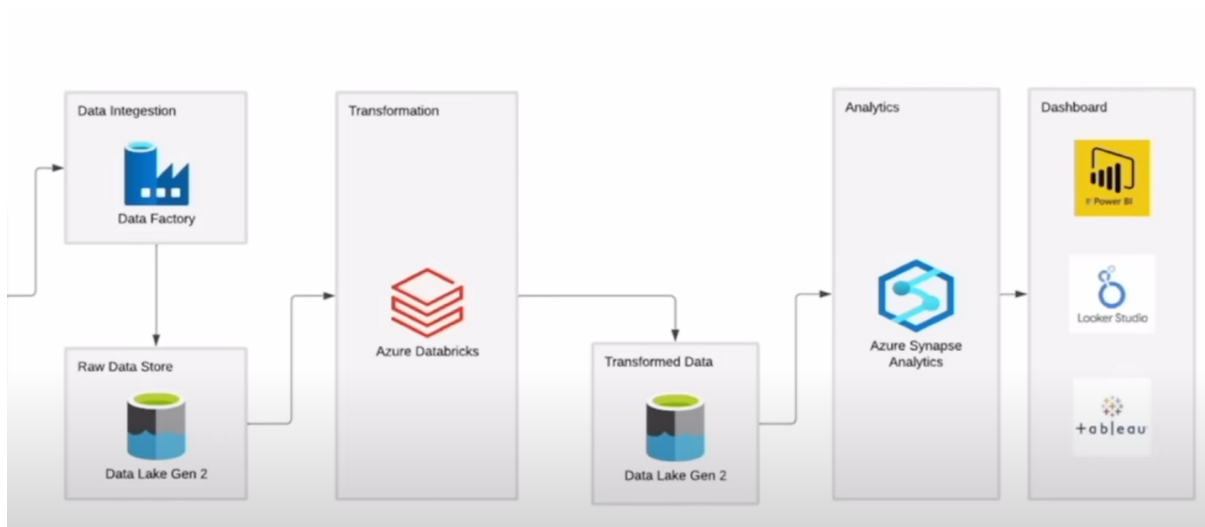Dataflow diagram:



Data Factory – to create ETL pipelines to ingest the data and load the data into other locations

Data Lake Gen 2 – to combine data lake features with Azure Blob Storage

Databricks – Analytics platform built on top of Apache Spark for big data and ML frameworks

