

Name : Parth Nandedkar

Date : 20 Feb 2024

Topics : Azure Data Factory

Batch : Data Engineering Batch-1

Handwritten Notes :

20/01/24

• **Azure Data Factory :**

Azure Data Factory is a cloud-based data integration service that allows you to create data-driven workflows in the cloud or orchestrating & automating data movement and data transformation.

• **Use cases :**

- Supporting data migrations
- Getting data from a client's server or online data to a Azure data lake
- Carrying out various data integration processes
- Integrating data from different ERP systems and loading it into Azure Synapse for Enterprise Resource Planning reporting.

• **How Working :**

Hosted on Azure

<div style="border: 1px solid black; padding: 5px; display: inline-block;">sql python pyspark SparkSQL</div>	⇒	<div style="border: 1px solid black; padding: 5px; display: inline-block;">Azure Data Factory</div> <div style="border: 1px solid black; padding: 5px; display: inline-block; margin-top: 10px;">Azure Synapse</div> <div style="border: 1px solid black; padding: 5px; display: inline-block; margin-top: 10px;">Azure DevOps</div>
--	---	--

Azure Data Factory :
Data driven pipeline
Move, transform, run data
into the pipeline..

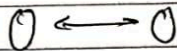
Cloud based integration service

Step 1 → Connect & collect

Step 2 → Transform & Enrich

Step 3 → Publish

Data Migration Activity



Copy Activity → Source to sink data store

(Blob, Cosmos DB)

Transformation activity

Hive, Map reduce, Spark

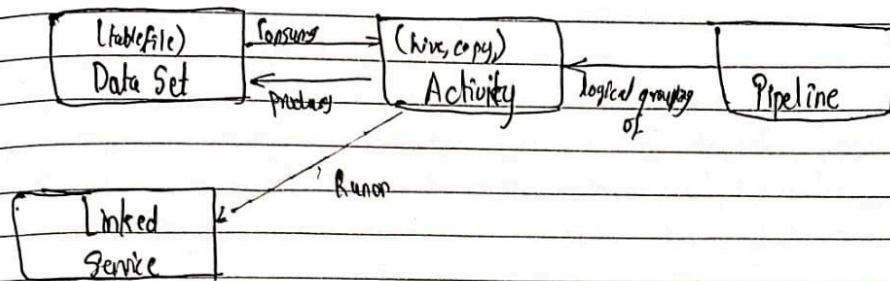
ADF pipeline →

- Dataset represent data structure within data stores:

= A pipeline is a graph activity

Activities define the actions to perform on your data

Linked services define the information needed for Azure Data Factory to connect to external resources.



You can use one of the following tools or APIs to create data pipeline in Azure Data Factory.

- Azure Portal
- Visual Studio
- PowerShell
- .NET API
- REST API
- Azure resource manager template.

Azure Data Factory :

Azure Data Factory (ADF) is a cloud-based data integration service provided by Microsoft Azure. It enables you to create, schedule, and manage data pipelines for ingesting, transforming, and processing data from various sources. Here's an overview of Azure Data Factory and its key features:

Data Orchestration: Azure Data Factory allows you to orchestrate and automate the movement and transformation of data across on-premises and cloud-based data stores. You can create pipelines to orchestrate complex data workflows involving data ingestion, data transformation, and data movement tasks.

Data Integration: ADF provides built-in connectors for integrating with various data sources and destinations, including Azure services (such as Azure Blob Storage, Azure SQL Database, Azure Data Lake Storage, Azure Synapse Analytics), on-premises data sources (using self-hosted integration runtimes), and external sources (such as SQL Server, Oracle, Salesforce, Amazon S3, and more).

Data Transformation: With Azure Data Factory, you can perform data transformations using Azure Data Flows, which provide a visual, code-free environment for building data transformation logic. Data Flows allow you to clean, enrich, aggregate, and transform data at scale, using familiar data manipulation techniques.

Data Movement: ADF enables efficient and scalable data movement between various data stores. You can use Copy Activity to move data between supported data sources, with support for parallelism, fault tolerance, and data compression to optimize data transfer performance.

Data Orchestration and Scheduling: Azure Data Factory allows you to schedule and orchestrate data pipelines using triggers. You can define triggers based on schedules (e.g., recurring intervals, specific times) or events (e.g., data arrival, external events) to automate the execution of your data pipelines.

