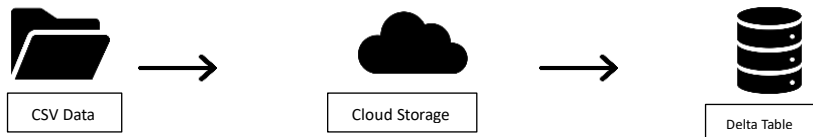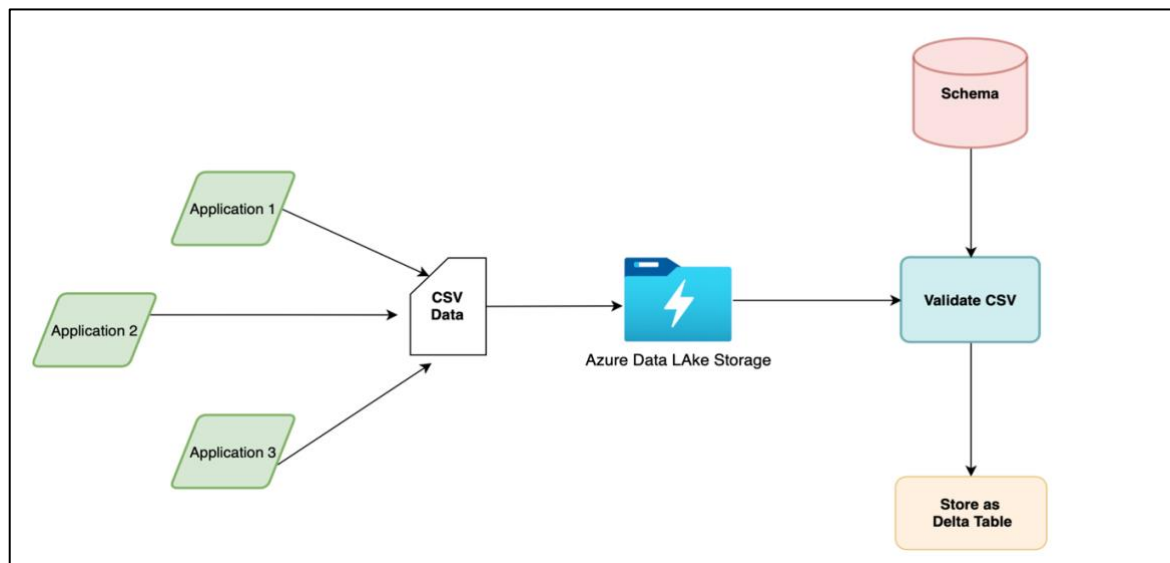# Wells Fargo – Data platform

## Business Case:

Wells Fargo has multiple internal applications that send huge amounts of data (*petabytes*) in CSV format **daily** in the company Azure cloud. We are required to perform **data/schema validation** on this incoming data and this needs to be stored in a **delta table** for downstream systems to make use for analysis.