Synapse Analytics – Data warehouse available on Azure
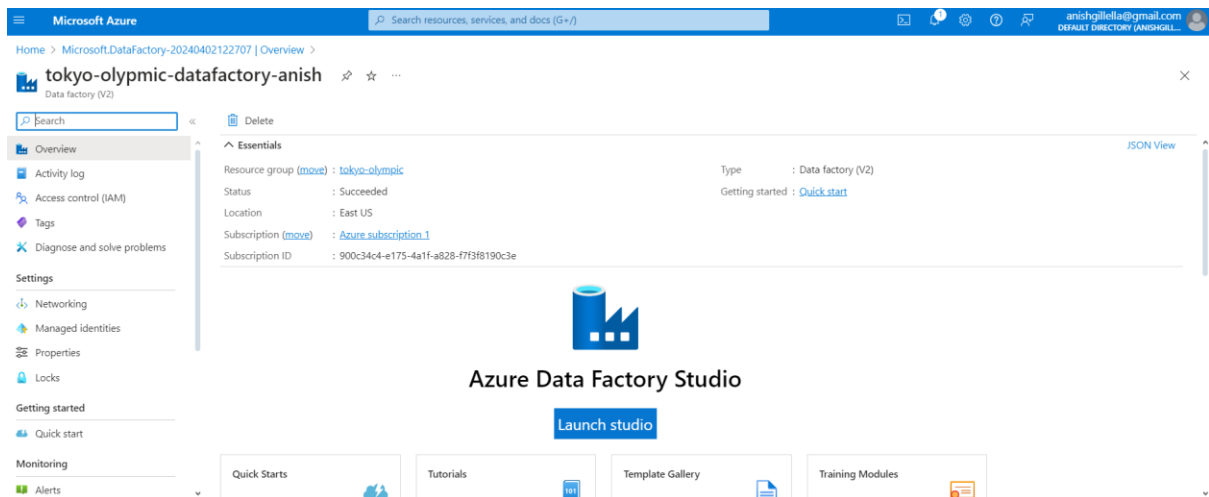
**Creating a storage account:**
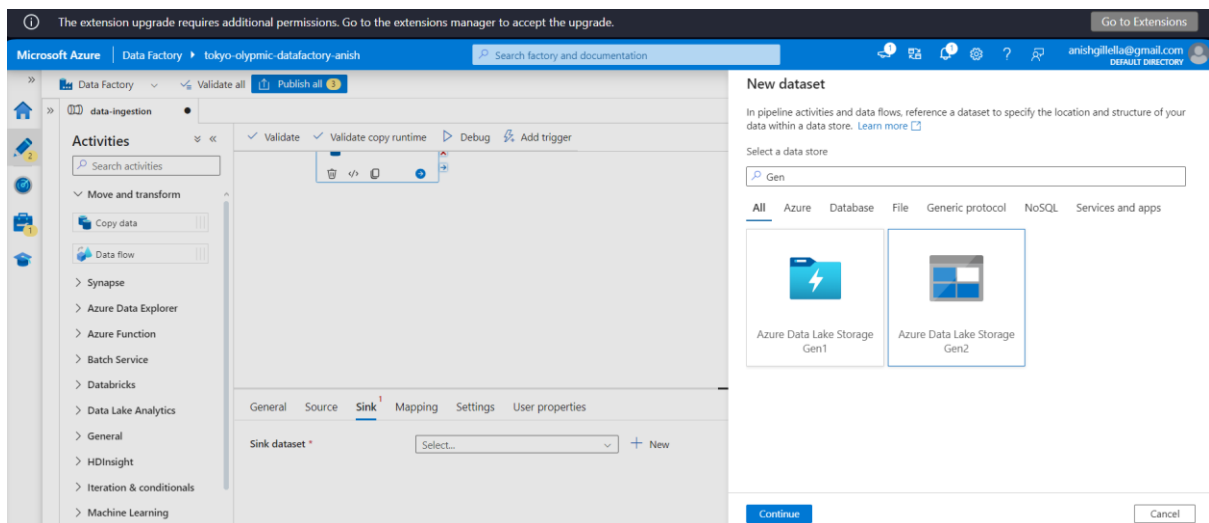


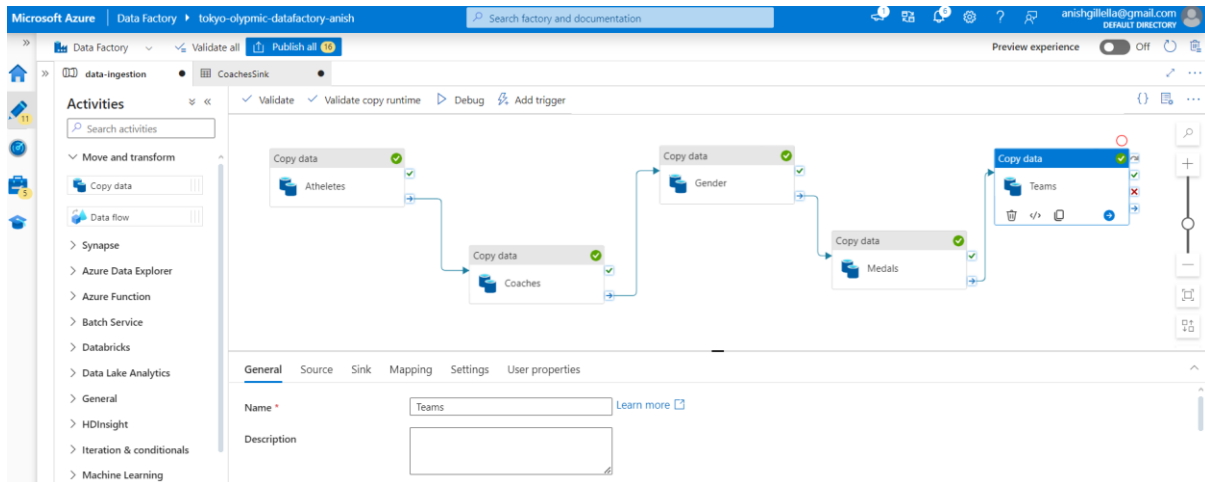**Storing the data in a container:**

**Creating a data factory:**



**Forking the dataset from github and ingesting the data connecting the GitHub repository to the data factory:**



**Sinking it into Gen 2 data lake storage:**

**Bringing in all the data regarding the teams, medals, gender, coaches and athletes to the data factory:**



**Running the pipeline to import the data:**



**INGESTION AND STORING THE DATA HAS BEEN COMPLETED UPTO THIS POINT**

**CREATING A DATABRICKS RESOURCE**

## CREATING A COMPUTE CLUSTER IN AZURE DATABRICKS



Mounting azure data lake storage to the data bricks azure data factory to easy access the data

