notebook named **ETL\_Processing**.

#### 2. Read Data from Curated (Silver) Container:

### 3.2 Transformation Logic

#### 1. Calculate Total Balance for Each Customer:

### 3.3 Save Data to Refined (Gold) Container



The screenshot shows a Databricks notebook interface with the title 'ETL Processing'. The notebook contains three code cells, each with a status bar indicating it was executed at 04:06 PM. The first cell sets the storage account name and key, and configures Spark to use the storage account key. The second cell sets up configuration for accessing ADLS with a SAS token. The third cell reads data from the silver container into two DataFrames: 'accounts\_df' and 'customers\_df'.

```
storage_account_name = "strgacc2831"
storage_account_key = "nk93TB+8Q47M4L1V45af/QN9W01RMS1NJYTe3ITBFE/Qt517jp2jAtmC/DsNenMgt5aIqGjDBxR+AStA8o5SA=="

# Configure Spark to use the storage account key
spark.conf.set(("fs.azure.account.key.strgacc2831.dfs.core.windows.net","nk93TB+8Q47M4L1V45af/QN9W01RMS1NJYTe3ITBFE/Qt517jp2jAtmC/DsNenMgt5aIqGjDBxR+AStA8o5SA=="))

# Set up configuration for accessing ADLS with a SAS token
spark.conf.set({
    "fs.azure.sas.silver.strgacc2831.dfs.core.windows.net",
    "sv=2022-11-02&ss=bfqt&srt=sco&sp=rwdlacupyx&se=2024-11-08T09:12:09Z&st=2024-11-08T01:12:09Z&spr=https&sig=93v7nCVWLv3xH6v%2F5KvVA1oHg83JJU5vBX9dcAR2Bjw%3D"
})

# Assuming both files are in CSV format
accounts_df = spark.read.format("csv").option("header", "true").load("abfss://silver@strgacc2831.dfs.core.windows.net/cleaned_accounts.csv/part-00000-tid-1480919646493366281-e08f047d-1733-4b59-8af7-4a1ba5f7a417-42-1-c000.csv")
customers_df = spark.read.format("csv").option("header", "true").load("abfss://silver@strgacc2831.dfs.core.windows.net/customers.csv")
```

ETL Processing Python

```
from pyspark.sql import functions as F

# Convert balance column to float if necessary
accounts_df = accounts_df.withColumn("balance", F.col("balance").cast("float"))

# Join customers and accounts DataFrames on customer_id and select specific columns to avoid ambiguity
joined_df = customers_df.alias("cust").join(accounts_df.alias("acct"), F.col("cust.customer_id") == F.col("acct.customer_id"), "inner") \
    .select(
        F.col("cust.customer_id").alias("customer_id"),
        "first_name",
        "last_name",
        "address",
        "city",
        "state",
        "zip",
        "account_type",
        "balance"
    )

# Aggregate to calculate the total balance for each customer
result_df = joined_df.groupBy(
    "customer_id",
    "first_name",
    "last_name",
    "address",
    "city",
    "state",
    "zip",
    "account_type"
)
```

ETL Processing Python

```
# Aggregate to calculate the total balance for each customer
result_df = joined_df.groupBy(
    "customer_id",
    "first_name",
    "last_name",
    "address",
    "city",
    "state",
    "zip",
    "account_type"
).agg(
    F.sum("balance").alias("total_balance")
)

# Display the result
result_df.show()
```

(3) Spark Jobs

- accounts\_df: pyspark.sql.dataframe.DataFrame = [customer\_id: string, account\_id: string ... 2 more fields]
- joined\_df: pyspark.sql.dataframe.DataFrame = [customer\_id: string, first\_name: string ... 7 more fields]
- result\_df: pyspark.sql.dataframe.DataFrame = [customer\_id: string, first\_name: string ... 7 more fields]

26	Abigail	Parker	2525 Poplar St	Barrie	ON	L4M0A1	Checking	6700.5
20	Nia	Nelson	1919 Birch Blvd	London	ON	N6A0A1	Checking	6100.0
5	David	Wilson	202 Birch Blvd	Vancouver	BC	V5K0A1	Checking	1600.5
53	James	Jenkins	5252 Willow Rd	Queensville	ON	L0G0A1	Savings	300.25
83	David	Fisher	8282 Ash Blvd	Vernor	ON	P0H0A1	Savings	275.75

