# Data Ingestion End-to-End Pipeline
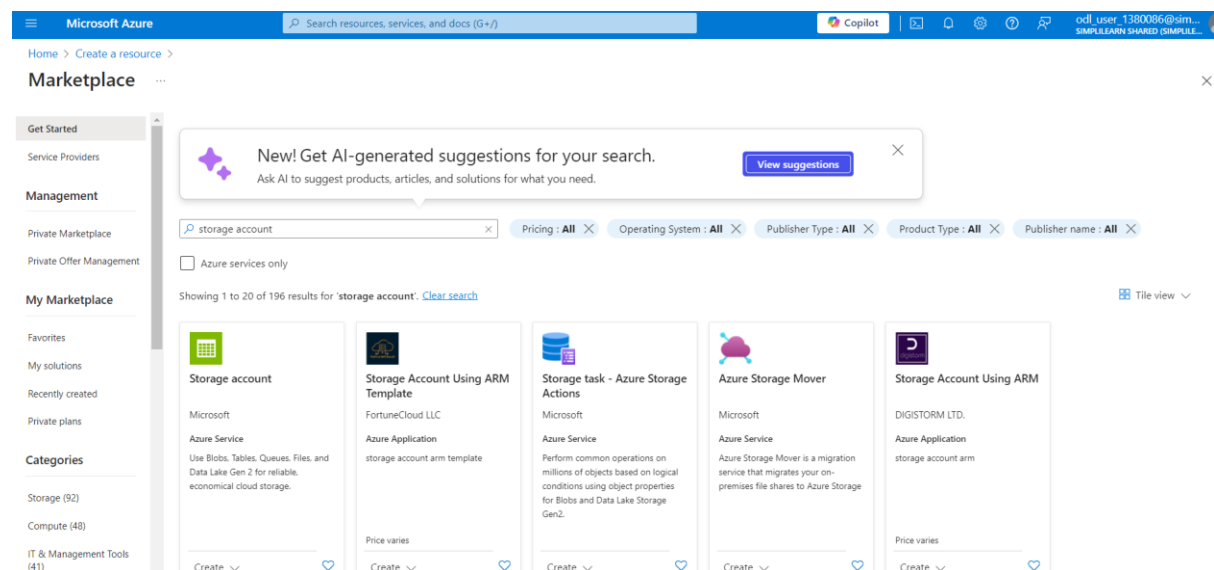
Course-end Project 1

## Description

Your company is looking for a data engineer and is inviting candidates to apply for this position by providing a portal where applicants can add their credentials.

As thousands of candidates have applied for this position, the company has a huge amount of data that it needs to upload to its website. This data is moved to Azure Data Lake Storage parallelly. The company wants to save the contents of all CSV files to Delta Lake of Azure Databricks so that these files can be retrieved and accessed from Azure Databricks when required.

### Step 1: Create a Landing Storage Account in Azure

1. **Log in to the Azure Portal**.
2. **Create a Storage Account**:
   - Go to **"Create a resource"** and search for **"Storage account"**.
   - Click on **"Create"**.
   - Fill in the required details (Subscription, Resource group, Storage account name, Region, Performance, and Replication).
   - Click on **"Review + create"** and then **"Create"**.

## Create a storage account

storage accounts

### Project details

Select the subscription in which to create the new storage account. Choose a new or existing resource group to organize and manage your storage account together with other resources.

Subscription *  Simplilearn SS - 009

Resource group *  databricks-1380086
Create new

### Instance details

Storage account name *  samirlanding

Region *  (US) East US
Deploy to an Azure Extended Zone

Performance *  ◉ Standard: Recommended for most scenarios (general-purpose v2 account)
○ Premium: Recommended for scenarios that require low latency.