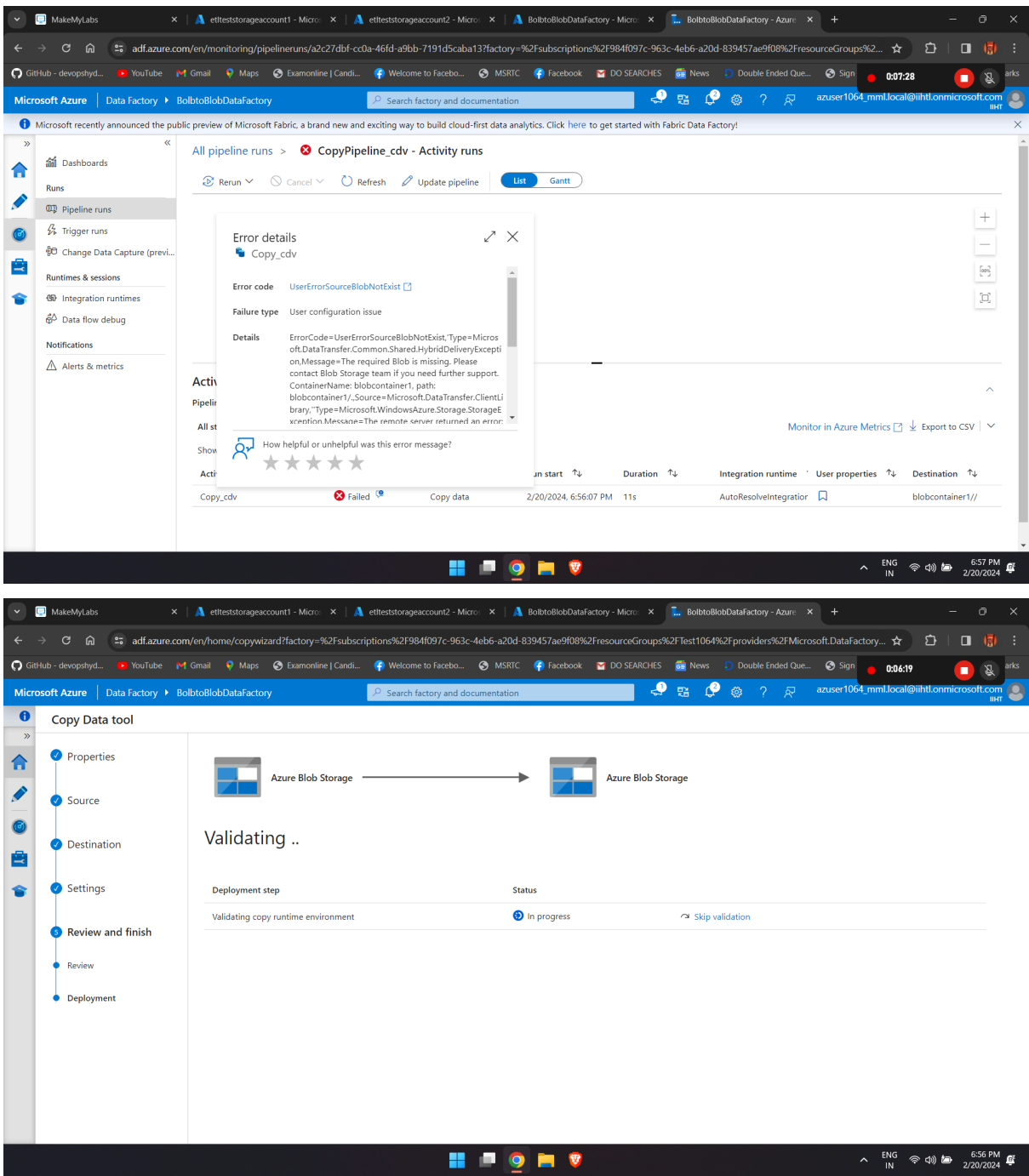
Monitoring and Management: ADF provides monitoring and management capabilities to track the performance and health of your data pipelines. You can monitor pipeline runs, track data lineage, view execution logs, and set up alerts to be notified of pipeline failures or performance issues.

Integration with Azure Services: Azure Data Factory integrates seamlessly with other Azure services, such as Azure Monitor, Azure Active Directory, Azure Key Vault, and Azure DevOps, enabling you to leverage additional capabilities for monitoring, security, and governance.

Security and Compliance: ADF provides features for securing your data pipelines and complying with data governance requirements. It supports encryption at rest and in transit, role-based access control (RBAC), audit logging, and data masking to protect sensitive data and ensure compliance with regulatory standards.

Overall, Azure Data Factory is a powerful data integration service that enables organisations to build scalable, reliable, and efficient data pipelines for ingesting, transforming, and processing data across hybrid and multi-cloud environments. It simplifies the process of managing complex data workflows and enables organisations to unlock insights from their data faster and more effectively.

