

Assessment Title
Azure Practical Scenario Evaluation

CCC602 | Assessment-3

The Student Name | Geni Bahar
Student ID: 27085022

Tutor/Accessor | Ovesh Vohra

Programme
Diploma in Cloud Engineering and Cyber Security (120 Credits)

Course
CCC602: Data Engineering in the Public Cloud
(Level 6, 30 Credits)

Date: 17th January 2026

Table of Contents

Task 1: Cloud Vendor Evaluation (Theoretical Evaluation).....	2
Question 1: Comparative Evaluation: Microsoft Azure vs. Amazon Web Services (AWS).....	2
Big Data and Analytics Capabilities.....	2
Integration with Business Intelligence (BI) Tools.....	2
Cost Structures and Pricing Comparison	2
Scalability and Automation.....	3
Vendor Selection Report.....	3
Task 2: Dataset Acquisition and Azure Resource Deployment Using IaC.....	4
Azure Resource Deployment Using ARM Template	4
Tokyo Olympic Dataset Acquisition	14
Task 3: Data Ingestion Using Azure Data Factory	18
Data Ingestion Pipeline Design	18
Source and Sink Configuration.....	24
Pipeline Validation and Verification	26
Pipeline Summary Report	28
Task 4: Create an Azure Databricks Workspace and perform data transformation using code to automate the process.	30
Azure Databricks Workspace and Cluster Setup	30
Secure Data Access and Authentication	35
Data Transformation Using PySpark.....	42
Databricks notebook with transformation code	42
Transformdata folder output.....	45
Task 5: Azure Synapse Analytics and Automated Data Integration	47
Azure Synapse Analytics Workspace Creation.....	47
Data Integration Pipeline	49
External Table Creation Using SQL.....	57
Data Analysis Using SQL.....	58
Data Visualisation	59

Task 1: Cloud Vendor Evaluation (Theoretical Evaluation)

Question 1: Comparative Evaluation: Microsoft Azure vs. Amazon Web Services (AWS)

Big Data and Analytics Capabilities

- Both Microsoft Azure and AWS offer mature big data and analytics services, but they differ in how seamlessly these services work together. From my learning and hands-on experience during this project, Azure provides a well-connected analytics ecosystem that includes Azure Data Factory, Azure Data Lake Storage Gen2, Azure Databricks, and Azure Synapse Analytics. I found that these services were easier to connect and manage together when building an end-to-end data pipeline. These services are designed to work together, which reduces setup complexity and learning overhead.
- AWS offers equivalent services such as Amazon S3, AWS Glue, Amazon EMR, and Amazon Redshift. While these services are powerful and widely used, they often require additional configuration and integration effort to achieve the same end-to-end workflow that Azure provides more naturally. From a learning and implementation perspective, Azure felt more structured and easier to manage for a complete data pipeline.

Integration with Business Intelligence (BI) Tools

- Integration with BI tools is a key requirement for analysing and presenting insights from the Tokyo Olympic dataset. Azure integrates natively with Power BI, which I found particularly useful as it allows simple connectivity and real-time data visualisation without requiring complex additional configuration. Since Power BI is part of the Microsoft ecosystem, it works smoothly with Azure services without additional connectors or complex configurations.
- AWS supports BI tools such as Amazon QuickSight and Tableau. While these tools are effective, integration with Power BI is less direct and often requires extra setup. For organisations already using Microsoft products, Azure offers a more seamless and user-friendly BI experience.

Cost Structures and Pricing Comparison

- Azure provides flexible pricing models such as pay-as-you-go, reserved capacity, and serverless options. Azure Data Lake Storage Gen2 offers low-cost storage for large datasets, and Azure Synapse serverless SQL allows querying data without provisioning dedicated resources, which helps control costs. Azure Databricks also allows scaling compute resources only when required.
- AWS offers similar pricing models for Amazon S3 and EMR; however, EMR clusters can become expensive if not carefully managed. In comparison, Azure's serverless and consumption-based services provide better cost transparency for smaller projects and student-level environments.

Scalability and Automation

- Both platforms support scalability and automation through Infrastructure as Code. Azure supports automation through ARM templates, which helped me understand how infrastructure can be deployed in a repeatable and consistent way without manually creating each resource through the Azure Portal. AWS uses CloudFormation, which provides similar functionality but has a steeper learning curve.
- Azure's automation tools felt easier to understand and implement, especially when deploying multiple resources such as storage, data factory, and analytics services together.

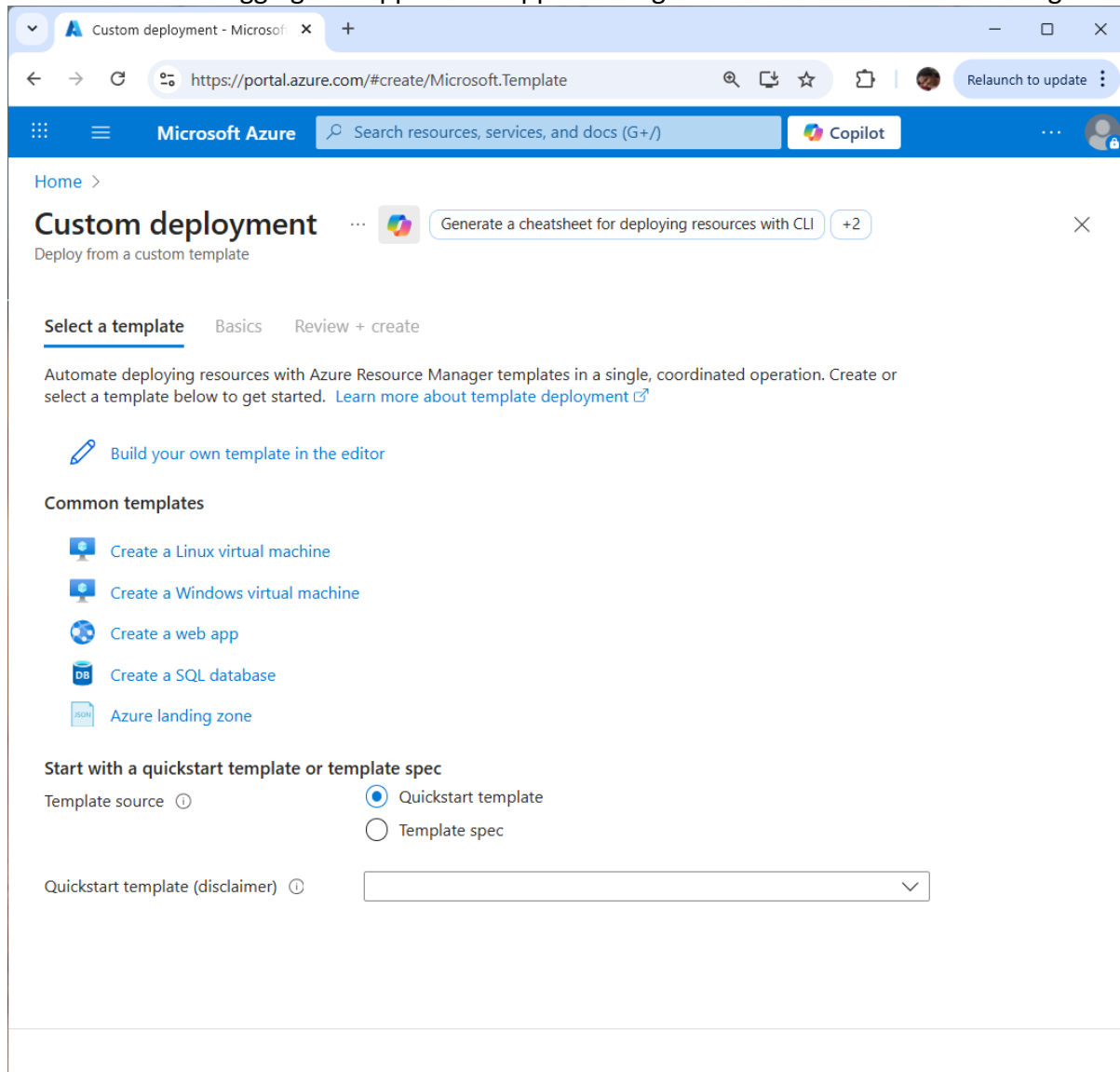
Vendor Selection Report

- Microsoft Azure was selected as the preferred cloud platform due to its strong integration across analytics services, native compatibility with Power BI, flexible cost structure, and robust automation capabilities using ARM templates. Azure aligns well with data analytics workloads and the Microsoft ecosystem, making it a practical and efficient choice for building a scalable and cost-effective data engineering solution.

Task 2: Dataset Acquisition and Azure Resource Deployment Using IaC

Azure Resource Deployment Using ARM Template

- In this task, core Azure infrastructure was deployed using an Azure Resource Manager (ARM) template to ensure consistency, repeatability, and cost control.
- A resource group named **geni_Cloud_RG** was created. Within this resource group, an Azure Data Lake Storage Gen2 account named **geniadls** was deployed. The storage account was configured with the StorageV2 type and hierarchical namespace enabled to support big data analytics.
- Resource tagging was applied to support cost governance and resource tracking.



Custom deployment - Microsoft

https://portal.azure.com/#create/Microsoft.Template

Microsoft Azure Search resources, services, and docs (G+/) Copilot

Home >

Custom deployment

Deploy from a custom template

New! Deployment Stacks let you manage the lifecycle of your deployments. Try it now →

Select a template **Basics** Review + create

Template

Customized template 1 resource

Edit template Edit parameters Visualize

Project details

Select the subscription to manage deployed resources and costs. Use resource groups like folders to organize and manage all your resources.

Subscription * Azure subscription 1

Resource group * (New) geniwaruni_Cloud_RG

Create new

Instance details

Region * Australia East

Storage Account Name geniwarunicloudstorage

Location australiacentral

Previous Next **Review + create**

Review + create

Custom deployment - Microsoft

https://portal.azure.com/#create/Microsoft.Template

Microsoft Azure Search resources, services, and docs (G+)

Home >

Custom deployment

Deploy from a custom template

I want to deploy resources using Bicep Automate Azure deployments with CLI +1

Select a template Basics **Review + create**

Summary

Customized template
1 resource

Terms

[Azure Marketplace Terms](#) | [Azure Marketplace](#)

By clicking "Create," I (a) agree to the applicable legal terms associated with the offering; (b) authorize Microsoft to charge or bill my current payment method for the fees associated the offering(s), including applicable taxes, with the same billing frequency as my Azure subscription, until I discontinue use of the offering(s); and (c) agree that, if the deployment involves 3rd party offerings, Microsoft may share my contact information and other details of such deployment with the publisher of that offering.

Microsoft assumes no responsibility for any actions performed by third-party templates and does not provide rights for third-party products or services. See the [Azure Marketplace Terms](#) for additional terms.

Deploying this template will create one or more Azure resources or Marketplace offerings. You acknowledge that you are responsible for reviewing the applicable pricing and legal terms associated with all resources and offerings deployed as part of this template. Prices and associated legal terms for any Marketplace offerings can be found in the [Azure Marketplace](#); both are subject to change at any time prior to deployment.

Neither subscription credits nor monetary commitment funds may be used to purchase non-Microsoft offerings. These purchases are billed separately.

If any Microsoft products are included in a Marketplace offering (e.g. Windows Server or SQL Server), such products are

Previous Next **Create** Create

Custom deployment - Microsoft

https://portal.azure.com/#create/Microsoft.Template

Microsoft Azure Search resources, services, and docs (G+/) Copilot

Home >

Custom deployment

Deploy from a custom template

I want to deploy resources using Bicep Automate Azure deployments with CLI +1

By clicking "Create," I (a) agree to the applicable legal terms associated with the offering; (b) authorize Microsoft to charge or bill my current payment method for the fees associated the offering(s), including applicable taxes, with the same billing frequency as my Azure subscription, until I discontinue use of the offering(s); and (c) agree that, if the deployment involves 3rd party offerings, Microsoft may share my contact information and other details of such deployment with the publisher of that offering.

Microsoft assumes no responsibility for any actions performed by third-party templates and does not provide rights for third-party products or services. See the [Azure Marketplace Terms](#) for additional terms.

Deploying this template will create one or more Azure resources or Marketplace offerings. You acknowledge that you are responsible for reviewing the applicable pricing and legal terms associated with all resources and offerings deployed as part of this template. Prices and associated legal terms for any Marketplace offerings can be found in the [Azure Marketplace](#); both are subject to change at any time prior to deployment.

Neither subscription credits nor monetary commitment funds may be used to purchase non-Microsoft offerings. These purchases are billed separately.

If any Microsoft products are included in a Marketplace offering (e.g. Windows Server or SQL Server), such products are licensed by Microsoft and not by any third party.

Basics

Subscription	Azure subscription 1
Resource group	geniwaruni_Cloud_RG
Region	Australia East
Storage Account Name	geniwarunicloudstorage
Location	australiacentral

Previous Next Create

Microsoft.Template-202601171 | Overview

Deployment

Search

Delete Cancel Redeploy Download Refresh

Deployment is in progress

Deployment name : Microsoft.Template-20260117110453
 Subscription : Azure subscription 1
 Resource group : geniwaruni_Cloud_RG
 Start time : 1/17/2026, 11:04:49 AM
 Correlation ID : 7f6eb1c0-cf4d-4201-bdcb-f94e9aff5a0

Deployment details

Resource	Type	Status
geniwaruniclou...	Storage account	Accepted

Give feedback
[Tell us about your experience with deployment](#)

Microsoft Defender for Cloud
 Secure your apps and infrastructure
[Go to Microsoft Defender for Cloud >](#)

Add or remove favorites by pressing Ctrl+Shift+F

The screenshot shows the Microsoft Azure portal interface. At the top, the browser address bar displays the URL: <https://portal.azure.com/#view/HubsExtension/DeploymentDetailsBla...>. The page title is "Microsoft.Template-20260117110453 | Overview".

The main content area features a green checkmark icon and the heading "Your deployment is complete". Below this, the following deployment details are listed:

- Deployment name : Microsoft.Template-20260117110453
- Subscription : Azure subscription 1
- Resource group : geniwaruni_Cloud_RG
- Start time : 1/17/2026, 11:04:55 AM
- Correlation ID : 7f6eb1c0-cf4d-4201-bdcb-f94e9aff5a0

Below the details, there are two expandable sections: "Deployment details" and "Next steps". A blue button labeled "Go to resource" is positioned below the "Next steps" section.

At the bottom of the main content area, there is a "Give feedback" section with a link to "Tell us about your experience with deployment". Below that is a "Cost management" section with a green circular icon containing a dollar sign and the text "Get notified to stay within your budget and prevent unexpected charges on your bill. Set up cost alerts >".

The left sidebar contains a search bar and a navigation menu with the following items: Overview (selected), Inputs, Outputs, and Template.

At the bottom left of the page, there is a small note: "Add or remove favorites by pressing Ctrl+Shift+F".

The screenshot displays the Azure portal interface for a deployment overview. At the top, the browser address bar shows the URL: `https://portal.azure.com/#view/HubsExtension/DeploymentDetailsBla...`. The Azure portal header includes the 'Microsoft Azure' logo, a search bar, and the 'Copilot' button. The main content area is titled 'Microsoft.Template-20260117110453 | Overview' and shows a 'Deployment' status of 'Complete'. A green checkmark icon indicates success. The deployment details listed are:

- Deployment name: Microsoft.Template-20260117110453
- Subscription: Azure subscription 1
- Resource group: geniwaruni_Cloud_RG
- Start time: 1/17/2026, 11:04:55 AM
- Correlation ID: 7f6eb1c0-cf4d-4201-bdcb-f94e9aff5a0

Below the details, there are sections for 'Deployment details', 'Next steps', and a 'Go to resource' button. A 'Give feedback' section with a 'Tell us about your experience with deployment' link is also present. At the bottom, a 'Cost management' section offers to 'Set up cost alerts'.

Home > Resource Manager | Resource groups >

geniwaruni_Cloud_RG

Resource group

Search

+ Create Manage view Delete resource group Group by none

Overview

- Activity log
- Access control (IAM)
- Tags
- Resource visualizer
- Events
- Settings
- Cost Management
- Monitoring
- Automation
- Help

Essentials JSON View

Resources Recommendations

Filter for any field...