Redundancy *  Locally-redundant storage (LRS)

Previous  Next  Review + create

---

Home >

### samirlanding_1718809337855 | Overview
Deployment

🗑 Delete  ⊘ Cancel  ⬆ Redeploy  ⬇ Download  ↻ Refresh

**Overview**
Inputs
Outputs
Template

✅ Your deployment is complete

Deployment name: samirlanding_1718809337855          Start time: 6/19/2024, 8:32:49 PM
Subscription: Simplilearn SS - 009                   Correlation ID: 8917cae2-1052-479c-9a99-956046bb5ca7
Resource group: databricks-1380086

∨ Deployment details

∧ Next steps

Go to resource

Give feedback

↗ Tell us about your experience with deployment

---

## Step 2: Store the CSV Files in the Landing Storage Account

1. **Upload CSV Files**:
   o Navigate to the created storage account.
   o Go to **"Containers"** and create a new container.
   o Upload the CSV files to this container.

## Step 3: Create a Staging Storage Account in Azure

Repeat the steps from Step 1 to create another storage account that will be used for staging.

## Step 4: Create an Azure Data Factory Resource

1. **Create Azure Data Factory**:
   - Go to **"Create a resource"** and search for **"Data Factory"**.
   - Click on **"Create"**.
   - Fill in the required details (Subscription, Resource group, Data Factory name, Version, and Region).
   - Click on **"Review + create"** and then **"Create"**.

## Step 5: Create Linked Services for the Storage Accounts

1. **In Azure Data Factory**:
   - Go to **"Manage"** on the left panel.
   - Under **"Connections"**, click on **"New"**.

o Choose **"Azure Blob Storage"** and configure the linked service for both landing and staging storage accounts by providing the necessary details

## Step 6: Use Azure Databricks as a Part of the ADF Pipeline

1. **Create a Databricks Workspace**:
    - o Go to **"Create a resource"** and search for **"Azure Databricks"**.
    - o Click on **"Create"**.
    - o Fill in the required details and create the workspace.
2. **Create a Databricks Cluster** within the workspace.

## Step 7: Create a Linked Service in ADF for Databricks

1. **In Azure Data Factory**:
   - o Go to **"Manage"** on the left panel.
   - o Under **"Connections"**, click on **"New"**.
   - o Choose **"Azure Databricks"** and configure the linked service by providing the necessary details (e.g., Databricks workspace URL, token).

## Step 8: Create an Azure Data Factory Pipeline

1. **Create a Pipeline**:
   - o Go to **"Author"** on the left panel.
   - o Click on the **"Pipeline"** icon and then **"New pipeline"**.
   - o Add a **"Data Flow"** activity from the activities pane.

The data has to be transposed. Since there's no way to directly transpose a table in ADF, we apply a series of Unpivot and Pivot to Transpose

**Unpivot settings**   Optimize   Inspect   Data preview ●

| | |
|---|---|
| Output stream name * | unpivot2 |
| Description | Unpivots columns into row values and ungroups columns |
| Incoming stream * | File2 ⌄ |

? Help      Learn more ⬈

↻ Reset

1. Ungroup by    **2. Unpivot key**    3. Unpivoted columns

| | |
|---|---|
| Unpivot column name * | candidateName |
| Unpivot column type * | abc  string ⌄ |
| Option * | ◉ Pick column names as values   ◯ Enter values |

---

**Unpivot settings**   Optimize   Inspect   Data preview ●

ungroups columns

| | |
|---|---|
| Incoming stream * | File2 ⌄ |

1. Ungroup by    2. Unpivot key    **3. Unpivoted columns**

| | |
|---|---|
| Column arrangement * | [Normal] [Lateral] |
| Drop rows with null ⓘ | ✓ |
| Columns * | |

| Column name | Type | |
|---|---|---|
| value ⌄ | abc  string ⌄ | ＋ 🗑 |

## Pivot settings  Optimize  Inspect  Data preview ●

| | |
|---|---|
| Output stream name * | pivot2 |
| Description | Pivots row values into columns, groups columns and aggregates data |
| Incoming stream * | unpivot2 |

