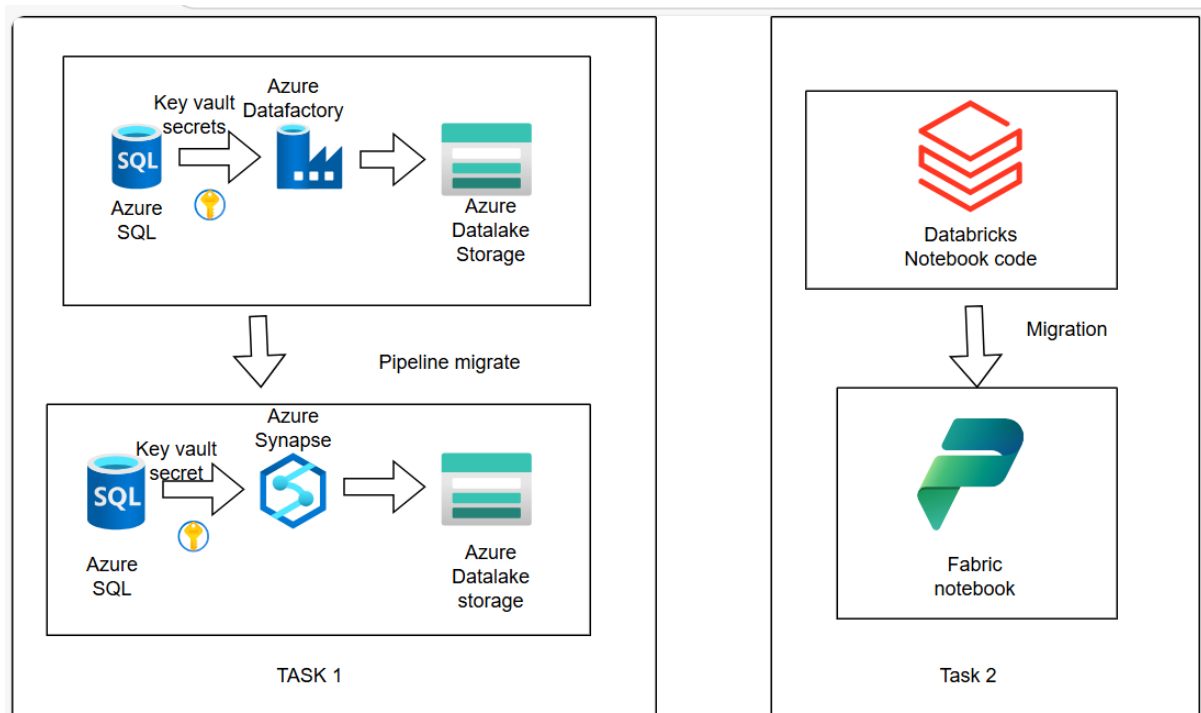# Bootcamp Project 5

## Project Title: Migrating pipelines from ADF to Synapse

## Problem Statement:

Develop ADF pipelines (copy activity, foreach loop, look up) and migrate the pipelines, datasets, linked services to Azure Synapse.
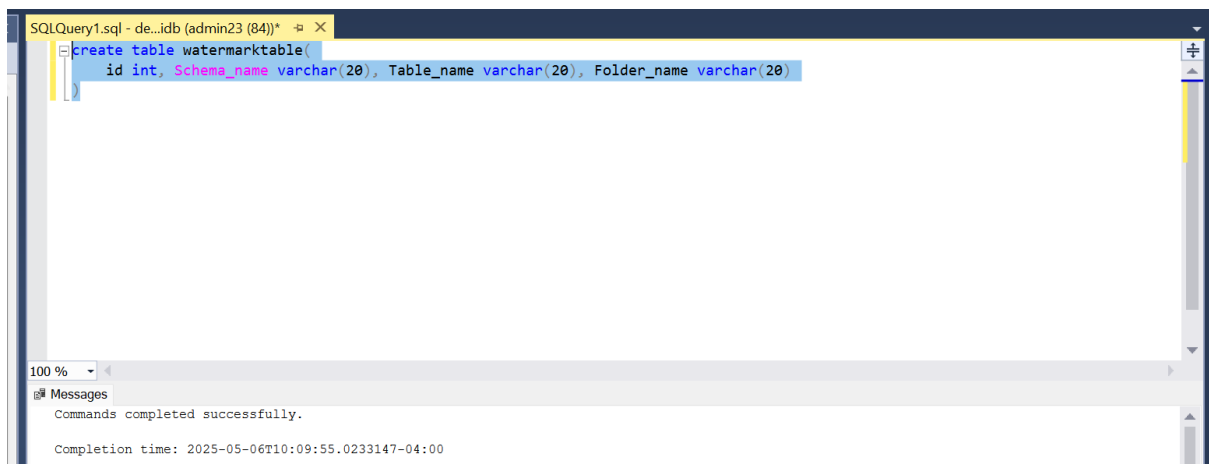
## Architecture Diagram:



## Tools & Technologies:

• Azure Data factory
• Azure Synapse
• Azure SQL
• Azure data lake Gen2
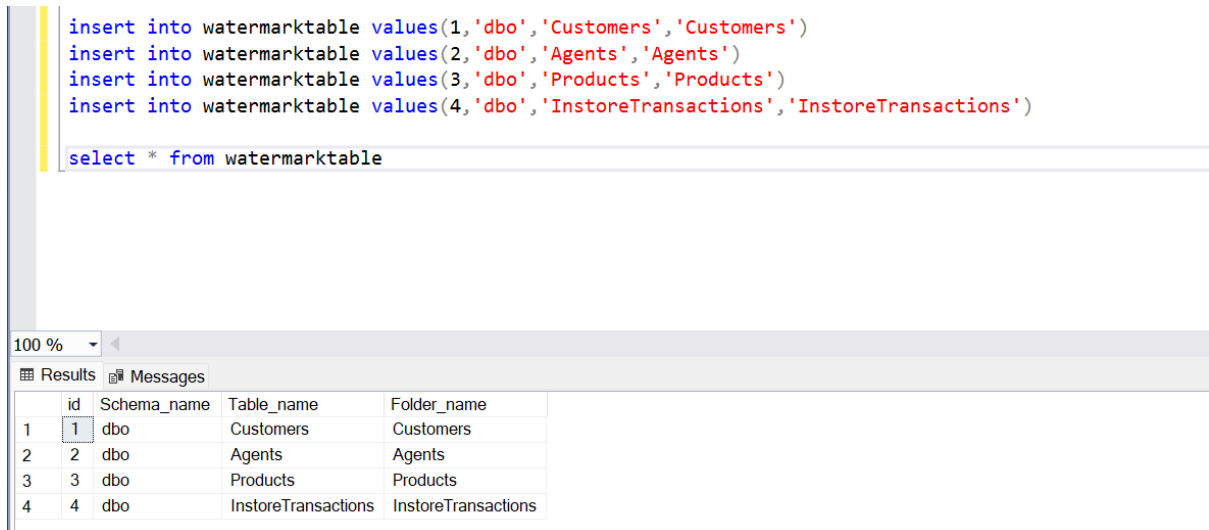• Azure key vault
• Draw.io
• Databricks
• Fabric