Type equals all Location equals all Add filter

<input type="checkbox"/>	Name ↑	Type	Location
<input type="checkbox"/>	geniwarunicloudstorage	Storage account	Australia Central

Showing 1 - 1 of 1. Display count: auto

Give feedback

Add or remove favorites by pressing Ctrl+Shift+F

The screenshot displays the Microsoft Azure portal interface for a storage account. The browser address bar shows the URL: <https://portal.azure.com/#@genistudent97gmail.onmicrosoft.com/res...>. The page title is "geniwarunicloudstorage" and it is identified as a "Storage account".

Navigation and Search: The top navigation bar includes "Microsoft Azure", a search bar for resources, services, and docs, and a "Copilot" button. The breadcrumb trail is: Home > Resource Manager | Resource groups > geniwaruni_Cloud_RG > geniwarunicloudstorage.

Storage Account Overview: The main content area shows the "Overview" tab selected in the left-hand navigation pane. The storage account name is "geniwarunicloudstorage". A notification banner at the top right states: "Storage is retiring support for TLS 1.0 and 1.1 starting Feb 3, 2026. If you are using any of these, please migrate to TLS 1.2. [Learn more](#)".

Essentials: A section titled "Essentials" provides key information:

- Resource group (move): [geniwaruni_Cloud_RG](#)
- Location: australiacentral
- Subscription (move): [Azure subscription 1](#)
- Subscription ID: 02ed530d-5bcf-442a-838c-b3f805b08091
- Disk state: Available
- Performance: Standard
- Replication: Locally-redundant storage (LRS)
- Account kind: StorageV2 (general purpose v2)
- Provisioning state: Succeeded
- Created: 1/17/2026, 11:04:56 AM

Tags: The "Tags" section shows two tags: "Department : Data" and "Project : Olympic2021".

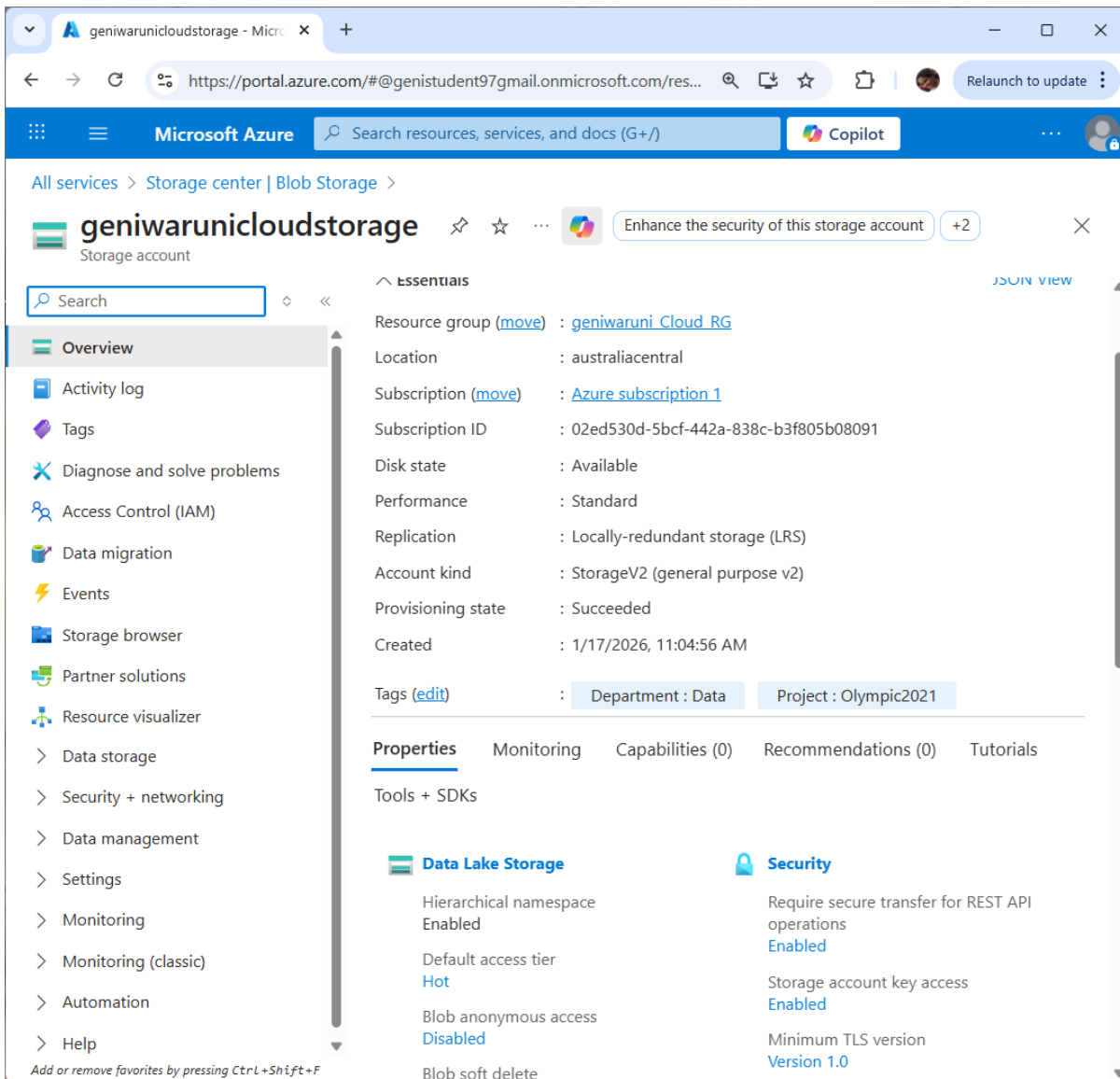
Properties and Tools: Below the essentials, there are tabs for "Properties", "Monitoring", "Capabilities (5)", "Recommendations (0)", and "Tutorials". Under the "Properties" tab, the "Data Lake Storage" section is expanded, showing "Hierarchical namespace Enabled". The "Security" section is also visible, with "Require secure transfer for REST API operations" listed.

Left Navigation Pane: The navigation pane includes the following items:

- Overview
- Activity log
- Tags
- Diagnose and solve problems
- Access Control (IAM)
- Data migration
- Events
- Storage browser
- Partner solutions
- Resource visualizer
- Data storage
- Security + networking
- Data management
- Settings
- Monitoring
- Monitoring (classic)
- Automation
- Help

At the bottom of the navigation pane, it says: "Add or remove favorites by pressing Ctrl+Shift+F".

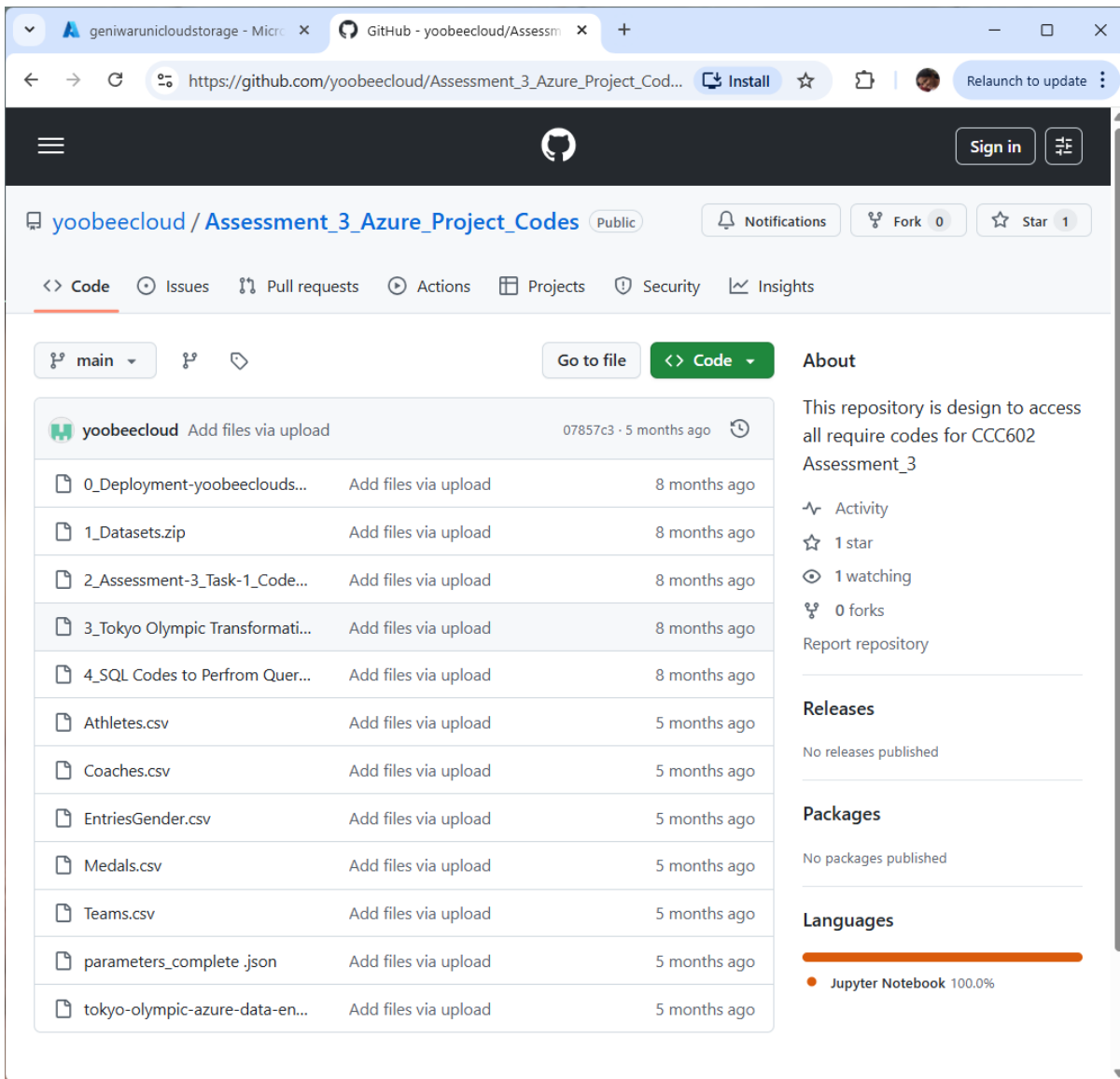
The screenshot displays the Microsoft Azure portal interface. At the top, the browser address bar shows the URL: `https://portal.azure.com/#@genistudent97gmail.onmicrosoft.com/res...`. The Azure portal header includes the 'Microsoft Azure' logo, a search bar, and the 'Copilot' icon. The main content area is titled 'All services > Storage center | Blob Storage >' and features the storage account name 'geniwarunicloudstorage' with a 'Storage account' subtitle. A notification banner at the top right states: 'Enhance the security of this storage account +2'. Below this, a search bar and a set of action buttons (Upload, Open in Explorer, Delete, Move, Refresh) are visible. A warning message indicates: 'Storage is retiring support for TLS 1.0 and 1.1 starting Feb 3, 2026. If you are using any of these, please migrate to TLS 1.2. Learn more'. The 'Essentials' section lists key properties: Resource group (geniwaruni_Cloud_RG), Location (australiacentral), Subscription (Azure subscription 1), Subscription ID (02ed530d-5bcf-442a-838c-b3f805b08091), Disk state (Available), Performance (Standard), Replication (Locally-redundant storage (LRS)), Account kind (StorageV2 (general purpose v2)), Provisioning state (Succeeded), and Created (1/17/2026, 11:04:56 AM). The 'Tags' section shows 'Department : Data' and 'Project : Olympic2021'. Below the Essentials section, there are tabs for 'Properties', 'Monitoring', 'Capabilities (0)', 'Recommendations (0)', and 'Tutorials'. The 'Data Lake Storage' section shows 'Hierarchical namespace Enabled', and the 'Security' section shows 'Require secure transfer for REST API operations'. A footer note at the bottom left reads: 'Add or remove favorites by pressing Ctrl+Shift+F'.



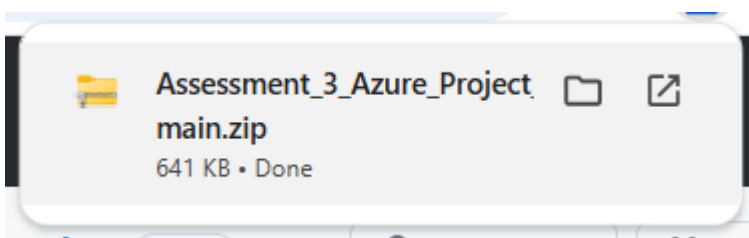
Using an ARM template for this task helped me better understand the importance of Infrastructure as Code, especially how it ensures consistent deployments and reduces manual configuration errors.

