

## Migrating Data from GCP to Azure SQL Managed Instance (Incremental Load)

Priya Akunuri

### GCP to Managed Instance (MI) Incremental Data Pipeline Using Azure Data Factory

#### Objective

This document outlines the steps for creating an incremental data pipeline to load the latest data from Google Cloud Big Query into an Azure Managed Instance (MI) via Blob Storage. The pipeline automates the process, ensuring that the latest file from Big Query is moved daily and loaded into MI while allowing schema drift during loading. The pipeline also includes optional file deletion from the staging area to save space.

#### Prerequisites

- Google Cloud Big Query account access.
- Azure Blob Storage set up as the staging area.
- Azure SQL Managed Instance (MI) for loading data.
- Azure Data Factory (ADF) with appropriate permissions.

#### Source and Destination Details

##### Source : Big Query

Project ID: **partner-portal-426309**

##### Destination : Azure Managed Instance

Server name: [smineuprdsq1s01.09ed840669e8.database.windows.net](https://smineuprdsq1s01.09ed840669e8.database.windows.net)

Database name: **DW**

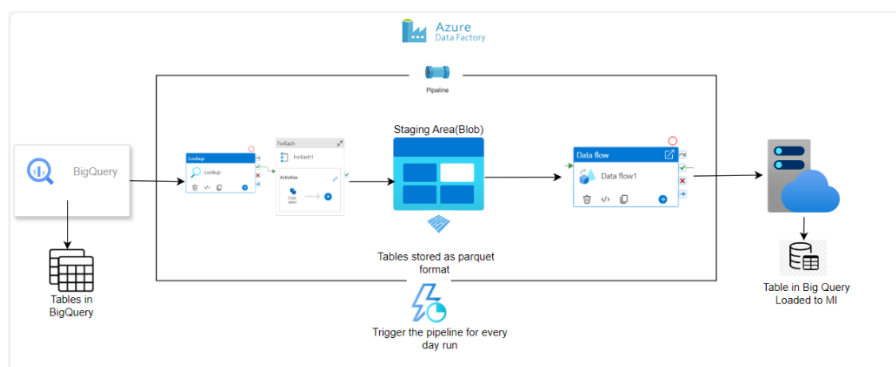
Table name : **BigQ**

#### ETL Details :

Data Factory Name : **Adfbigqueryprod**

Storage account Name : **stoweuprdbigq01**

#### Architecture Diagram



## Step-by-Step Guide

### Step 1: Create Linked Service for Big Query in Azure Data Factory

Go to Azure Data Factory.

- In Manage > Linked Services, click New.
- Select Google Big Query as the service.
- Configure the authentication settings to connect to your Big Query environment.

**Usage:** This linked service will be used in the Lookup and Copy activities to fetch the latest file from Big Query.

The screenshot shows the 'Edit linked service' configuration page for Google BigQuery. The page has a title 'Edit linked service' with a sub-header 'Google BigQuery' and a 'Learn more' link. The configuration fields are as follows:

- Name \***: LS\_GCP\_LU\_INC
- Description**: A large empty text area.
- Connect via integration runtime \***: A dropdown menu with a green checkmark icon and the text 'AutoResolveIntegrationRuntime'.
- Project ID \***: A text input field containing 'partner-portal-426309'.
- Authentication type \***: A dropdown menu with the text 'Service authentication'.
- Key file \***: A section with two tabs: 'Key file' (active) and 'Azure Key Vault'. Below the tabs is a text input field with a dotted pattern, a 'Choose File' button, and the text 'No file chosen'.
- Annotations**: A section with a '+ New' button.
- Parameters**: A section with a '> Parameters' link.



At the bottom of the page, there are three buttons: 'Apply' (blue), 'Cancel' (white), and 'Test connection' (blue with a test icon).

### Step 2: Create Linked Service for Azure Blob Storage (Staging Area)

- Go to Manage > Linked Services > New.
- Select Azure Blob Storage as the service.
- Set up the connection to your Blob Storage account.

**Usage:** This linked service will store the file temporarily before loading it to MI using a Copy activity.

**Edit linked service**


 Azure Blob Storage [Learn more](#) 

**Name \***

LC\_BLB\_SRC\_INC

**Description**

**Connect via integration runtime \*** ⓘ

 AutoResolveIntegrationRuntime

**Authentication type**

Account key

**Connection string** **Azure Key Vault**

**Account selection method** ⓘ

☐ From Azure subscription ☒ Enter manually

**Storage account name \***


stoweuprdbigq01

**Storage account key** **Azure Key Vault**

**Storage account key \***

.....