## Task:
1. Validation needs to be applied as follows:
    1.1. Need to check for **duplicate rows**. If file contains duplicate rows, it's rejected.
    1.2. Need to **validate the data format** for all the data fields, the date column names, and the date format is stored in the **Azure SQL server**, if validation fails the file is rejected.
2. Need to move all rejected files to Rejected folder.
3. Need to move all valid files to Staging folder.
4. Write all the validated files as a delta table in **Azure Databricks**
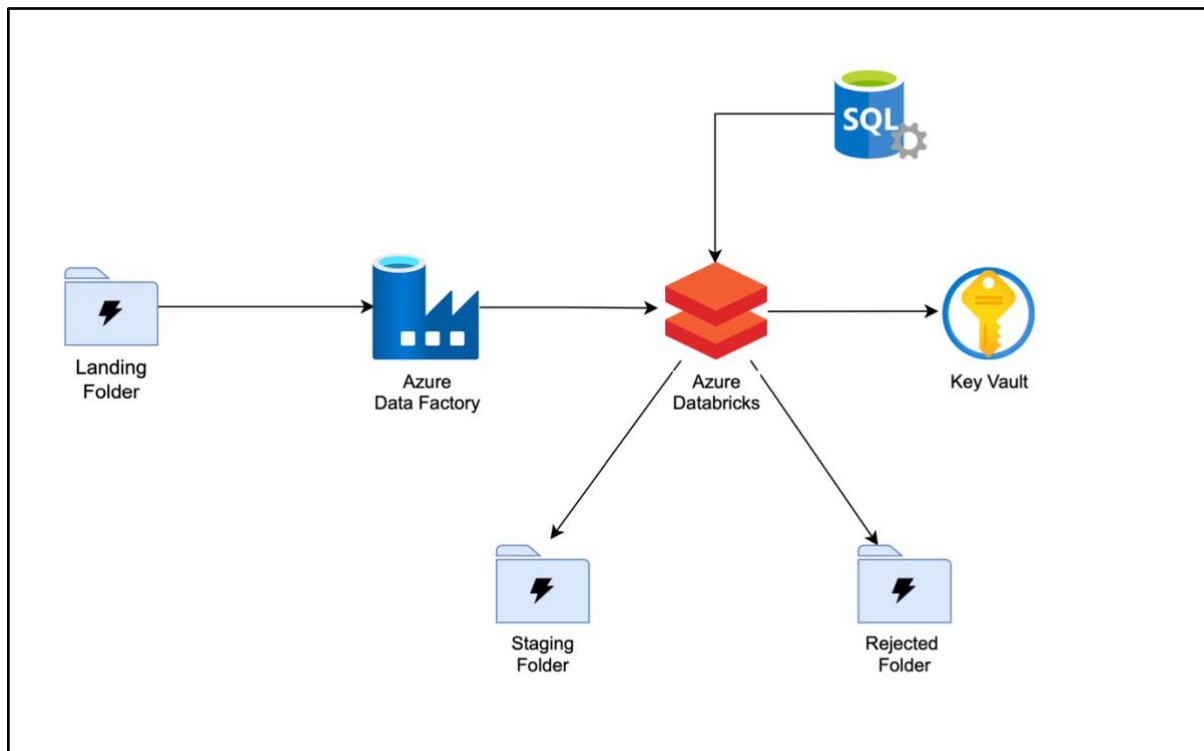
## High-level Scope:



(*A diagrammatic representation of the high-level workflow of the architecture*)

## Architecture Diagram:

1. Input data will come into the **Azure Data Lake Storage** account in a landing folder.
2. Then the data will be validated, by task scheduling in **azure data factory**.
3. ADF will pull the input data and invoke **Databricks** (spark) and create a delta table.

4. To apply validation rules, the **schema** will be fetched by the **SQL server**.
5. The file in landing folder will be validated against obtained schema and categorized.
6. credentials required to connect these services to each other will be stored in **azure key vault.**



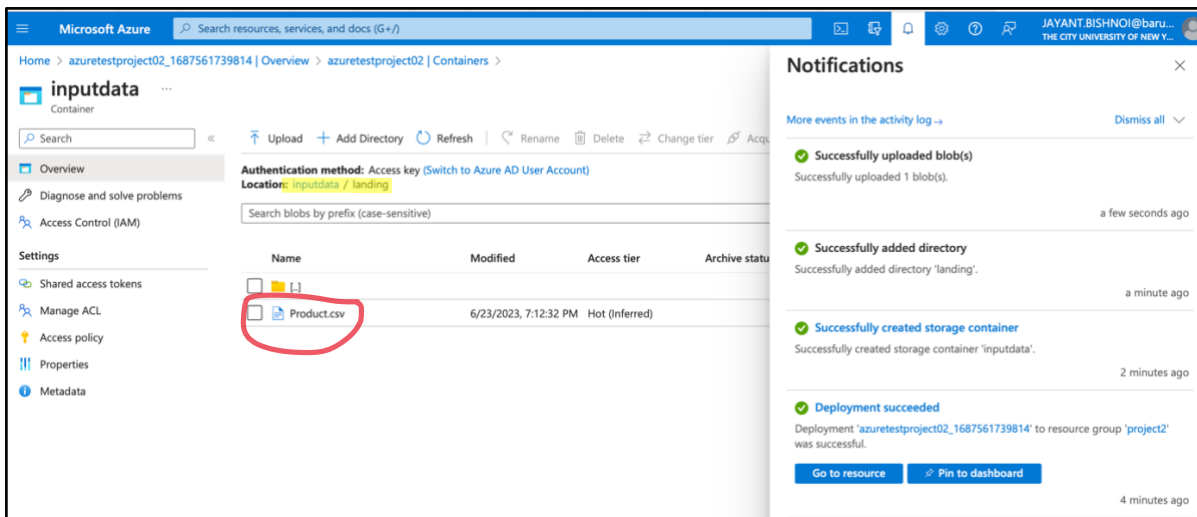(*A diagrammatic representation of the architecture*)