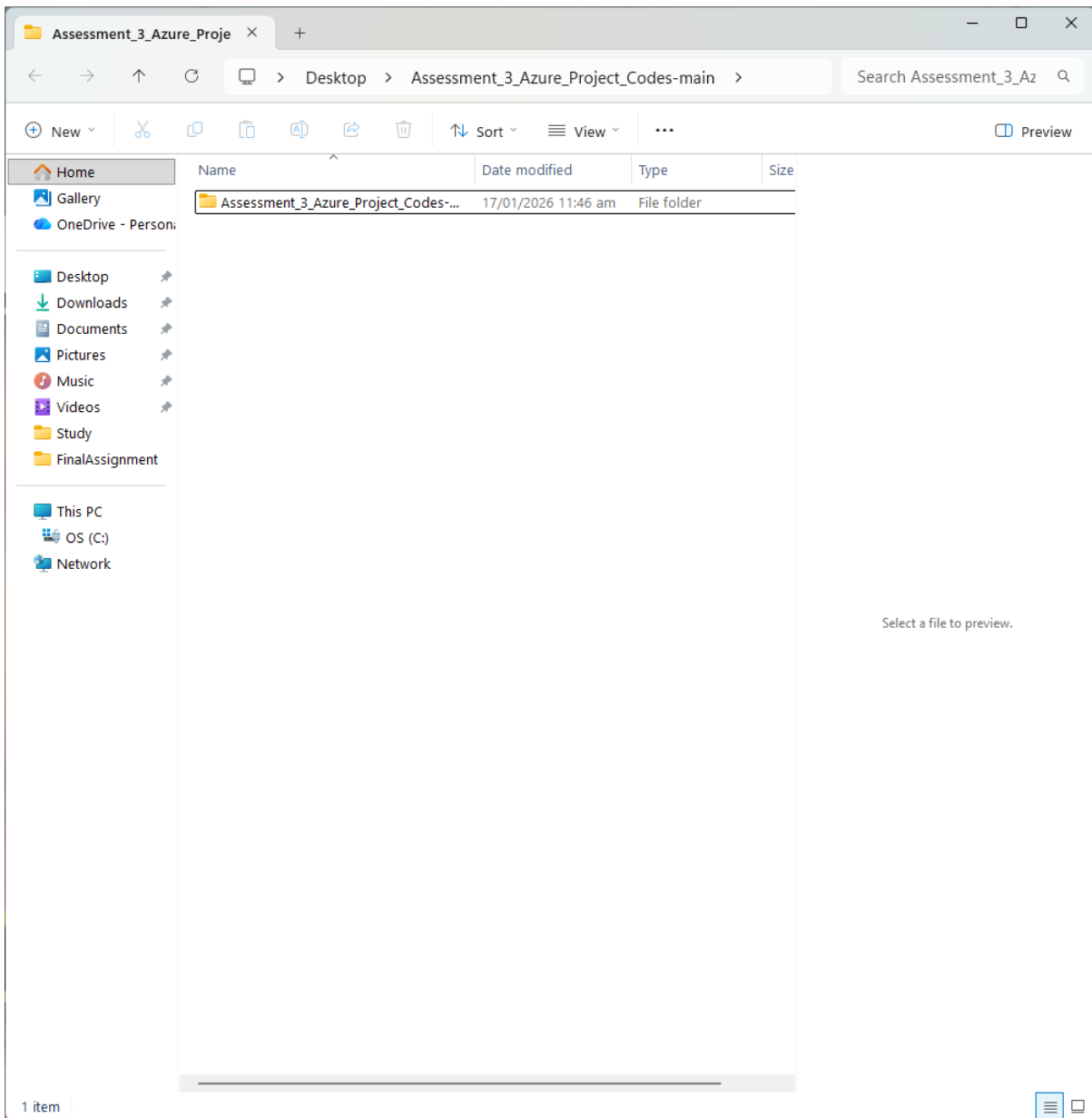
Tokyo Olympic Dataset Acquisition

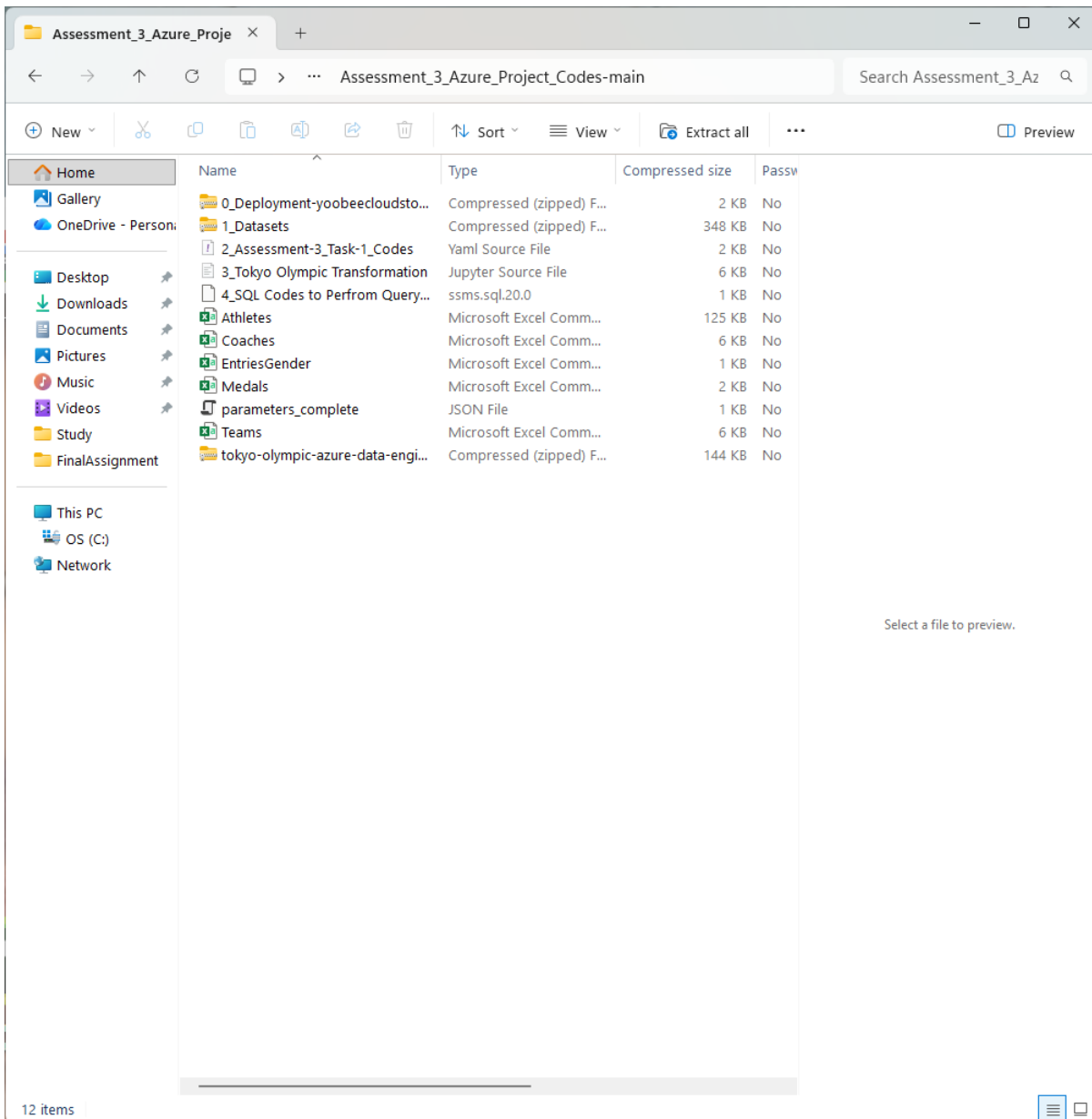
- The Tokyo Olympic dataset was downloaded from the GitHub repository provided by the tutor. The dataset contains multiple files related to athletes, medals, teams, coaches, and gender participation.
- The files were reviewed to understand their structure, file formats, and key fields before ingestion into Azure.



File Name	File Type	Description
athletes	CSV	Athlete details and country information
medals	CSV	Medal counts by country
teams	CSV	Team and discipline details
coaches	CSV	Coach information
entriesgender	CSV	Gender participation by discipline







Task 3: Data Ingestion Using Azure Data Factory

Data Ingestion Pipeline Design

- Azure Data Factory was used to automate the ingestion of data from Azure Data Lake Storage into the /raw folder. The Copy Data activity was configured using ADLS Gen2 as both the source and sink to preserve the original data format.
- The pipeline ensured that data ingestion was repeatable and scalable.

The screenshot shows the 'Create Data Factory' wizard in the Microsoft Azure portal. The 'Basics' tab is active, and the following information is visible:

- Project details:**
 - Subscription: Azure subscription 1
 - Resource group: (New) geni_Cloud_RG (with a 'Create new' link below)
- Instance details:**
 - Name: geniADF (with a green checkmark)
 - Region: Australia East
 - Version: V2
- Notification:** A blue banner recommends upgrading to Fabric Data Factory for a modern, unified data integration experience, with a link to 'Try Fabric instead!'.
- Navigation:** At the bottom, there are buttons for 'Previous', 'Next', and 'Review + create'.

The screenshot shows a web browser window with the URL <https://portal.azure.com/#create/Microsoft.DataFactory>. The page title is "Create Data Factory" and the current step is "Git configuration". The navigation tabs include "Basics", "Git configuration" (selected), "Networking", "Advanced", "Tags", and "Review + create". The main content area contains the following text: "Azure Data Factory allows you to configure a Git repository with either Azure DevOps or GitHub. Git is a version control system that allows for easier change tracking and collaboration. [Learn more about Git integration in Azure Data Factory](#)". Below this is a checkbox labeled "Configure Git later" which is checked. At the bottom of the page, there are three buttons: "Previous", "Next", and "Review + create" (highlighted in blue). A "Give feedback" link is also present in the bottom right corner.

The screenshot shows the 'Create Data Factory' wizard in the Microsoft Azure portal, specifically the 'Review + create' step. The browser address bar shows the URL <https://portal.azure.com/#create/Microsoft.DataFactory>. The page title is 'Create Data Factory'. The navigation tabs include 'Basics', 'Git configuration', 'Networking', 'Advanced', 'Tags', and 'Review + create' (which is the active tab). A link for 'View automation template' is visible. The configuration details are as follows:

Basics	
Subscription	Azure subscription 1
Resource group	geni_Cloud_RG
Name	geniADF
Region	Australia East
Version	V2

Networking	
Connect via	Public endpoint

At the bottom, there are three buttons: 'Previous', 'Next', and 'Create'. The 'Create' button is highlighted in blue. A 'Give feedback' link is located in the bottom right corner.

Microsoft Azure portal interface for creating a Data Factory. The page title is "Create Data Factory". The breadcrumb navigation shows "Home >". The main heading is "Create Data Factory". Below the heading are tabs for "Basics", "Git configuration", "Networking", "Advanced", "Tags", and "Review + create" (which is the active tab). A link "View automation template" is present. The configuration details are as follows:

Basics	
Subscription	Azure subscription 1
Resource group	geni_Cloud_RG
Name	geniADF
Region	Australia East
Version	V2

Networking	
Connect via	Public endpoint

At the bottom, there are buttons for "Previous", "Next", and "Create". A "Give feedback" link is also visible. A notification bubble in the top right corner reads: "Initializing deployment... Initializing template deployment to resource group 'geni_Cloud_RG'."

The screenshot displays the Microsoft Azure portal interface. At the top, the browser address bar shows the URL: `https://portal.azure.com/#view/HubsExtension/DeploymentDetailsBlade...`. The Azure portal header includes the 'Microsoft Azure' logo, a search bar, and the 'Copilot' icon.

The main content area is titled 'Microsoft.DataFactory-20260117121925 | Overview' under the 'Deployment' section. A notification banner at the top right states: 'Deployment succeeded' with a green checkmark icon, indicating that the deployment of 'Microsoft.DataFactory-20260117121925' to resource group 'geni_Cloud_RG' was successful. Below this banner are buttons for 'Pin to dashbo...' and 'Go to resource gr...'.

The left sidebar contains navigation options: 'Overview' (selected), 'Inputs', 'Outputs', and 'Template'. The main content area shows the deployment status as 'Deployment is in progress' (despite the success message above). The deployment details are as follows:

- Deployment name: Microsoft.DataFactory-20260117121925
- Subscription: Azure subscription 1
- Resource group: geni_Cloud_RG
- Start time: 1/17/2026, 12:27:40 PM
- Correlation ID: 26f4183a-10ab-4daa-9195-f3d601c3904a

Below the deployment details is a table with the following data:

Resource	Type	Status	Op
geniADF	Data factory (V2)	OK	Op

At the bottom of the page, there is a 'Give feedback' section with a link to 'Tell us about your experience with deployment' and a 'Microsoft Defender for Cloud' section with a link to 'Go to Microsoft Defender for Cloud >'. A footer note at the bottom left reads: 'Add or remove favorites by pressing Ctr+L+Shift+F'.

The screenshot displays the Microsoft Azure portal interface. At the top, the browser address bar shows the URL: <https://portal.azure.com/#@genistudent97gmail.onmicrosoft.com/resou...>. The page title is "geniADF - Microsoft Azure".

The main content area shows the configuration for a "Data factory (V2)" named "geniADF". The "Essentials" section provides the following details:

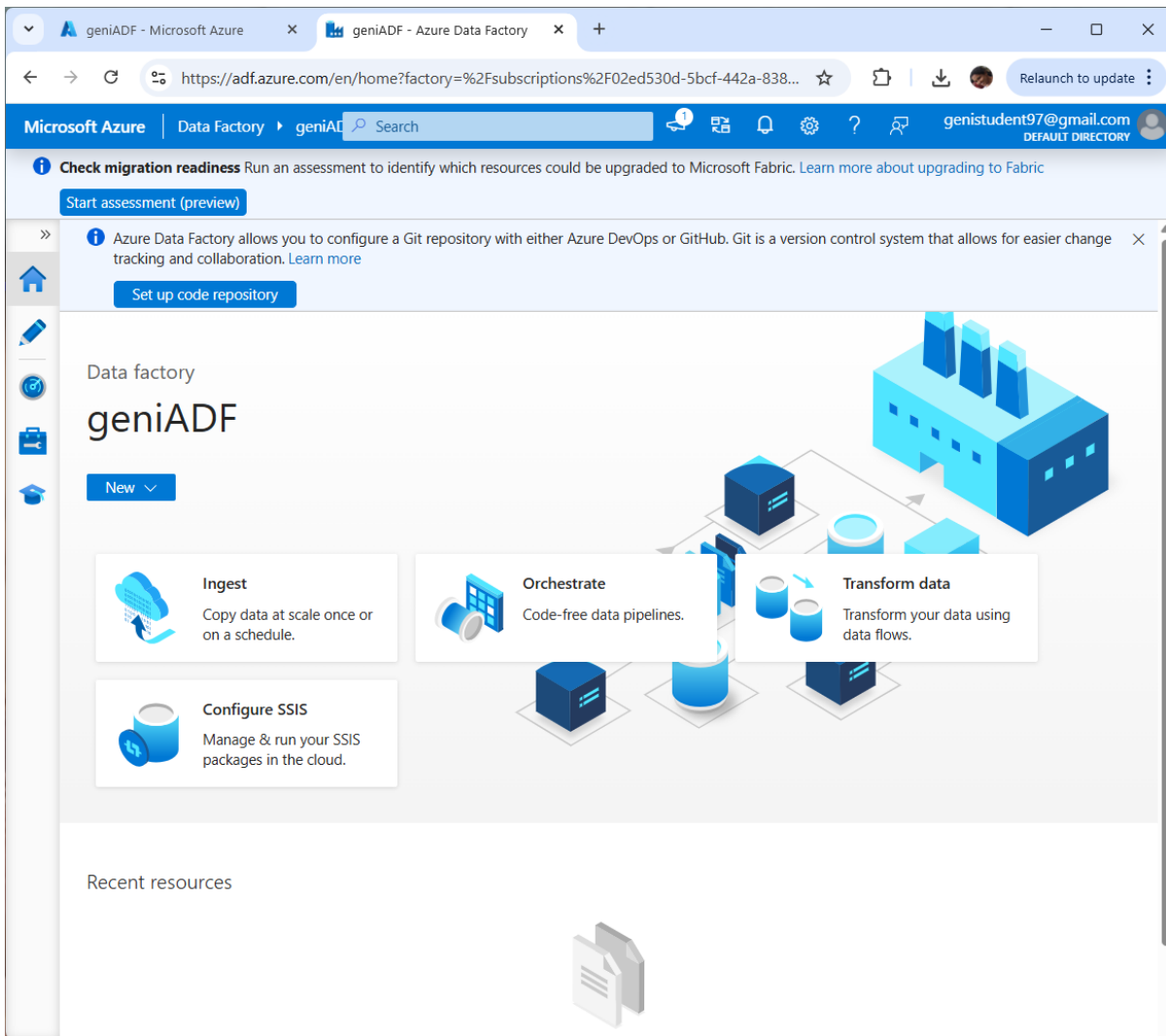
- Resource group: [geni_Cloud_RG](#) (move)
- Type: Data factory (V2)
- Status: Succeeded
- Location: Australia East
- Subscription: [Azure subscription 1](#) (move)
- Subscription ID: 02ed530d-5bcf-442a-838c-b3f805b08091

Below the essentials, there is a large blue icon representing a factory and the text "Azure Data Factory Studio". A prominent blue button labeled "Launch studio" is centered below this text.

At the bottom of the page, there are four tiles for additional resources:

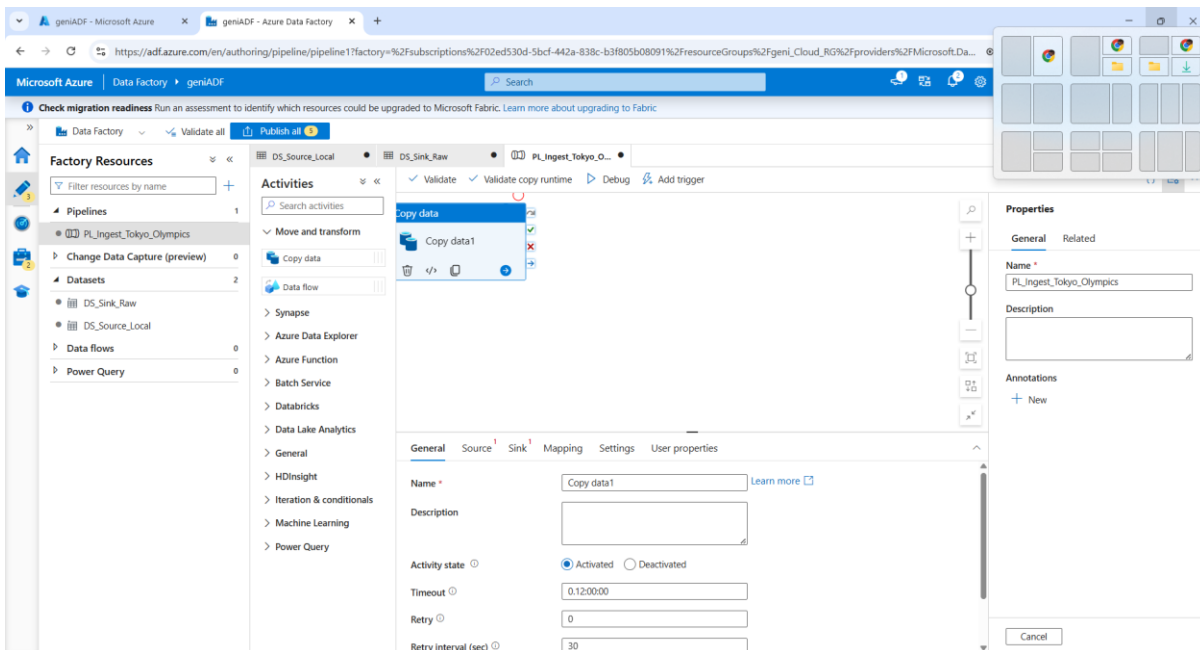
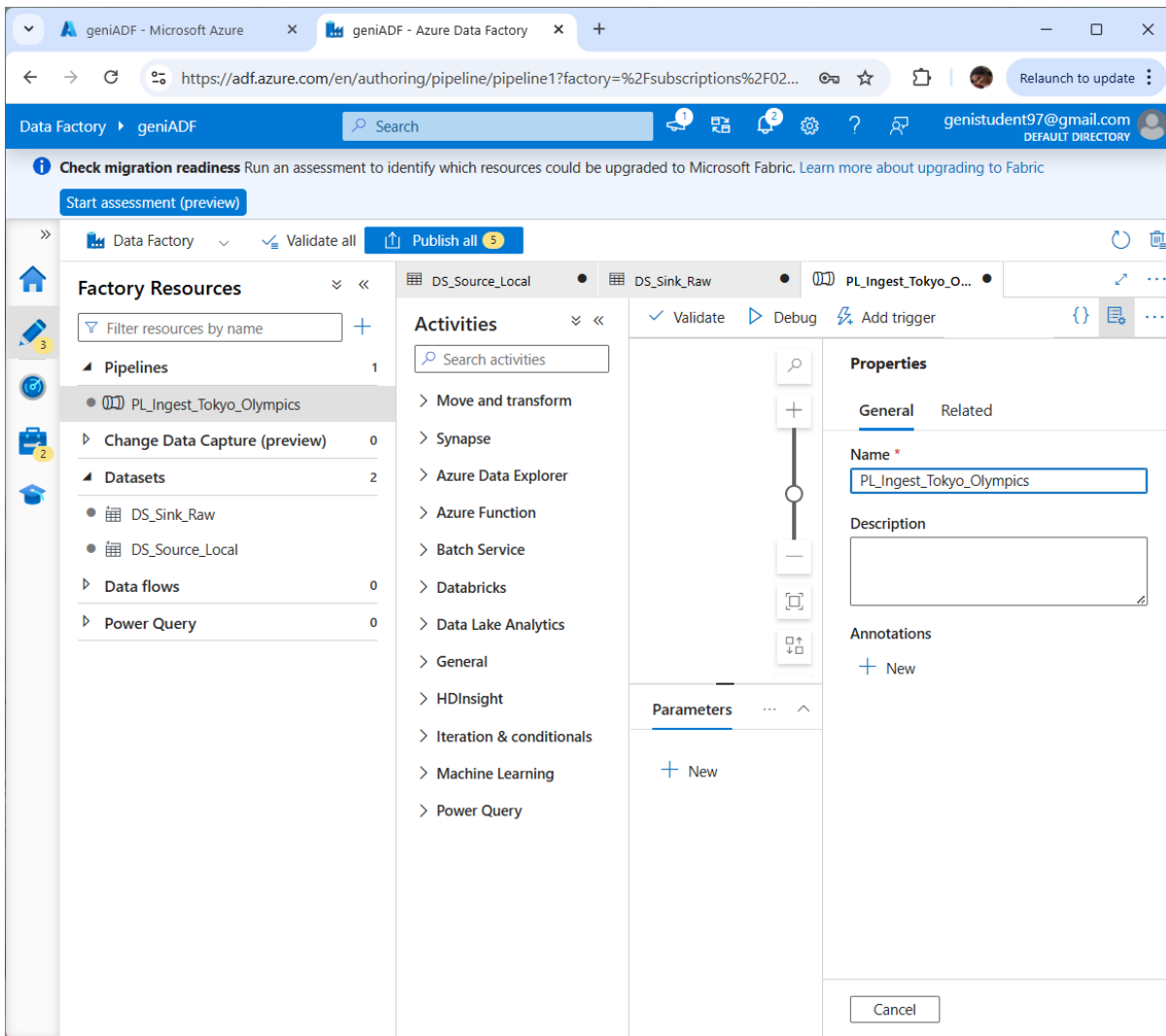
- Quick Starts (with a lightning bolt icon)
- Tutorials (with a document icon labeled "101")
- Template Gallery (with a document icon)
- Training Modules (with a document icon and a ribbon)

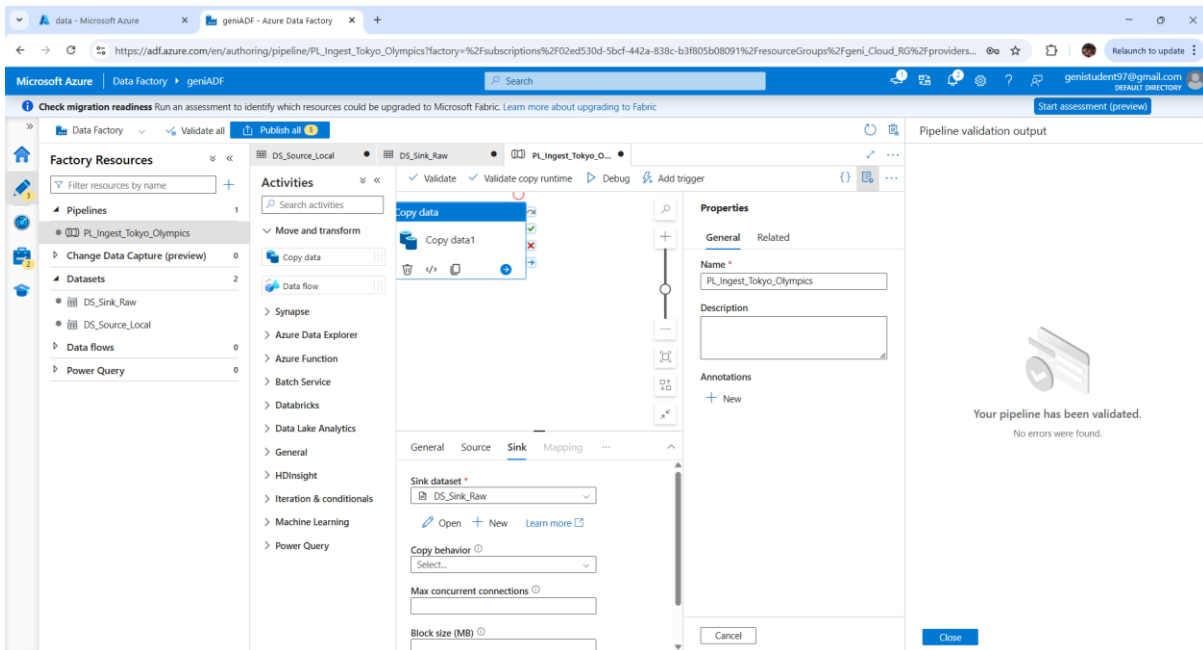
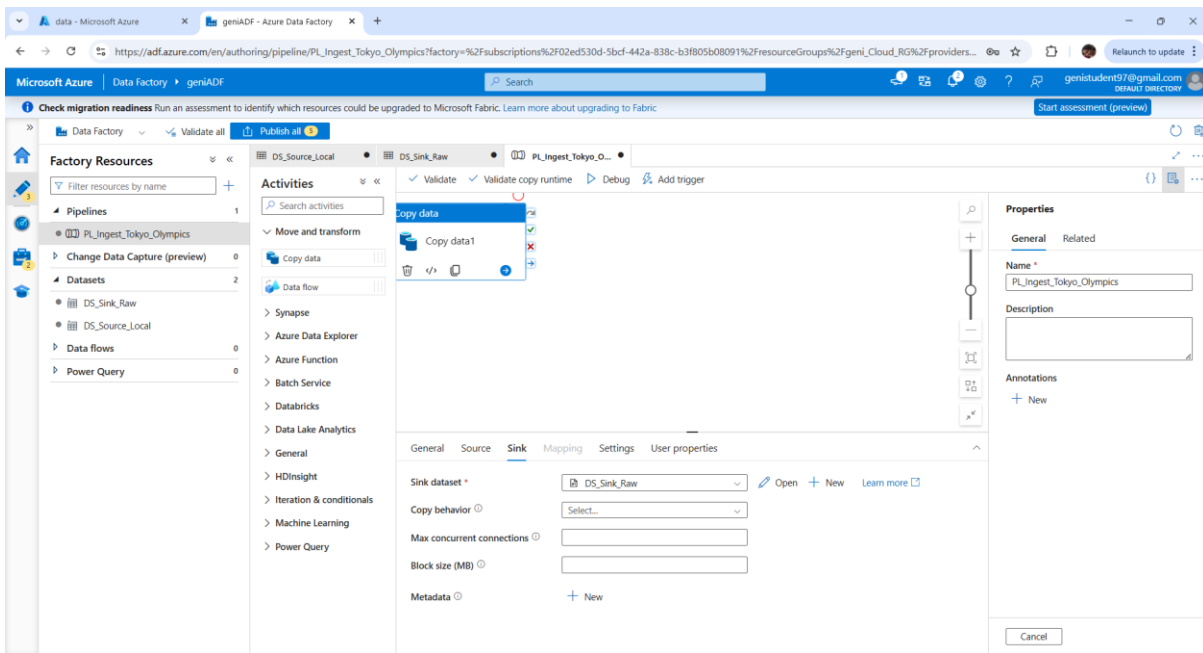
A footer note at the bottom left reads: "Add or remove favorites by pressing Ctrl+L+Shift+F".



Source and Sink Configuration

- The source dataset pointed to the /source folder in the Data Lake, while the sink dataset wrote the ingested data to the /raw folder. Linked services were configured using managed identity authentication.





Pipeline Validation and Verification

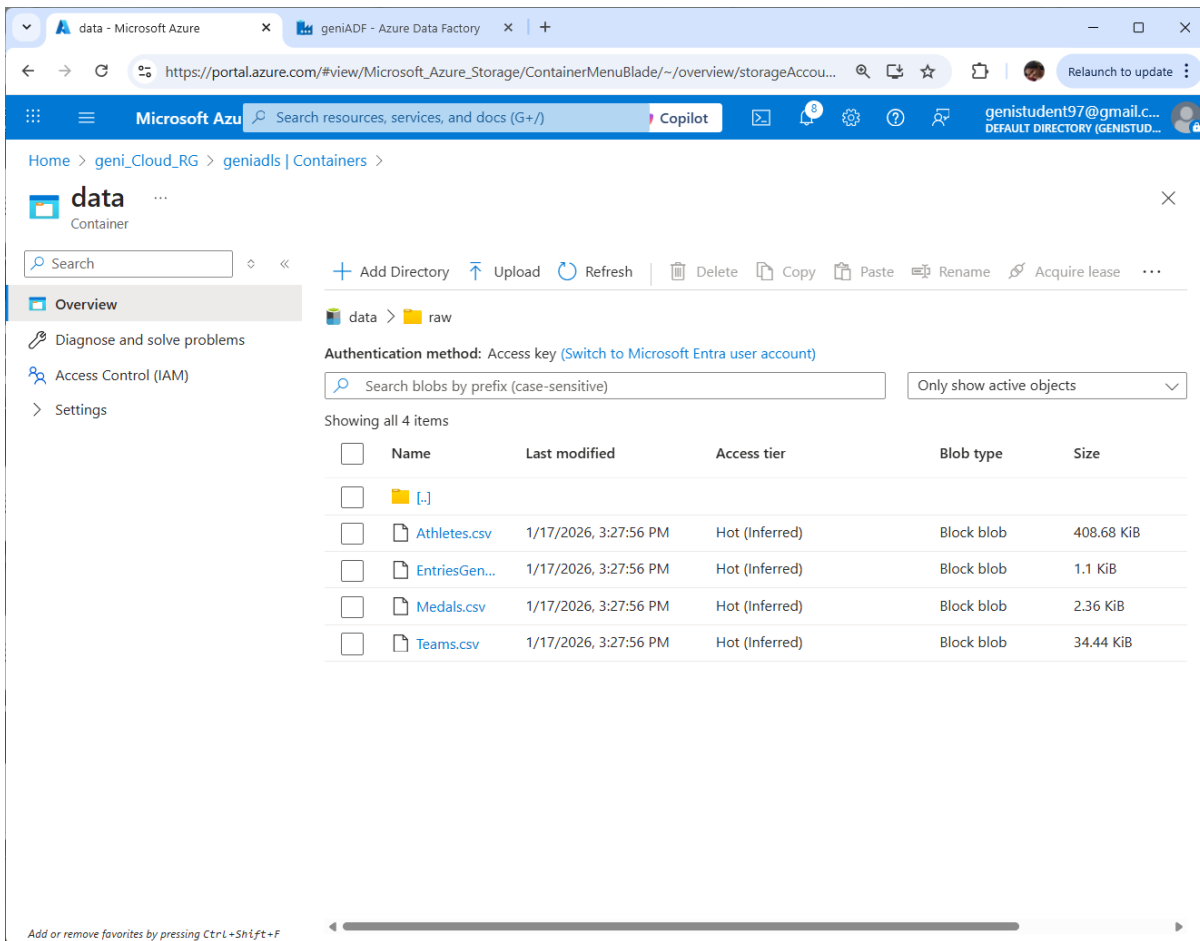
- The pipeline was validated and executed successfully. After execution, the /raw folder was checked to confirm that all files were ingested correctly.

The screenshot shows the Microsoft Azure Data Factory 'Authoring' view for a pipeline named 'PL_Ingest_Tokyo_Olympics'. The 'Activities' pane shows a 'Copy data' activity with a green checkmark, indicating it has completed successfully. The 'Properties' pane on the right shows the activity's name and description. Below the activity, the 'Output' tab displays the pipeline run ID and a table of activity results.

Activity name	Activity st...	Activit...	Run start	Duration	Integration runtime
Copy data1	✓ Succeeded	Copy data	1/17/2026, 3:27:42 PM	16s	AutoResolveIntegra

The screenshot shows the 'Monitoring' view for pipeline runs in the same Data Factory. The 'Pipeline runs' section displays a table with columns for Pipeline name, Run start, Run end, Duration, Status, Triggered by, and Run ID. The table shows one successful run for the 'PL_Ingest_Tokyo_Olympics' pipeline.

Pipeline name	Run start	Run end	Duration	Status	Triggered by	Run ID
PL_Ingest_Tokyo_Olympics	1/17/2026, 3:27:36 PM	1/17/2026, 3:27:58 PM	23s	✓ Succeeded	Manual trigger	1ae6083c-355b-40cd-bcad...



Pipeline Summary Report

Purpose of the Pipeline

- The main purpose of the Azure Data Factory pipeline was to automate the ingestion of the Tokyo Olympic dataset into Azure Data Lake Storage Gen2 in a structured and repeatable way. Instead of manually uploading files every time, the pipeline ensures that raw data is consistently moved into the /raw folder, where it can later be used for transformation and analysis. This approach follows real-world data engineering practices, where raw data must be preserved in its original format before any processing is applied.
- The pipeline also supports scalability. If more Olympic datasets or updated files are added in the future, the same pipeline can be reused without redesigning the entire process. This reduces manual effort and helps maintain data integrity across the project.

Data Flow (Source to Sink)

- The data flow begins with the Tokyo Olympic dataset stored in Azure Data Lake Storage Gen2. Initially, the files were uploaded manually into a /source folder in the data container. Azure Data Factory was then configured to use this folder as the source location.
- Within the pipeline, a **Copy Data** activity was used. The source dataset pointed to the ADLS Gen2 /source directory, and the sink dataset wrote the data into the /raw directory within the same storage account. Using ADLS Gen2 for both source and sink ensured that the data

transfer was secure, efficient, and compatible with downstream services such as Azure Databricks and Azure Synapse Analytics.

- The pipeline was validated and tested using the Debug option. After successful execution, the /raw folder was checked in the Azure Portal to confirm that all files were transferred correctly. This verification step was important to ensure that the data ingestion process worked as expected before moving to the transformation stage.

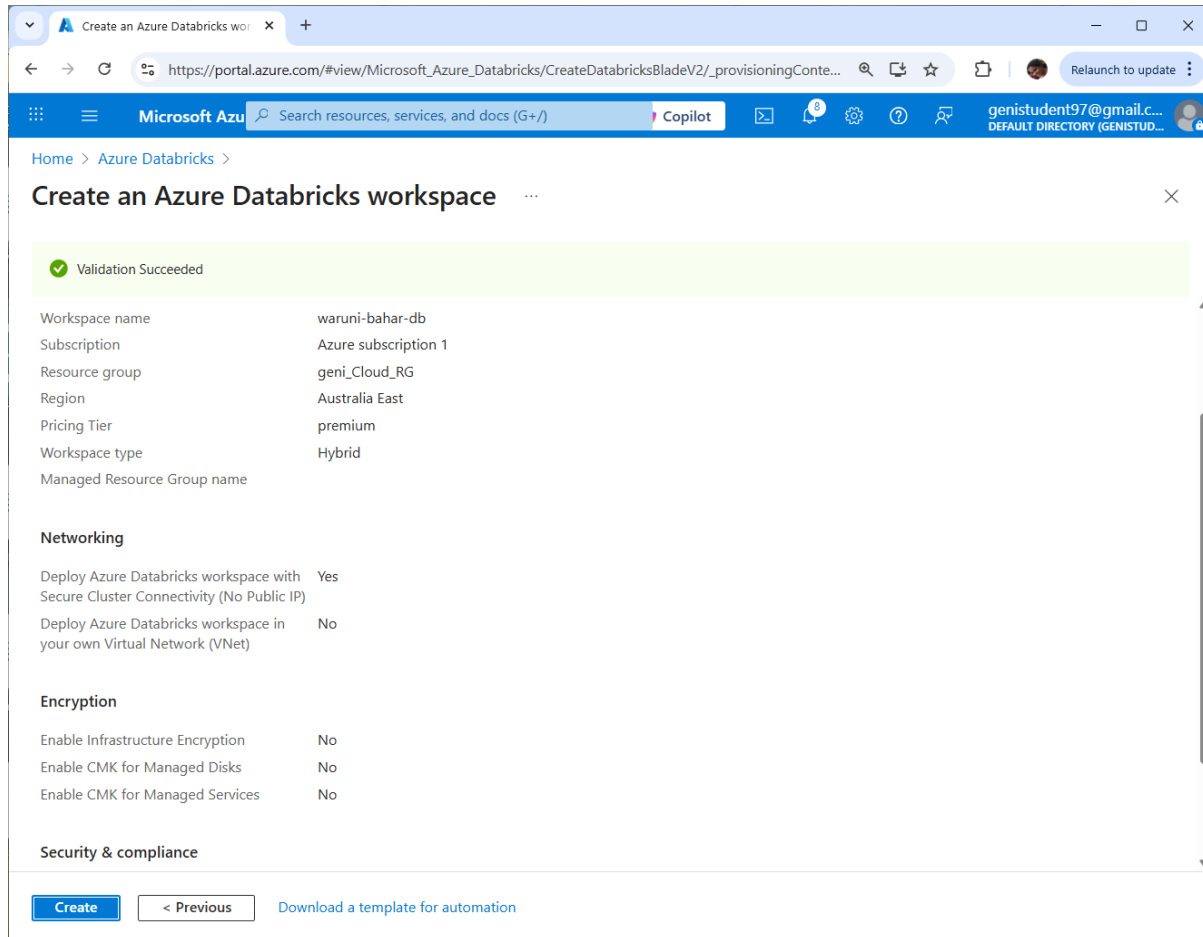
Challenges and Observations During Implementation

- Several challenges were encountered during the implementation of the pipeline, which helped improve understanding of real-world cloud data engineering issues.
- One of the main challenges was related to **permissions and authentication**. Initially, the pipeline failed with “Forbidden” and “AuthorizationPermissionMismatch” errors. This occurred because Azure Data Factory did not have sufficient permissions to access the Data Lake Storage account. The issue was resolved by enabling the system-assigned managed identity for Data Factory and assigning the **Storage Blob Data Contributor** role at the storage account level.
- Another challenge was confusion around using a **local file system dataset**. At first, a local (File System) dataset was configured as the source, which resulted in connection errors because Azure Data Factory cannot directly access files from a personal computer without a self-hosted integration runtime. This was resolved by uploading the files to ADLS Gen2 and using an ADLS-based dataset instead.
- Additionally, sink configuration errors occurred due to missing parameters and file paths. In some cases, leaving the file name empty in the sink dataset caused null reference errors during pipeline execution. Carefully reviewing dataset settings and explicitly defining directory paths helped resolve this issue.
- Overall, these challenges highlighted the importance of correct identity configuration, understanding how cloud services access data, and carefully validating pipeline settings. Resolving these issues provided valuable hands-on experience that closely reflects real industry scenarios.

Task 4: Create an Azure Databricks Workspace and perform data transformation using code to automate the process.

Azure Databricks Workspace and Cluster Setup

- An Azure Databricks workspace was created within the same resource group. A single-node cluster was configured with auto-termination enabled to minimise costs.



The screenshot shows the Microsoft Azure portal interface. The browser address bar displays the URL: <https://portal.azure.com/#view/HubsExtension/DeploymentDetailsBlade/~/overview/id/%2Fsubscriptions...>. The user is logged in as 'genistudent97@gmail.c...' with the role 'DEFAULT DIRECTORY (GENISTUD...)'.

The main content area is titled 'geni_Cloud_RG_waruni-bahar-db | Overview' and is categorized as a 'Deployment'. A navigation sidebar on the left includes 'Overview' (selected), 'Inputs', 'Outputs', and 'Template'. At the top of the main area, there are action buttons: 'Delete', 'Cancel', 'Redeploy', 'Download', and 'Refresh'.

The primary message is 'Deployment is in progress'. Below this, the following deployment metadata is listed:

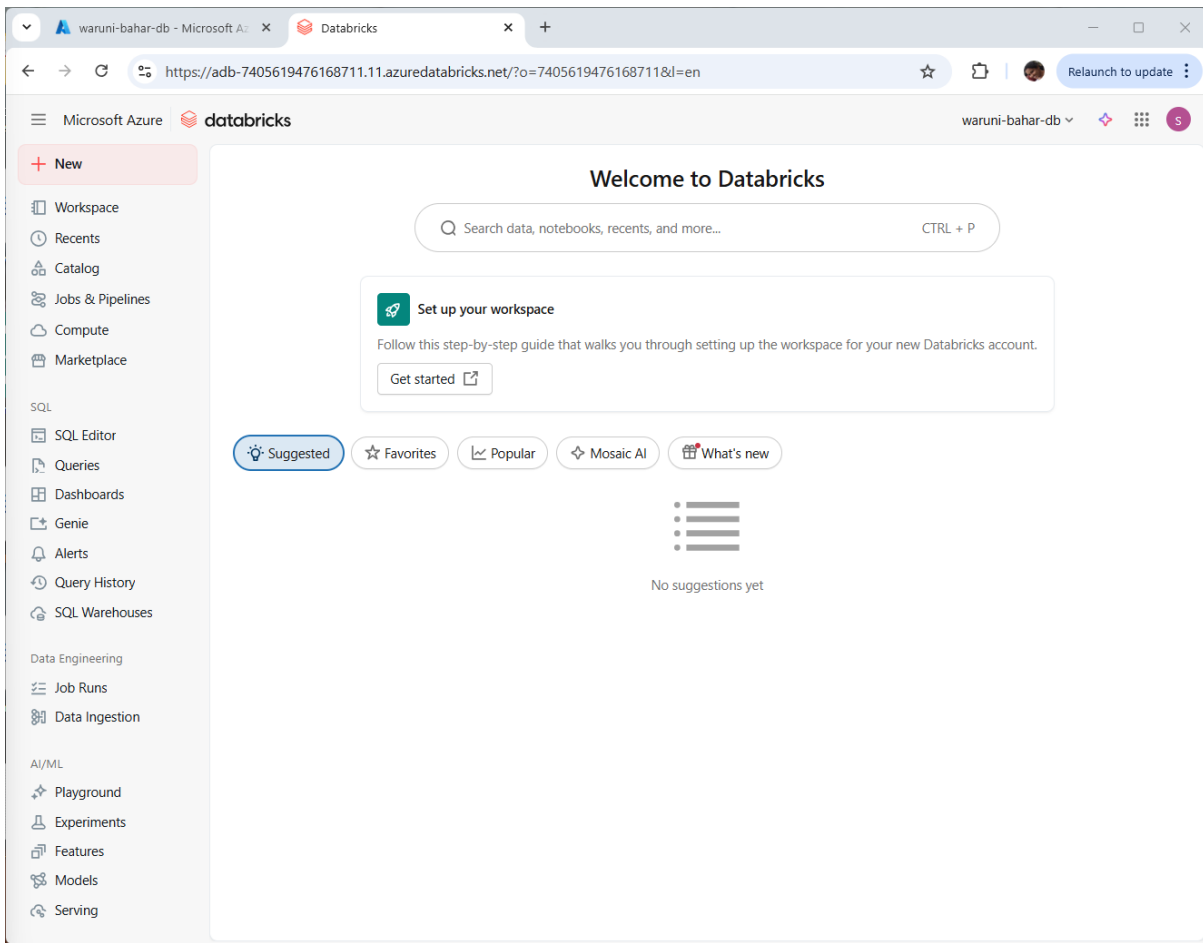
- Deployment name : geni_Cloud_RG_waruni-bahar-db
- Subscription : Azure subscription 1
- Resource group : geni_Cloud_RG
- Start time : 1/17/2026, 3:37:29 PM
- Correlation ID : e9633b2e-d0c9-48cb-8ae8-cbfc94aba20

Under the 'Deployment details' section, a table provides a summary of the resources:

Resource	Type	Status	Operation details
waruni-bahar-db	Azure Databricks Service	Created	Operation details

At the bottom of the page, there are two informational sections: 'Microsoft Defender for Cloud' with the text 'Secure your apps and infrastructure' and a link to 'Go to Microsoft Defender for Cloud >', and 'Free Microsoft tutorials'.

The screenshot shows the Microsoft Azure portal interface. At the top, the browser address bar displays the URL: <https://portal.azure.com/#view/HubsExtension/DeploymentDetailsBlade/~/overview/id/%2Fsubscriptions...>. The page title is "geni_Cloud_RG_waruni-bahar-db | Overview". The main content area features a large green checkmark icon and the text "Your deployment is complete". Below this, the deployment details are listed: Deployment name: geni_Cloud_RG_waruni-bahar-db, Subscription: Azure subscription 1, Resource group: geni_Cloud_RG, Start time: 1/17/2026, 3:37:36 PM, and Correlation ID: e9633b2e-d0c9-48cb-8ae8-cbfc94aba20. A "Go to resource" button is visible. On the left, a navigation pane includes "Overview", "Inputs", "Outputs", and "Template". At the bottom, there is a "Cost management" section with a "Set up cost alerts" link. The footer of the page contains the text: "Add or remove favorites by pressing Ctrl+Shift+F".



The screenshot displays the 'Create new compute' page in the Databricks interface. The browser address bar shows the URL: `https://adb-7405619476168711.11.azuredatabricks.net/compute/clusters/new?o=7405619476168711`. The page title is 'Create new compute'.

Summary: 16 GB Memory, 4 Cores, 1.5 DBU/h. Data access is set to 'Unity Catalog'.

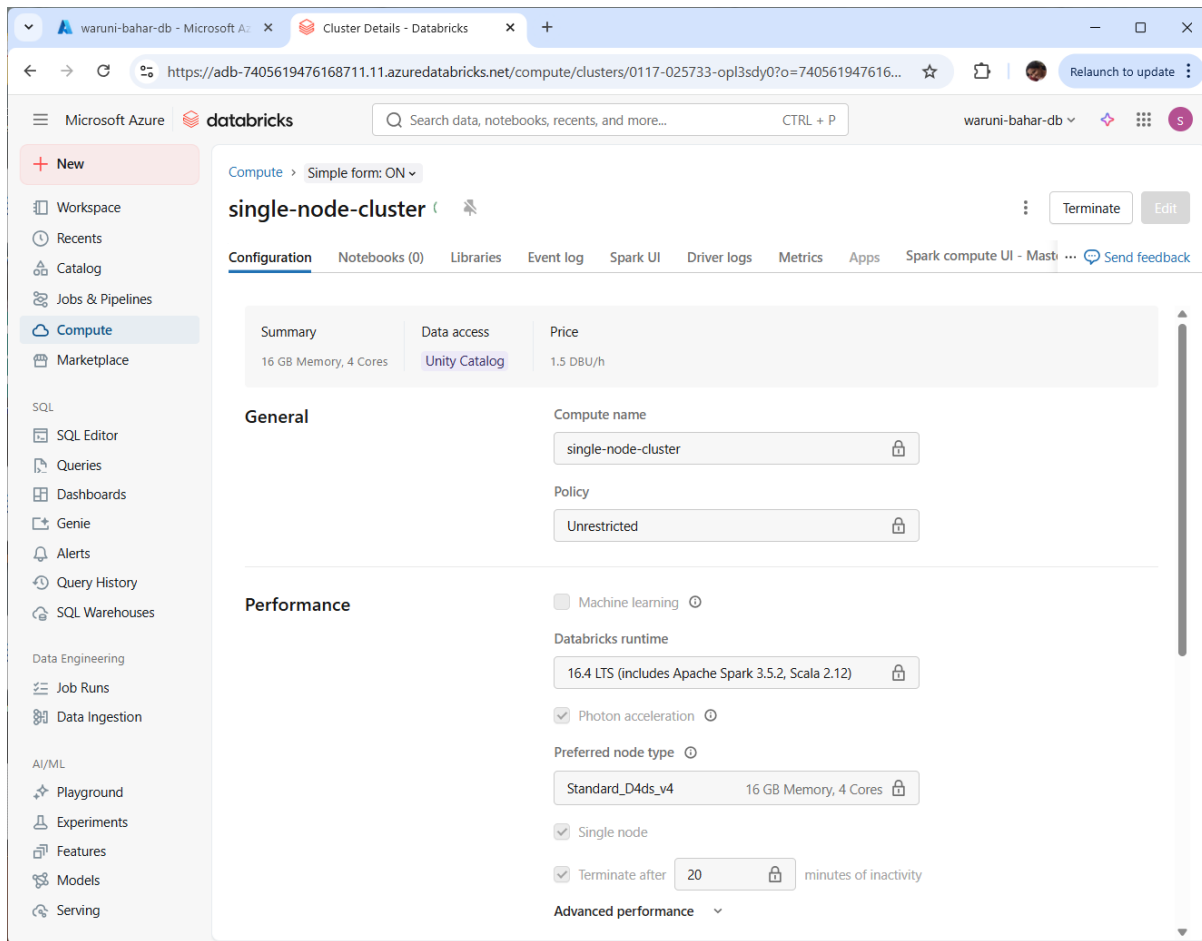
General:

- Compute name:
- Policy:

Performance:

- Machine learning
- Databricks runtime: Scala 2.12, Spark 3.5.2
- Photon acceleration
- Preferred node type: 16 GB Memory, 4 Cores
- Single node
- Terminate after minutes of inactivity
- Advanced performance:

Buttons:



Secure Data Access and Authentication

- An Azure App Registration was created to allow Databricks to securely access Azure Data Lake Storage Gen2. The app was granted the Storage Blob Data Contributor role.

The screenshot shows the 'Register an application' page in the Azure portal. The browser address bar shows the URL: `https://portal.azure.com/#view/Microsoft_AAD_RegisteredApps/CreateApplicationBlade/isMSAApp~/false`. The page title is 'Register an application'. Below the title, there is a text input field for the user-facing display name, which contains 'geni-app'. Underneath, there is a section for 'Supported account types' with four radio button options: 'Accounts in this organizational directory only (Default Directory only - Single tenant)', 'Accounts in any organizational directory (Any Microsoft Entra ID tenant - Multitenant)', 'Accounts in any organizational directory (Any Microsoft Entra ID tenant - Multitenant) and personal Microsoft accounts (e.g. Skype, Xbox)', and 'Personal Microsoft accounts only'. The first option is selected. Below this is a 'Redirect URI (optional)' section with a text input field containing 'e.g. https://example.com/auth'. At the bottom, there is a 'Register' button and a link to 'Microsoft Platform Policies'.

The screenshot shows the Microsoft Azure portal interface. At the top, there are browser tabs for 'geni-app - Microsoft Azure' and 'Cluster Details - Databricks'. The address bar shows the URL: https://portal.azure.com/#view/Microsoft_AAD_RegisteredApps/ApplicationMenuBlade/~/Overview/appId.... The Microsoft Azure logo and search bar are visible. The user's profile is 'genistudent97@gmail.c...' with the role 'DEFAULT DIRECTORY (GENISTUD...)'.

The main content area is titled 'geni-app' and includes a search bar and action buttons: 'Delete', 'Endpoints', and 'Preview features'. A notification banner reads: 'Got a second? We would love your feedback on Microsoft identity platform (previously Azure AD for developer). →'. Below this is an 'Essentials' section with the following details:

Display name geni-app	Client credentials Add a certificate or secret
Application (client) ID c9e3a213-3eef-4c0e-ad92-397b8680afc0	Redirect URIs Add a Redirect URI
Object ID dff9fdb7-c495-48c2-88fb-8291c6bfff733	Application ID URI Add an Application ID URI
Directory (tenant) ID 051bc474-1bf0-4155-8e6c-11ae24bf070b	Managed application in local directory geni-app
Supported account types My organization only	

Below the essentials section are two informational messages:

- Welcome to the new and improved App registrations. Looking to learn how it's changed from App registrations (Legacy)? [Learn more](#)
- Starting June 30th, 2020 we will no longer add any new features to Azure Active Directory Authentication Library (ADAL) and Azure Active Directory Graph. We will continue to provide technical support and security updates but we will no longer provide feature updates. Applications will need to be upgraded to Microsoft Authentication Library (MSAL) and Microsoft Graph. [Learn more](#)

At the bottom, there are links for 'Get Started' and 'Documentation'. A footer note says: 'Add or remove favorites by pressing Ctrl+Shift+F'.

The screenshot shows the Microsoft Azure portal interface. The main content area displays the 'Add a client secret' dialog for the application 'geni-app'. The dialog has two input fields: 'Description' with the value 'databricks-secret' and 'Expires' with a dropdown menu set to 'Recommended: 180 days (6 months)'. At the bottom of the dialog are 'Add' and 'Cancel' buttons.

The left-hand navigation pane is visible, showing the 'Certificates & secrets' section selected under the 'Manage' category. The main content area also shows a list of client secrets, which is currently empty, with a 'Description' header and a message stating 'No client secrets have been created'.

Home > Storage center | Blob Storage > dbstoragekrxrmwosxhkw | Access Control (IAM) >

Add role assignment

Role **Members** Conditions Review + assign

Selected role Storage Blob Data Contributor

Assign access to

- User, group, or service principal
- Managed identity

Members [+ Select members](#)

Name	Object ID	Type
geni-app	c6bd6402-903b-4019-bd8e-68c24...	App

Description

Optional

[Review + assign](#) [Previous](#) [Next](#) [Feedback](#)

Browser tabs: Add role assignment - Microsoft | Cluster Details - Databricks

URL: https://portal.azure.com/#view/Microsoft_Azure_AD/AddRoleAssignmentsLandingBlade/scope/%2Fsubscr...

Microsoft Azure Search resources, services, and docs (G+/) Copilot genistudent97@gmail.com DEFAULT DIRECTORY (GENISTUD...

Home > Storage center | Blob Storage > dbstoragekrxrmwosxhkw | Access Control (IAM) >

Add role assignment

Role Members Conditions **Review + assign**

Role Storage Blob Data Contributor

Scope /subscriptions/02ed530d-5bcf-442a-838c-b3f805b08091/resourceGroups/databricks-rg-waruni-bahar-db-xbpfnrfvq5r6/providers/Microsoft.Storage/storageAccounts/dbstoragekrxrmwosxhkw

Members

Name	Object ID	Type
geni-app	c6bd6402-903b-4019-bd8e-68c241a0c660	App

Description No description

Condition None

Buttons: Review + assign Previous Next Feedback

The screenshot shows the Microsoft Azure portal interface. At the top, there are browser tabs for 'geni-app - Microsoft Azure' and 'Cluster Details - Databricks'. The address bar shows the URL: https://portal.azure.com/#view/Microsoft_AAD_RegisteredApps/ApplicationMenuBlade/~/Overview/quick.... The Microsoft Azure logo and search bar are visible, along with the user profile 'genistudent97@gmail.c...' and 'DEFAULT DIRECTORY (GENISTUD...'.

The main content area is titled 'Home > Default Directory | App registrations > geni-app'. A left-hand navigation pane includes 'Overview' (selected), 'Quickstart', 'Integration assistant', 'Diagnose and solve problems', 'Manage', and 'Support + Troubleshooting'. The main content area features a search bar and action buttons: 'Delete', 'Endpoints', and 'Preview features'. A blue notification banner at the top reads: 'Got a second? We would love your feedback on Microsoft identity platform (previously Azure AD for developer). →'.

Under the 'Essentials' section, the following details are listed:

- Display name:** [geni-app](#)
- Application (client) ID:** acd4cbfe-5045-4db2-bb82-568a0414fae8
- Object ID:** aa3035d0-a98b-4345-99f3-5d87421f2f56
- Directory (tenant) ID:** 051bc474-1bf0-4155-8e6c-11ae24bf070b
- Supported account types:** [My organization only](#)

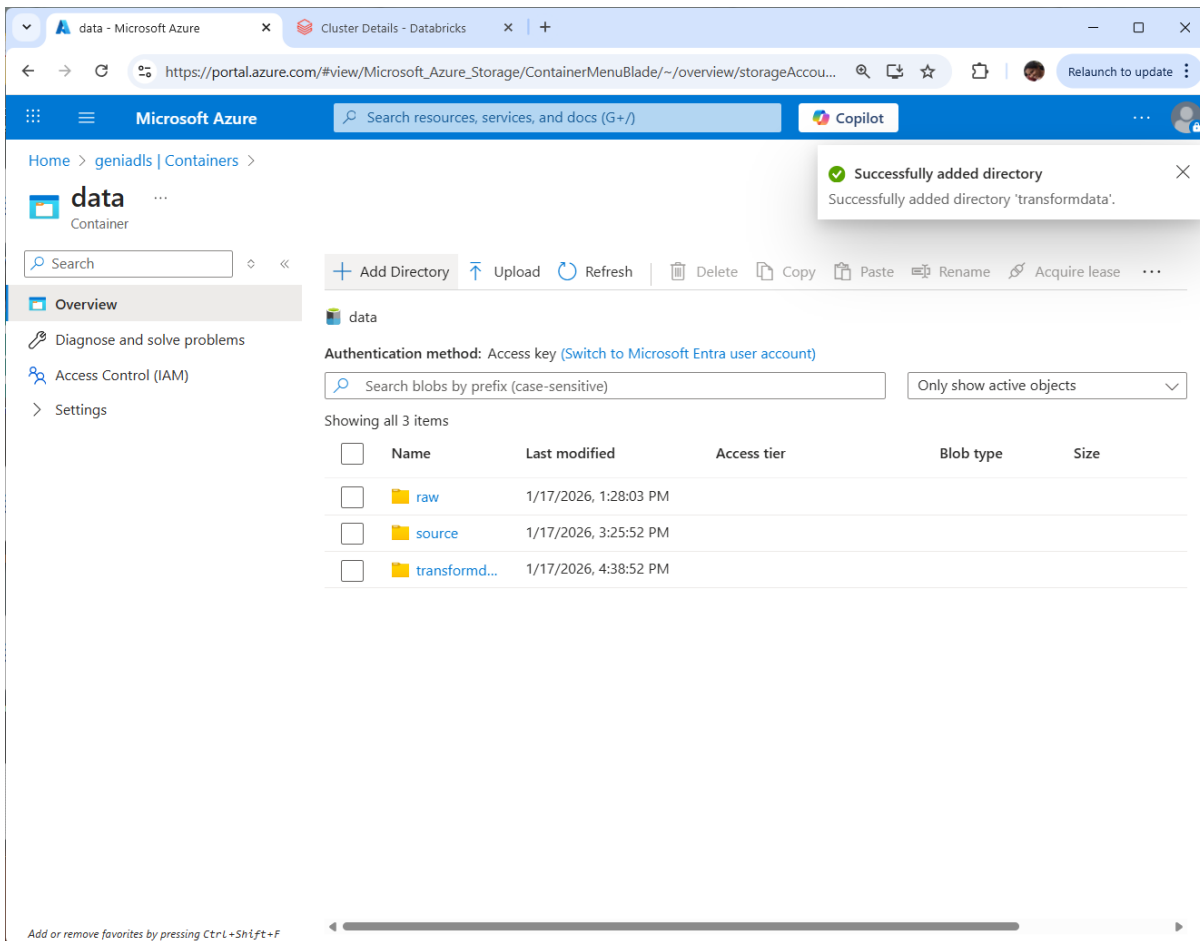
Client credentials and other links are listed on the right:

- Client credentials:** [Add a certificate or secret](#)
- Redirect URIs:** [Add a Redirect URI](#)
- Application ID URI:** [Add an Application ID URI](#)
- Managed application in local directory:** [geni-app](#)

Two informational messages are displayed in light blue boxes:

- Message 1: 'Welcome to the new and improved App registrations. Looking to learn how it's changed from App registrations (Legacy)? [Learn more](#)'
- Message 2: 'Starting June 30th, 2020 we will no longer add any new features to Azure Active Directory Authentication Library (ADAL) and Azure Active Directory Graph. We will continue to provide technical support and security updates but we will no longer provide feature updates. Applications will need to be upgraded to Microsoft Authentication Library (MSAL) and Microsoft Graph. [Learn more](#)'

At the bottom, there are links for 'Get Started' and 'Documentation'. A footer note at the very bottom reads: 'Add or remove favorites by pressing Ctrl+Shift+F'.



Data Transformation Using PySpark

- PySpark code was used to read data from the /raw folder, perform transformations such as data type casting and aggregation, and write the processed data to the /transformdata folder.
- This step ensured that data was cleaned and optimised for analytics.

Databricks notebook with transformation code

```

from pyspark.sql.functions import col
from pyspark.sql.types import IntegerType

storage_account = "geniadls"
container = "data"

account_key = "<storage-account-key-masked>"

```

```

try:
    dbutils.fs.mount(
        source=f"abfss://{container}@{storage_account}.dfs.core.windows.net/",
        mount_point=mount_point,
        extra_configs={f"fs.azure.account.key.{storage_account}.dfs.core.windows.net": account_key}
    )
    display(dbutils.fs.ls(mount_point))
except Exception as e:
    print(str(e))

base_path = f"abfss://{container}@{storage_account}.dfs.core.windows.net"
raw_path = f"{base_path}/raw"
transform_path = f"{base_path}/transformdata"

def read_csv(path):
    return (spark.read
            .format("csv")
            .option("header", "true")
            .option("inferSchema", "true")
            .option(f"fs.azure.account.key.{storage_account}.dfs.core.windows.net", account_key)
            .load(path))

def write_parquet(df, path, mode="overwrite"):
    (df.write
     .format("parquet")
     .mode(mode)
     .option(f"fs.azure.account.key.{storage_account}.dfs.core.windows.net", account_key)
     .save(path))

athletes = read_csv(f"{raw_path}/Athletes.csv")

```

```

coaches = read_csv(f"{raw_path}/Coaches.csv")
entriesgender = read_csv(f"{raw_path}/EntriesGender.csv")
medals = read_csv(f"{raw_path}/Medals.csv")
teams = read_csv(f"{raw_path}/Teams.csv")

print(f"Athletes: {athletes.count()} rows")
print(f"Coaches: {coaches.count()} rows")
print(f"Entries Gender: {entriesgender.count()} rows")
print(f"Medals: {medals.count()} rows")
print(f"Teams: {teams.count()} rows")

display(athletes.limit(5))
display(medals.limit(5))

entriesgender_transformed = (entriesgender
    .withColumn("Female", col("Female").cast(IntegerType()))
    .withColumn("Male", col("Male").cast(IntegerType()))
    .withColumn("Total", col("Total").cast(IntegerType()))
)

top_gold = medals.orderBy(col("Gold").desc()).select("TeamCountry", "Gold")
display(top_gold.limit(10))

average_entries_by_gender = (entriesgender_transformed
    .withColumn("Avg_Female", col("Female") / col("Total"))
    .withColumn("Avg_Male", col("Male") / col("Total"))
)

display(average_entries_by_gender.limit(10))

```

```
write_parquet(athletes, f"{transform_path}/athletes")
write_parquet(coaches, f"{transform_path}/coaches")
write_parquet(entriesgender_transformed, f"{transform_path}/entriesgender")
write_parquet(medals, f"{transform_path}/medals")
write_parquet(teams, f"{transform_path}/teams")

print("Data written to /transformdata/")

medals_out = (spark.read
              .format("parquet")
              .option(f"fs.azure.account.key.{storage_account}.dfs.core.windows.net", account_key)
              .load(f"{transform_path}/medals"))

display(medals_out.limit(5))
```

Code output

Microsoft Azure databricks

Search data, notebooks, recent, and more... CTRL + P

waruni-bahar-db

File Edit View Run Help Python Table ON

```

1 from pyspark.sql.functions import col
2 from pyspark.sql.types import IntegerType
3
4 storage_account = "gmiadls"
5 container = "data"
6
7 account_key = "*****"
8
9 try:
10     dbutils.fs.mount(
11         source="abfss://{container}@{storage_account}.dfs.core.windows.net",
12         mount_point=mount_point,
13         extra_configs={"fs.azure.account.key.{storage_account}.dfs.core.windows.net": account_key}
14     )
15     display(dbutils.fs.ls(mount_point))
16 except Exception as e:
17     print(str(e))
18
19 base_path = "abfss://{container}@{storage_account}.dfs.core.windows.net"
20 raw_path = f"{base_path}/raw"
21 transform_path = f"{base_path}/transformdata"
22
23 def read_csv(path):
24     return (spark.read
25             .format("csv")
26             .option("header", "true")
27             .option("inferSchema", "true")
28             .option("fs.azure.account.key.{storage_account}.dfs.core.windows.net", account_key)
29             .load(path))
30
31 def write_parquet(df, path, mode="overwrite"):

```

Output Terminal Debug console pip logs

- athletes: pyspark.sql.connect.dataframe.DataFrame = [PersonName: string, Country: string ... 1 more field]
- average_entries_by_gender: pyspark.sql.connect.dataframe.DataFrame = [Discipline: string, Female: integer ... 4 more fields]
- coaches: pyspark.sql.connect.dataframe.DataFrame = [Name: string, Country: string ... 2 more fields]
- entriesgender: pyspark.sql.connect.dataframe.DataFrame = [Discipline: string, Female: integer ... 2 more fields]
- entriesgender_transformed: pyspark.sql.connect.dataframe.DataFrame = [Discipline: string, Female: integer ... 2 more fields]
- medals: pyspark.sql.connect.dataframe.DataFrame = [Rank: integer, TeamCountry: string ... 5 more fields]
- medals_joined: pyspark.sql.connect.dataframe.DataFrame = [Rank: integer, TeamCountry: string ... 5 more fields]
- teams: pyspark.sql.connect.dataframe.DataFrame = [TeamName: string, Discipline: string ... 2 more fields]
- top_gold: pyspark.sql.connect.dataframe.DataFrame = [TeamCountry: string, Gold: integer]

Method public com.databricks.backend.daemon.dbtutils.DBUtilsCoreResult com.databricks.backend.daemon.dbtutils.DBUtilsCore.mount()

Athletes: 11985 rows
Coaches: 304 rows
Entries Gender: 46 rows
Medals: 93 rows
Teams: 743 rows

Table - +

PersonName	Country	Discipline
1 AALERUD Katrine	Norway	Cycling Road
2 ABAD Nestor	Spain	Artistic Gymnastics
3 ABAGNALE Giovanni	Italy	Rowing
4 ABALDI Alberto	Spain	Basketball
5 ABALDI Tamara	Spain	Basketball

5 rows | 12.86s runtime

Table - +

Rank	TeamCountry	Gold	Silver	Bronze	Total	RankByTotal
1	United States of America	39	41	33	113	1
2	Peoples Republic of China	38	32	18	88	2
3	Japan	27	14	17	58	5
4	Great Britain	22	21	22	65	4
5	ROC	20	28	23	71	3

5 rows | 12.86s runtime

Table - +

TeamCountry	Gold
1 United States of America	39
2 Peoples Republic of China	38
3 Japan	27
4 Great Britain	22
5 ROC	20
6 Australia	17
7 Netherlands	10
8 France	10
9 Germany	10
10 Italy	10

10 rows | 12.86s runtime

Table - +

Discipline	Female	Male	Total	1.2 Avg. Female	1.2 Avg. Male
1 3x3 Basketball	32	32	64	0.5	0.5
2 Archery	64	64	128	0.5	0.5
3 Artistic Gymnastics	98	98	196	0.5	0.5
4 Artistic Swimming	105	0	105	1	0
5 Athletics	969	1072	2041	0.4747672709451640	0.5252327290548360
6 Badminton	86	87	173	0.49710962658939315	0.50289037341060685
7 Baseball/Softball	90	144	234	0.38451308461538464	0.61548691846153846
8 Basketball	144	144	288	0.5	0.5
9 Beach Volleyball	48	48	96	0.5	0.5
10 Boxing	102	187	289	0.35294119647058826	0.6470588235294119

10 rows | 12.86s runtime

Data written to /transformdata/

Table - +

Rank	TeamCountry	Gold	Silver	Bronze	Total	RankByTotal
1	United States of America	39	41	33	113	1
2	Peoples Republic of China	38	32	18	88	2
3	Japan	27	14	17	58	5
4	Great Britain	22	21	22	65	4
5	ROC	20	28	23	71	3

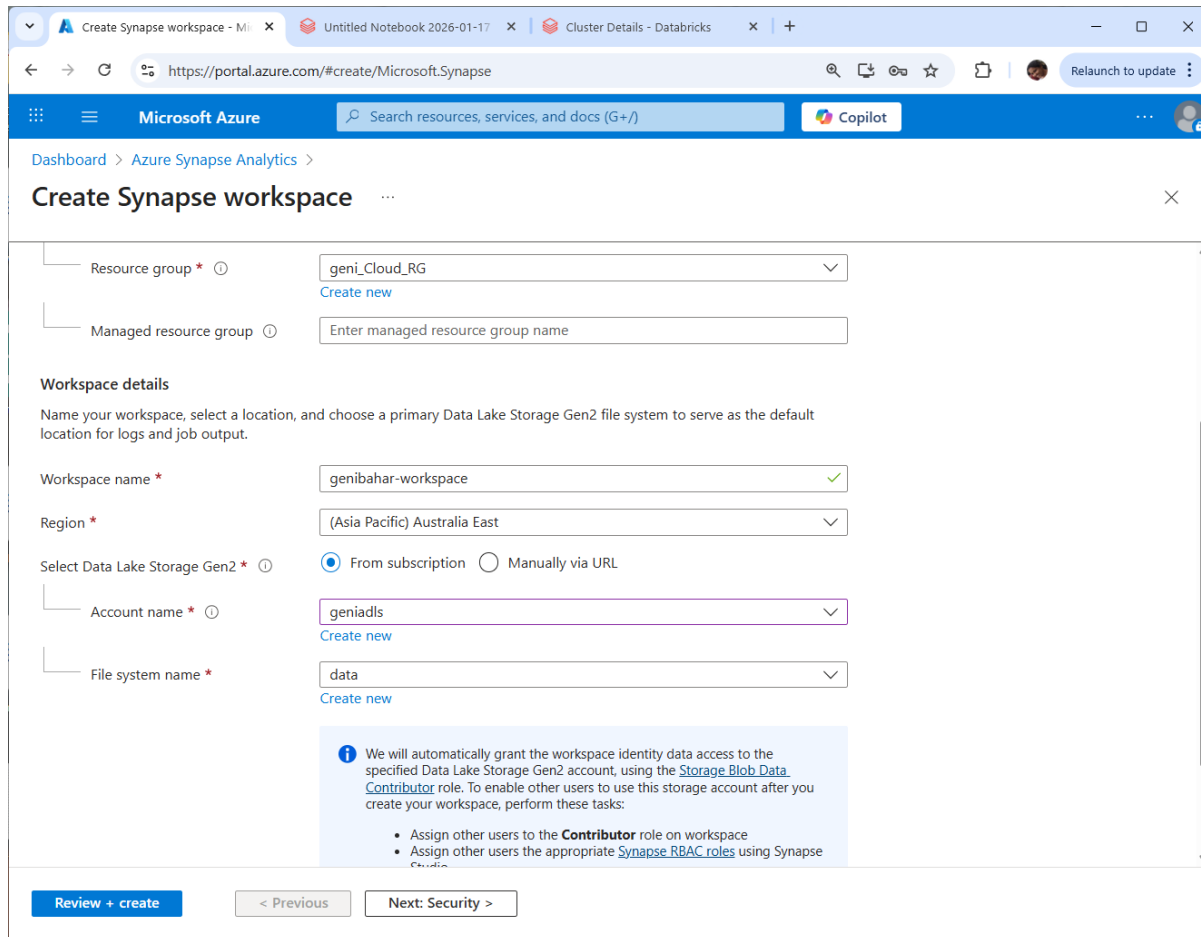
5 rows | 12.86s runtime

See performance (15)

Task 5: Azure Synapse Analytics and Automated Data Integration

Azure Synapse Analytics Workspace Creation

- An Azure Synapse Analytics workspace was created in the same resource group and region. The workspace was linked to the Azure Data Lake Storage Gen2 account to enable direct access to transformed datasets.



Dashboard > Azure Synapse Analytics >

Create Synapse workspace

Validation succeeded

* Basics * Security Networking Tags **Review + create**

Product Details

Azure Synapse Analytics workspace by Microsoft
 Serverless SQL est. cost/TB **5.00 USD**
[Terms of use](#) | [Privacy policy](#)

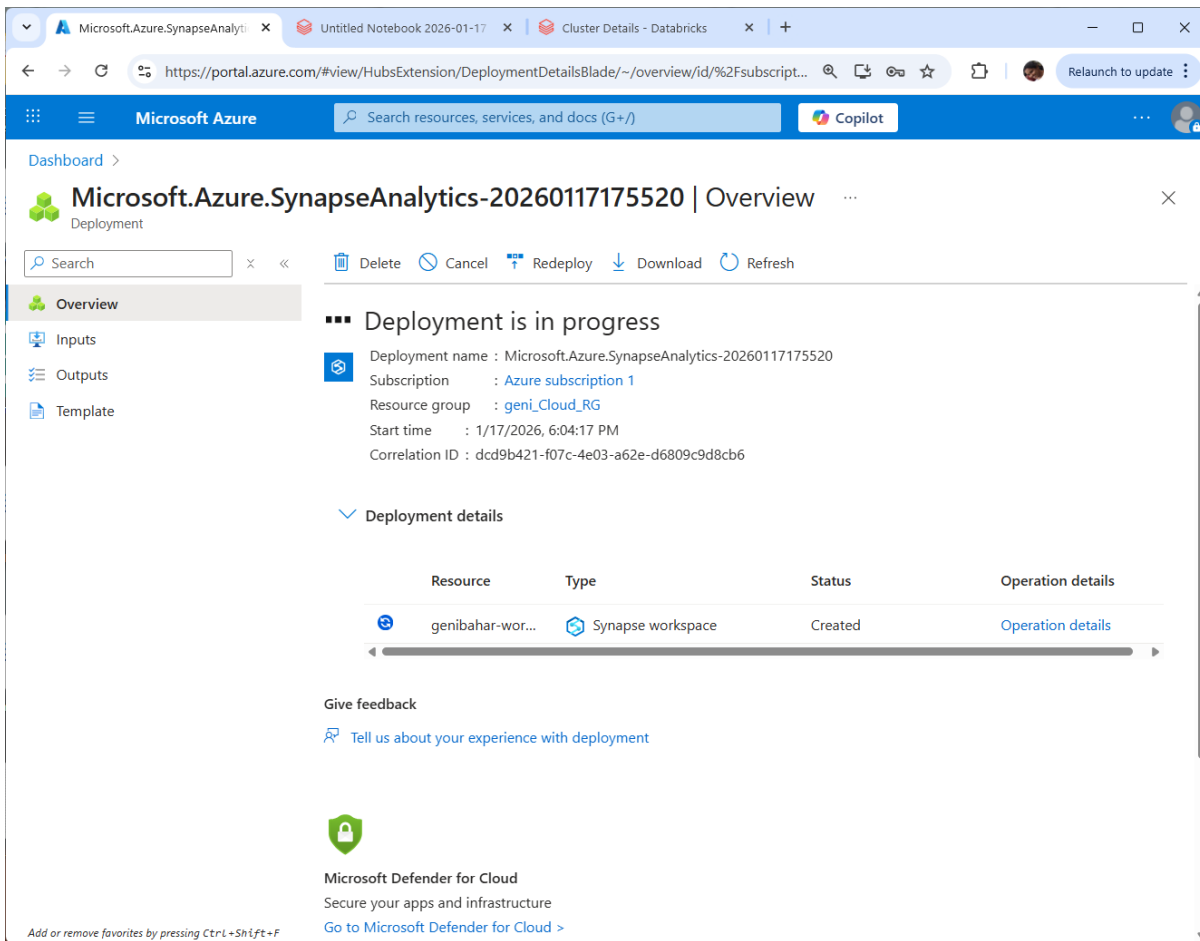
Terms

By clicking Create, I (a) agree to the legal terms and privacy statement(s) associated with the Marketplace offering(s) listed above; (b) authorize Microsoft to bill my current payment method for the fees associated with the offering(s), with the same billing frequency as my Azure subscription; and (c) agree that Microsoft may share my contact, usage and transactional information with the provider(s) of the offering(s) for support, billing and other transactional activities. Microsoft does not provide rights for third-party offerings. For additional details see [Azure Marketplace Terms](#).

Basics

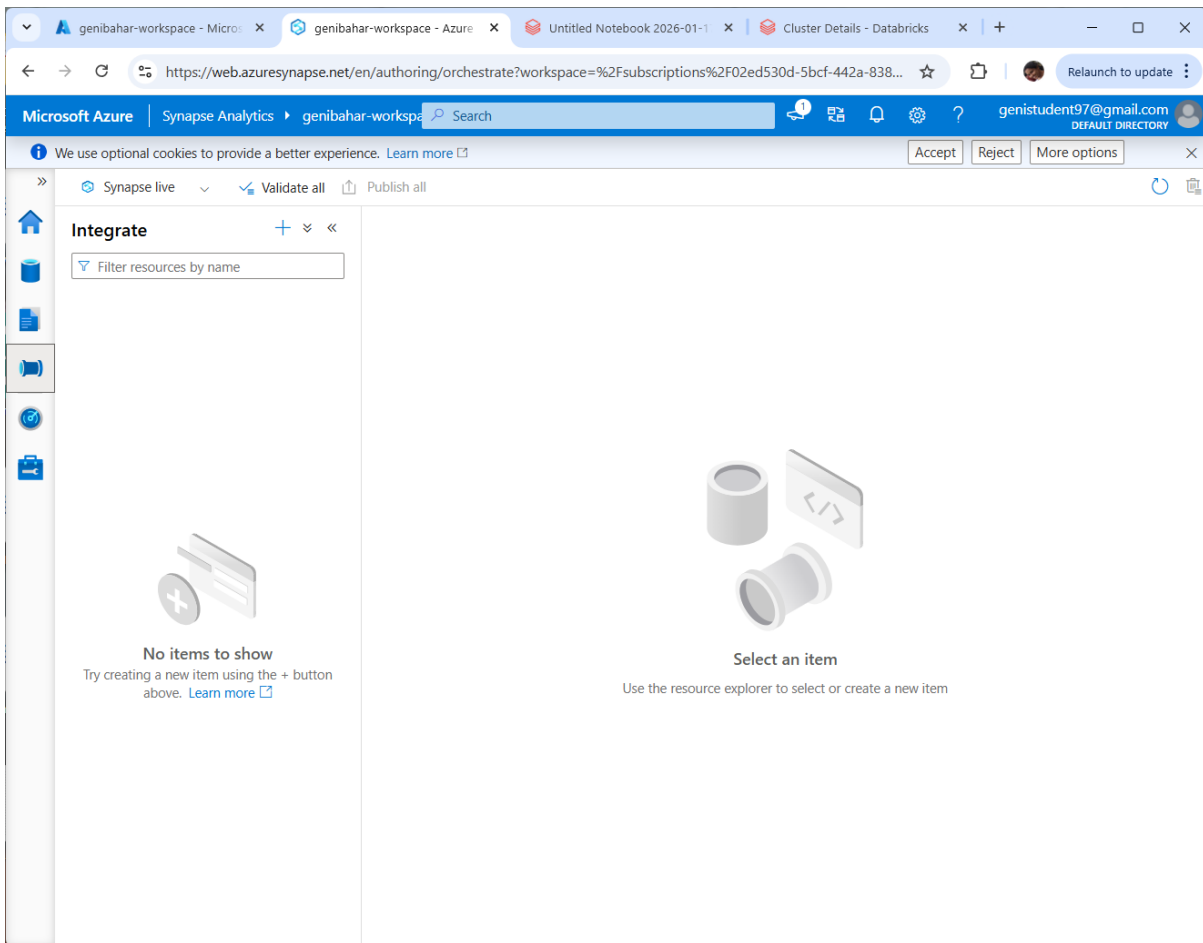
Subscription	Azure subscription 1
Resource group	geni_Cloud_RG
Region	Australia East
Workspace name	(new) genibahar-workspace
Data Lake Storage Gen2 account	https://geniadl.dfs.core.windows.net

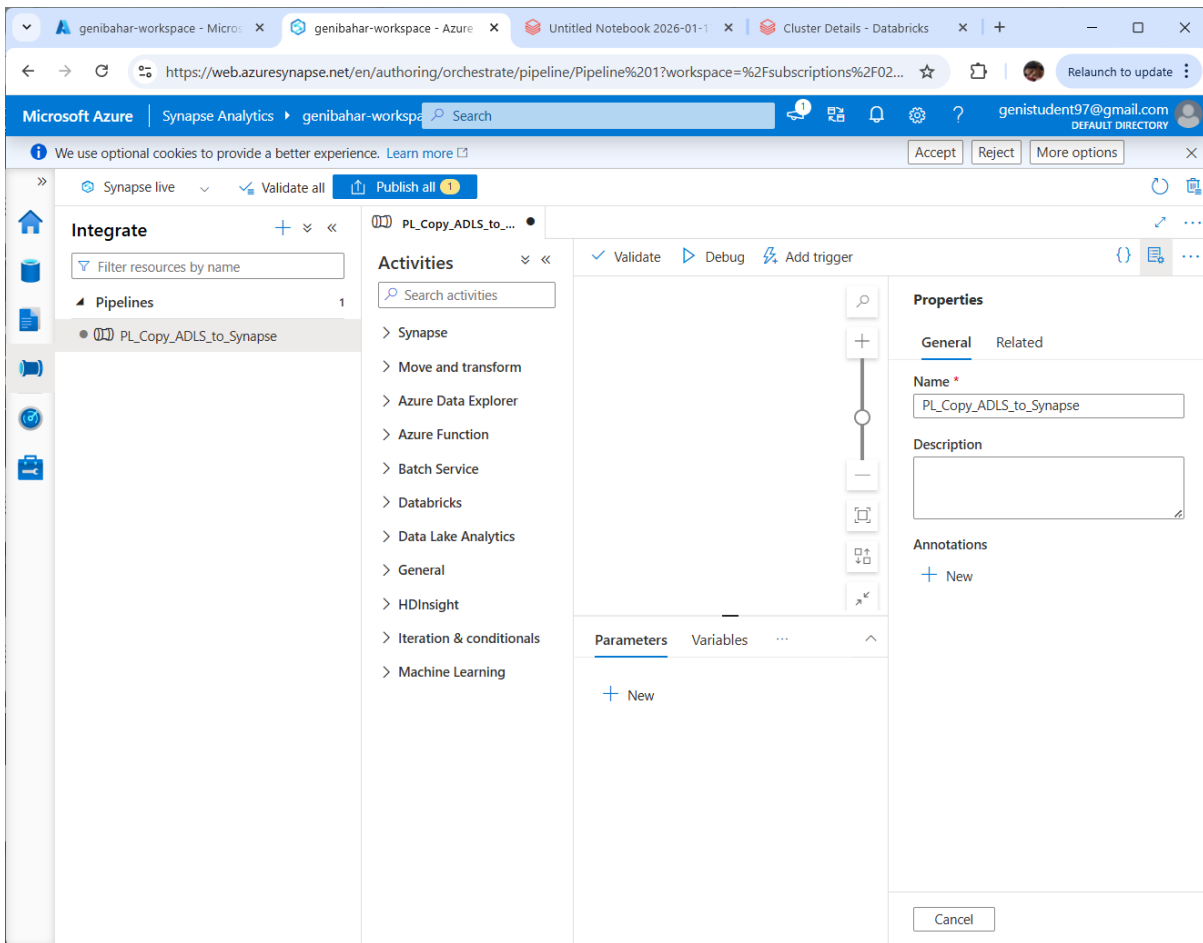
Create < Previous Next > Download a template for automation



Data Integration Pipeline

- A Synapse pipeline was designed to automate the data integration process between Azure Data Lake Storage Gen2 and Synapse Analytics. The pipeline was validated and executed successfully.





