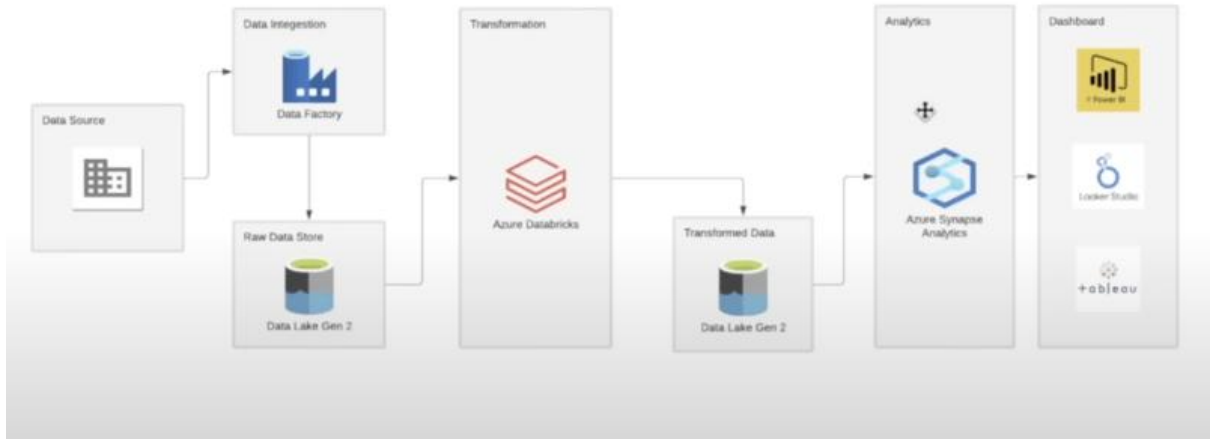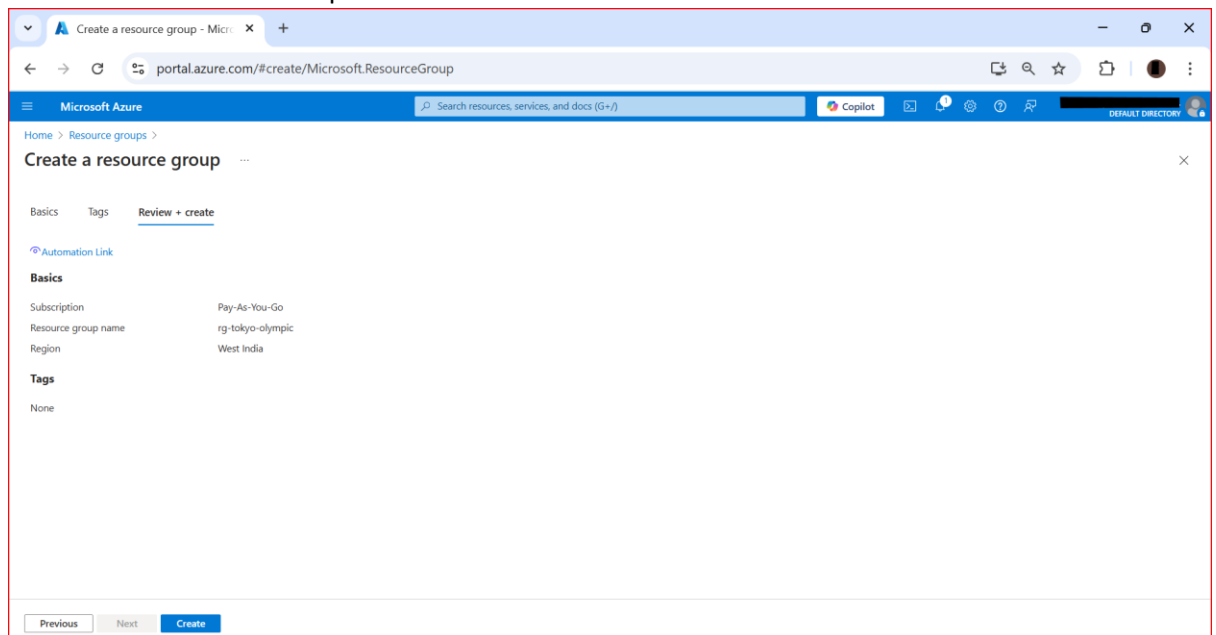# Azure End to End Data Engineering Project