Configuring and Mounting the Notebook to the container in data storage account using the tenant ID and secret key:
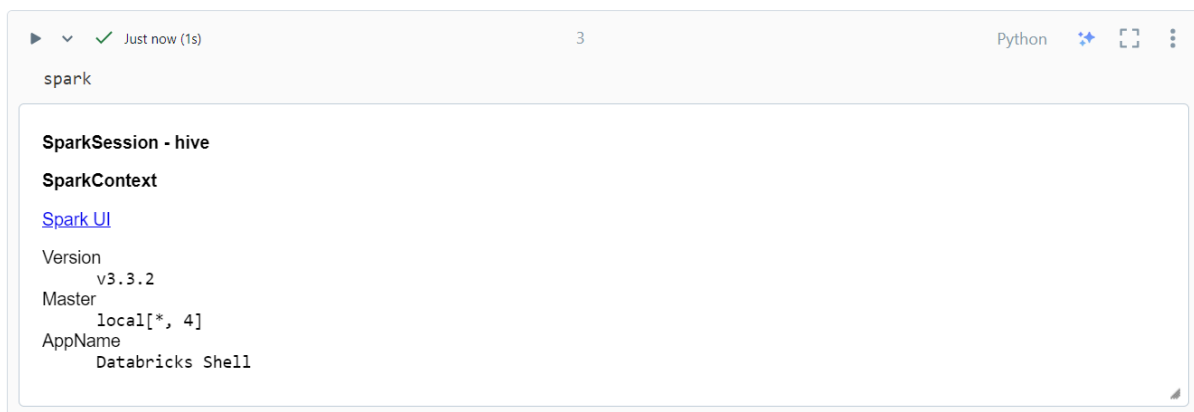


