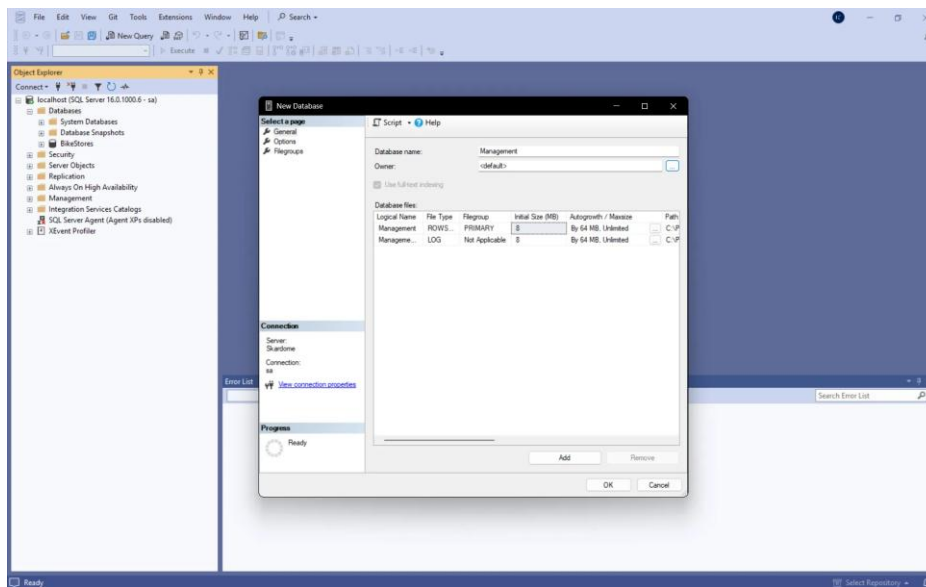# CDC with Databricks

## Project Documentation

<u>Concepts to cover</u>: Focusing on implementing Change Data Capture (CDC) using Databricks, a unified data analytics platform. The primary objective is to efficiently identify and process data changes (inserts, updates, deletes) in source systems and reflect them in the target data lake or data warehouse in near real-time. Leveraging Apache Spark, Delta Lake, and streaming capabilities within Databricks, the project aims to build a scalable and reliable CDC pipeline that supports incremental data ingestion, ensures data consistency, and minimizes processing overhead.

### Step 1

Initially, the project involves setting up a reliable database environment by installing and configuring MS SQL Server in a local environment, including necessary network protocols and user permissions. Automated SQL scripts are developed to create databases and tables, specifically designed to optimize CDC operations. Additionally, scripts are prepared to populate tables with initial datasets, with scheduled jobs ensuring continuous data insertion and updates to simulate a dynamic source environment.