Hands On :



MakeMyLabs

etteststorageaccount1 - Micro...etteststorageaccount2 - Micro...BolbtoblobDataFactory - Micro...BolbtoblobDataFactory - Azure...

adf.azure.com/en/home/copywizard?factory=%2Fsubscriptions%2F984f097c-963c-4eb6-a20d-839457ae9f08%2FresourceGroups%2FTest1064%2Fproviders%2FMicrosoft.DataFactory...

GitHub - devopshyd...YouTubeGmailMapsExamination | Candi...Welcome to Facebo...MSRTCFacebookDO SEARCHESNewsDouble Ended Que...Sign0:06:16arks

Microsoft AzureData FactoryBolbtoblobDataFactorySearch factory and documentationazuser1064_mml.local@ihl.onmicrosoft.com

Copy Data tool

PropertiesSourceDestinationSettingsReview and finishReviewDeployment

Summary

You are running pipeline to copy data from Azure Blob Storage to Azure Blob Storage.

Destination

Connection nameetlstorageaccount1Dataset nameDestinationDataset_cdvColumn delimiter,Escape character \Quote char "First row as headertrue

Copy settings

Timeout0.12:00:00Retry0Retry interval (sec)30Secure outputfalseSecure inputfalse

< PreviousNext >Cancel

MakeMyLabs

etteststorageaccount1 - Micro...etteststorageaccount2 - Micro...BolbtoblobDataFactory - Micro...BolbtoblobDataFactory - Azure...

adf.azure.com/en/home/copywizard?factory=%2Fsubscriptions%2F984f097c-963c-4eb6-a20d-839457ae9f08%2FresourceGroups%2FTest1064%2Fproviders%2FMicrosoft.DataFactory...

GitHub - devopshyd...YouTubeGmailMapsExamination | Candi...Welcome to Facebo...MSRTCFacebookDO SEARCHESNewsDouble Ended Que...Sign0:06:13arks

Microsoft AzureData FactoryBolbtoblobDataFactorySearch factory and documentationazuser1064_mml.local@ihl.onmicrosoft.com

Copy Data tool

PropertiesSourceDestinationSettingsReview and finishReviewDeployment

Summary

You are running pipeline to copy data from Azure Blob Storage to Azure Blob Storage.

