NAME-MALLU SRAVANA ANDHYA

MAIL.ID-mallusravanasandhya9948@gmail.com

# DATA PROCESSING

## Project Statement

Implement a data processing pipeline where Azure Data Factory orchestrates data workflows, and Azure Databricks is used as a processing engine for on-demand analytics and transformations.

## Project Overview

This project can implement a data processing pipeline using Azure Data Factory (ADF) and Azure Databricks. This setup allows you to orchestrate, schedule, and automate data workflows using ADF while leveraging the power of Databricks for on-demand analytics and transformations.

## Project Requirements

**Azure Databricks:**

Create a Databricks workspace from the Azure portal. This workspace will act as the processing engine for your data pipeline.

**Azure Data Factory:**

Create a Databricks Linked Service in ADF. This linked service will enable ADF to connect to your Databricks workspace and submit jobs for execution.

**Data processing:**

Define the necessary data inputs and outputs for your data processing pipeline. This could involve creating datasets in ADF.

## Execution Overview

- Azure Databricks is responsible for creating a notebook and loading, processing and analysis the data
- Azure Data Factory is used to linked service in ADF.This linked service will enable ADF to connect to Databricks Workspace and submit jobs for execution in Pipelines

## Tools/Technology used in project

- ✓ Azure Databricks
- ✓ Azure Data Factory
- ✓ Linked services
- ✓ Pipelines
- ✓ Data Processing
- ✓ Python

## Prerequisites

**Azure Subscription:** Ensure you have an active Azure Subscription to provision the required services.

**Access to Azure Portal:** You'll need to Azure portal to create and manage resources.

## Source Data Files:

We are referring to the source file which is used to run the code in the notebook.

| Filename | Type |
| --- | --- |
| Diamond | csv |

## Steps to implement:

To implement this data processing pipeline, you can follow these steps:

Set up Azure Data Factory (ADF) to create a new pipeline to define the data workflow.

Create an Azure Databricks workspace.

Install and configure the Databricks connectors for both Azure Blob Storage and Azure Data Lake Storage Gen2.

Upload the Databricks Notebook to the Azure Databricks workspace.

Create a Databricks linked service in ADF.

Configure the Databricks Notebook activity in ADF.

Connect the Copy activity and the Databricks Notebook activity in the ADF pipeline.

Trigger the pipeline execution to initiate the data workflow.

Monitor the progress and view the execution details.

**IMPLEMENTING PRACTICALLY IN AZURE PORTAL:**

Creating Azure Databricks Workspace:



**Databricks Workspace is successfully created:**

Create A Data Factory in azure portal:



Lanching Azure databricks workspace:

To create a cluster go to new and select cluster:



Give a name to the cluster and select the policy as personal compute, select theruntime version from the dropdown then click on create compute

Cluster is successfully created:



Create a notebook as shown in the figure:

Launch studio in azure data factory:



