

Customer Churn Azure Project

Note: The dataset provided for this assignment was already clean, with no missing values or Personally Identifiable Information (PII). However, for demonstration purposes, a filter condition was implemented to show how missing values can be handled using an Azure Data Factory (ADF) pipeline.

Project Objective

The goal of this project is to migrate data from an API source to an Azure SQL Database and establish an automated orchestration workflow using Azure Data Factory.

Technology Stack

- Resource Group
- Storage Accounts
- Azure Data Factory (ADF)
- Azure Databricks
- Azure SQL Database
- Azure Key Vault

ADF Activities Used

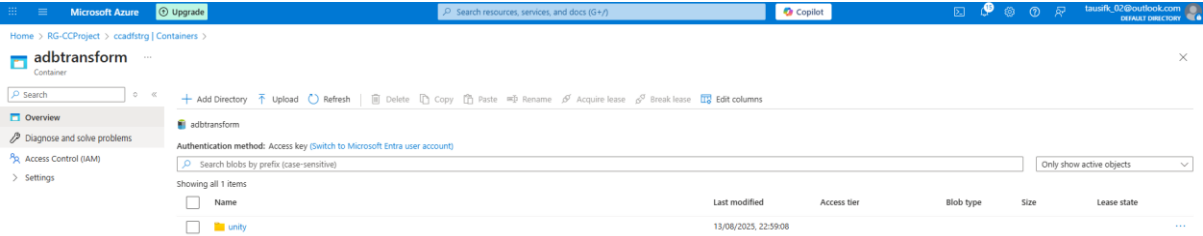
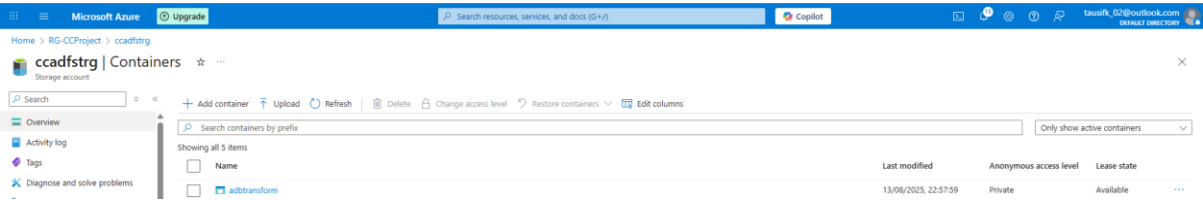
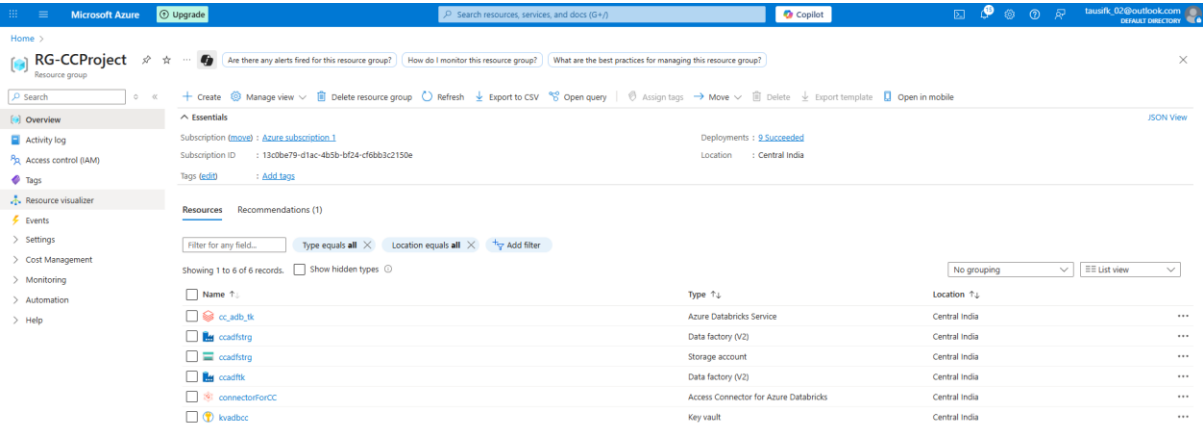
- Lookup Activity
- Copy Data Activity
- Azure Databricks Notebook
- Data Flow