**1. Group by**    2. Pivot key    3. Pivoted columns

| Columns | Name as | | |
|---|---|---|---|
| abc  candidateName | candidateName | + | 🗑 |

## Pivot settings  Optimize  Inspect  Data preview ●

| | |
|---|---|
| Output stream name * | pivot2 |
| Description | Pivots row values into columns, groups columns and aggregates data |
| Incoming stream * | unpivot2 |

1. Group by    **2. Pivot key**    3. Pivoted columns

| | |
|---|---|
| Pivot key * | abc  Candidate Name |

**Value**

| | | |
|---|---|---|
| Enter value (optional)... | + | 🗑 |

☐ Null value

One of files has Null Values in one of the columns. Use a Derived Column to replace null values with a string

The final Data Flow looks like this:-

## Step 9: Convert CSV Files to Parquet Files in Staging Storage

- o  The Column Names have spaces in between, and therefore can't be directly converted to parquet. Use "replace" function is Select to convert the columns
- o  Publish the pipeline and put a manual trigger to run the pipeline

dataflow1 ● | File5 ● | File4 ● | File3 ● | File2 ● | File1 ●

✓ Validate | ⬤ Data flow debug ✓ | Debug Settings

0
0
8

| File1 | unpivot1 | pivot1 | union1 | union2 | union3 | union4 | unpivot6 | derivedColumn3 | 1 Columns |

pivot2 | pivot3 | pivot4 | pivot5

| File2 | select1 | unpivot2 | pivot2 |

| File3 | derivedColumn1 | unpivot3 | pivot3 |

| File4 | unpivot4 | pivot4 |

| File5 | derivedColumn2 | unpivot5 | pivot5 |

1
0

**Properties**

Genera

Name
dataflo

Descrip

---

Pivot settings | Optimize | Inspect | Data preview ●

Number of rows | ✦ INSERT 5 | ✦ UPDATE 0 | ✕ DELETE 0 | ✦ UPSERT 0 | LOOKUP 0 | ❌ ERROR 0 | TOTAL 5

🔄 Refresh ⌄ | Typecast ⌄ | Modify ⌄ | ▤ Map drifted | ▤ Statistics | ✕ Remove | ⬇ Export to CSV ⌄

| candidateName | CertificationDe... | Education | PhoneNo | PrimarySkill | SecondarySkill |
|---|---|---|---|---|---|
| ✦ Vinod Kumar | null | Bachelors of Engin... | 6221990099 | Javascript, ASP.NET | WCF, Azure |
| ✦ Amit Kumar | DP 203, AZ 104 | Bachelors of Engin... | 7821990099 | .NET, Azure | AWS |
| ✦ Vivek Kumar | Azure Data Scientist | Bachelors of Engin... | 7829999999 | Machine Learning | Data Engineering |
| ✦ Sandip Kumar | Google Cloud Prof... | Bachelors of Engin... | 7999990099 | Java | Google Cloud |
| ✦ Gunjan Kumar | null | MCA | 7821990099 | C# | Java |

---

📊 Data Factory ⌄ | ✓ Validate all | ⬆ Publish all 1

**Factory Resources** ⌄ «

🔍 Filter resources by name   +

▲ Pipelines   1
   ◉ ⬡ pipeline1
▷ Change Data Capture (preview)   0
▲ Datasets   8
   ▦ DelimitedText1
   ▦ File1
   ▦ File2
   ▦ File3
   ▦ File4
   ▦ File5
   ▦ File6
   ▦ Parquet1
▲ Data flows   1
   ⬡ dataflow1
▷ Power Query   0

**Activities** ⌄ «

🔍 data

⌄ Move and transform
   ▤ Copy data
   🔷 Data flow

⌄ Azure Data Explorer
   Azure Data Explorer C...