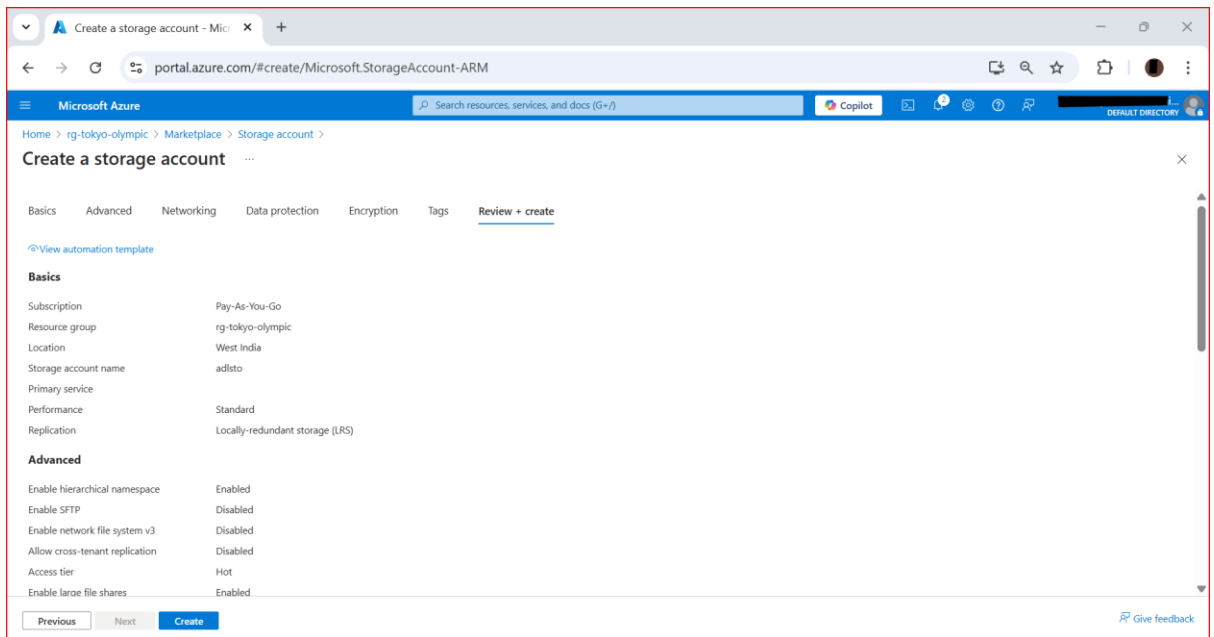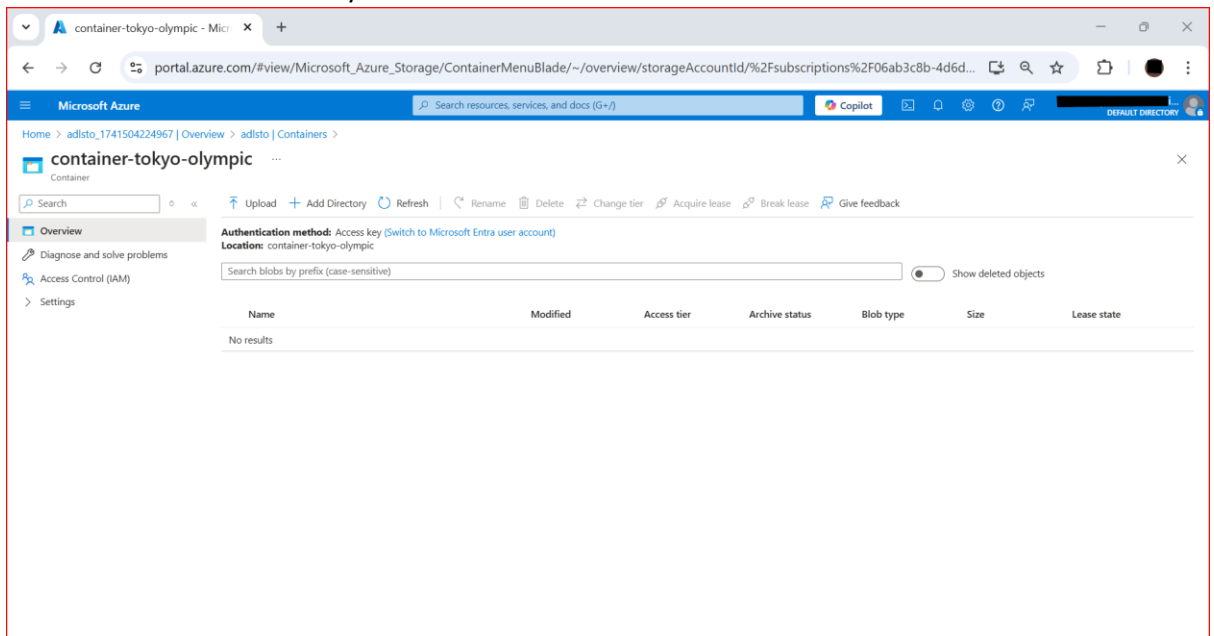
## Project Flow