ETL Processing Python

```
# Write the sorted DataFrame to Azure Data Lake Storage in CSV format
result_df.write.format("csv").mode("overwrite").option("header", "true").save("abfss://silver@strgacc2831.dfs.core.windows.net/total_balance_per_customer.csv")
```

(3) Spark Jobs

```
# Assuming 'total_balance_df' is the DataFrame with the aggregated total balances
gold_container_path = "abfss://gold@strgacc2831.dfs.core.windows.net/total_balance_per_customer.csv"

result_df.write.format("csv").mode("overwrite").option("header", "true").save("abfss://gold@strgacc2831.dfs.core.windows.net/total_balance_per_customer.csv")

result_df.show()
```

(6) Spark Jobs

26	Abigail	Parker	2525 Poplar St	Barrie	ON	L4M0A1	Checking	6700.5
20	Nia	Nelson	1919 Birch Blvd	London	ON	N6A0A1	Checking	6100.0
5	David	Wilson	202 Birch Blvd	Vancouver	BC	V5K0A1	Checking	1600.5
53	James	Jenkins	5252 Willow Rd	Queensville	ON	L0G0A1	Savings	300.25
83	David	Fisher	8282 Ash Blvd	Vernon	ON	P0M0A1	Savings	275.75
40	Sophia	Rivera	3939 Poplar St	Milton	ON	L9T0A1	Checking	8500.0
57	David	Patterson	5656 Cedar Ln	King City	ON	L7B0A1	Savings	350.25
82	Elizabeth	Reynolds	8181 Poplar St	Sturgeon Falls	ON	P2B0A1	Savings	775.75
65	Daniel	Bryant	6464 Redwood Dr	Elmvale	ON	L0L0A1	Savings	800.25
21	Andrew	Mitchell	2020 Spruce Ln	Hamilton	ON	L8P0A1	Checking	10700.5
71	Christopher	Mueller	3070 Cedar Ln	Goldthorn	ON	L0P0A1	Savings	425.75

## Step 4: Azure Synapse Analytics

### Create External Tables in Synapse

Connect to Synapse Studio and create a SQL Database in 'Data' tab

The screenshot shows the Microsoft Azure Synapse Analytics interface. The left sidebar displays the 'Data' section with a tree view of resources. The main area shows a SQL script being executed. The script is as follows:

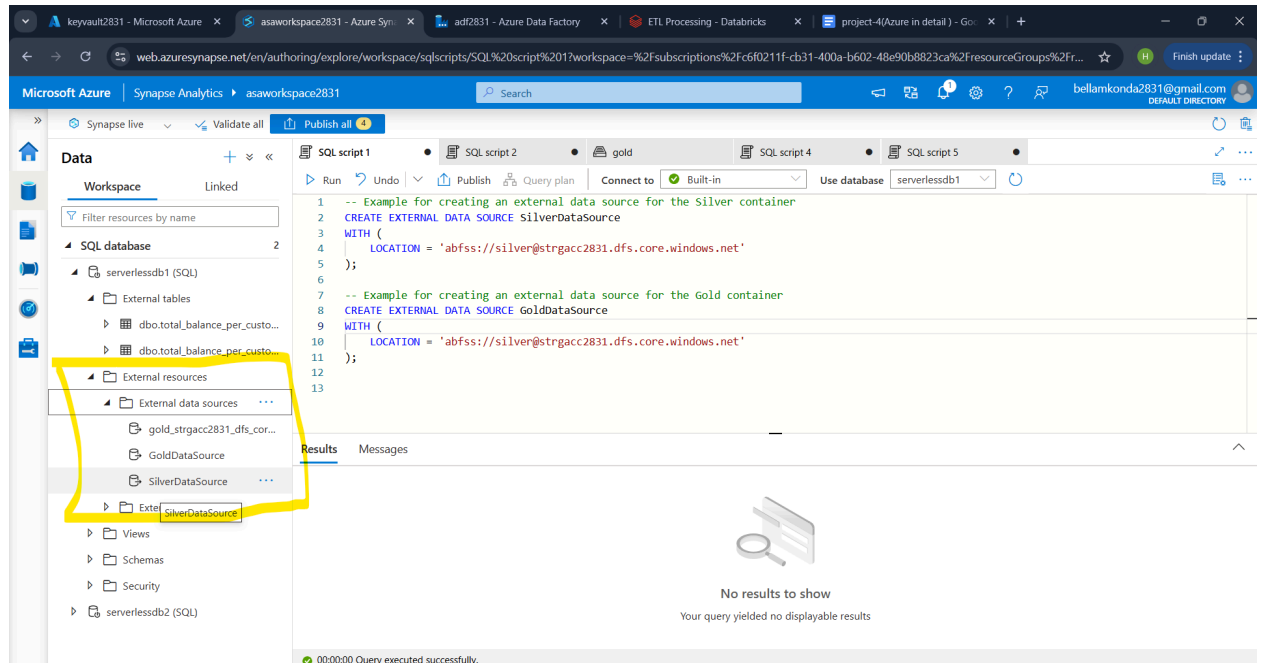
```
1 -- Example for creating an external data source for the Silver container
2 CREATE EXTERNAL DATA SOURCE SilverDataSource
3 WITH (
4     LOCATION = 'abfss://silver@strgacc2831.dfs.core.windows.net'
5 );
6
7 -- Example for creating an external data source for the Gold container
8 CREATE EXTERNAL DATA SOURCE GoldDataSource
9 WITH (
10     LOCATION = 'abfss://silver@strgacc2831.dfs.core.windows.net'
11 );
12
13
```