Project Workflow

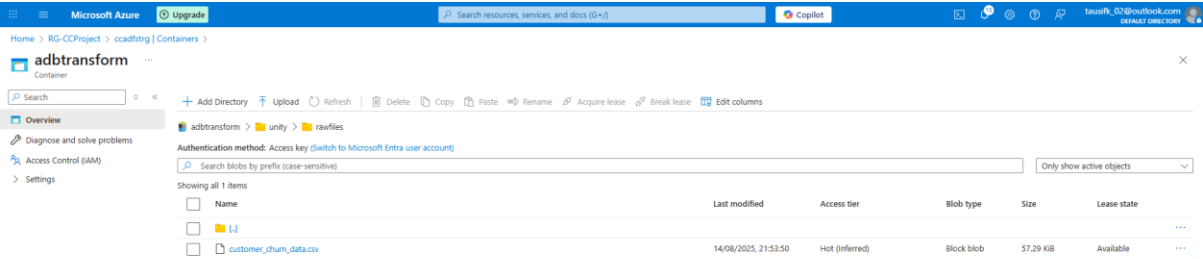
Azure Resources Setup

1. Created a Resource Group in Azure.
2. Added a Storage Account and Azure Data Factory instance under this resource group.
3. Within the Storage Account, created one container '**adbtransform**' and '**unity**' folder under the container.

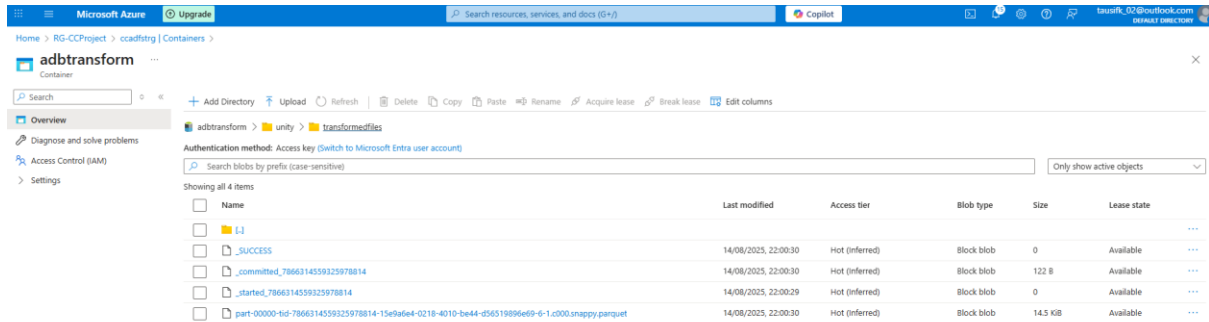
The following screenshot shows all the resources created under the Resource Group.



CSV file is loaded using lookup activity dynamically using ADF pipeline



Parquet file is created using Azure Databricks notebook.

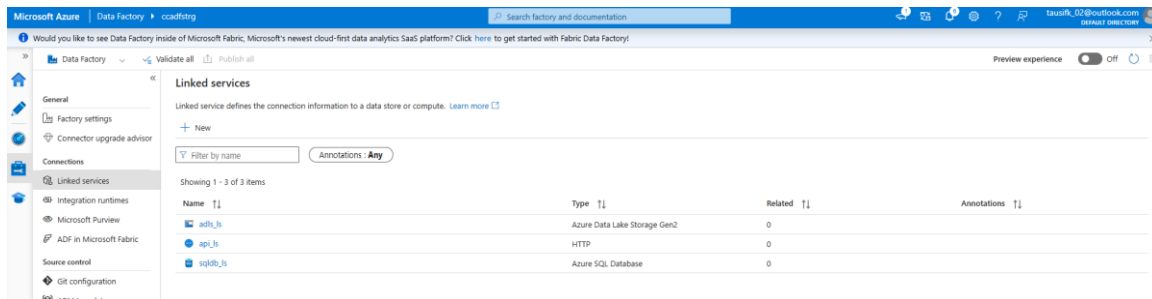


Name	Last modified	Access tier	Blob type	Size	Lease state
1-1					...
_SUCCESS	14/08/2025, 22:00:30	Hot (Inferred)	Block blob	0	Available
_committed_7866314559325978814	14/08/2025, 22:00:30	Hot (Inferred)	Block blob	122 B	Available
_started_7866314559325978814	14/08/2025, 22:00:29	Hot (Inferred)	Block blob	0	Available
part-00000-tid-7866314559325978814-15e9a6e4-0218-4010-be44-d56519896e69-6-1.c000.snappy.parquet	14/08/2025, 22:00:30	Hot (Inferred)	Block blob	14.5 KiB	Available

Linked Services Configuration

Configured Linked Services in ADF for connections to:

- Azure Storage Account
- Github Account
- Azure SQL Database



Name	Type	Related	Annotations
adls_ls	Azure Data Lake Storage Gen2	0	
api_ls	HTTP	0	
sqldb_ls	Azure SQL Database	0	

Dataset Creation

- The config file in the Storage Account – config_ds
- Loading file from GitHub – api_ds
- Create CSV file in ADLS Gen 2 – adls_adb_csv_ds
- Creating Parquet file using Notebook – adb_parquet_ds
- Create table in SQL Database – adb_Sql_ds

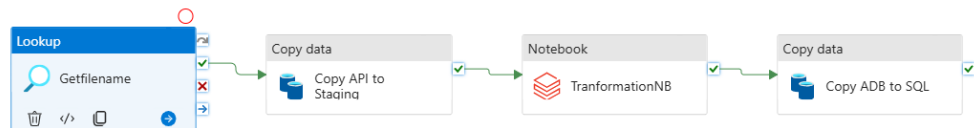
▲ Datasets 8

- adb_parquet_ds
- adb_sql_ds
- adls_adb_csv_ds
- adls_ds
- api_ds
- ccparquet_ds
- config_ds
- sqldb_ds

ETL Pipeline Development

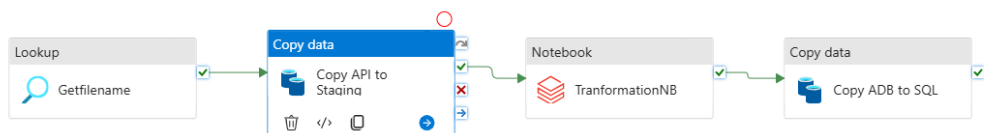
Pipeline Workflow:

1. The config file containing the file name is passed through a Lookup Activity to the Copy Data Activity.
2. In **Azure Databricks notebook**, data is filtered/transformed by applying filter condition as “InternetService” is not equal to "None".
3. Saved the filtered data as a Parquet file in the transformedfiles folder.
4. Used another **Copy Data Activity** to load the Parquet file into the SQL Database.



General Settings User properties

Source dataset *	config_ds	Open	New	Preview data	Learn more
First row only	<input type="checkbox"/>				
File path type	<input checked="" type="radio"/> File path in dataset <input type="radio"/> Wildcard file path <input type="radio"/> List of files				
Filter by last modified	Start time (UTC)	End time (UTC)			
Recursively	<input checked="" type="checkbox"/>				
Enable partitions discovery	<input type="checkbox"/>				
Max concurrent connections					
Skip line count					



General **Source** Sink Mapping Settings User properties

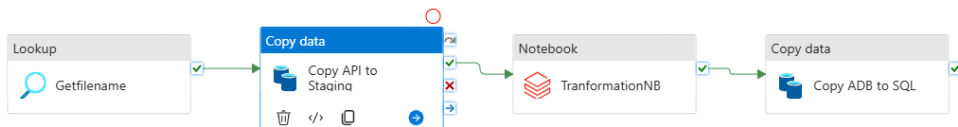
Source dataset * [Open](#) [+ New](#) [Preview data](#) [Learn more](#)

Dataset properties ⓘ

Name	Value
load	@activity('Getfilename').output.value...

Request method ⓘ

Additional headers ⓘ



General Source **Sink** Mapping Settings User properties

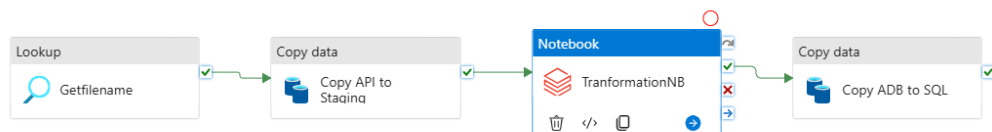
Sink dataset * [Open](#) [+ New](#) [Learn more](#)

Dataset properties ⓘ

Name	Value
filename	@activity('Getfilename').output.value...

Copy behavior ⓘ

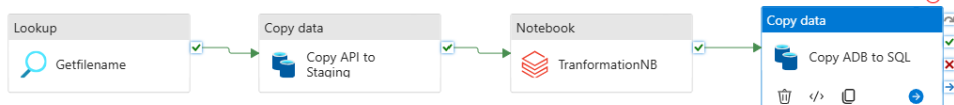
Max concurrent connections ⓘ



[Expand toolbox pane](#)

General **Azure Databricks** Settings User properties

Databricks linked service * [Test connection](#) [Edit](#) [+ New](#)



General **Source** Sink Mapping Settings User properties

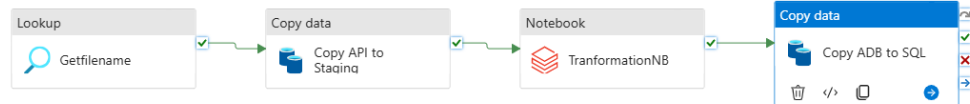
Source dataset * adb_parquet_ds [Open](#) [New](#) [Preview data](#) [Learn more](#)

File path type ☐ File path in dataset ☒ Wildcard file path ☐ List of files

Wildcard paths adbtransform / Wildcard folder path / *.parquet

Filter by last modified ☐ Start time (UTC) End time (UTC)

Recursively ☒



General Source **Sink** Mapping Settings User properties

Sink dataset * adb_sql_ds [Open](#) [New](#) [Learn more](#)

Dataset properties

Name	Value
filename	customers

Write behavior ☒ Insert ☐ Upsert ☐ Stored procedure

Microsoft Azure **Upgrade** Search resources, services, and docs (G+7) Copilot tmsuRt_02@outlook.com

Home > RG-CCProject > kvadbcc | Secrets >

adb Versions

Search [New Version](#) [Refresh](#) [Delete](#) [Download Backup](#)

Overview Access control (IAM)

Version	Status	Activation date	Expiration date
CURRENT VERSION			
395572556dcf4a67b72eb14865adcbf3	✓ Enabled		

Microsoft Azure **Upgrade** Search resources, services, and docs (G+7) Copilot tmsuRt_02@outlook.com

Home > RG-CCProject >

cc_adb_tk Azure Databricks Service

Search [Delete](#)

Overview [Essentials](#) [JSON View](#)

Status	: Active	Managed Resource Group	: databricks-rg-cc_adb_tk-vmw5em4dsk
Resource group	: RG-CCProject	URL	: https://adb-202354061623116.azuredatabricks.net
Location	: Central India	Pricing Tier	: Premium (+ Role-based access control) (Click to change)
Subscription	: Azure subscription 1	Enable No Public IP	: Yes
Subscription ID	: 13cd8e79-d1ac-4b5b-bf24-cf6b63c2150e		
Tags	(edit) Add tags		

[Launch Workspace](#)

08:16 PM (5s) 5 Python

```
1 df_filter.write.mode("overwrite").format("parquet").save("abfss://adbtransform@ccadfstrg.dfs.core.windows.net/transformedfiles")
```

See performance (1) Optimize

03:09 PM (1s) 6

```
1 df_filter.createOrReplaceTempView("customers")
2
3
```

See performance (1) Optimize

03:10 PM (5s) 7

```
1 %sql
2
3 CREATE TABLE customerchurn.default.customers
4 AS
5 select * from customers
```

See performance (1) Optimize

Table +

num_affected_rows	num_inserted_rows
-------------------	-------------------

03:11 PM (7s) 8 SQL

```
1 %sql
2 select * from customerchurn.default.customers
```

See performance (1) Optimize

_sqlidf: pyspark.sql.connect.dataframe.DataFrame = [CustomerID: integer, Age: integer ... 8 more fields]

Table +

	CustomerID	Age	Gender	Tenure	MonthlyCharges	ContractType	InternetService	TotalCharges	TechSupport
1	1	49	Male	4	88.35	Month-to-Month	Fiber Optic	353.4	Yes
2	2	43	Male	0	36.67	Month-to-Month	Fiber Optic	0	Yes
3	3	51	Female	2	63.79	Month-to-Month	Fiber Optic	127.58	No
4	4	60	Female	8	102.34	One-Year	DSL	818.72	Yes
5	6	42	Female	16	119.75	Two-Year	DSL	1916	Yes
6	9	40	Female	53	49.81	Two-Year	Fiber Optic	2639.9300000000003	Yes
7	10	50	Female	10	61.55	Month-to-Month	Fiber Optic	615.5	Yes
8	11	40	Female	1	63.53	Month-to-Month	Fiber Optic	63.53	Yes
9	12	40	Female	19	36.18	One-Year	Fiber Optic	687.42	Yes
10	15	27	Male	14	95.05	Month-to-Month	Fiber Optic	1330.7	Yes
11	16	39	Female	41	89.11	Two-Year	Fiber Optic	3653.5099999999998	Yes
12	19	35	Male	98	49.59	Two-Year	Fiber Optic	4859.8200000000001	Yes
13	20	30	Female	2	89.55	Month-to-Month	Fiber Optic	179.1	No
14	21	59	Female	6	73.56	One-Year	DSL	441.36	No

703 rows | 7.34s runtime Refreshed 7 hours ago

This result is stored as _sqlidf and can be used in other Python and SQL cells.