## Steps to follow

1. Create Azure Portal Account and Login
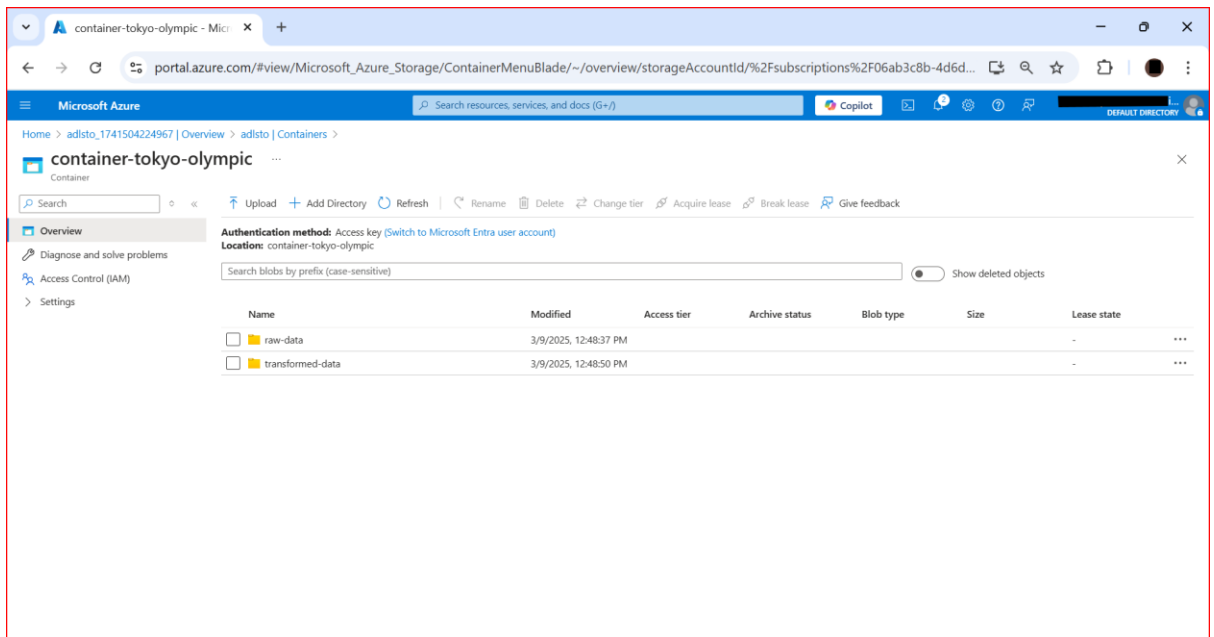
2. Create New Resource Group



3. Goto new resource group, click create, and create Azure Data Lake Storage Gen2 Account
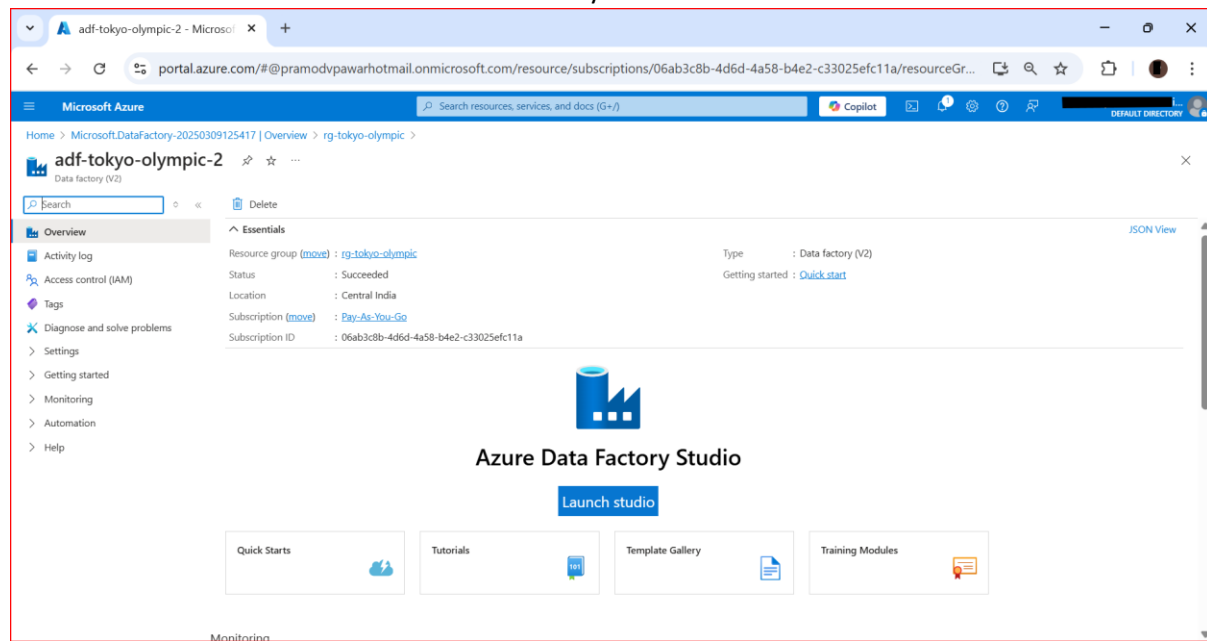   Note: - Enable Hierarchical Namespace
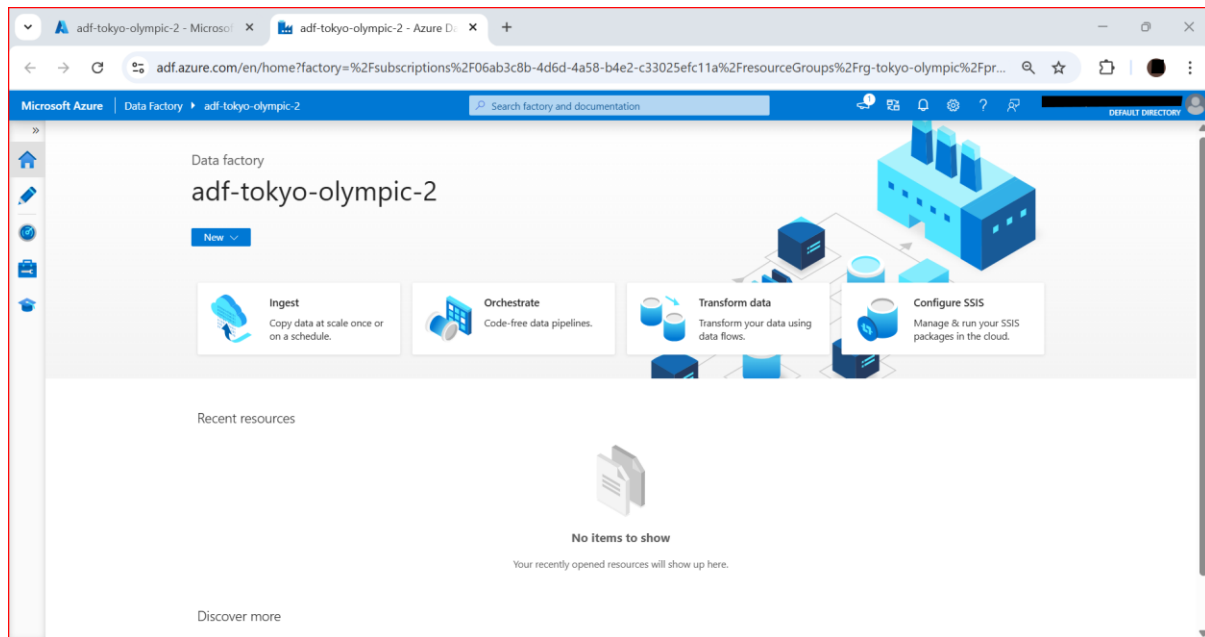
4. Create container under newly created ADLS Gen2 account