The 'Use database' dropdown is set to 'master'. The 'Create SQL database' dialog is open on the right, showing an error message: 'A database with the name 'serverlessdb1' already exists. Please provide a different name.'

**Define External Data Sources for Silver and Gold containers by selecting appropriate database in 'use database'**

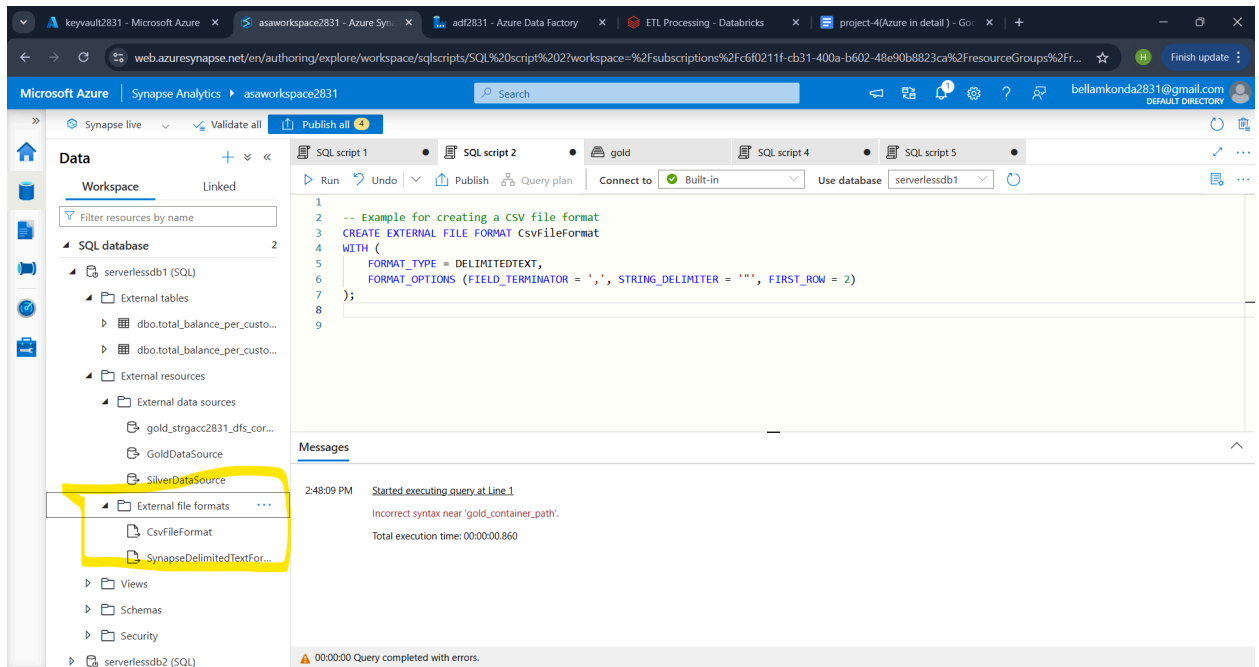
The screenshot shows the Microsoft Azure Synapse Analytics interface. The left sidebar displays the 'Develop' section with a tree view of resources. The main area shows the same SQL script being executed. The 'Use database' dropdown is now set to 'serverlessdb1'. The 'Results' section shows 'No results to show'.