⌄ Databricks
   Notebook
   Jar
   Python

⌄ Data Lake Analytics
   U-SQL

⌄ General
   ℹ Get Metadata

dataflow1 | File5 | File4 | File3 | File2 | File1

✓ Validate | ▷ Debug ⌄ | ⚡ Add trigger | ⬤ Data flow debug ✓

Data flow ⬚
🔷 Data flow1

General | **Settings** | Parameters 1 | User properties

Data flow *    [ dataflow1 ⌄ ]   ✏ Open   + New
Run on (Azure IR) * ⓘ   [ AutoResolveIntegrationRuntime ⌄ ]
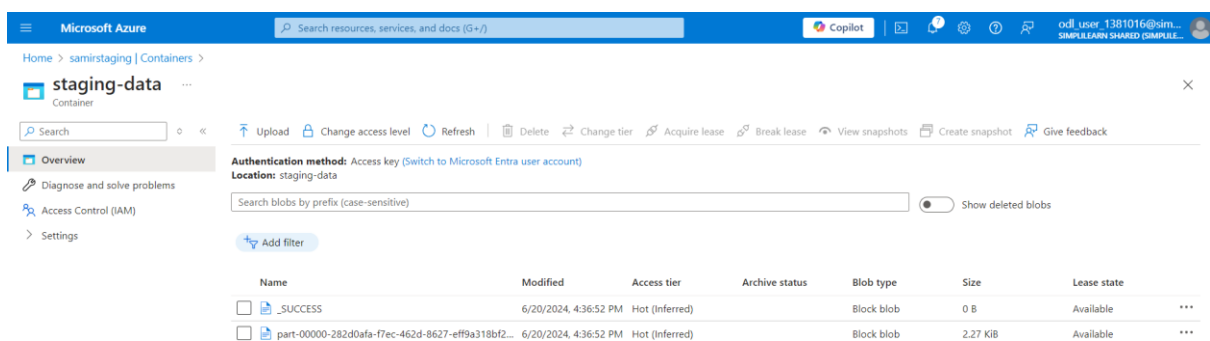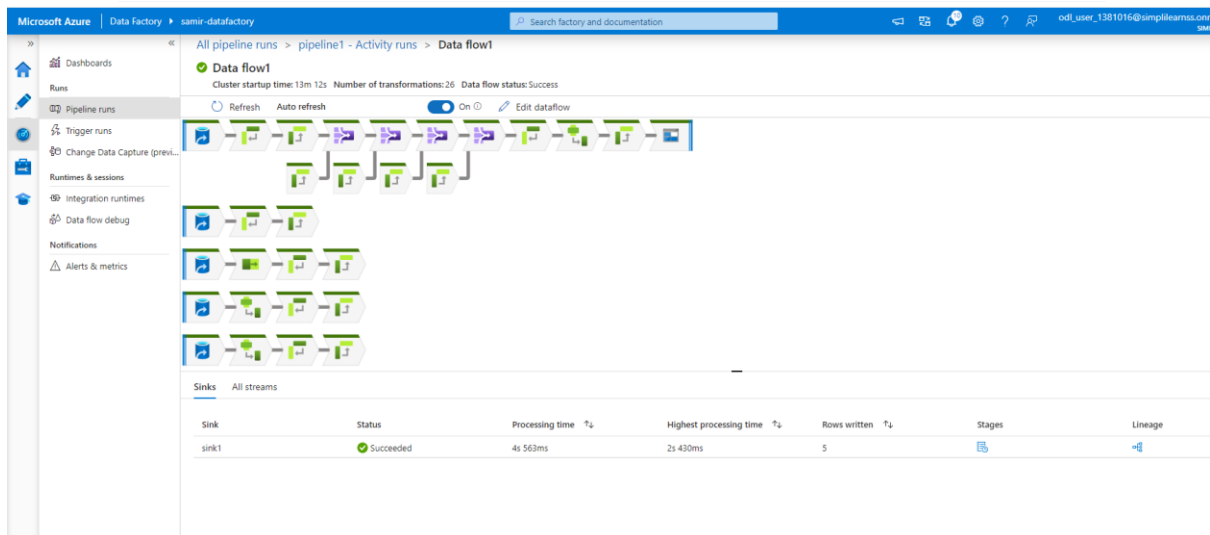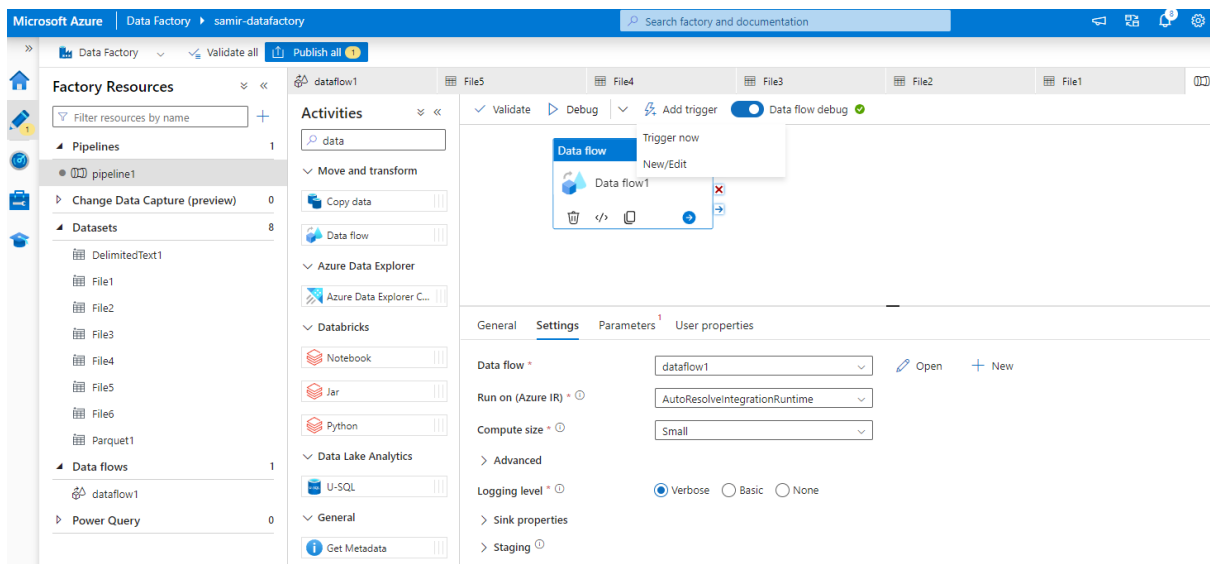Compute size * ⓘ   [ Small ⌄ ]
▷ Advanced
Logging level * ⓘ   ◉ Verbose   ○ Basic   ○ None
▷ Sink properties
▷ Staging ⓘ