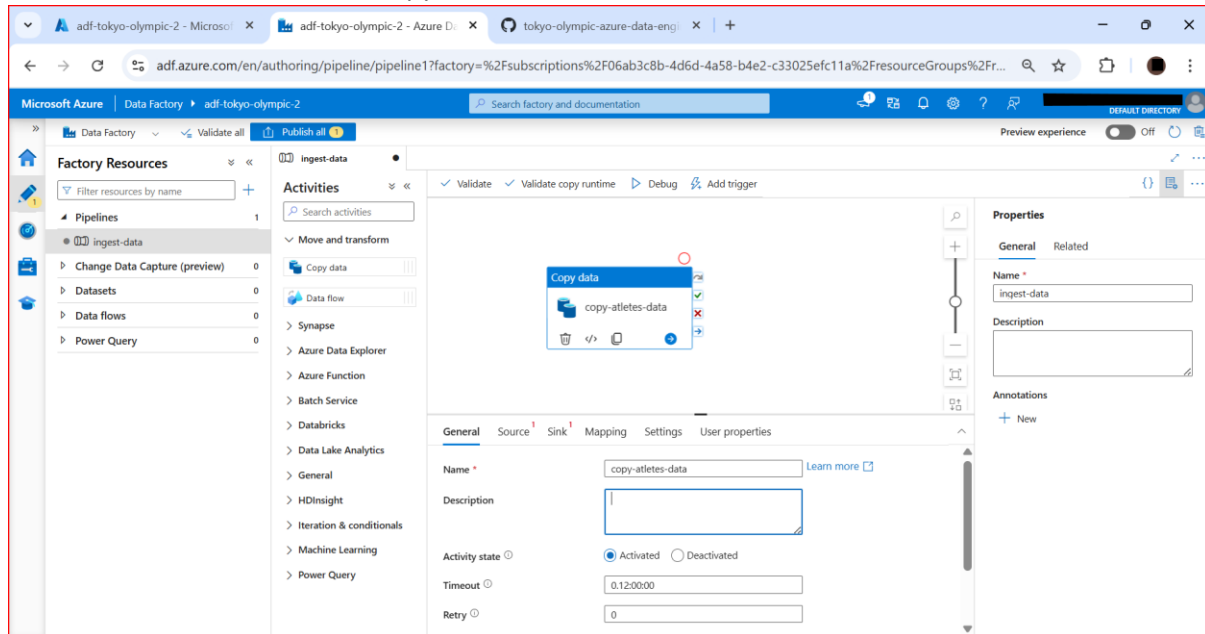5. Create Directories under container, raw-data, transformed-data