Verify whether the data sources are created for both the containers in the Data>workspace.

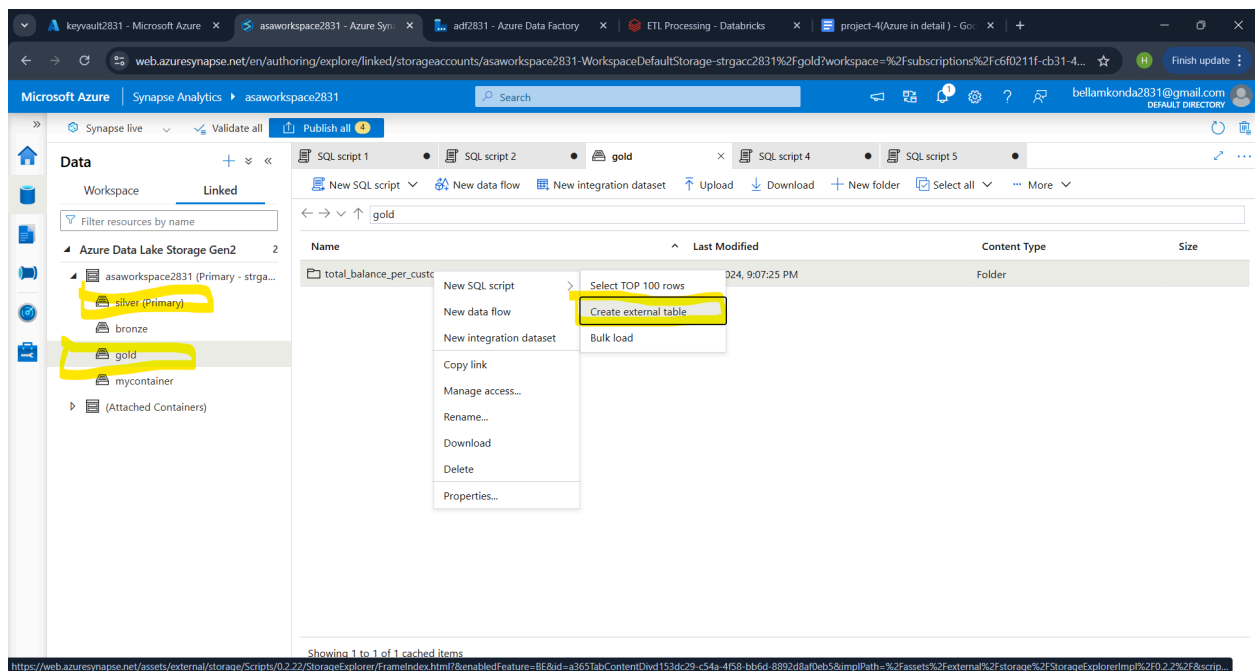


Create an External file format in which you want to share it with the Analyst or scientist teams.

