# Data Engineering in the Cloud

## Azure Data Lake Storage Gen2

Xuemao Zhang
East Stroudsburg University

January 18, 2025

# Outline

- Introduction to Azure Data Lake Storage Gen2
- Lab: Create an Azure Database

# Introduction to Azure Data Lake

- A **data lake** is a single, centralized repository where you can store all your data, both structured and unstructured.
- Key Benefits:
  - ▶ Centralized Storage: Quickly and more easily store, access, and analyze a wide variety of data in a single location.
  - ▶ Flexible Data Storage: Store your data in its raw or native format, usually as files or binary large objects (blobs).
  - ▶ No Predefined Structure: No need to conform your data to fit an existing structure.
- Advantages:
  - ▶ Unified Data Management: Simplifies data storage and management across the organization.
  - ▶ Enhanced Data Analysis: Facilitates advanced analytics and insights from diverse data types.
  - ▶ Scalability: Accommodates growing data volumes efficiently.

# Introduction to Azure Data Lake

- Azure Data Lake Storage Gen2 is a set of capabilities dedicated to **big data analytics**, built on Azure Blob Storage.
  - The data that you ingest persist as blobs in the storage account.
- It combines the scalable and cost-effective storage capabilities of Blob Storage with advanced data management, security, and performance features required for large-scale analytics.
- *Azure Data Lake Storage Gen2* refers to the current implementation of Azure's Data Lake Storage solution. The previous implementation, Azure Data Lake Storage Gen1 was retired on February 29, 2024.

# Introduction to Azure Data Lake

- Data Lake Storage Gen2 includes the following capabilities.
  - Hadoop-compatible access
  - Hierarchical directory structure
  - Optimized cost and performance
  - Finer grain security model
  - Massive scalability
- The hierarchical namespace is a key feature that enables Azure Data Lake Storage Gen2 to provide high-performance data access at object storage scale and price.

# Introduction to Azure Data Lake

- You can use the `hierarchical namespace` feature to organize all the objects and files within your storage account into a hierarchy of directories and nested subdirectories.
  - Data is organized in much the same way that files are organized on your computer.
- To unlock `Data Lake Storage Gen2` capabilities, we need to enable the hierarchical namespace setting when we creat a storage account.
- When you use Azure Data Lake Storage to store data, the data is stored using a hierarchical file system that is compatible with Hadoop Distributed File System (HDFS).
  - We'll discuss HDFS later.
- For more information, read Introduction to Azure Data Lake Storage Gen2

# Lab

- Create a storage account with `Enable hierarchical namespace` checked.

# Lab

- Home $\rightarrow$ SQL databases and click `Create SQL database`
- We need to create a server
    - Authentication method: **Use SQL authentication**
    - Set Username: sqladmin and pwd: EsuMath23
- SQL elastic pool: No
- Workload environment: Production or Development(more Cost-Effective)
- `Compute + storage`:
    - Service tier: General purpose

# Lab

- Networking tab:
  - Connectivity method: Public endpoint
  - Firewall rules: Yes and Yes
- Security tab:
  - Enable Microsoft Defender for SQL: Not now
- Additional settings: No changes
- Tags: No changes

# Lab

- Review + create

# Lab

# Lab

- Click `Go to resource`

# Lab

- To access Query editor, you need the username and pwd you created before.

# Lab

- Now go back to Home

# Lab

- Then click on the storage account and add a `container` container1

# Lab

# Lab

- Open container1, and upload the data set `population.csv` to the container
  - We did not upload any data to a storage account in the last lab

# License