6. Open Azure Data Factory to ingest data
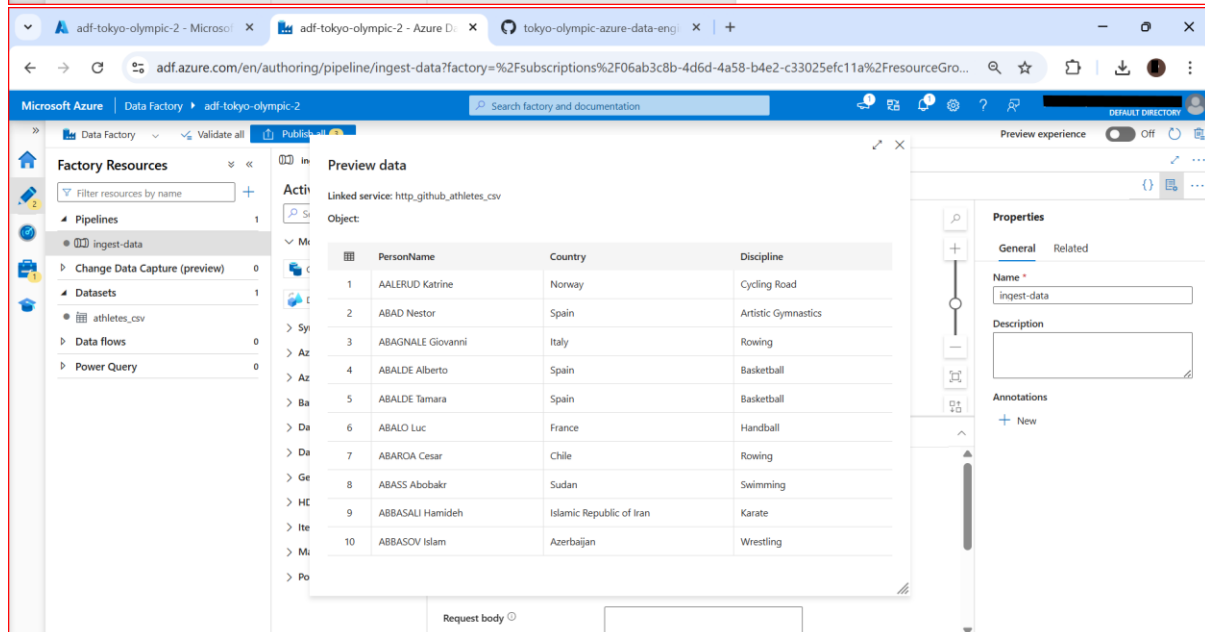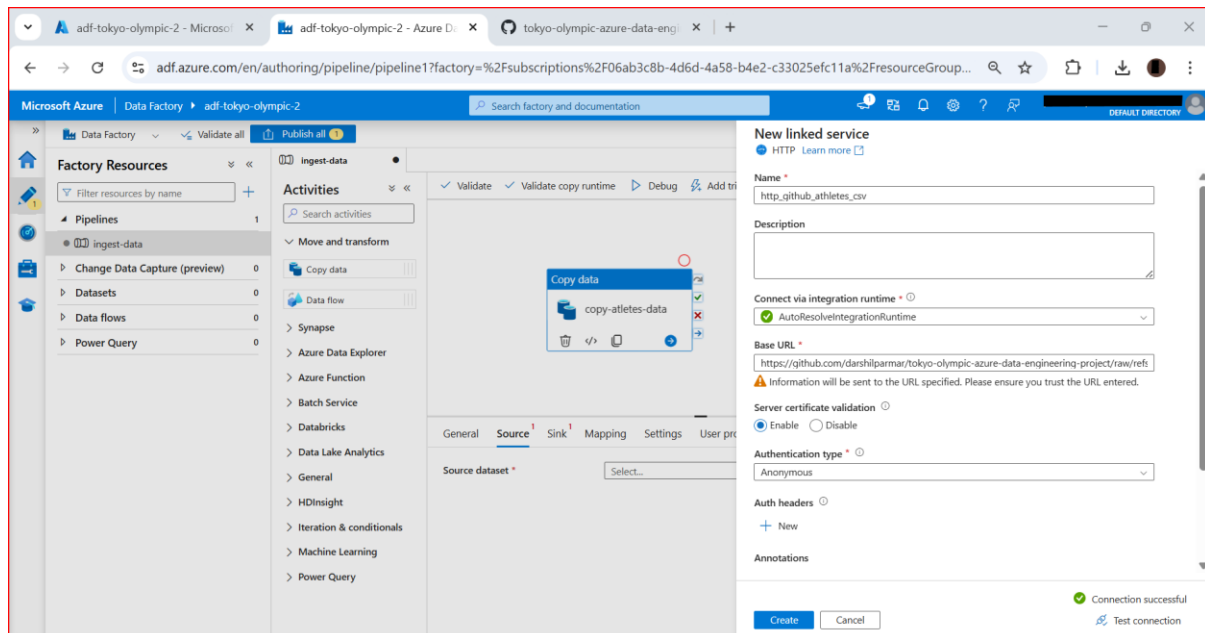