Click on Create a Linked service:
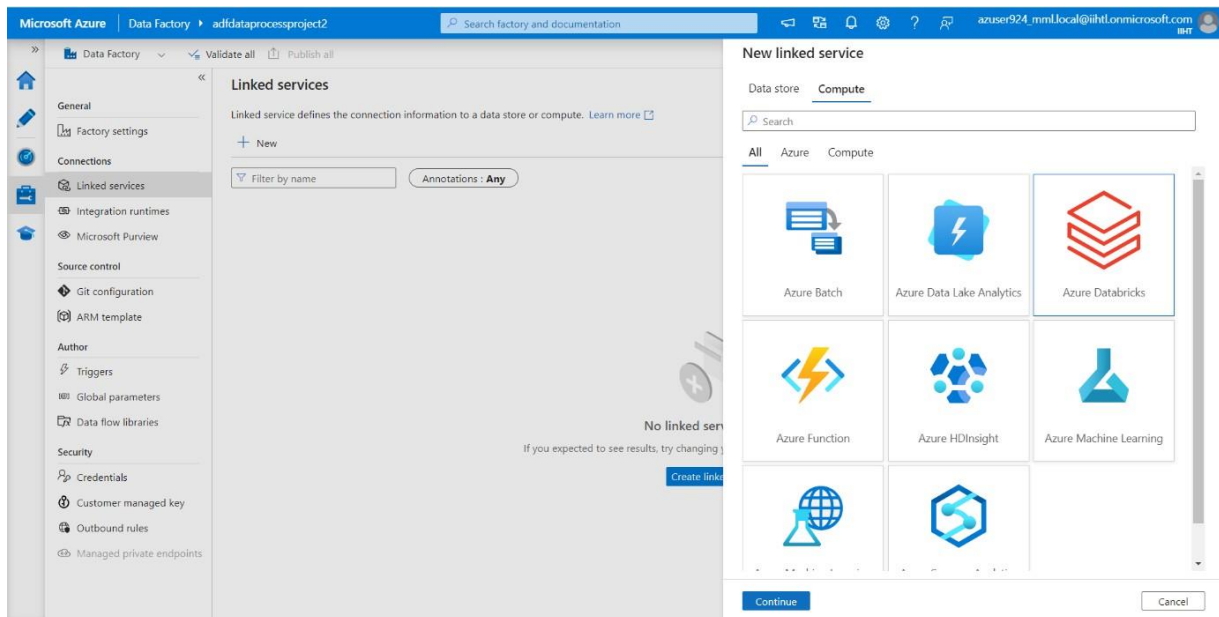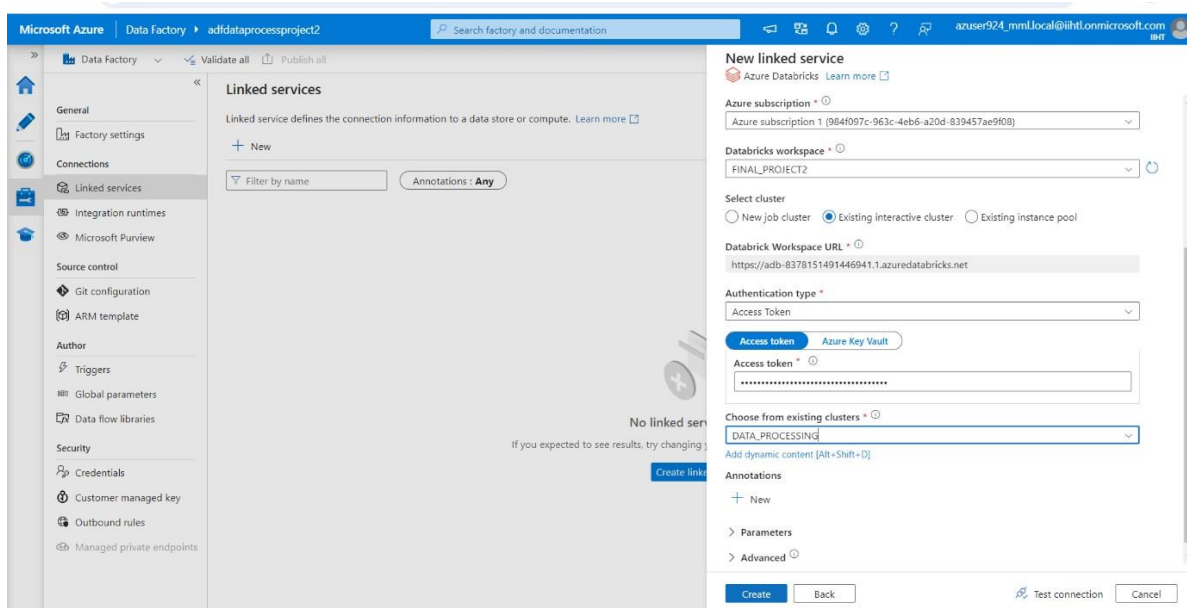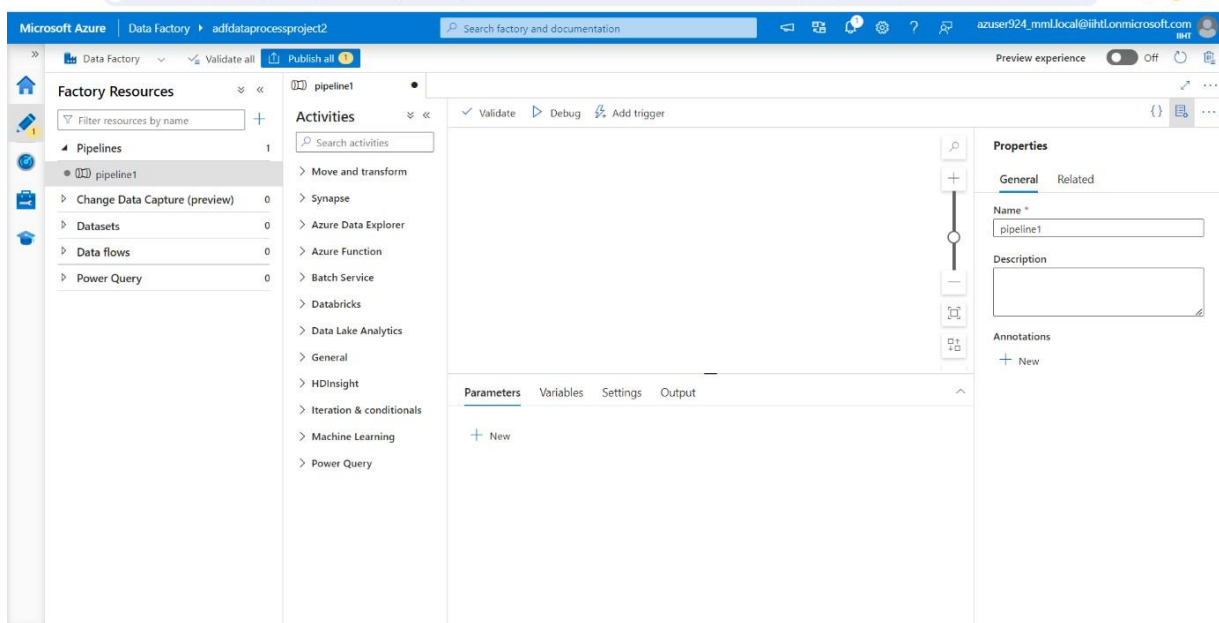