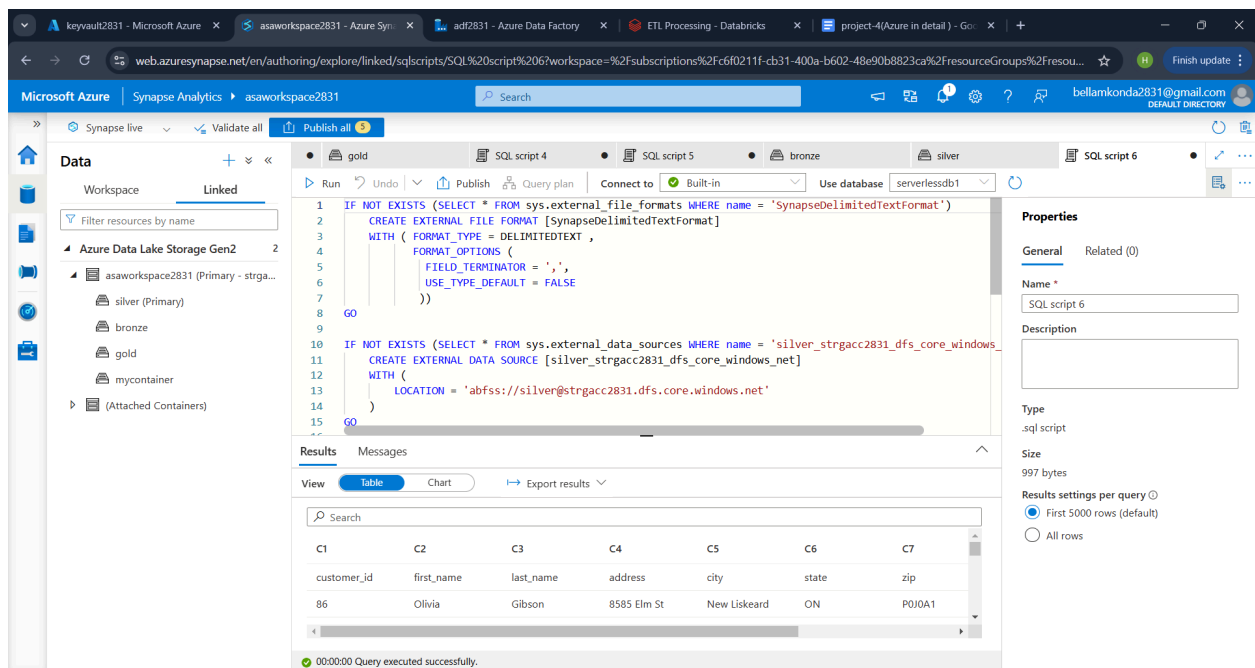
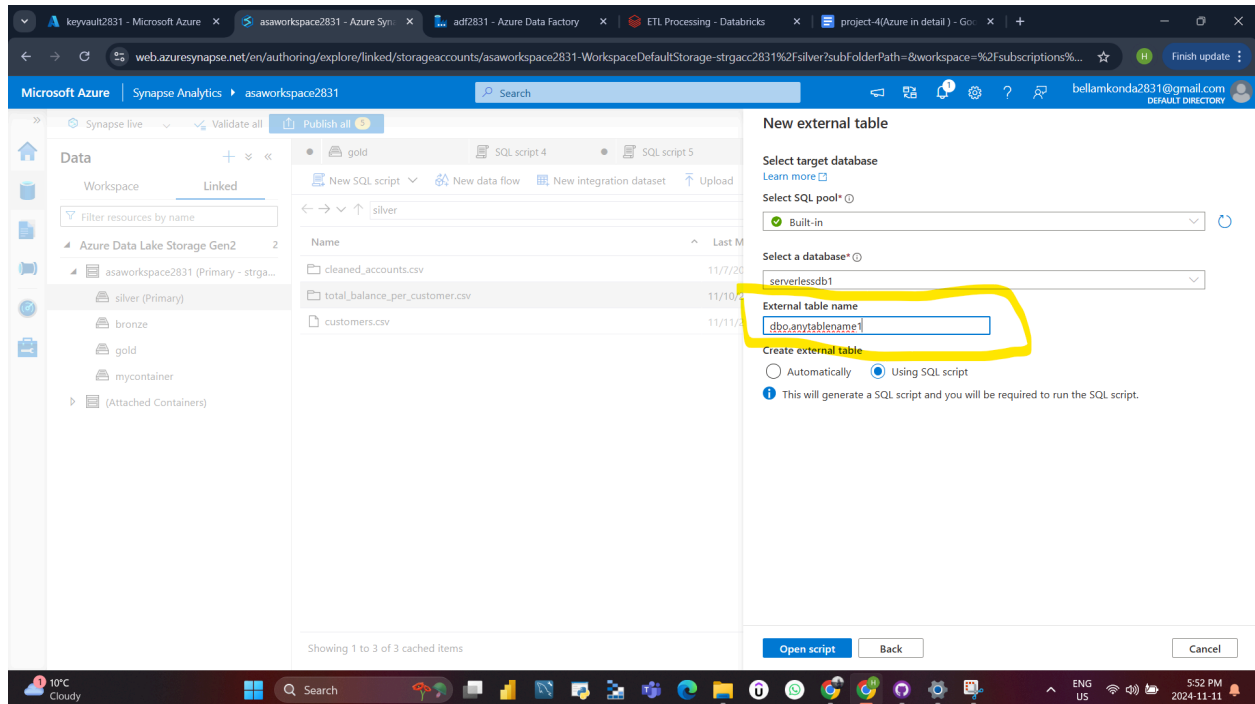




**Create External Tables** in Synapse for both the curated and refined data. Choose the container in Data>linked for which you wanted to create an external table.



Name the external table name a unique table name followed by schema  
dbo(default schema name).



This is how we can create an External table with an Automatic code generator.

This allows data analysts and business intelligence teams to access and query the data directly using tools like Synapse Studio or notebooks.