## Task 1:
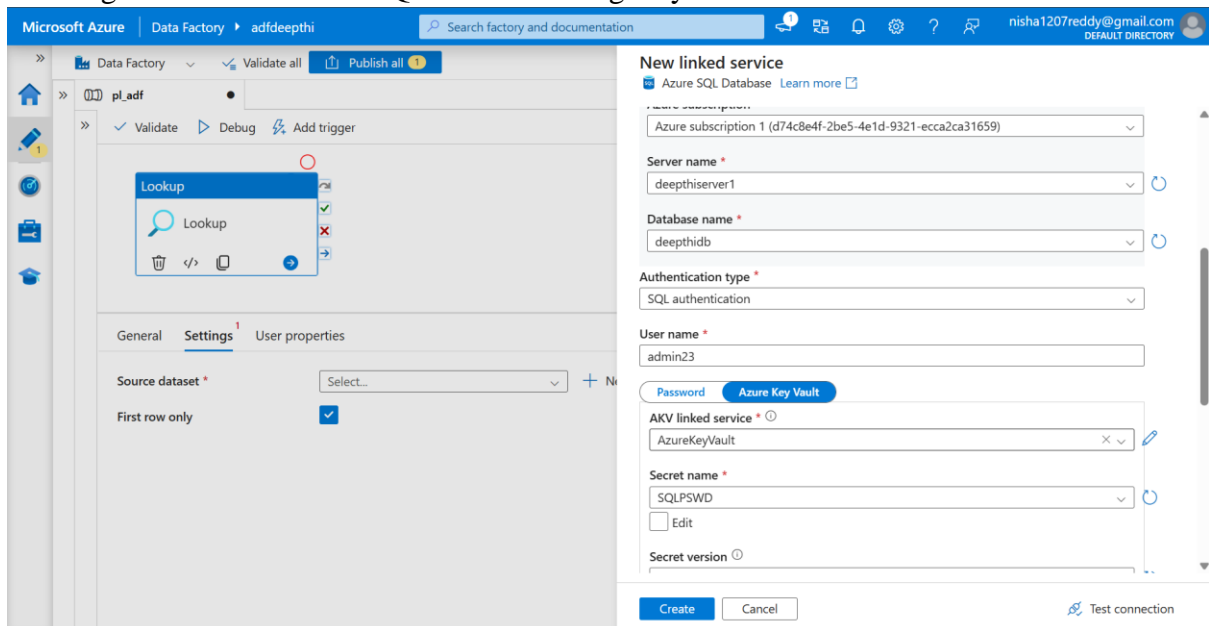**Migrating pipeline from ADF to Azure synapse.**
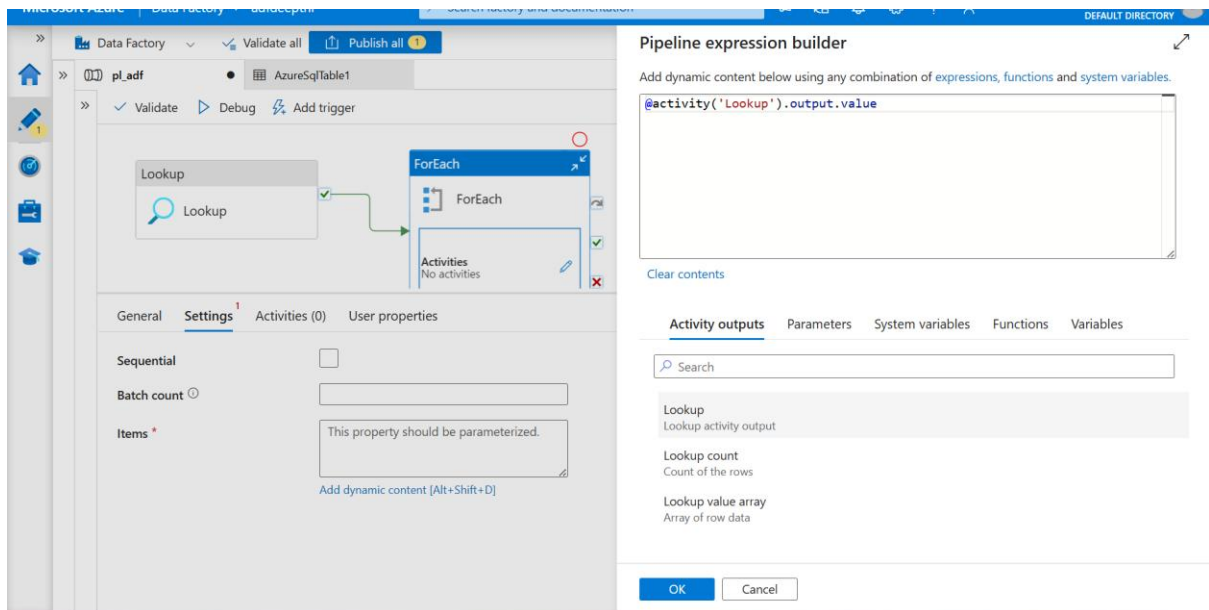
Creating a meta_table

```sql
create table watermarktable(
    id int, Schema_name varchar(20), Table_name varchar(20), Folder_name varchar(20)
)
```

Commands completed successfully.