We do not have permission since the keys used in app registry do not have permission to access the contents of data storage account. So we assign the IAM in the containers to let the app registry access it:
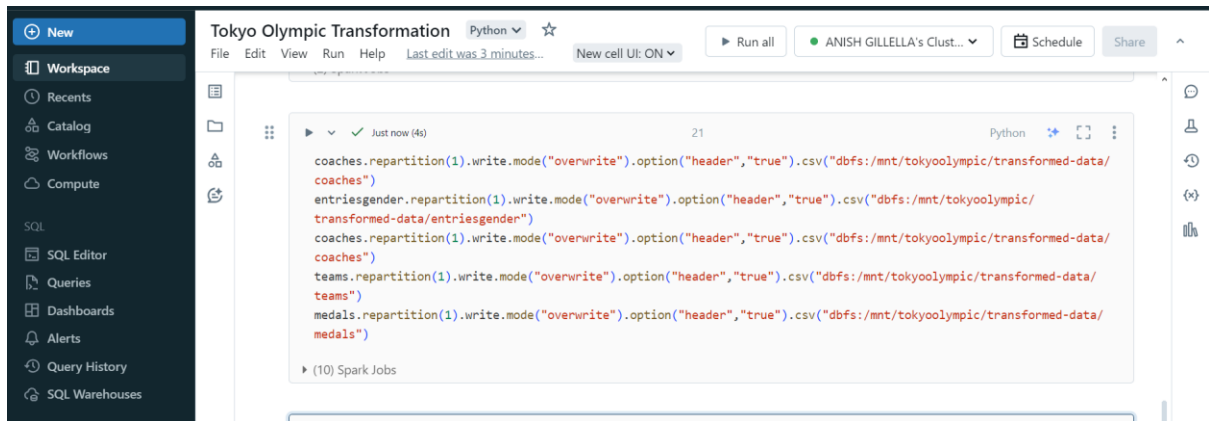
**Since it is azure databricks we do not have to create a spark session as it is already inbuilt like we usually do when using spark**
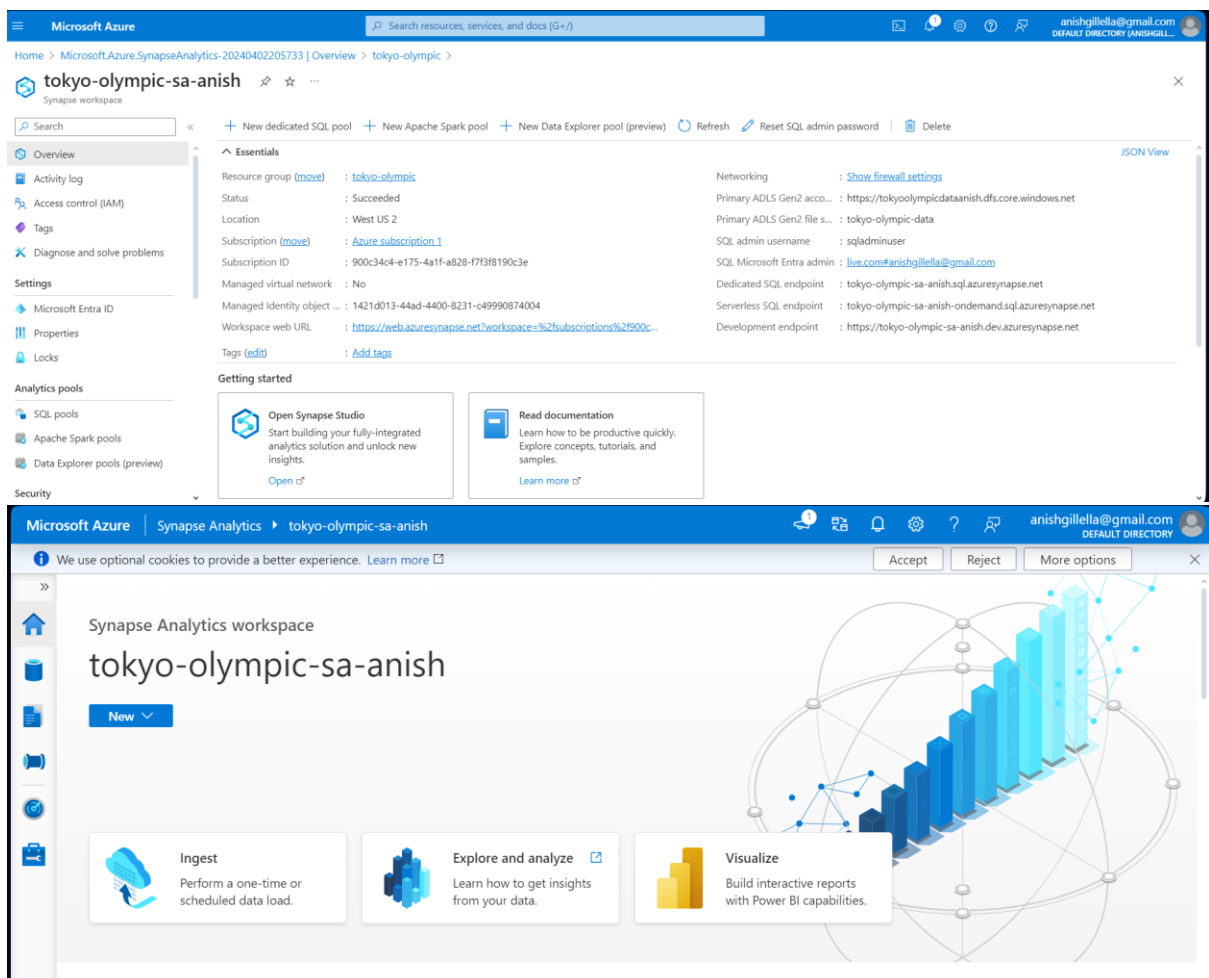


**Apache Spark uses lazy evaluation wherein spark does not perform any loading or transformation, it only performs the action when you call the particular data which was loaded with a function like we did using show() here**

**Writing the transformed data back to the data storage factory**



## CREATING AN AZUR SYNAPSE ANALYTICS WORKSPACE

## CREATING A TABLE IN SYNAPSE ANALYTICS FROM THE DATA LAKE



## QUERYING THE DATA IN SYNAPSE ANALYTICS: