# Data Engineering in the Cloud
## Azure Data Factory

Xuemao Zhang
East Stroudsburg University

January 18, 2025

# Outline

- Overview of Azure Data Factory
- Data Integration
- Lab

# Overview of Azure Data Factory

- Once the data is loaded, it is important to perform serverless transformation and integration of the data.
- For data transformation, you can take the help of Azure Data Factory to create data transformation workflows.
- Azure Data Factory is an extract load transform (ETL) service that creates data driven workflows on cloud.
- The following is the working of Data Factory by a series of interconnected system:

# Overview of Azure Data Factory

**Connect and Collect**

- Connect Data Sources: All important data and processing sources, such as SaaS providers, databases, and other systems, are connected to build an information production system.
  - For example, Connect to relational databases.
- Centralize Data: Move the data to a central location for further processing.
  - Data Lake: Use Azure Data Lake Storage to store raw data in its native format.
- Data Movement with Data Factory: Use Azure Data Factory to create a copy activity in a pipeline.
  - Create a pipeline in Azure Data Factory to copy data from the source to the central location.
    - ⋆ The source and central location in an Azure Data Factory (ADF) pipeline can both be containers in Azure Data Lake Storage.
- Data Analysis: After centralizing the data, perform further analysis as needed.
  - For example, use Azure Data Factory to transform data as needed (e.g., cleaning, aggregating).

# Overview of Azure Data Factory

**Transform and Enrich:** Azure Data Factory (ADF) offers robust capabilities for data transformation using mapping data flows, as well as integration with services like Apache Spark and Azure Machine Learning

- Azure Data Factory Mapping Data Flows: Mapping data flows in ADF allow you to design visually-driven transformations.
- We can use the drag-and-drop interface to build complex transformation logic.
- Services such as Spark and Machine Learning can transform the data with the help of the Azure Data Factory.
  - ▶ Azure Synapse Analytics integrates with Apache Spark, enabling large-scale data processing and transformation.
  - ▶ Azure Data Factory can orchestrate data processing and machine learning workflows by integrating with Azure Machine Learning.

# What is Data Integration?

- Data integration involves combining data from different sources to provide a unified view, making it easier to manage and analyze.
- Data integration can be performed with Azure Data Factory, which enables for the development of data driven processes for data transformation at scale.
- It may also be used to build and organize data driven processes that allow for the consumption of data from multiple sources.
  - ADF supports integration with various data sources such as databases, file systems, APIs, and cloud services.
- In this process, complex ETL processes to modify data visually may also be created with the use of data flows or computing services like Azure Databricks, Azure Synapse, and so on.

# Transformation Data Using Mapping Data Flow

- Mapping Data Flows provides a visual platform for creating a variety of data transformations without needing to use code.
- The data flows are built and then executed on scaled out Apache Spark clusters that are deployed automatically when the Mapping Data Flow is executed.
- It also enables the monitoring of the transformation's execution to observe their progress and comprehend any difficulties that may develop.

# Transformation Data Using Compute Resources

- Azure Data Factory can transform data by using the computation resources and a platform service that is better suited for the task.
- For example, Azure Data Factory may build a pipeline to an analytical data platform, such as Spark pools on an Azure Synapse Analytics instance, to perform a complex python computation.
- Another example is sending data to an Azure SQL Database instance in order to run a Transact SQL stored procedure.

# Types of Azure Data Factory Transformation

A range of transformation types are available in Mapping Data Flows to allow data to be changed. They are divided into the following categories:

- Schema Modifier Transformations allow you to change the structure of your data by adding, removing, or modifying columns.
- Row Modifier Transformations are used to modify the data within each row. These transformations enable you to clean, aggregate, and manipulate your data on a row-by-row basis, like sort, filter, aggregate
- Multiple Inputs/Outputs Transformations allow you to handle data from multiple sources or direct data to multiple destinations. These transformations are useful for complex data processing scenarios involving multiple datasets.

# Prerequisites

- Prerequisites (5 minutes):
  - Create a datalake (storage account with hierarchical namespace enabled)
    - To lower the cost, you may choose Redundancy as LRS
  - Create a container in the storage: raw
    - Upload the data set CoursesData.csv to the container
  - Create an Azure Synapse Workspace
    - Create a dedicated SQL pool

# Lab

- Click on `create a resource` and search for a `Data Factory`
- Go to `Data Factory` and click on `create`
- Set the `Basics` as shown in the below image

# Lab

- Click on "Next: Git configuration" and "Enable configure Git later"

## Create Data Factory ⋯

Basics **Git configuration** Networking Advanced Tags Review + create

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

Configure Git later ⓘ      ☑

# Lab

- Now leave the other configurations as it is and go to `review + create`
- Once the configuration is validated, click on the "create" button

# Lab

- Click `Go to resource` go to the `Data Factory` and lunch `Azure Data Factory Studio`



- Now go to the home section and click on `Ingest`

# Lab

- Choose Task type as `Built-in copy task`
- Choose "Run once now" and click on "next"

# Lab

- Click on "+ New connection"

## Source data store

Specify the source data store for the copy task. You can use an existing data store connection or specify a new data store.

Source type          All                                        ∨

Connection *         Select...                                  ∨      + New connection

# Lab

- Select **Azure Data Lake Storage Gen-2** and click on `continue`. You are choosing this because we have a dataset stored in the Gen-2 already.

# Lab

- Enter the details, *test connection*, and click on the `create` button

# Lab

- Select the source type as `Azure Data Lake Storage Gen-2`, now click on the `Browse` button

# Lab

- Untick the `recursively` option
- Click on `Next` and click on "Preview data" to check the connection

# Lab

- Click on `Next`
- Now let us specify our "new target". Click on the "new connection"

# Lab

- Select `Azure Synapse Analytics` and click on the `Continue` button

# Lab

- Enter the details, *test connection*, and click on the Create button

# Lab

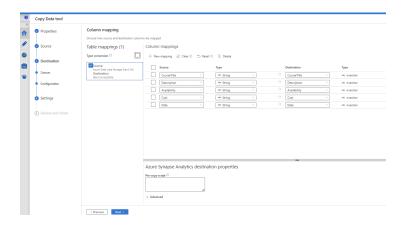- Click on "Use existing table" and click `Auto-create a destination table with the source schema`
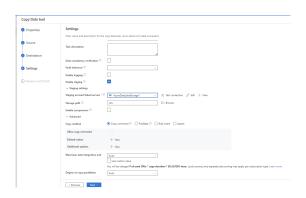
# Lab

- Click on "Next"

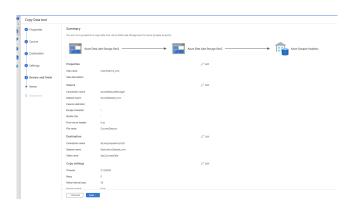# Lab

- Keep the "conversion settings".

# Lab

- Have a look at the "Column mappings", make sure **Untick** the type "Type conversion" option, and click on Next
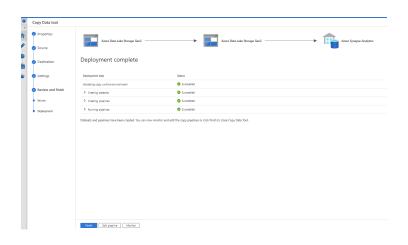
## Lab

- Give a name to the task and the storage path. Specify the staging account linked service as well

# Lab
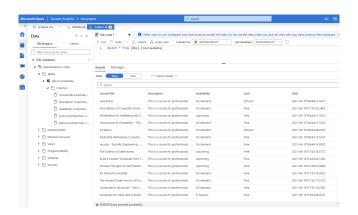


- Click on Next again

# Lab



- Click on the "Finish" button

# Lab

- Now go to the "monitor" section and you should be able to see the pipeline in succeeded state or progress state

## Lab

- Now go to or open the "Azure Synapse studio",
  - and check the database or
  - create a SQL script and run the following command and you must be able to see the data fetched in Azure Synapse

```
SELECT * from [dbo].[CoursesData]
```

# lab

- Delete everything we have done in Azure

# License