Microsoft Azure Upgrade

Home > tk-azure-sql-server > CCpresentationlayer (tk-azure-sql-server/CCpresentationlayer)

CCpresentationlayer (tk-azure-sql-server/CCpresentationlayer) | Query editor (preview)

SQL database

Search Login + New Query Open query Feedback Getting started

Overview Activity log Tags

Diagnose and solve problems

Query editor (preview)

Mirror database in Fabric (preview)

Resource visualizer

Settings

Data management

Integrations

Power Platform

Security

Intelligent performance

Monitoring

Automation

Help

CCpresentationlayer (tausir_02@outlo...

Showing limited object explorer here. For full capability please click here to open Azure Data Studio.

Tables

- dbo.Customer_churn
- dbo.customers

Views

Stored Procedures

Query 1 X Query 2 X

Run Cancel query Save query Export data as Show only Editor

```
1 SELECT TOP (1000) * FROM [dbo].[customers]
```

Results Messages

Search to filter items...

CustomerID	Age	Gender	Tenure	MonthlyCharges	ContractType	InternetServ
1	49	Male	4	88.35	Month-to-Month	Fiber Optic
2	43	Male	0	36.67	Month-to-Month	Fiber Optic
3	51	Female	2	63.79	Month-to-Month	Fiber Optic
4	60	Female	8	102.34	One-Year	DSL
6	42	Female	16	119.75	Two-Year	DSL
9	40	Female	53	49.81	Two-Year	Fiber Optic
10	50	Female	10	61.55	Month-to-Month	Fiber Optic

Orchestration and Automation

- Set a trigger so the pipeline runs automatically at scheduled times. *As per assignment it is to be scheduled for every hour.*
- Screenshots of the setup are included for reference.

Microsoft Azure | Data Factory | ccaadfrg | Search

Dashboard | Pipeline runs | Trigger runs | Runtimes & sessions | Integration runtimes | Data flow debug | Notifications | Alerts & metrics

Pipeline runs

Triggered | Debug | Run | Cancel options | Refresh | Edit columns | List | Gantt

Filter by run ID or name | Chennai, Kolkata, Mu... | Last 24 hours | Pipeline name: All | Status: All | Runs: Latest runs | Triggered by: All | Add filter | Copy filters | Export to CSV

Showing 1 - 10 items | Last refreshed 0 minutes ago

Pipeline name	Run start	Run end	Duration	Triggered by	Status	Run	Parameters	Annotations	Run ID
copy from api to adfs_adb	8/14/2025, 9:53:17 PM	8/14/2025, 10:01:18 PM	8m 1s	Manual trigger	Succeeded	Original			eeae1cd48-8321-4b2b-b134-dca50eb0b
copy from api to gen2	8/14/2025, 12:32:26 AM	8/14/2025, 12:38:49 AM	6m 23s	Manual trigger	Succeeded	Original			bf2a3efe-b253-43f7-9aba-47b579e945
copy from api to gen2	8/14/2025, 12:19:05 AM	8/14/2025, 12:23:59 AM	4m 54s	Manual trigger	Canceled	Original			6dfb0e9a-9c29-4c1b-b4c9-b62e388b79
copy from api to gen2	8/14/2025, 12:09:43 AM	8/14/2025, 12:15:02 AM	5m 20s	Manual trigger	Canceled	Original			9a56b572-31d4-4bcb-b927-fd885a2e0e
copy from api to gen2	8/14/2025, 12:07:32 AM	8/14/2025, 12:07:33 AM	2s	Manual trigger	Failed	Original			5c814873-9719-4c7b-aa05-260878cd1c
copy from api to gen2	8/13/2025, 11:40:15 PM	8/13/2025, 11:40:17 PM	2s	Manual trigger	Failed	Original			4dc7dbdc-69cb-40aa-b75f-63f457e356e
copy from api to gen2	8/13/2025, 11:21:20 PM	8/13/2025, 11:26:46 PM	5m 26s	Manual trigger	Succeeded	Original			66852149-12f1-4a3c-987a-ae8738668e
copy from api to gen2	8/13/2025, 11:12:58 PM	8/13/2025, 11:19:10 PM	6m 12s	Manual trigger	Succeeded	Original			a19067cd-a9e1-4b6c-9a7c-ba3e5988ea
copy from api to gen2	8/13/2025, 11:10:21 PM	8/13/2025, 11:10:50 PM	29s	Manual trigger	Failed	Original			4229dae2-7131-47af-aa21-063ec68a47
copy from api to gen2	8/13/2025, 11:06:11 PM	8/13/2025, 11:06:41 PM	30s	Manual trigger	Failed	Original			0ed747af-5d57-49f1-b0f1-bbbed3f667c