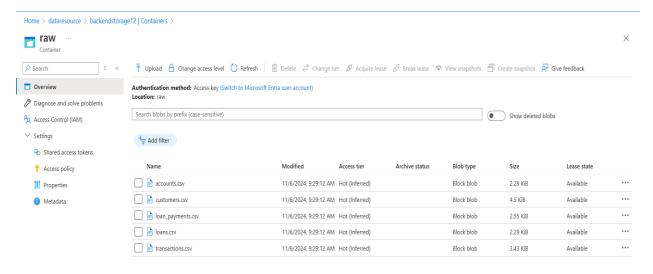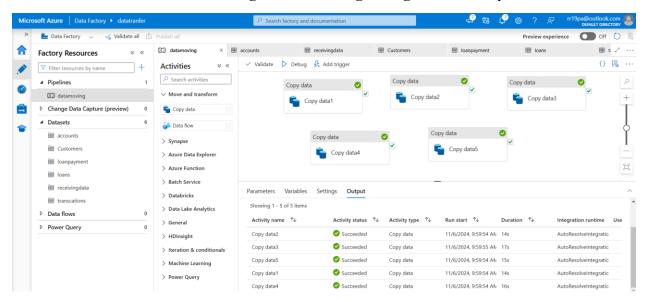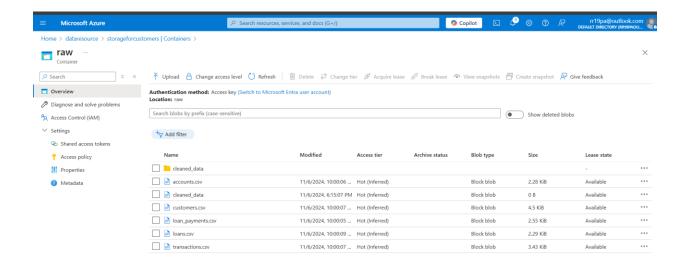# Data Pipeline for Customer

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



**We moved files from backendstorage  to new storage using Data Factory**

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

Here we read data and clean all the data and move them to cleaned_data



also used key vault to access id secret and you can see it in cleaning. Ipynb

## Step 3: Databricks Activity (ETL Processing)

We created another databricks notebook for ETL Processing here we read all the cleaned data we created in step 2 and wrote the query to calculate the total balance across all accounts for each customer, ensuring that all

columns from the accounts and customers tables are selected and included and stored it in gold container.





## Step 4: Azure Synapse Analytics

**Create external tables in Azure Synapse Analytics to map to the data stored in the Curated(silver) and Refined(gold) containers of your data lake. This allows data analysts and business intelligence teams to access and query the data directly using tools like Synapse Studio or notebooks.**

# Integrate

## Pipelines — Pipeline 1

### Activities

- Synapse
- Move and transform
- Azure Data Explorer
- Azure Function
- Batch Service
- Databricks
- Data Lake Analytics
- General
- HDInsight
- Iteration & conditionals
- Machine Learning

Validate · Debug · Add trigger

Copy data — Copy data1
Copy data — Copy data2
Copy data — Copy data3
Copy data — Copy data4
Copy data — Copy data5

Parameters | Variables | Settings | **Output**

Pipeline run ID: ea631a88-256a-4eb8-811b-465191259b2b

Pipeline status ✔ Succeeded — View debug run consumption

All status

Monitor in Azure Metrics — Export to CSV

Showing 1 - 5 of 5 items

| Activity name | Activity status | Activity type | Run start | Duration | Integration runtime |
|---|---|---|---|---|---|
| Copy data5 | ✔ Succeeded | Copy data | 11/11/2024, 7:54:12 PM | 13s | AutoResolveIntegration |
| Copy data4 | ✔ Succeeded | Copy data | 11/11/2024, 7:53:58 PM | 14s | AutoResolveIntegration |
| Copy data3 | ✔ Succeeded | Copy data | 11/11/2024, 7:53:42 PM | 15s | AutoResolveIntegration |
| Copy data2 | ✔ Succeeded | Copy data | 11/11/2024, 7:53:28 PM | 13s | AutoResolveIntegration |
| Copy data1 | ✔ Succeeded | Copy data | 11/11/2024, 7:53:13 PM | 15s | AutoResolveIntegration |

---

# Data

Workspace | Linked

- SQL database
  - accounts (SQL)
    - External tables
      - dbo.accounts
      - dbo.customer_total_balances
      - dbo.customers
      - dbo.loan_payment
      - dbo.loans
      - dbo.transactions
    - External resources
    - Views
    - Schemas
    - Security
  - acounts (SQL)

Run · Undo · Publish · Query plan — Connect to: Built-in — Use database: accounts

```sql
SELECT TOP (100) [customer_id]
    ,[first_name]
    ,[last_name]
    ,[address]
    ,[city]
    ,[state]
    ,[zip]
```

**Results** | Messages

View: Table | Chart — Export results

| customer_id | first_name | last_name | address | city | state | zip |
|---|---|---|---|---|---|---|
| 42 | Charlotte | Richardson | 4141 Beech Dr | Newmarket | ON | L3Y0A1 |
| 43 | Joseph | Cox | 4242 Cedar Ln | Aurora | ON | L4G0A1 |
| 44 | Amelia | Howard | 4343 Elm St | Bradford | ON | L3Z0A1 |
| 45 | Christopher | Ward | 4444 Maple Ave | Keswick | ON | L4P0A1 |
| 46 | Mia | Brooks | 4545 Oak Dr | Stouffville | ON | L4A0A1 |
| 47 | Andrew | Gray | 4646 Pine Rd | Uxbridge | ON | L9P0A1 |

00:00:00 Query executed successfully.

## Properties

**General** | Related (0)

Name *
SQL script 9

Description

Type
.sql script

Size
188 bytes

Results settings per query
● First 5000 rows (default)
○ All rows