**Partitioned DNS enabled** ⓘ ☐



**Apply** **Cancel**  **Test connection**

### Step 3: Create Linked Service for Azure SQL Managed Instance

- Go to Manage > Linked Services > New.
- Select Azure SQL Managed Instance.
- Configure the connection details for the MI database.
- Create new integration run time because we are using Managed instance so auto integration runtime is not possible to use.
- Create the managed instance private end points and you can approve the end points from the managed instance
- Enable the interactive authoring

**Usage:** This linked service will be used in the Data Flow or Copy activity to load data into the MI table, allowing schema drift.

**Edit linked service**




 Azure SQL Database Managed Instance [Learn more](#) 


**Name \***

LS\_BLB\_BLB2MI\_INC

**Description**

**Connect via integration runtime \*** ⓘ

 integrationRuntime1 (Managed Virtual Network)  

 Interactive authoring disabled ⓘ

**Version**

☐ Recommended ☒ Legacy

**Connection string** **Azure Key Vault**

**Account selection method** ⓘ

☐ From Azure subscription ☒ Enter manually


**Fully qualified domain name \***

smineuprdsqsls01.09ed840669e8.database.windows.net

**Database name \***

DW

**Authentication type \***

**Apply** **Cancel**  Test connection

#### Step 4: Create Datasets for Google Cloud BigQuery (Source)

- Create two datasets for BigQuery:

##### For the Copy Activity:

This dataset will point to the BigQuery source, from which data will be copied to the staging area (Blob).

##### For the Lookup Activity:

This dataset will be used to fetch the latest file using a query.

#### Step 5: Create Dataset for Azure Blob Storage (Staging Area)

Create a dataset in ADF for the Blob Storage.

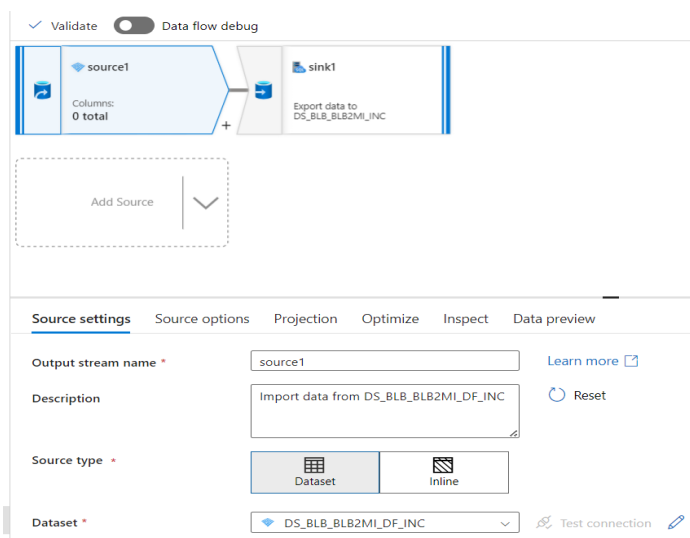
- This dataset will be used as the destination for the Copy Activity, which temporarily stores the data before transferring it to MI.
- I am using the parquet format.

#### Step 6: Create Dataset for Azure SQL Managed Instance (MI)

- Create a dataset for the SQL Managed Instance that will be used in the Data Flow activity to load the data from the Blob storage to MI.

#### Step 7: Create a Data Flow for Loading Data from Blob to MI

- In ADF, create a new Data Flow.
  - Set Blob Storage as the Source.
  - Set SQL Managed Instance as the Sink.
  - Allow schema drift to accommodate any changes in the schema from the source.
- Enable Debug Mode to troubleshoot during development.



#### Step 8: Create the Pipeline

##### Add Lookup Activity:

- In your ADF pipeline, add a Lookup Activity.
- Use the dataset created for BigQuery Lookup.
- Write the following query to fetch the latest file from BigQuery:

## SQL

```
SELECT table_name, table_schema
```

```
FROM `partner-portal-426309.analytics_393043220.INFORMATION_SCHEMA.TABLES`
```