Completion time: 2025-05-06T10:09:55.0233147-04:00

Values inserted into table

```sql
insert into watermarktable values(1,'dbo','Customers','Customers')
insert into watermarktable values(2,'dbo','Agents','Agents')
insert into watermarktable values(3,'dbo','Products','Products')
insert into watermarktable values(4,'dbo','InstoreTransactions','InstoreTransactions')

select * from watermarktable
```

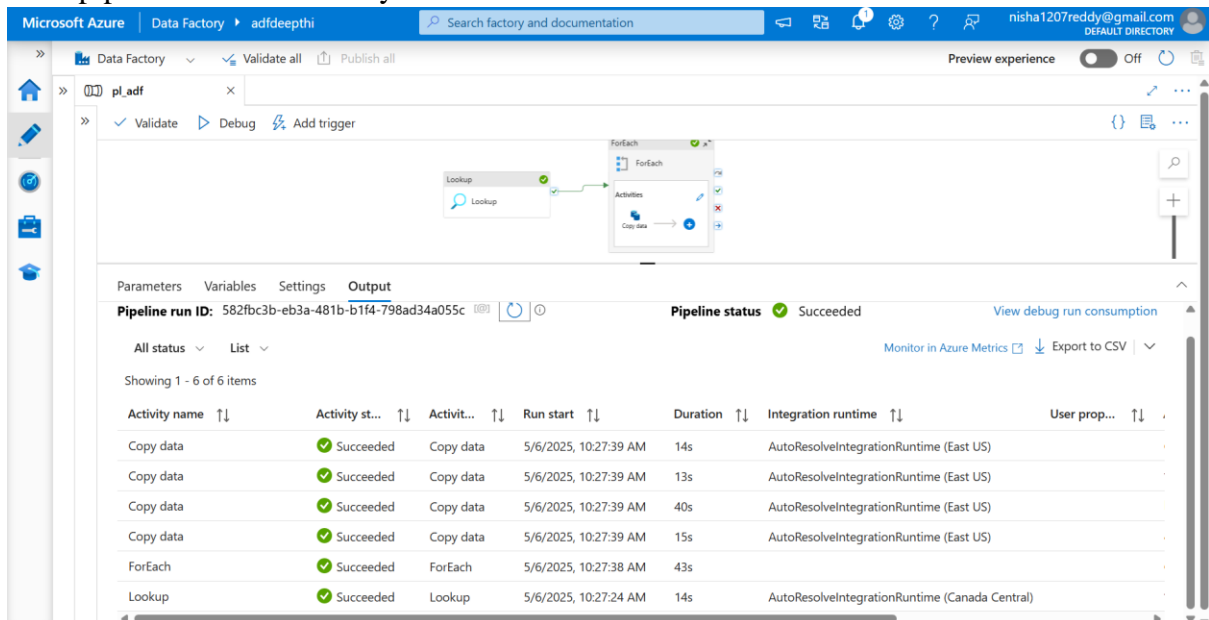| | id | Schema_name | Table_name | Folder_name |
|---|---|---|---|---|
| 1 | 1 | dbo | Customers | Customers |
| 2 | 2 | dbo | Agents | Agents |
| 3 | 3 | dbo | Products | Products |
| 4 | 4 | dbo | InstoreTransactions | InstoreTransactions |

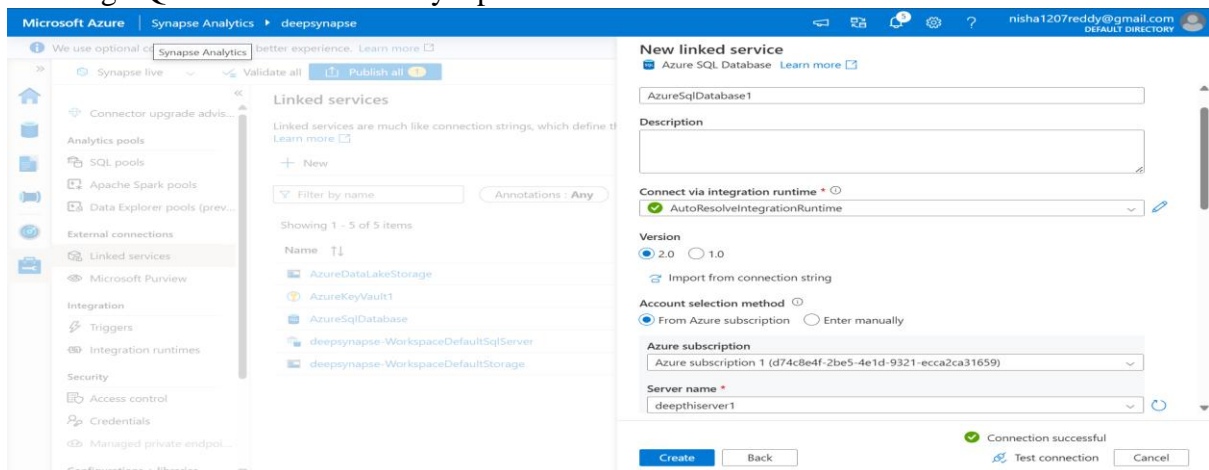Creating linked services for SQL database using Key vault.



Connecting lookup to foreach activity
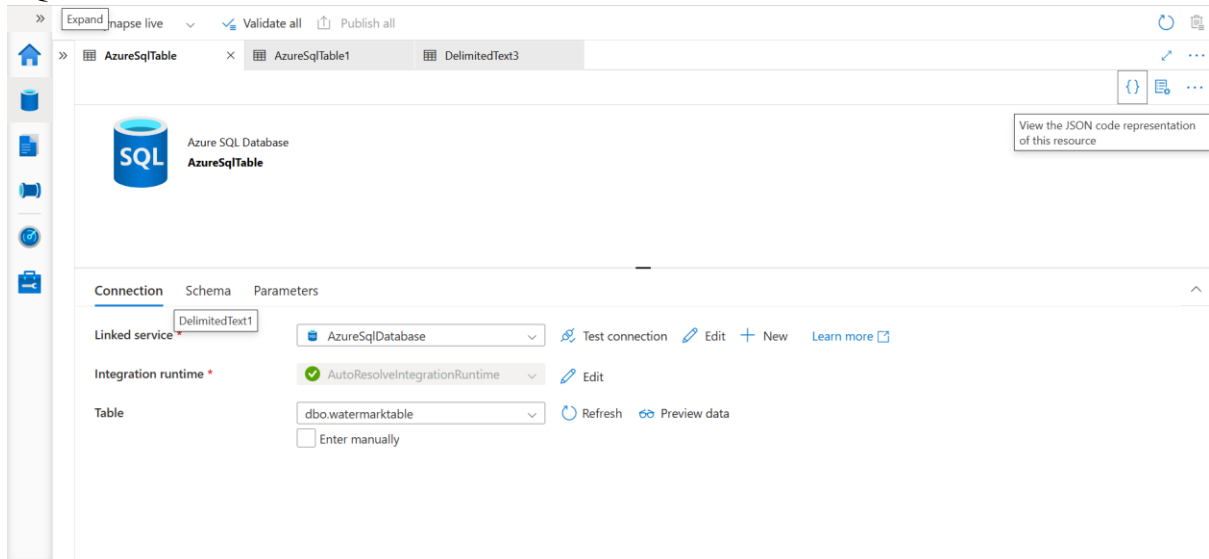
ADF pipeline ran successfully