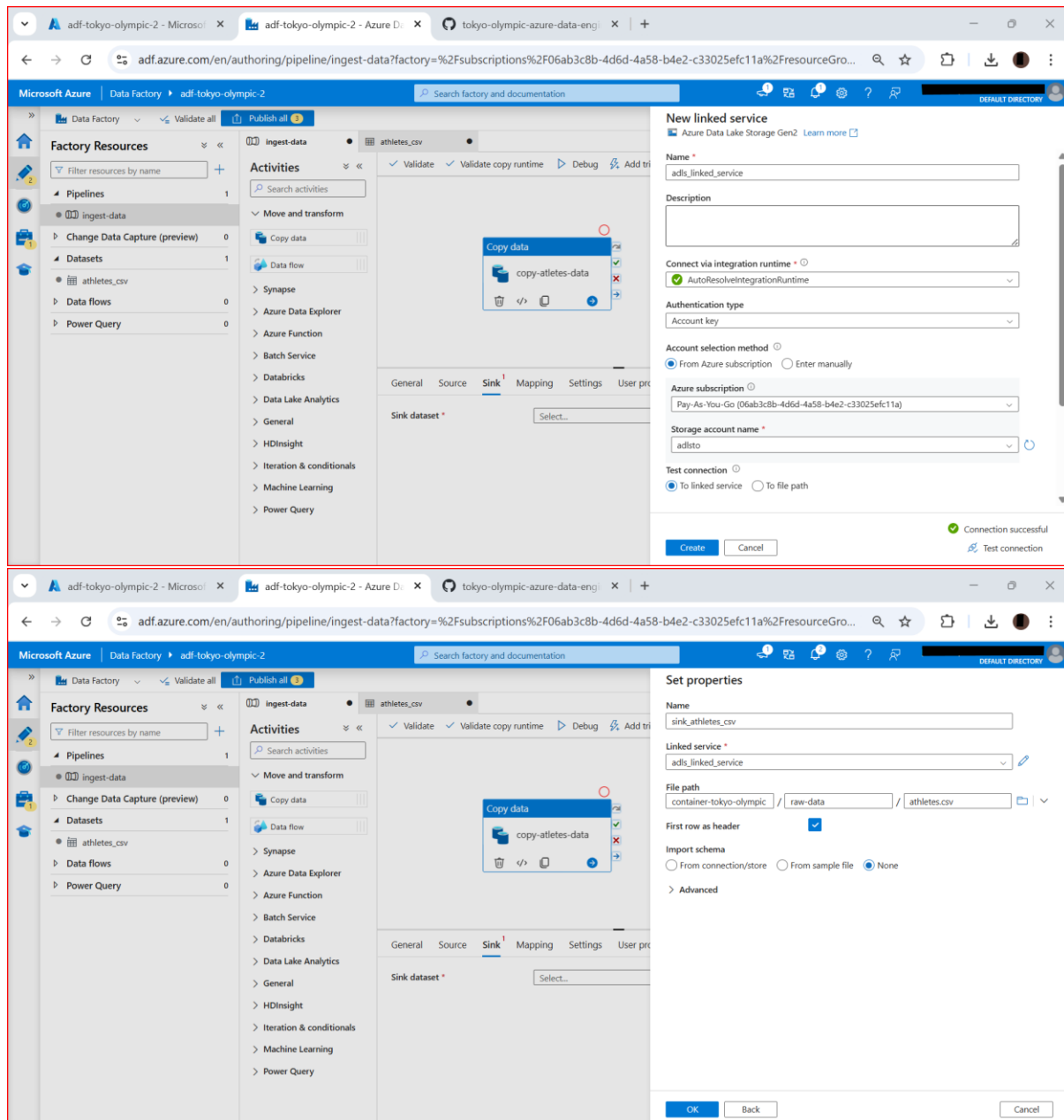   a. Search Data factories and create new data factory