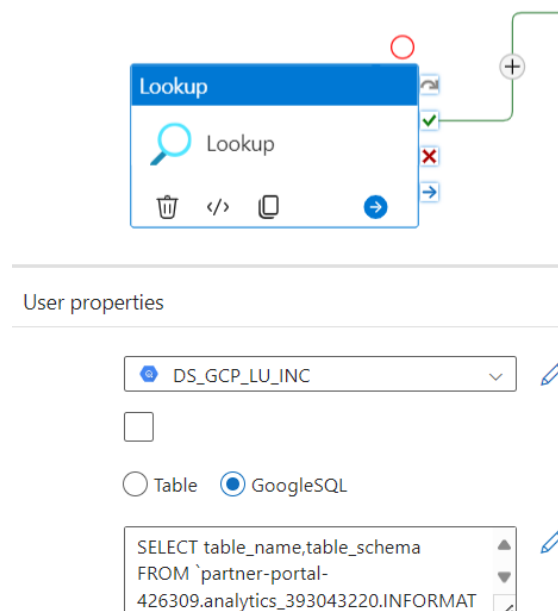
```
WHERE table_type = 'BASE TABLE'
```

```
AND table_name LIKE 'events_%'
```

```
ORDER BY creation_time DESC
```

```
LIMIT 1;
```

This query will return the latest table.



### Add For Each Activity:

- Add a For Each Activity to loop through the results from the Lookup.  
[Use @activity('Lookup').output.value]  
to fetch the file name dynamically.
- Configure parameters for table\_name and table\_schema to fetch the latest file.  
Set the number of files to process, e.g., LIMIT 1 for the latest file.



- In the Sink properties, use the Blob dataset and specify the file format (e.g., .parquet):

The screenshot shows the 'Sink properties' dialog in Azure Data Factory. The 'Connection' tab is active, displaying a dropdown for 'Linked service' set to 'LS\_GCP\_CD\_INC'. To the right are buttons for 'Test connection', 'Edit', 'New', and 'Learn more'. The 'Parameters' tab is also visible, showing two input fields for '@dataset().table\_name' and '@dataset().table\_name', with a 'Preview data' button. A checkbox labeled 'Enter manually' is checked.

```
@concat(item().table_schema,'_',item().table_name,'.pqt')
```

### Add Data Flow Activity:

- After the Copy Activity, add the Data Flow Activity.
  - Connect the For Each loop to the Data Flow.
- This will load the data from Blob Storage to MI, allowing schema drift.

### Optional: Add Delete Activity:

- Add a Delete Activity to remove the file from Blob Storage after the data is loaded into MI.
- This helps to save storage space by deleting the file after processing.

**Note:** Skip this step if you want to retain the file in Blob Storage.

### Step 9: Set Up Daily Trigger

- In Triggers, create a new trigger that runs the pipeline every day.
- Attach this trigger to the pipeline for automatic execution without manual



intervention.

The screenshot shows the 'Trigger' configuration window in Azure Data Factory. The 'Name' field is set to 'trigger\_INC'. The 'Description' field is empty. The 'Type' is set to 'ScheduleTrigger'. The 'Start date' is '9/30/2024, 1:00:00 AM'. The 'Time zone' is 'Dublin, Edinburgh, Lisbon, London (UTC+1)'. A note indicates that this time zone observes daylight savings. The 'Recurrence' is set to 'Every 1 Day(s)'. Under 'Advanced recurrence options', the 'Execute at these times' section is expanded, showing 'Hours' and 'Minutes' fields, both of which are empty. At the bottom, there are 'OK' and 'Cancel' buttons.

Name \*

trigger\_INC

Description

Type \*

ScheduleTrigger

Start date \* ⓘ

9/30/2024, 1:00:00 AM

Time zone \* ⓘ

Dublin, Edinburgh, Lisbon, London (UTC+1)

ⓘ This time zone observes daylight savings. Trigger will auto-adjust for one hour difference.

Recurrence \* ⓘ

Every 1 Day(s)

Advanced recurrence options

Execute at these times ⓘ

Hours

Minutes

Schedule execution times

OK Cancel

## Step 10: Debug and Run the Pipeline

- Enable Debug Mode in ADF to test the pipeline before deployment.
- After successful debugging, publish the pipeline.  
The pipeline will now run daily, fetching the latest file from BigQuery, loading it to Blob Storage, then transferring it to MI.

## Conclusion:

By following the steps outlined above, we have successfully created an incremental data pipeline in Azure Data Factory that automates the process of fetching the latest data from Google Cloud Big Query and loading it into an Azure SQL Managed Instance (MI). This solution efficiently moves the latest files daily, utilizes Blob Storage as a staging area, and incorporates flexible data handling by allowing schema drift. Additionally, the optional file deletion step helps optimize storage usage.