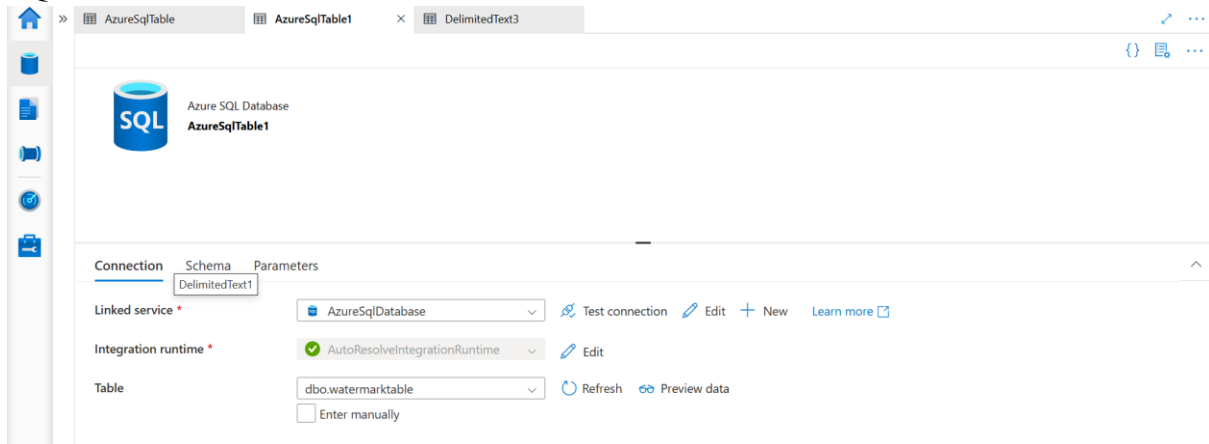Synapse Pipeline:

Creating SQL linked service in Synapse

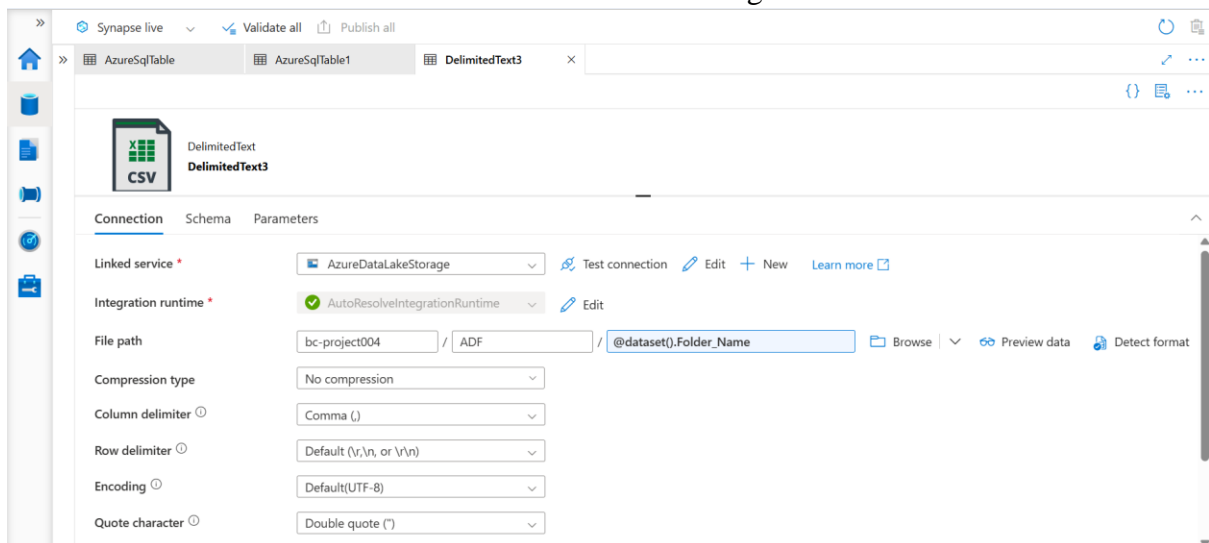Creating Data sets which are in Data factory in Synapse as well.
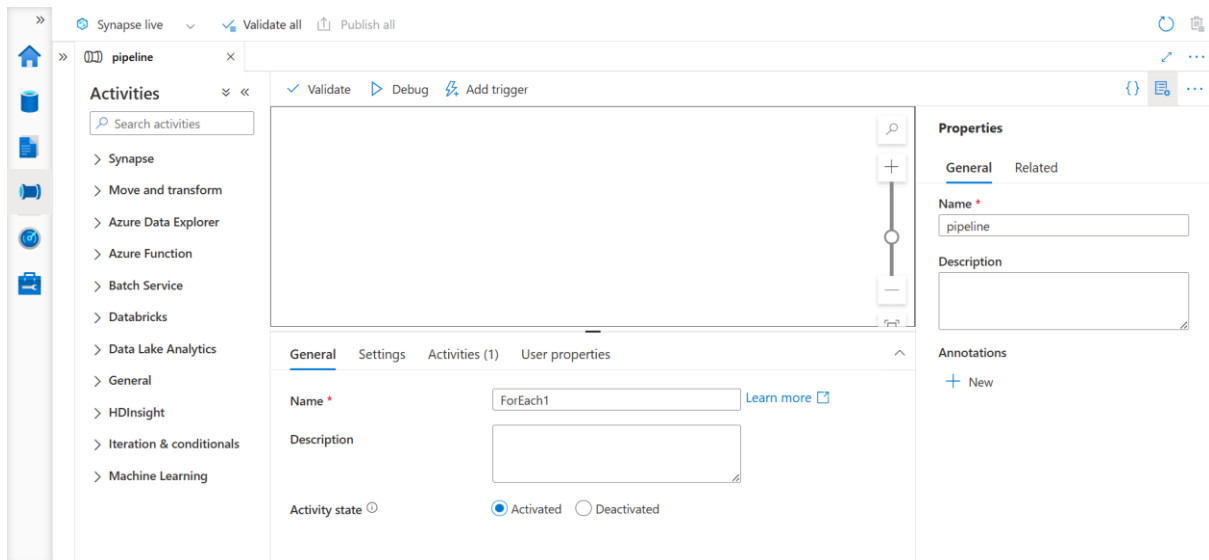
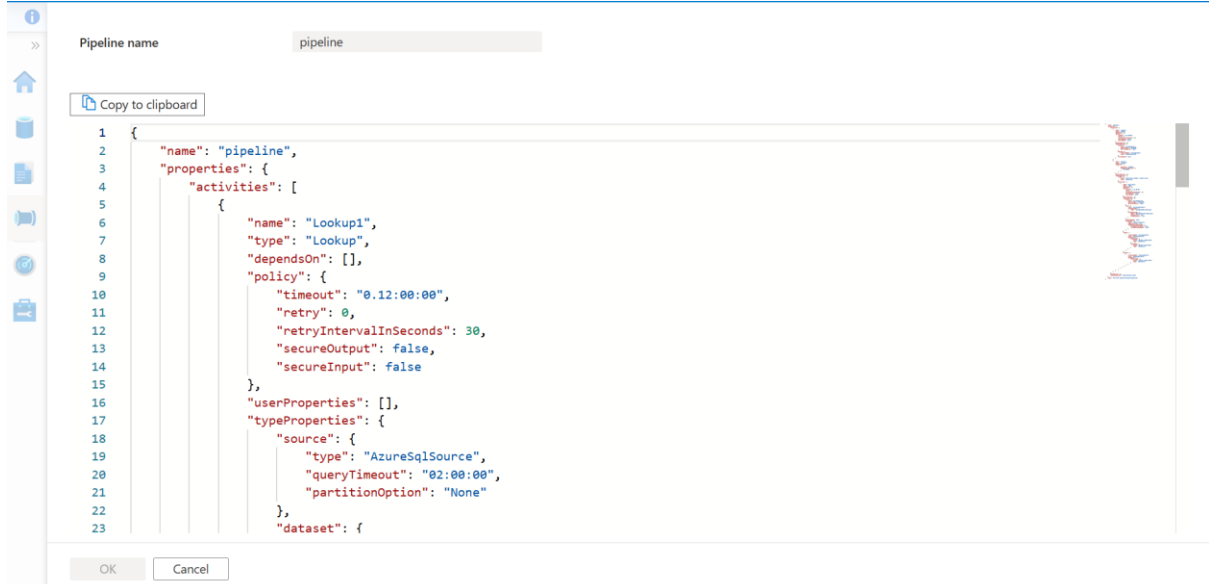SQL Dataset 1



SQL Dataset 2



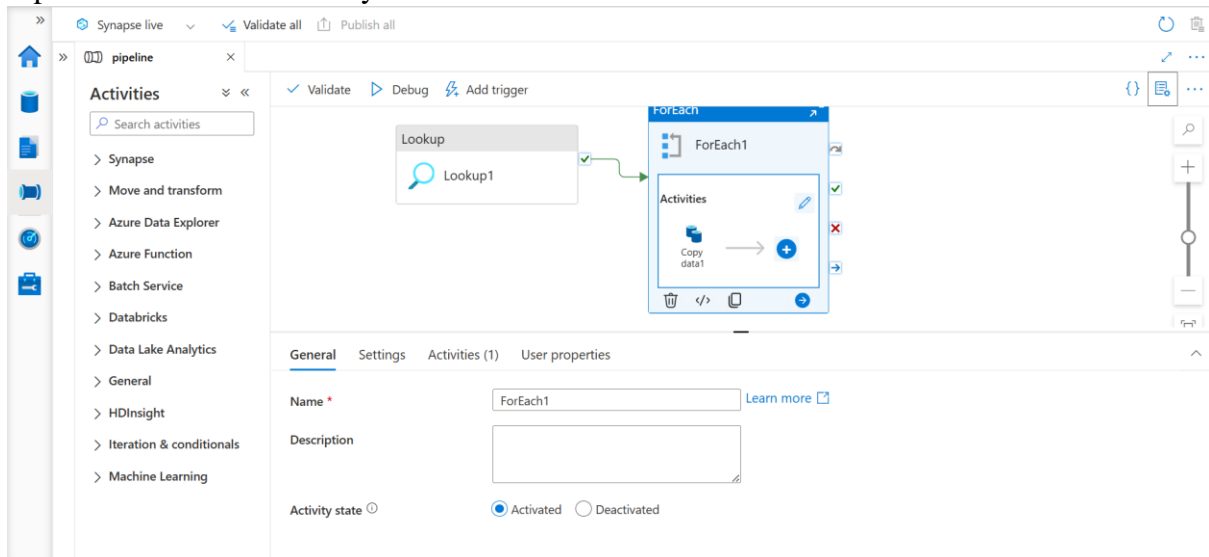Delimited Text 3 dataset for Azure Data Lake Gen2 Storage.



Now created a pipeline in Synapse and named it same as pipeline in ADF.

Copied JSON code from ADF to synapse



Pipeline created successfully.



Published pipeline and ran pipeline.
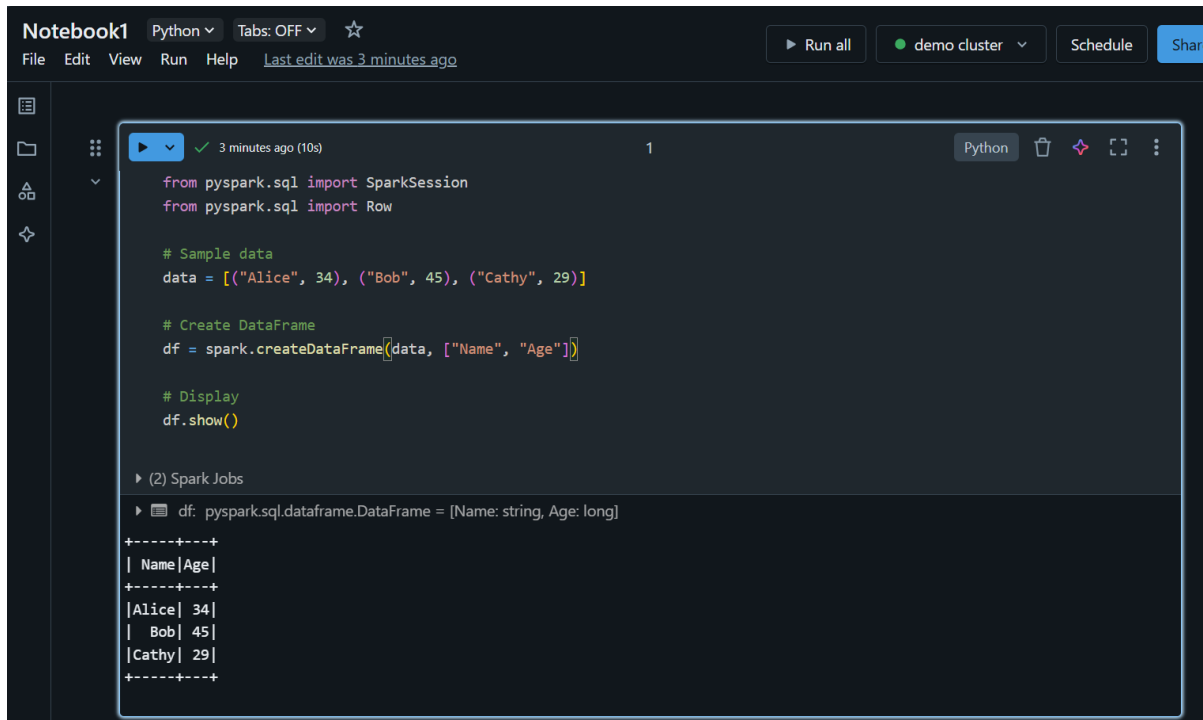
Pipeline ran successfully in Synapse.



Task 2:
Migrating Databricks notebook to Fabric notebook.
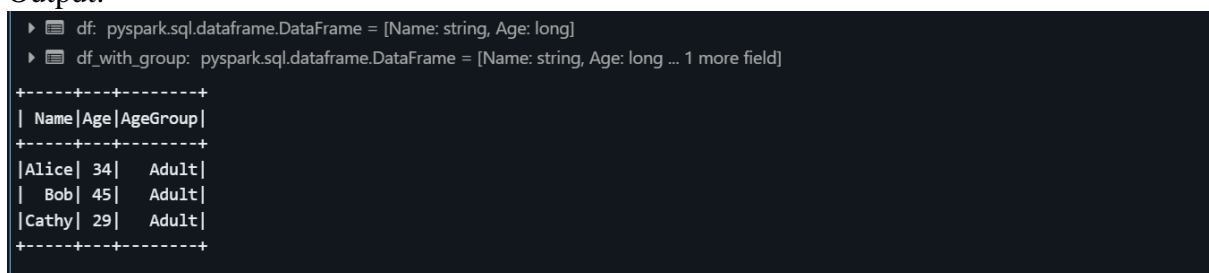Creating notebooks in Databricks
Notebook 1.



Importing this Notebook as ipynb file.

Creating Notebook 2



Output:



```
+-----+---+--------+
| Name|Age|AgeGroup|
+-----+---+--------+
|Alice| 34|   Adult|
|  Bob| 45|   Adult|
|Cathy| 29|   Adult|
+-----+---+--------+
```
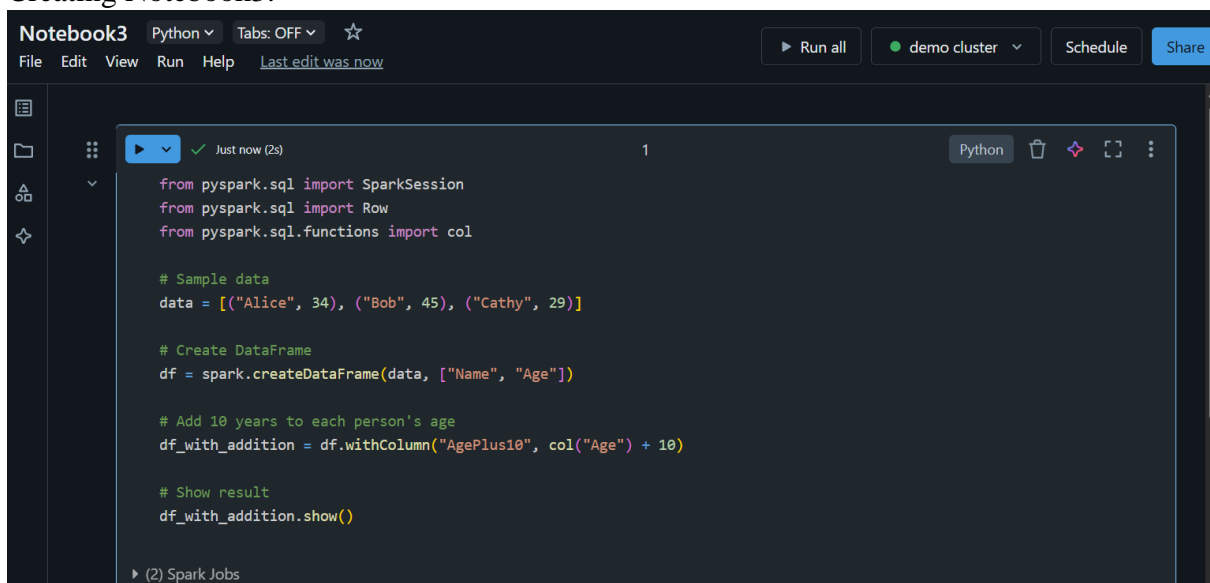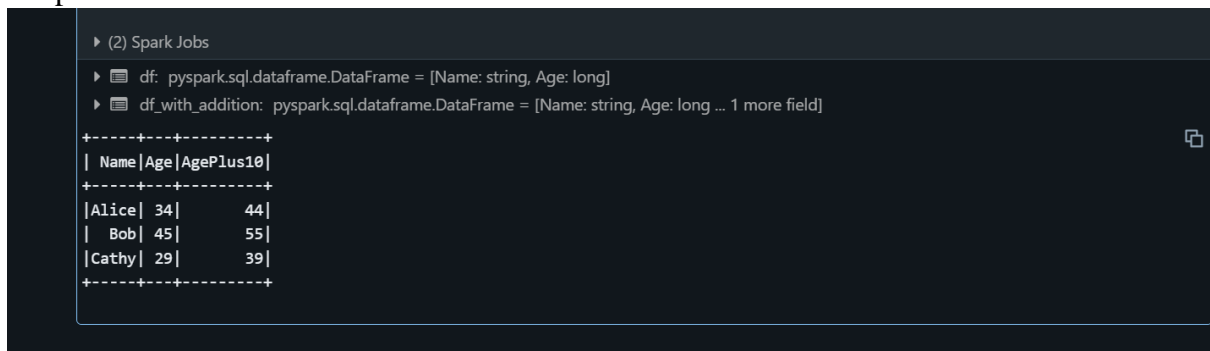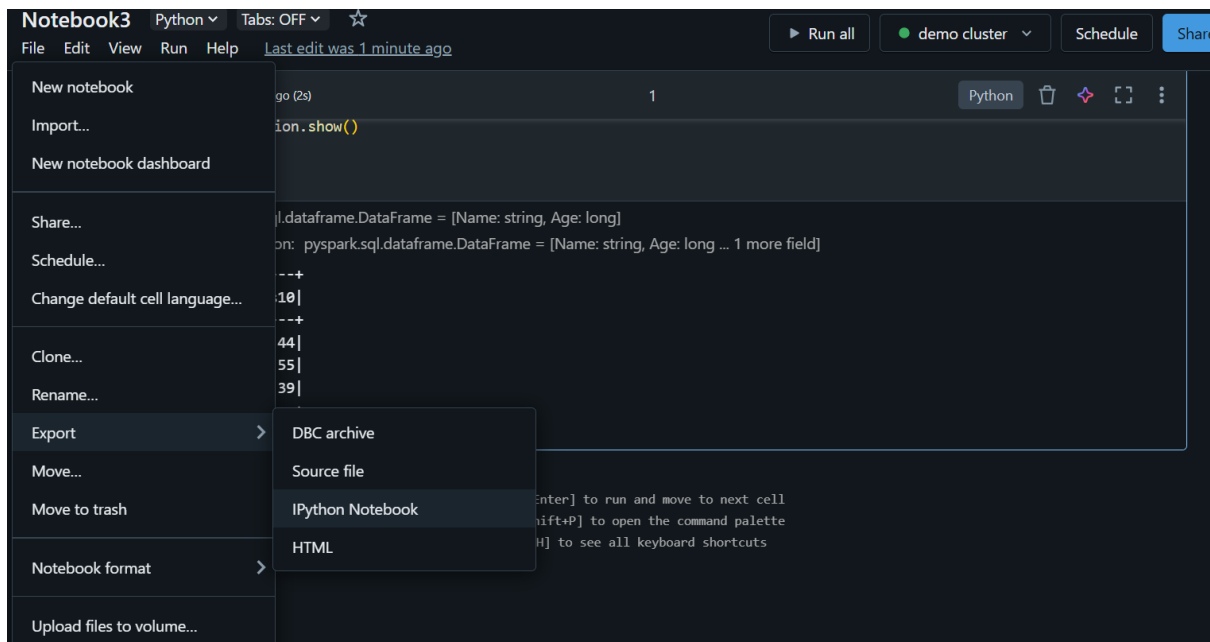
Exporting Notebook2.

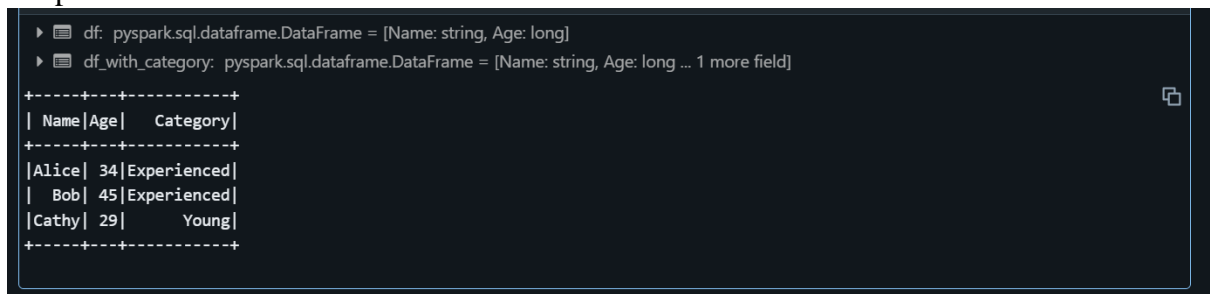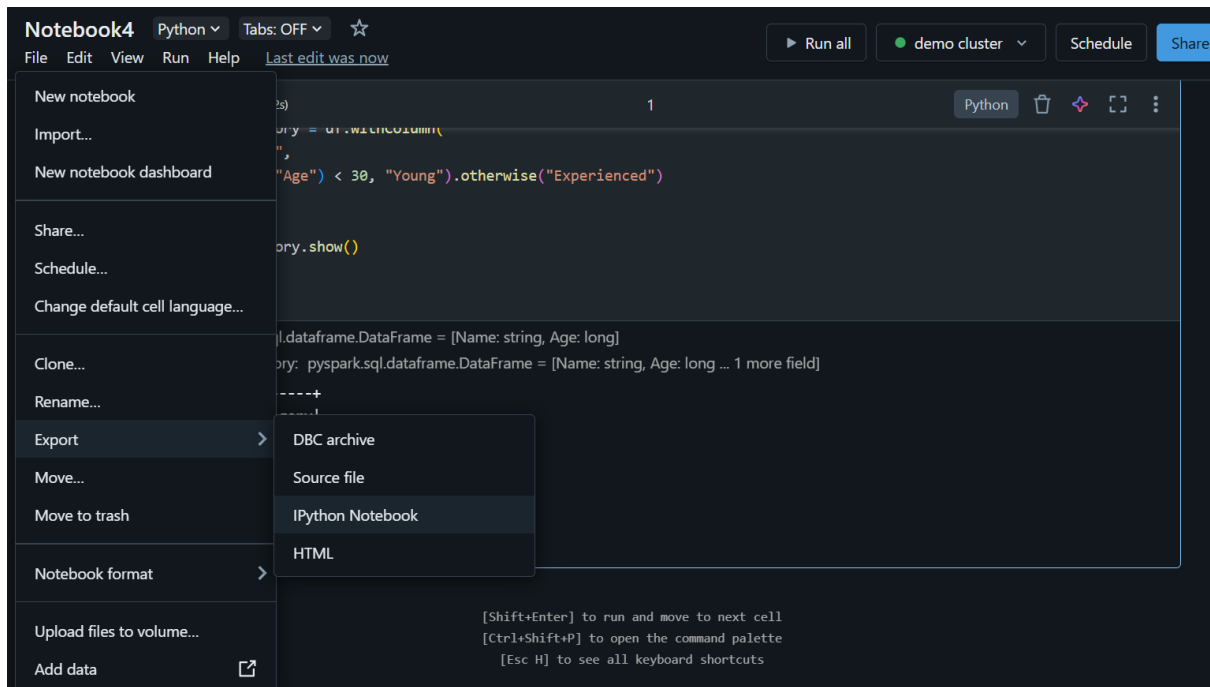Creating Notebook3.



Output:



Exporting this Notebook3.

**Notebook3** Python ∨ Tabs: OFF ∨ ☆

File Edit View Run Help Last edit was 1 minute ago

Run all demo cluster ∨ Schedule Share

New notebook

Import...

New notebook dashboard

Share...

Schedule...

Change default cell language...

Clone...

Rename...

Export ›

Move...

Move to trash

Notebook format ›

Upload files to volume...

DBC archive

Source file

IPython Notebook

HTML

ion.show()

1

sql.dataframe.DataFrame = [Name: string, Age: long]
on: pyspark.sql.dataframe.DataFrame = [Name: string, Age: long ... 1 more field]

--+
10|
--+
44|
55|
39|

Enter] to run and move to next cell
ift+P] to open the command palette
H] to see all keyboard shortcuts

Creating notebook 4



**Notebook4** Python ∨ Tabs: OFF ∨ ☆

File Edit View Run Help Last edit was now

Run all demo cluster ∨ Schedule Share

Just now (2s)

1

Python

```
from pyspark.sql import SparkSession
from pyspark.sql import Row
from pyspark.sql.functions import col
from pyspark.sql.functions import when

# Sample data
data = [("Alice", 34), ("Bob", 45), ("Cathy", 29)]

# Create DataFrame
df = spark.createDataFrame(data, ["Name", "Age"])

# Categorize age using when/otherwise
df_with_category = df.withColumn(
    "Category",
    when(col("Age") < 30, "Young").otherwise("Experienced")
)

df_with_category.show()
```

▶ (2) Spark Jobs

Output:



```
df: pyspark.sql.dataframe.DataFrame = [Name: string, Age: long]
df_with_category: pyspark.sql.dataframe.DataFrame = [Name: string, Age: long ... 1 more field]

+-----+---+-----------+
| Name|Age|   Category|
+-----+---+-----------+
|Alice| 34|Experienced|
|  Bob| 45|Experienced|
|Cathy| 29|      Young|
+-----+---+-----------+
```
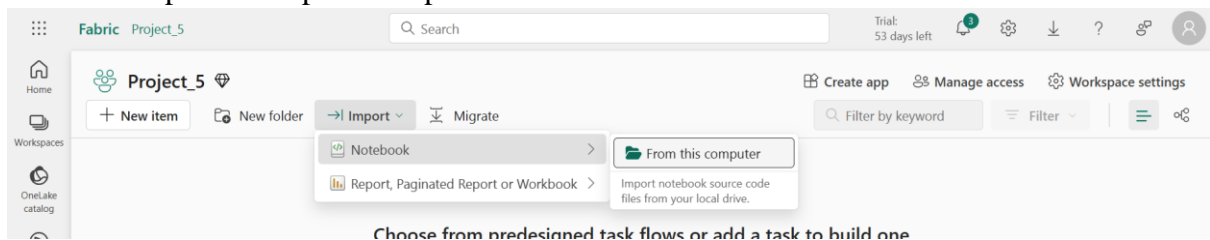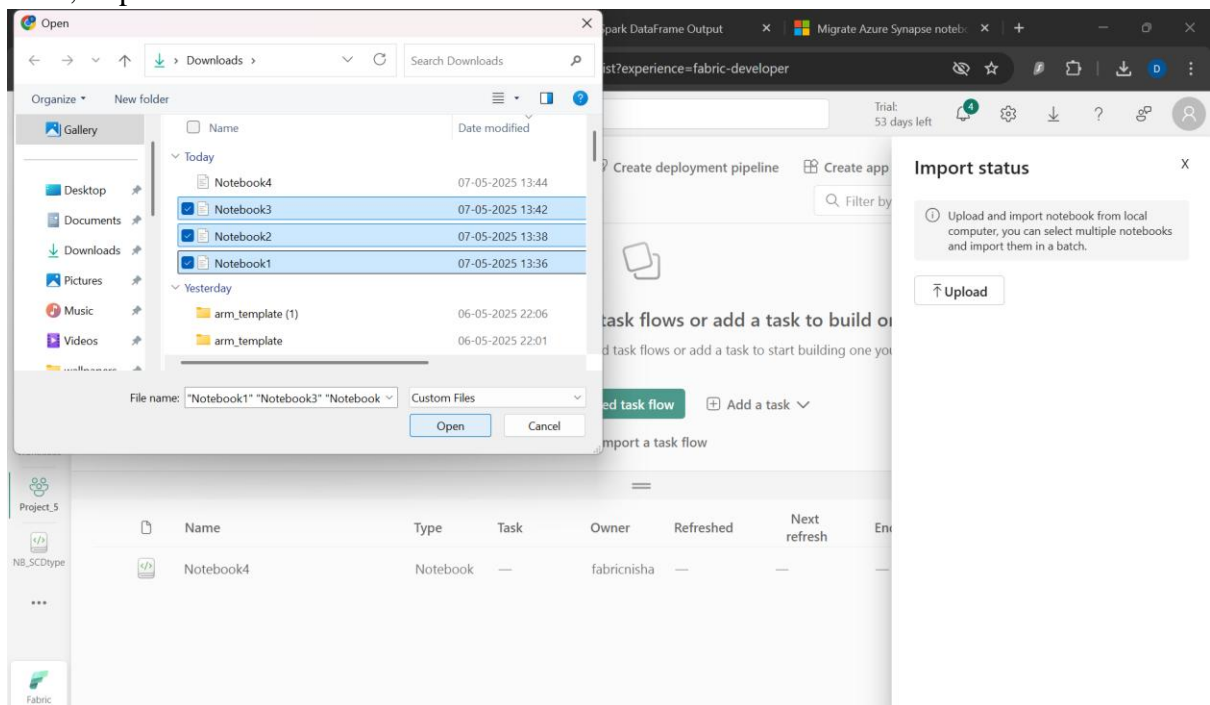
Exporting Notebook4

Now will import exported notebooks from Databricks to Fabric.
Go to workspace -> import-> import notebook.



Now, import all Notebooks at once



Now open Notebooks and run them.
Notebook 1 ran successfully.

Run Notebook2



```
[1]    ✓  10 sec - Command executed in 4 sec 596 ms by fabricnisha on 2:27:41 PM, 5/07/25        PySpark (Python) ⌄

    ⟩     Spark jobs (3 of 3 succeeded)   📊 Resources   📋 Log                                         ⋯

...  +-----+---+--------+
     | Name|Age|AgeGroup|
     +-----+---+--------+
     |Alice| 34|   Adult|
     |  Bob| 45|   Adult|
     |Cathy| 29|   Adult|
     +-----+---+--------+
```

Notebook 2 ran successfully.

## Run Notebook 3



Notebook3 ran successfully.

## Run Notebook4



Notebook 4 ran successfully.