The parquet file is stored in the staging area.

## Step 10: Access Parquet Files from the Staging Account in Azure Databricks

1. **In Azure Databricks**:
   - o Mount the staging storage account if needed.
   - o Read the Parquet files into Databricks.

🔑 **samirstaging | Access keys** ☆ ⋯
Storage account

🔍 access ✕ «

👤 Access Control (IAM)

**Security + networking**

🔑 Access keys

🔗 Shared access signature

🔒 Encryption

**Data management**

◎ Lifecycle management

**Settings**

🗄 Configuration

🕐 Set rotation reminder   ↻ Refresh   👥 Give feedback

Access keys authenticate your applications' requests to this storage account. Keep your keys in a secure location like Azure Key Vault, and replace them often with new keys. The two keys allow you to replace one while still using the other.

Remember to update the keys with any Azure resources and apps that use this storage account.
Learn more about managing storage account access keys ⧉

Storage account name

| samirstaging | 📋 |

**key1** ↻ Rotate key

Last rotated: 6/20/2024 (0 days ago)

Key

| ●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●● | Show |

Connection string

| ●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●... | Show |

**key2** ↻ Rotate key

Last rotated: 6/20/2024 (0 days ago)

Key

| ●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●● | Show |

Connection string

| ●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●●... | Show |

---

**Project1-PGPDE** [Python ∨] ☆
File  Edit  View  Run  Help     Last edit was 4 minutes ago   💬 Provide feedback

▶ Run all   ● Samir's Cluster ∨   📅 Schedule

▶  ✓ 07:56 PM (32s)                                          1

```python
# Define your storage account name and key
storage_account_name = "samirstaging"
storage_account_key = 

# Define the container name and mount point
container_name = "staging-data"
mount_point = "/mnt/staging"

# Mount the storage account
dbutils.fs.mount(
    source = f"wasbs://{container_name}@{storage_account_name}.blob.core.windows.net/",
    mount_point = mount_point,
    extra_configs = {f"fs.azure.account.key.{storage_account_name}.blob.core.windows.net": storage_account_key}
)

# List the files in the mounted directory to verify
display(dbutils.fs.ls(mount_point))
```

▶ (3) Spark Jobs

Table ∨  +                                          🔍 ▽ ▢

| | name | size | modificationTime |
|---|---|---|---|
| 1 | -282d0afa-f7ec-462d-8627-eff9a318bf29-c000.snappy.parqu... | part-00000-282d0afa-f7ec-462d-8627-eff9a318bf29-c000.snappy.parqu... | 2322 | 1718890308000 |

# Step 11: Convert the Parquet Files to Azure Databricks Delta Tables

1. **In Azure Databricks**:
   - Convert Parquet to Delta.

```
# part-00000-282d0afa-f7ec-462d-8627-eff9a318bf29-c000.snappy.parquet

df = spark.read.parquet("/mnt/staging/part-00000-282d0afa-f7ec-462d-8627-eff9a318bf29-c000.snappy.parquet")
```

▶ (1) Spark Jobs

▼ ▦ df: pyspark.sql.dataframe.DataFrame
        candidateName: string
        CertificationDetails: string
        Education: string
        PhoneNo: string
        PrimarySkill: string
        SecondarySkill: string

```
df.write.format("delta").save("/mnt/delta/samir-delta-table")
```

▶ (4) Spark Jobs

## Step 12: Store and Visualize the Data from Azure Databricks Delta Tables

1. **Create Delta Tables**:
    o  Register the Delta table in the Databricks metastore.

2. **Visualize Data**:
    o  Use Databricks notebooks to visualize the data

```
spark.sql("CREATE TABLE candidates USING DELTA LOCATION '/mnt/delta/samir-delta-table'")

res = spark.sql("SELECT * FROM candidates")
display(res)
```

▶ (1) Spark Jobs

▶ ▦ res: pyspark.sql.dataframe.DataFrame = [candidateName: string, CertificationDetails: string ... 4 more fields]

Table ▾  +

| | candidateName | CertificationDetails | Education | PhoneNo | PrimarySkill | SecondarySkill |
|---|---|---|---|---|---|---|
| 1 | Amit Kumar | DP 203, AZ 104 | Bachelors of Engineeri... | 7821990099 | .NET, Azure | AWS |
| 2 | Gunjan Kumar | null | MCA | 7821990099 | C# | Java |
| 3 | Sandip Kumar | Google Cloud Professional Cloud Architect | Bachelors of Engineeri... | 7999990099 | Java | Google Cloud |
| 4 | Vinod Kumar | null | Bachelors of Engineeri... | 6221990099 | Javascript, ASP.NET | WCF, Azure |
| 5 | Vivek Kumar | Azure Data Scientist | Bachelors of Engineeri... | 7829999999 | Machine Learning | Data Engineering |

⬇ 5 rows | 5.80 seconds runtime                                    Refreshed 6 minutes ago