# Azure Databricks Coding Challenge

Pratik Wani

# Question 3

❖ **Azure Data factory and Copy Activity:**

- o Azure Data Factory is a cloud-based data integration service that allows us to manage data pipelines for data transformation tasks.

- o It enables us to build workflows to transform and publish data from various sources at scale.

- o One of the main features of Azure Data Factory is the Copy Activity, which is used to copy data between different data stores.

- o We perform copy activity between Source and Sink Data stores

- o We can perform copy activity between different data sources like

    - Blob – Blob

    - Blob – ADLS

    - ADLS – Blob

    - ADLS – ADLS

- o I perform the copy activity by creating two Blob storage '1056hexadeb1' and '1056hexadeb2'

- o '1056hexadeb1' is source storage and '1056hexadeb2' is sink storage

❖ Copy Activity:

   o Created Storage Account and Container

o   Create Data Factory

o Copy Data From one Blob Storage to Another Blob Storage

## Copy Data tool

### Source data store

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

**Properties** ✔

**Source** ②

Source type — All

Connection * — Select...   + New connection

**Dataset**

**Configuration** ○

**Destination** ③

**Settings** ④

**Review and finish** ⑤

[ < Previous ]  [ Next > ]

---

**New connection**

▦ Azure Blob Storage   Learn more ↗

Authentication type
Account key

[ Connection string ] [ Azure Key Vault ]

Account selection method ⓘ
◉ From Azure subscription   ○ Enter manually

Azure subscription ⓘ
Select all

Storage account name *
firsthexadeb1056                          ↻

Additional connection properties
+ New

Test connection ⓘ
◉ To linked service   ○ To file path

Annotations
+ New

> Parameters
> Advanced ⓘ

[ Create ]  [ Back ]            ⚡ Test connection   [ Cancel ]

---

## Copy Data tool

### 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.

**Properties** ✔

**Source** ②

Source type — All

Connection * — ▦ AzureBlobStorage1   ✎ Edit   + New connection

**Dataset**

**Configuration** ○

**Destination** ③

**Settings** ④

**Review and finish** ⑤

**File or folder** *
If the identity you use to access the data store only has permission to subdirectory instead of the entire account, specify the path to browse. Append a slash (/) at the end if the path refers to a folder.

firstcontainer/Azure Databricks Assignment 1.pdf          📁 Browse

**Options**

☐ Binary copy ⓘ

☑ Recursively ⓘ

☐ Enable partitions discovery ⓘ

Max concurrent connections ⓘ

**Filter by last modified**

Start time (UTC)                    End time (UTC)

[ < Previous ]  [ Next > ]                                    [ Cancel ]

o Data Copied Successfully