   b. Launch Azure Studio

c. Create new pipeline: data-ingestion
d. Select Move and Transform → Copy Data



e. Source → New Dataset → http → file format as CSV → Create Linked Service → Give Base URL (raw URL from GitHub) → First Row as Header → Click OK → Preview Data

f.  Go to Sink → New Dataset → New Linked Service → ADLS Gen 2 (Select ADLS Gen 2 created before) → Select Folder raw-data

g.  Click Validate and Debug to see if data is getting loaded from http source to ADLS Gen2 target location

h. Repeat the activity for other CSV files
i. Connect different activity with arrows to run them one after the another

Check if all the source files are copied to ADLS location



7. Create an Azure Databricks Workspace.

8. Once deployment is complete; Goto resource and click Launch workspace

9. Click Compute → Create Compute
   a. Select Single node
   b. Select runtime
   c. Select node type
   d. Create compute



Important Note: - I received Quota Exceeded error while creating compute. I raised the Quota increase requested as suggested in error message. After some time, the request got approved and I was able to create compute.

10. Create notebook: New → Notebook → Rename notebook. Select Spark Cluster that is newly created

11. Create new App (to be used for connecting Azure Databricks with ADLS Gen2)
    a. Go to Azure Portal → Search for App Registry
    b. Register App → Give some name → Click Register
    c. Copy Client ID, Tenant ID
    d. Click on Manage → Certificates and Secrets
    e. New Client Secrets → Give some name → Copy secret value

12. Go back to Databricks → Create new cell → Copy standard config from Azure documentation
    (https://learn.microsoft.com/en-us/azure/databricks/connect/storage/tutorial-azure-storage)
    configs = {"fs.azure.account.auth.type": "OAuth",
    "fs.azure.account.oauth.provider.type":
    "org.apache.hadoop.fs.azurebfs.oauth2.ClientCredsTokenProvider",
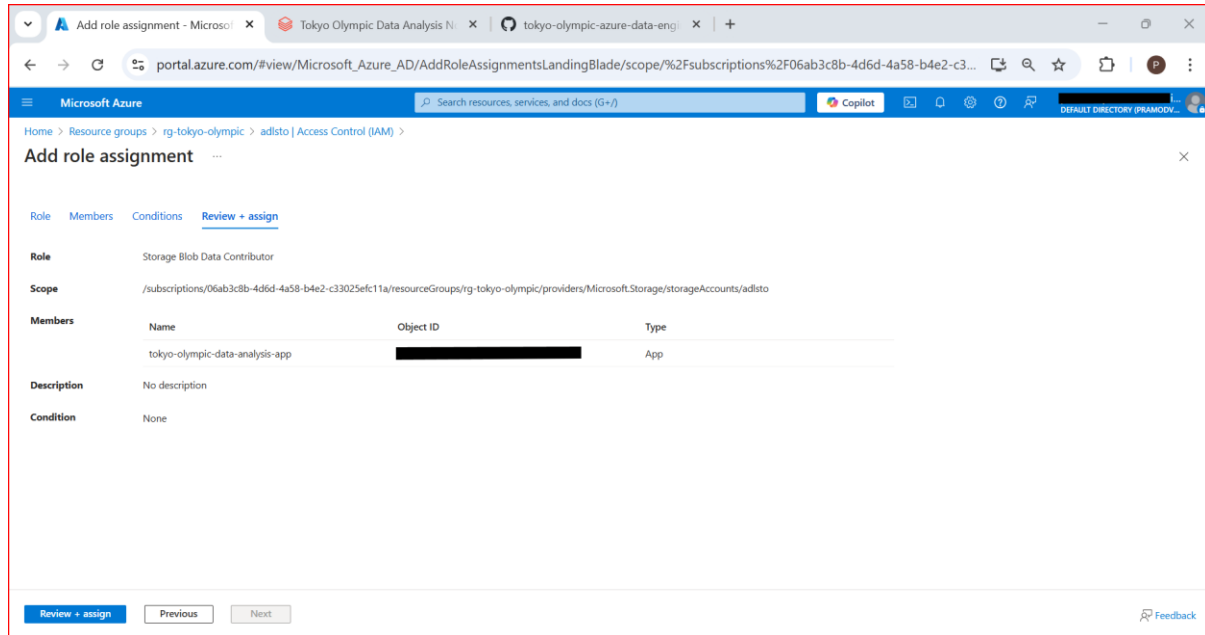    "fs.azure.account.oauth2.client.id": "<YOUR CLIENT ID>",
    "fs.azure.account.oauth2.client.secret": 'YOUR APP SECRET VALUE',
    "fs.azure.account.oauth2.client.endpoint": "https://login.microsoftonline.com/<YOUR
    TENANT ID>/oauth2/token"}

13. Replace the tenant id and other details and ensure the connection is successful.

14. Copy the mount point code given below. Replace the container name and ADLS Gen2 Storage
    account name
    dbutils.fs.mount(
    source = "abfss://your_container_name@your_adls_storage_name.dfs.core.windows.net", #
    contrainer@storageacc
    mount_point = "/mnt/tokyoolymic",
    extra_configs = configs)

15. We need to give permission to the app we created in order to access ADLS Gen2 storage
    from Azure databricks
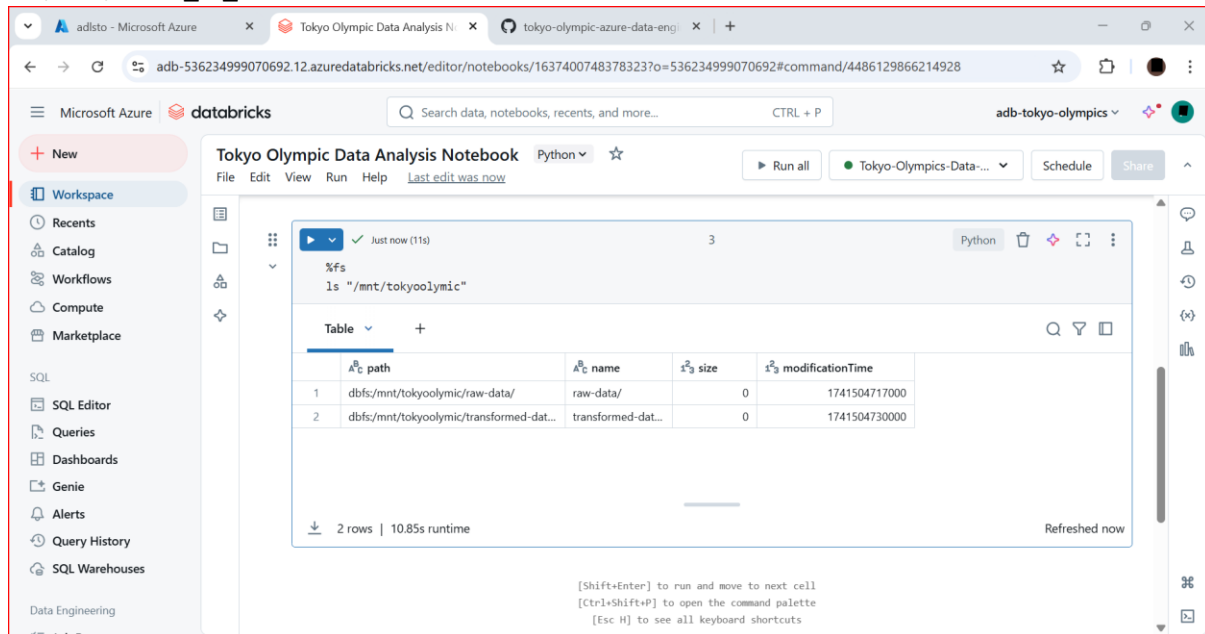
a. Goto ADLS Gen2 → select container → Click Access Control (IAM) → Click Add Role Assignment → Select "Storage Blob Data Contributor" → Click Next → Select Member → Type App Name (The App we created before) → Click Next → Review and assign → Wait for some time and revisit Azure Data Bricks



16. Test the Azure Databricks and ADLS Gen2 connection again
    %fs
    ls "/mnt/name_of_container"



17. We can now start reading CSV files in Spark from ADLS Gen2
    athletes =
    spark.read.format("csv").option("header","true").load("/mnt/container/folder/filename.csv"
    )
    We can use printSchema to check the data type

We can use withColumn and col operator to convert string column to Integer. We will need to add import from pyspark.sql to use these operators
gender = gender.withColumn("Female"), col("Female").cast(IntegerType))



18. We can use data type manually as mentioned above or we can use inferSchema for conversion by Spark

19. After doing some transformations, we can write data back to ADLS Gen2 in transformed-data folder
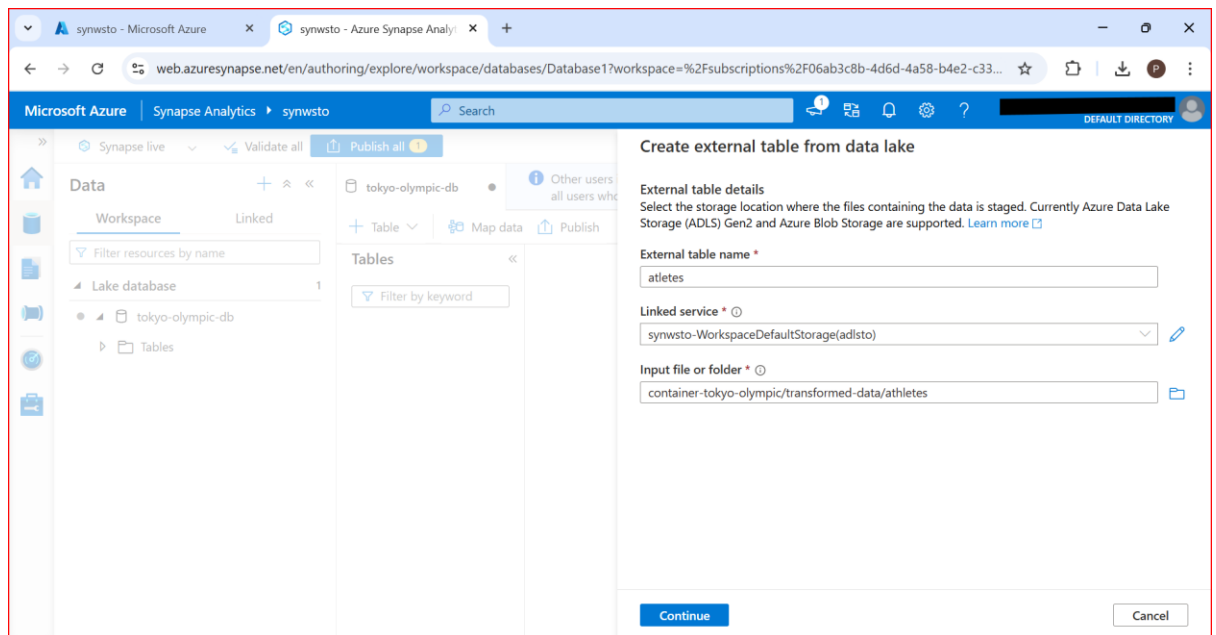


20. We can now open Azure Synapse Analytics. Create Workspace → Give unique workspace name → Select ADLS Gen 2 Account which is created before → Next → Review and Create. Wait for few minutes for workspace to be created

21. Open Synapse Studio → Click Data → Click Lake Database → Give DB Name

22. Create Table → From Data Lake → Give Table Name → Select Storage Account → Select
    Container → Folder → transformed-data → Click Validate and Publish. Create tables for
    other files.



23. We can run queries against newly created tables. Perform some basic analytics on newly
    created table

24. Create Dashboards using any visualization tool

25. Cleanup – Remove the resource group