**Setting up MS SQL Server**



**Adding Queries for Table Creation:**
https://github.com/Skaelum/CDC_Databricks_Celebal/blob/main/SQL%20Queries/Table%20Creation%20.sql
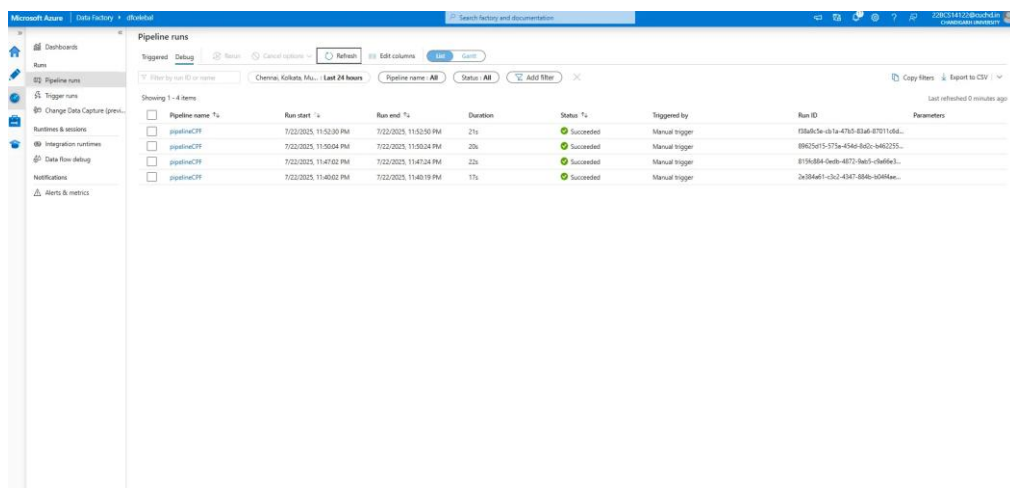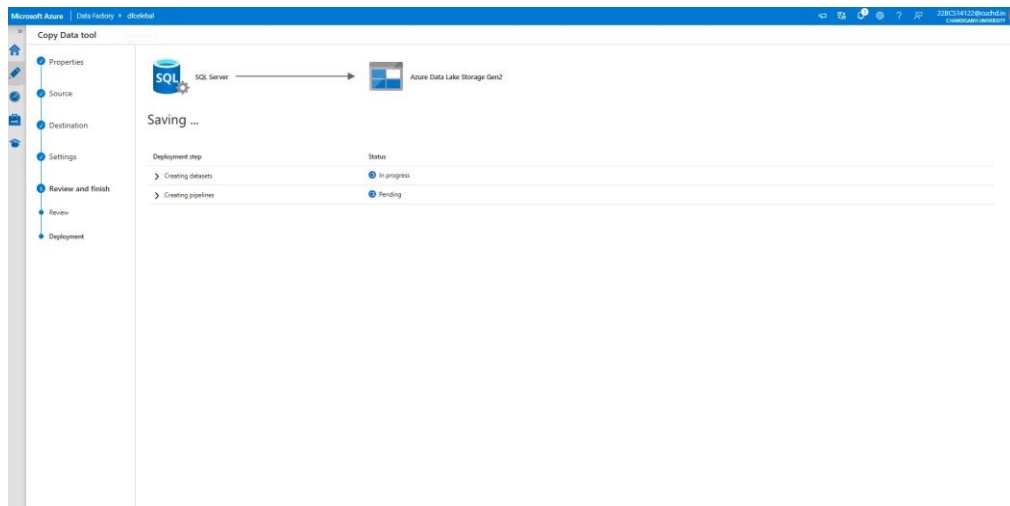
**Automating Table Population:**
https://github.com/Skaelum/CDC_Databricks_Celebal/blob/main/SQL%20Queries/Automate%20Table%20Population.sql
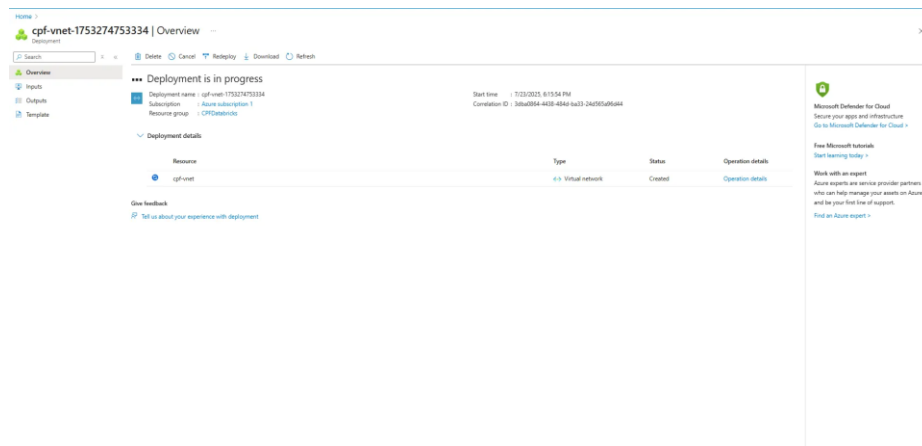
## Step 2

Azure Data Factory (ADF) plays a central role in data migration and ingestion, starting with its provisioning and configuration. A secure Self-Hosted Integration Runtime (SHIR) is set up to facilitate secure data transfers from the local SQL server environment to the Azure cloud. Rigorous connectivity tests are conducted to validate data movement and integration between the local infrastructure and Azure Data Factory.

Setting up a self-hosted runtime integrated (SHIR) connection for migrating the database from local to cloud server.

**Configured a Subnet for Databricks**



## Step 3

The ingestion process involves creating pipelines that transfer data from the local SQL server to Azure Storage Account (ADLS Gen2). Databricks is then utilized to configure Change Data Capture, leveraging Delta Lake tables to effectively track and identify row-level changes such as INSERT, UPDATE, and DELETE operations. Using Spark Structured Streaming, the incremental ingestion of data changes is streamlined, ensuring efficiency and near real-time processing. Additionally, Azure Logic Apps are implemented to create notification pipelines, which send operational alerts via email after each successful execution.

1. Creating a pipeline to ingest data into the storage account from the localhost.
2. Setting up change data capture in databricks (CDC) to detect row-level changes that reflect INSERT, UPDATE, DELETE queries in the source table.
3. Creating an email notification pipeline to notify operations after each run.

## Step 4

An incremental ETL pipeline is established to run at regular hourly intervals, capturing incremental data changes automatically. Robust error handling, logging, and retry mechanisms are implemented to enhance pipeline resilience. Furthermore, Databricks workspace is seamlessly integrated within Azure Resource Groups, ensuring smooth authentication and authorization processes. Apache Spark's powerful computing capabilities are harnessed through Azure Logic Apps to manage complex CDC processes, while optimization efforts are focused on resource allocation for performance efficiency and cost management.

1. Create an ETL pipeline to run every hour to capture the changes
2. Integrated the databricks workspace with azure resource group
3. Using Azure Logic App to implement apache spark architecture

## Step 5

The final phase involves configuring automated email alerts using Azure Logic Apps, triggered after successful pipeline executions. These alerts contain essential information including records inserted, updated, deleted, and various runtime metrics to facilitate operational monitoring and oversight.

1. Trigger email sending function to trigger after running every successful pipeline.