The screenshot displays the Microsoft Azure Synapse Analytics web interface. The browser address bar shows the URL: `https://web.azuresynapse.net/en/authoring/orchestrate/pipeline/Pipeline%201?workspace=%2Fsubscriptions%2F02...`. The user is logged in as `genistudent97@gmail.com`.

The interface is divided into several sections:

- Integrate:** A sidebar on the left with a search bar and a list of pipelines. The selected pipeline is `PL_Copy_ADLS_to_Synapse`.
- Activities:** A central pane showing a search bar and a list of activity categories: Synapse, Move and transform (with `Copy data` and `Data flow` sub-items), Azure Data Explorer, Azure Function, Batch Service, Databricks, Data Lake Analytics, General, HDInsight, Iteration & conditionals, and Machine Learning.
- Canvas:** The main workspace showing a pipeline diagram with a single activity named `Copy data1`.
- Properties:** A right-hand pane with tabs for `General` and `Related`. The `General` tab is active, showing:
 - Name:** `PL_Copy_ADLS_to_Synapse`
 - Description:** (Empty text area)
 - Annotations:** `+ New`
 - Activity state:** Activated Deactivated
 - Timeout:** (Partially visible)

The screenshot shows the 'Set properties' dialog for a pipeline activity in the Microsoft Azure Synapse Analytics interface. The activity is named 'DS_Source_TransformData'. The 'Linked service' is set to 'genibahar-workspace-WorkspaceDefaultStorage'. The 'Connect via integration runtime' is set to 'AutoResolveIntegrationRuntime'. The 'File path' is 'data / transformdata / File name'. The 'First row as header' checkbox is checked. The 'Import schema' options are 'From connection/store' (selected), 'From sample file', and 'None'. There is an 'Advanced' section that is currently collapsed. The dialog has 'OK', 'Back', and 'Cancel' buttons at the bottom.

The screenshot shows the Microsoft Azure Synapse Analytics interface. The browser address bar indicates the URL: <https://web.azuresynapse.net/en/authoring/orchestrate/pipeline/Pipeline%201?workspace=%2Fsubscriptions%2F02...>. The user is logged in as `genistudent97@gmail.com`. The main content area is titled "Set properties" and contains the following fields:

- Name:** `DS_Sink_Synapse`
- Linked service *:** `genibahar-workspace-WorkspaceDefaultSqlServer`
- Connect via integration runtime *:** `AutoResolveIntegrationRuntime`
- Table name:** `Select...` (with an "Enter manually" checkbox that is unchecked)
- Import schema:** `None` (selected via radio button)

At the bottom of the dialog, there are three buttons: `OK`, `Back`, and `Cancel`. The left sidebar shows the "Integrate" section with a search filter and a list of activities, including "PL_Copy_ADLS_to_Synapse".

Home > Azure Synapse Analytics > Create Synapse workspace

Workspace details
Name your workspace, select a location, and choose a primary Data Lake Storage Gen2 file system to serve as the default location for logs and job output.

Workspace name *

Region *

Select Data Lake Storage Gen2 * From subscription Manually via URL

Account name * [Create new](#)

File system name * [Create new](#)

i We will automatically grant the workspace identity data access to the specified Data Lake Storage Gen2 account, using the [Storage Blob Data Contributor](#) role. To enable other users to use this storage account after you create your workspace, perform these tasks:

- Assign other users to the **Contributor** role on workspace
- Assign other users the appropriate [Synapse RBAC roles](#) using Synapse Studio
- Assign yourself and other users to the **Storage Blob Data Contributor** role on the storage account

[Learn more](#)

[Review + create](#) [< Previous](#) [Next: Security >](#)

The screenshot shows the Microsoft Azure portal interface. At the top, there are browser tabs for 'Microsoft.Azure.SynapseAnalyti...', 'Untitled Notebook 2026-01-17', and 'Cluster Details - Databricks'. The address bar shows the URL: 'https://portal.azure.com/#view/HubsExtension/DeploymentDetailsBlade/~/overview/id/%2Fsubscriptions...'. The main header includes the 'Microsoft Azure' logo, a search bar, and the 'Copilot' button.

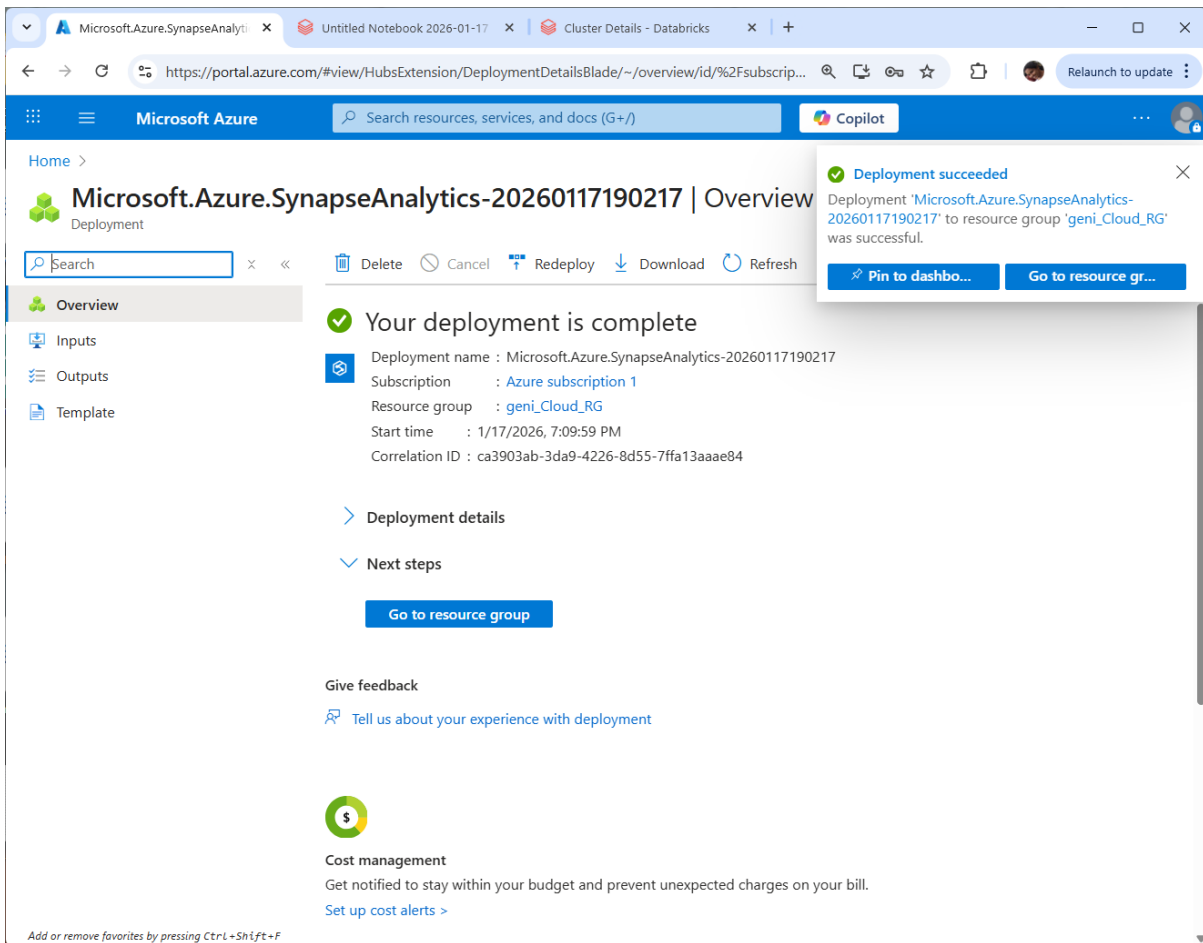
The main content area is titled 'Microsoft.Azure.SynapseAnalytics-20260117190217 | Overview'. Below the title is a 'Deployment' section with a search bar and action buttons: 'Delete', 'Cancel', 'Redeploy', 'Download', and 'Refresh'. A notification box in the top right corner states: 'Deployment in progress... Deployment to resource group 'geni_Cloud_RG' is in progress.'

The 'Deployment is in progress' section displays the following details:

- Deployment name : Microsoft.Azure.SynapseAnalytics-20260117190217
- Subscription : Azure subscription 1
- Resource group : geni_Cloud_RG
- Start time : 1/17/2026, 7:09:58 PM
- Correlation ID : ca3903ab-3da9-4226-8d55-7ffa13aaae84

Below these details is a 'Deployment details' section with a table. The table has four columns: 'Resource', 'Type', 'Status', and 'Operation details'. The table is currently empty, with a message below it stating: 'There are no resources to display.'

At the bottom of the page, there is a 'Give feedback' link and a 'Microsoft Defender for Cloud' section with the text: 'Secure your apps and infrastructure'.



External Table Creation Using SQL

- External tables were created in Synapse using SQL scripts. These tables reference data stored in the /transformdata folder, allowing serverless querying without data duplication.

```

1 CREATE EXTERNAL DATA SOURCE OlympicDataLake
2 WITH (
3     LOCATION = 'https://genieis-dfs.core.windows.net/data'
4 );
5
6 CREATE EXTERNAL FILE FORMAT CsvFormat
7 WITH (
8     FORMAT_TYPE = DELIMITEDTEXT,
9     FORMAT_OPTIONS (
10        FIELD_TERMINATOR = ',',
11        FIRST_ROW = 2
12    );
13 );
14
15 CREATE EXTERNAL TABLE Athletes (
16     Name NVARCHAR(100),
17     Country NVARCHAR(50),
18     Discipline NVARCHAR(50)
19 )
20 WITH (
21     LOCATION = 'transformdata/athletes',
22     DATA_SOURCE = OlympicDataLake,
23     FILE_FORMAT = CsvFormat
24 );
25
26 CREATE EXTERNAL TABLE Medals (
27     Rank INT,
28     Team_Country NVARCHAR(50),
29     Gold INT,
30     Silver INT,
31     Bronze INT,
32     Total INT,
33     TotalMedals INT
34 );
35 WITH (
36     LOCATION = 'transformdata/medals',
37     DATA_SOURCE = OlympicDataLake,
38     FILE_FORMAT = CsvFormat
39 );
40
41 CREATE EXTERNAL TABLE dbo_EntriesGender (
42     Discipline NVARCHAR(100),
43     Female INT,
44     Male INT
45 )
46 WITH (
47     LOCATION = 'transformdata/entriesgender',
48     DATA_SOURCE = OlympicDataLake,
49     FILE_FORMAT = CsvFormat
50 );
51 GO
    
```

Data Analysis Using SQL

Several SQL queries were executed to analyse the dataset:

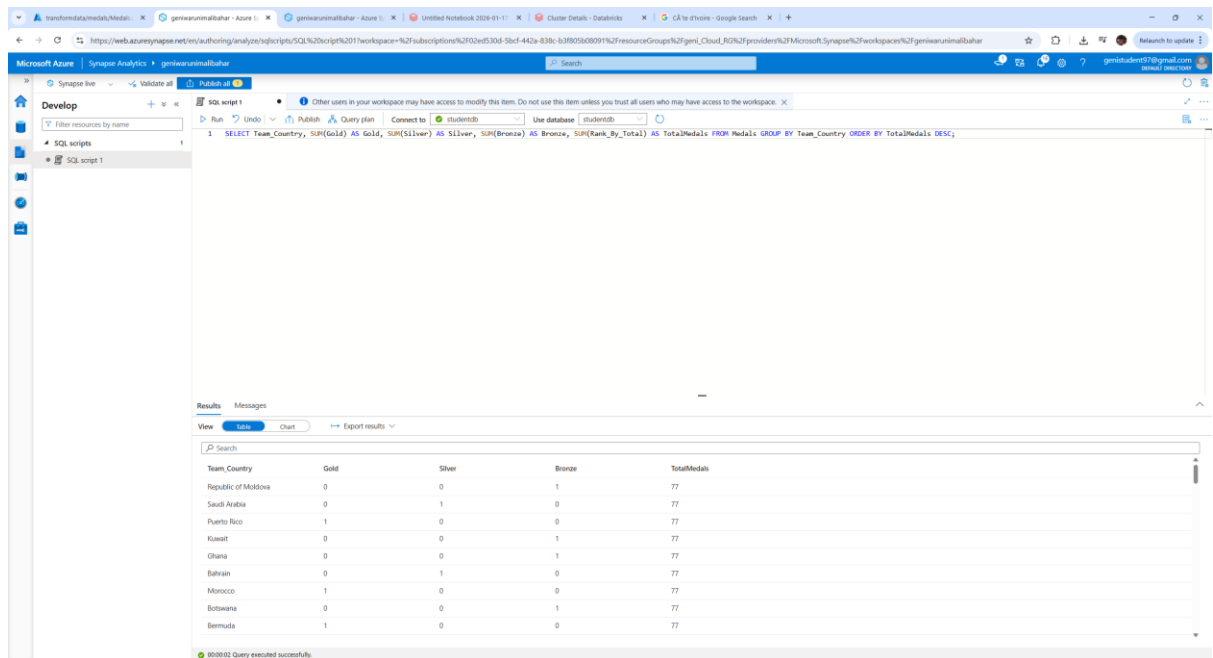
- Number of athletes by country

```

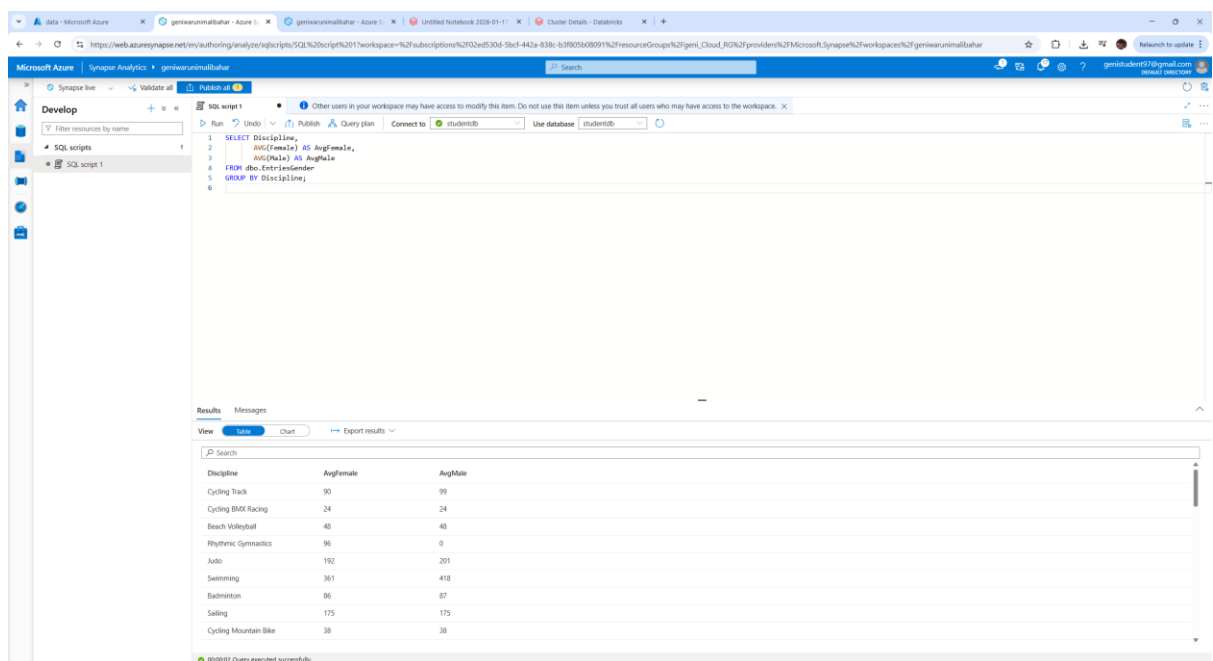
1 SELECT Country, COUNT(*) AS TotalAthletes
2 FROM Athletes
3 GROUP BY Country
4 ORDER BY TotalAthletes DESC;
5
    
```

Country	TotalAthletes
United States of America	615
Japan	586
Australia	470
People's Republic of China	401
Germany	400
France	377
Canada	368
Great Britain	366
Italy	356

- Total medals by country



- Average gender participation by discipline



These queries provided insights into Olympic performance and participation trends.

Working with external tables in Azure Synapse helped me understand how large datasets can be analysed directly from the data lake without moving or duplicating data. The visualisations made it easier to identify trends, such as medal distribution across countries, which would be harder to see using raw query results alone.

Data Visualisation

Charts were created within Synapse Studio to visually represent the query results. These visualisations help interpret the data more effectively.