Azure Blob Storage

Azure Blob Storage

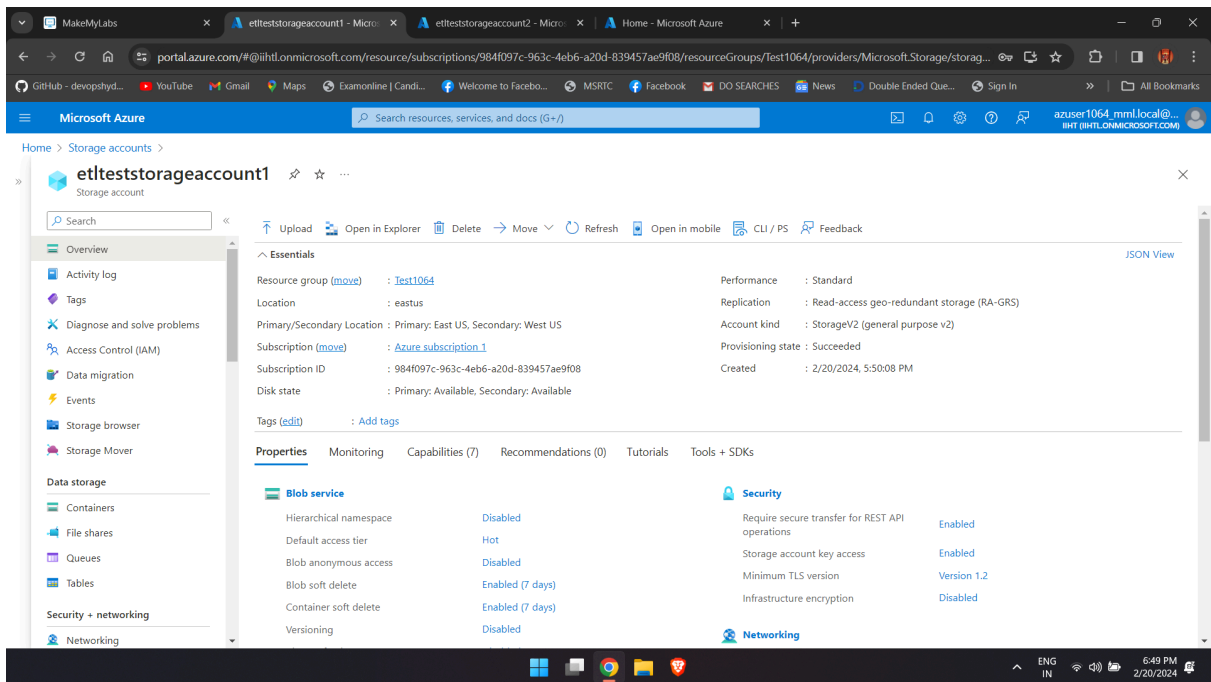
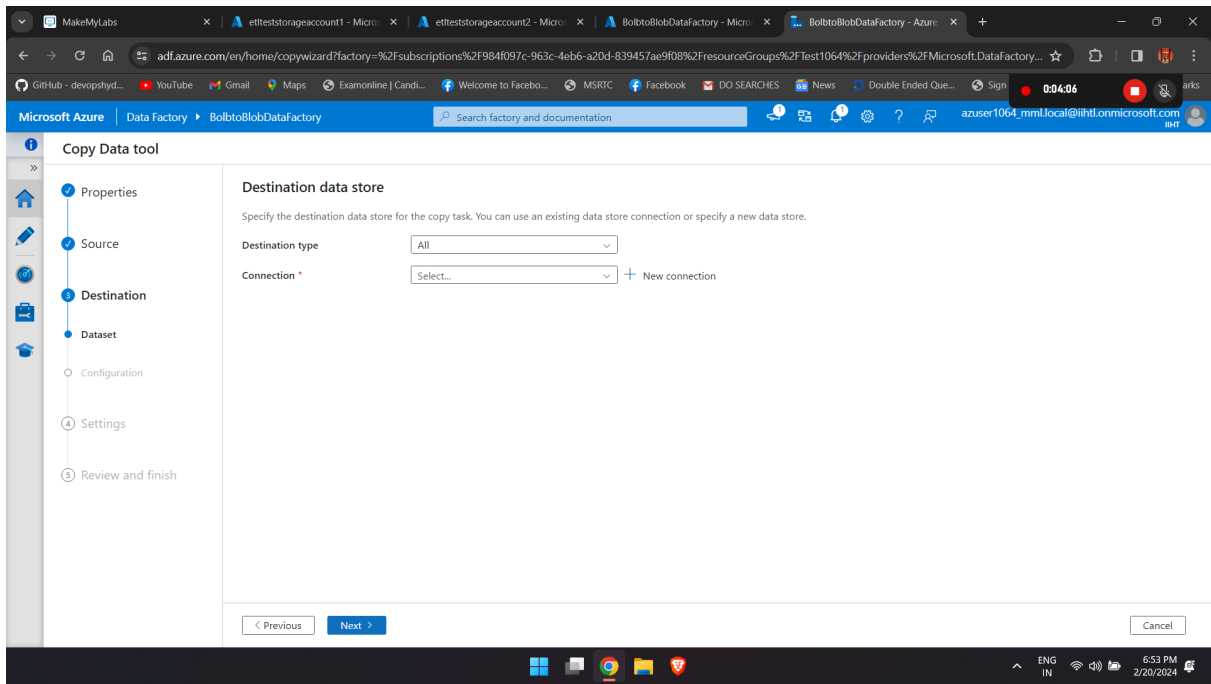
Properties

Task nameCopyPipeline_cdvTask description

Source

Connection nameetlstorageaccount1Dataset nameSourceDataset_cdvColumn delimiter,Escape character \Quote char "First row as headertrue

< PreviousNext >Cancel



Microsoft Azure portal showing the overview of the storage account **etteststorageaccount1**. The account is located in the **Test1064** resource group, East US region, and is associated with **Azure subscription 1**. The account is a **StorageV2 (general purpose v2)** type, created on 2/20/2024 at 5:50:08 PM. The account is in a **Succeeded** provisioning state.

Essentials

- Resource group (move): [Test1064](#)
- Location: **eastus**
- Primary/Secondary Location: **Primary: East US, Secondary: West US**
- Subscription (move): [Azure subscription 1](#)
- Subscription ID: **984f097c-963c-4eb6-a20d-839457ae9f08**
- Disk state: **Primary: Available, Secondary: Available**
- Tags (edit): [Add tags](#)

Properties

Blob service

- Hierarchical namespace: **Disabled**
- Default access tier: **Hot**
- Blob anonymous access: **Disabled**
- Blob soft delete: **Enabled (7 days)**
- Container soft delete: **Enabled (7 days)**
- Versioning: **Disabled**

Security

- Require secure transfer for REST API operations: **Enabled**
- Storage account key access: **Enabled**
- Minimum TLS version: **Version 1.2**
- Infrastructure encryption: **Disabled**

Networking

Microsoft Azure portal showing the **Copy Data tool** configuration page for **Datafactory1064**. The tool is configured to copy data from a **Source data store** to a **Destination**.

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](#)

Destination: **Configuration**

Settings

Review and finish

[Previous](#) [Next](#) [Cancel](#)