Click on Azure Databricks and click continue:
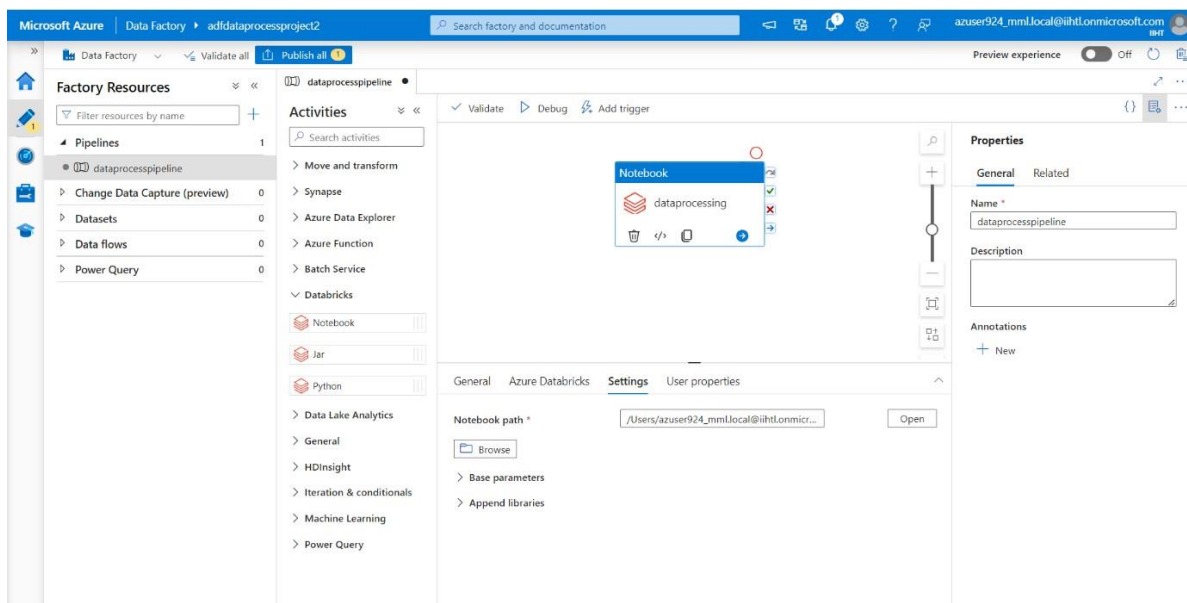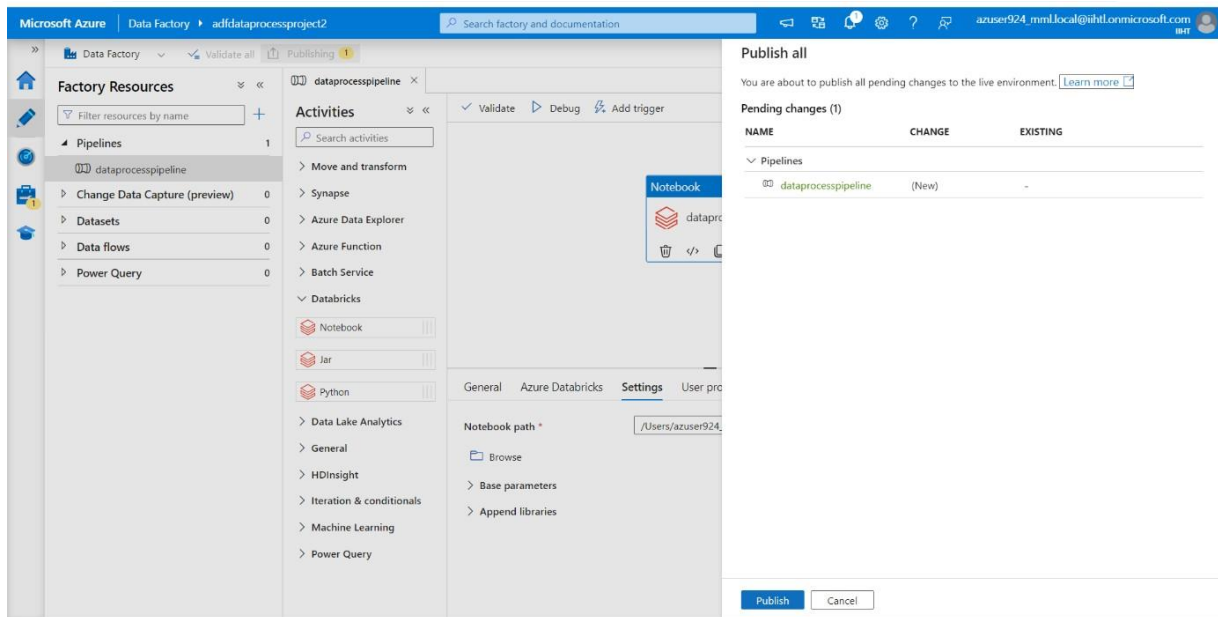


Create a new linked service by clicking create

Create a new pipline in Azure data factory:
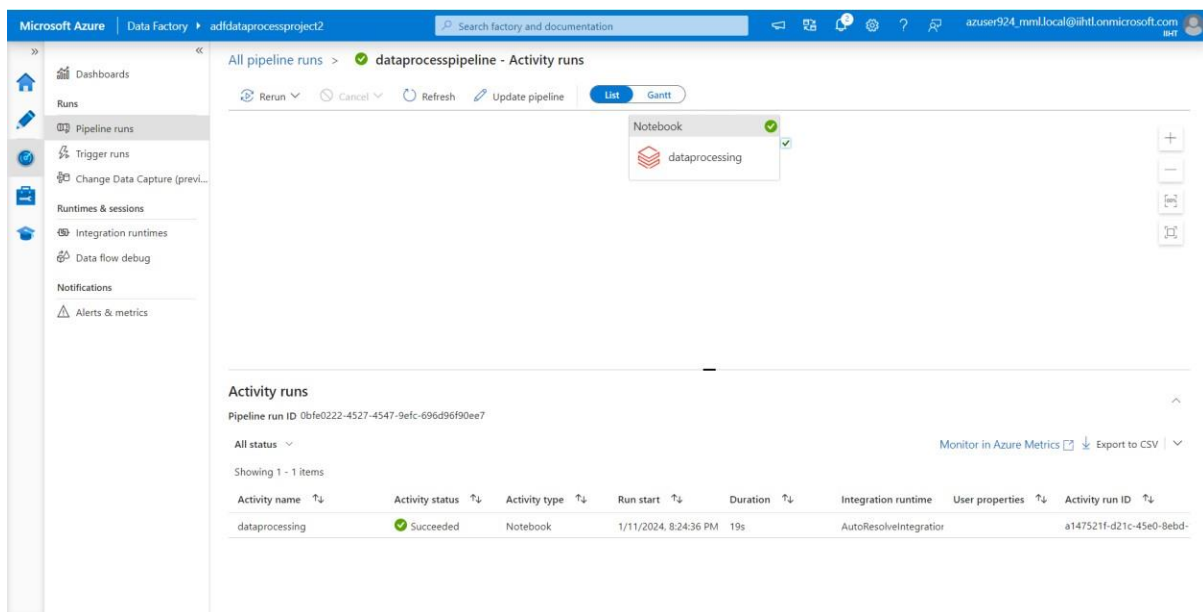


In pipeline trigger the notebook created in the azuredatabricks:

Publish pipeline and trigger the notebook



Check the notebook is successfully run in the azure data factory



By following these steps, you will have implemented a data processing pipeline where Azure Data Factory orchestrates data workflows, and Azure Databricks is used as a processing engine for on-demand analytics and transformations.