# Phase 1 - Setting up Framework

## Create Azure Data Lake Storage account:
(*In a real-world scenario, the files in landing folder will be provided by third party apps, but since this is a synthesized business use case, we are assuming that data is already present in the landing folder*)

### Storage account & Resource Group:

1. Search for storage account in azure portal and click on create new.
2. Create new resource group called 'project 2'
3. Enable '**hierarchical namespace'**, to make it into an **Azure data lake storage gen2** account.
4. Review+Create to create the storage account.
5. Inside the storage account click in '**container**' to create a storage container and call it 'input'
6. Inside 'input' container we will create our **landing folder** where we want out s3 data to show.
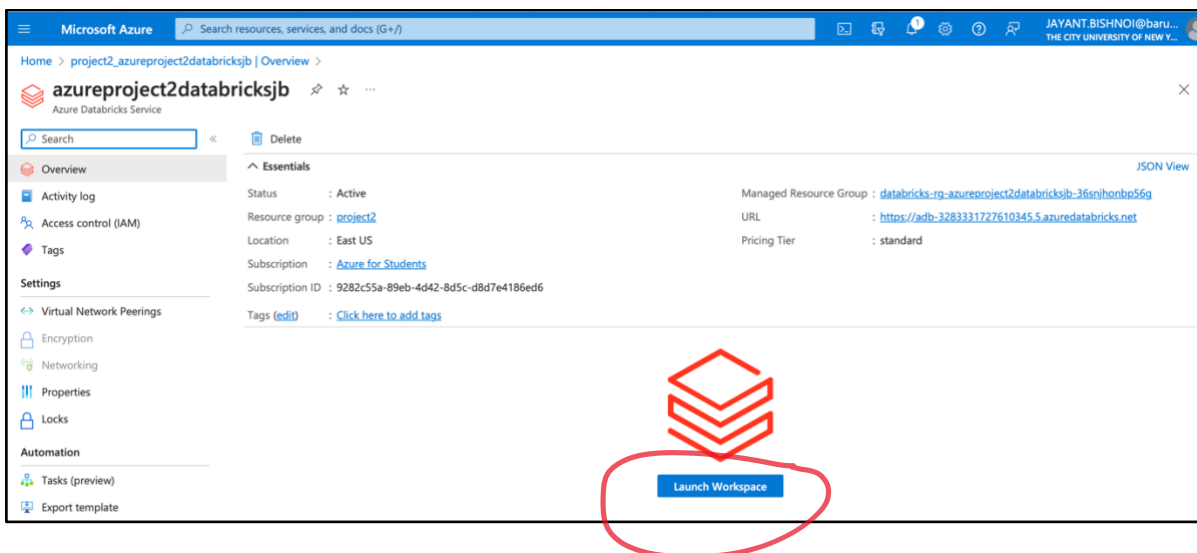
*(Storage account with the input/landing folders and relevant data)*

According to the architecture diagram, we need to set up an azure data factory account, but without a destination to send the data to it is not very practical, so we will set up the databricks account as destination first...

## Create Azure Databricks Account Workspace:

Azure Databricks:
1. Search for 'databricks' in azure portal and click on create button.
2. Give a suitable name and assign a subscription and resource group (project2)
3. Keep the pricing tier as 'standard' for the purpose of the project.
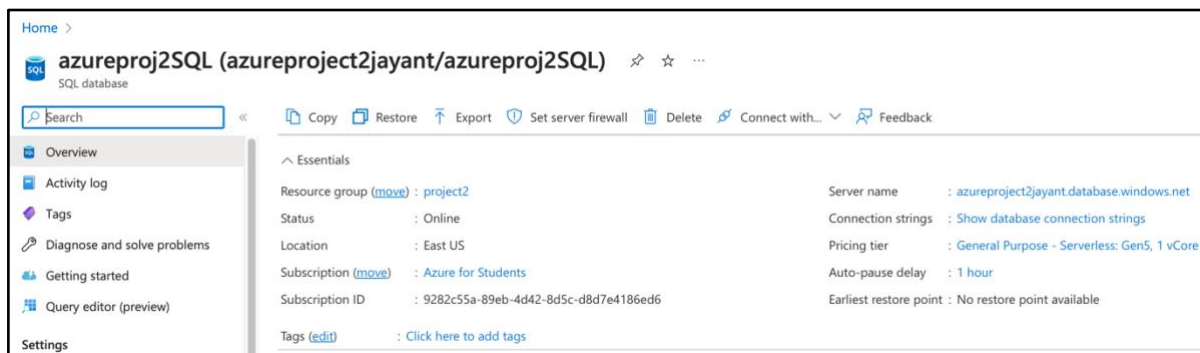4. Review and create



Before using ADF to call the databricks account we need to create a 'notebook' inside it, so that we can write the code for the validation and get that notebook called by ADF.

## Create Azure SQL Server:
1. Create > subscription > resource group > database name.
2. Create new server and give name > SQL authentication > define username and password > compute + storage (basic tier)> LRS> public endpoint> allow access to DB> review +create.

(Go to resource and go to query editor, give login and password)
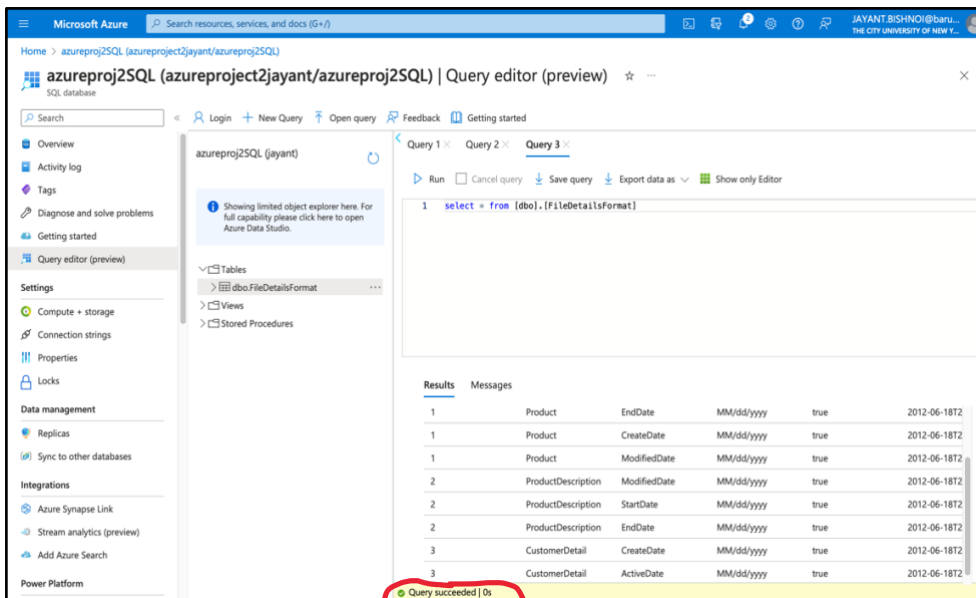


## Create schema table:

**SSM – SQL Server Management System is an IDE, is an editor where you can run queries and operations related to a database.**
**If working on MAC, the alternative is to use AzureSQL- Edge container on docker and connect it to azure data studio.**
**I am choosing to use the AzureDB to make the data table and load the data.**

In Azure SQl DB, go to the query editor and write the following code to make the schema table for the incoming files.

```sql
CREATE TABLE [FileDetailsFormat](
    [FileNo] [int] NOT NULL,
    [FileName] [nvarchar](100) NOT NULL,
    [ColumnName] [nvarchar](100) NULL,
    [ColumnDateFormat] [nvarchar](108) NULL,
    [ColumnIsNull] [nvarchar](100) NULL,
    [ModifiedDate] [datetime] NOT NULL
)
INSERT [dbo].[FileDetailsFormat] ([FileNo], [FileName], [ColumnName], [ColumnDateFormat], [ColumnIsNull], [ModifiedDate]) VALUES (1, N'Product', N'StartDate', N'MM-dd-yyyy', N'true', CAST(N'2012-06-18T22:34:09.000' AS DateTime))
INSERT [dbo].[FileDetailsFormat] ([FileNo], [FileName], [ColumnName], [ColumnDateFormat], [ColumnIsNull], [ModifiedDate]) VALUES (1, N'Product', N'EndDate', N'MM/dd/yyyy', N'true', CAST(N'2012-06-18T22:34:09.000' AS DateTime))
INSERT [dbo].[FileDetailsFormat] ([FileNo], [FileName], [ColumnName], [ColumnDateFormat], [ColumnIsNull], [ModifiedDate]) VALUES (1, N'Product', N'CreateDate', N'MM/dd/yyyy', N'true', CAST(N'2012-06-18T22:34:09.000' AS DateTime))
INSERT [dbo].[FileDetailsFormat] ([FileNo], [FileName], [ColumnName], [ColumnDateFormat], [ColumnIsNull], [ModifiedDate]) VALUES (1, N'Product', N'ModifiedDate', N'MM/dd/yyyy', N'true', CAST(N'2012-06-18T22:34:09.000' AS DateTime))
INSERT [dbo].[FileDetailsFormat] ([FileNo], [FileName], [ColumnName], [ColumnDateFormat], [ColumnIsNull], [ModifiedDate]) VALUES (2, N'ProductDescription', N'ModifiedDate', N'MM/dd/yyyy', N'true', CAST(N'2012-06-18T22:34:09.000' AS DateTime))
INSERT [dbo].[FileDetailsFormat] ([FileNo], [FileName], [ColumnName], [ColumnDateFormat], [ColumnIsNull], [ModifiedDate]) VALUES (2, N'ProductDescription', N'StartDate', N'MM/dd/yyyy', N'true', CAST(N'2012-06-18T22:34:09.000' AS DateTime))
INSERT [dbo].[FileDetailsFormat] ([FileNo], [FileName], [ColumnName], [ColumnDateFormat], [ColumnIsNull], [ModifiedDate]) VALUES (2, N'ProductDescription', N'EndDate', N'MM/dd/yyyy', N'true', CAST(N'2012-06-18T22:34:09.000' AS DateTime))
INSERT [dbo].[FileDetailsFormat] ([FileNo], [FileName], [ColumnName], [ColumnDateFormat], [ColumnIsNull], [ModifiedDate]) VALUES (3, N'CustomerDetail', N'CreateDate', N'MM/dd/yyyy', N'true', CAST(N'2012-06-18T22:34:09.000' AS DateTime))
INSERT [dbo].[FileDetailsFormat] ([FileNo], [FileName], [ColumnName], [ColumnDateFormat], [ColumnIsNull], [ModifiedDate]) VALUES (3, N'CustomerDetail', N'ActiveDate', N'MM/dd/yyyy', N'true', CAST(N'2012-06-18T22:34:09.000' AS DateTime))
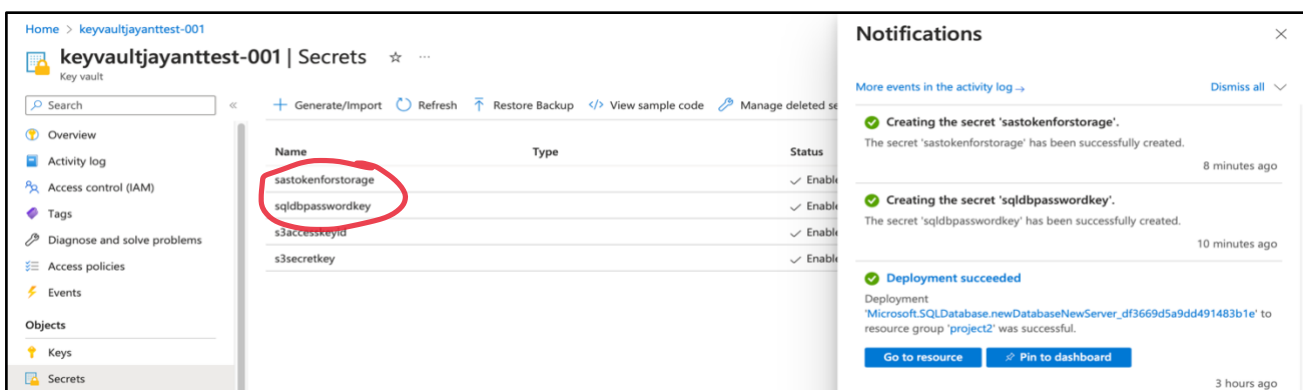```

# Create Key vault and store secrets:

*"We require SQL DB password & SAS token of storage account to able to connect to Databricks, we will store these secrets in key vault"*

## Azure key vault:
1. Go to azure portal, search for 'Key vault' and hit create, provide resource group (project
2. Give a suitable name, review +create.

## Secrets:
1. Go to 'Secrets' and hit 'generate' to store the SQL 'DB password'
2. Create another secret for the 'shared account signature (sas) token'
3. Go to storage account and click on shared access signature and give all permissions for files and containers and objects and click create.
4. Copy the SAS token and paste it as a secret value when you are creating the secret in key vault and hit create.
5. Now Databricks will be able to access the SQL DB and the landing and staging folders with the SAS token by using these secrets in the key vault.
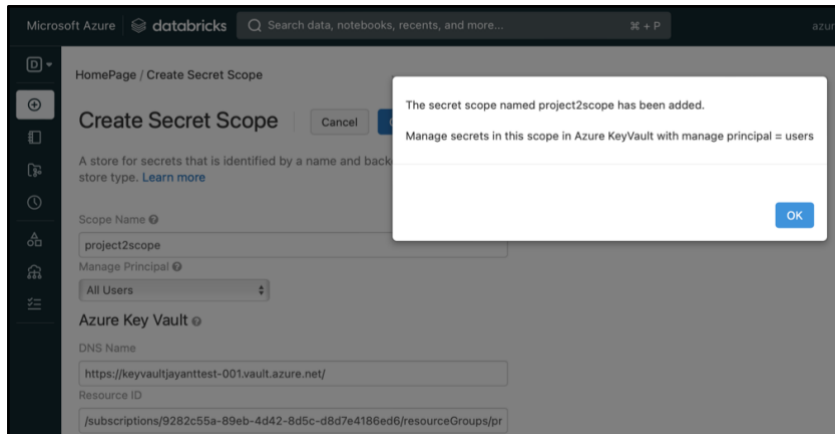


# Create a secrets scope in Databricks:

*"We require Databricks to communicate with SQL DB to get schema, but it cannot go to the key vault to fetch the secrets directly. So, we define a secret scope in databricks which in turn communicates with key vault by deploying a notebook"*
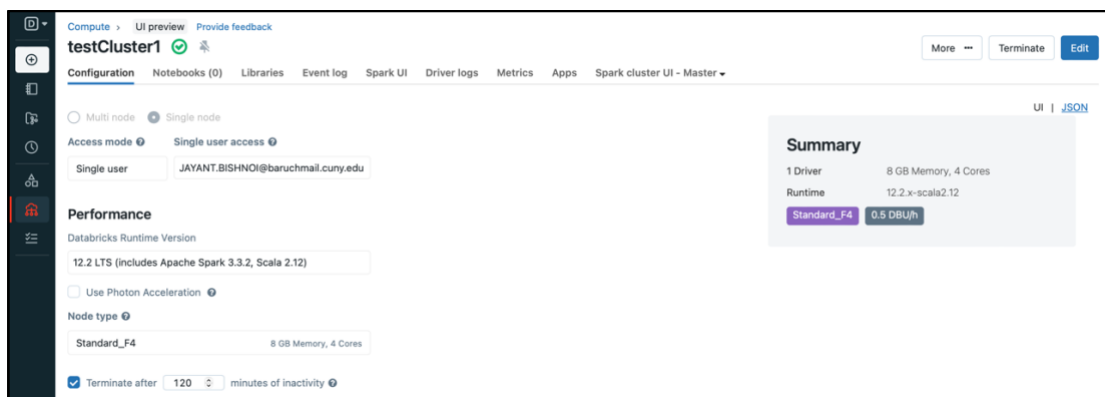
Databricks → Secret Scope → Key Vault

Secret Scope:

1. Go to Databricks account that you have created and 'Launch Workspace'
2. At the end of your URL in search bar add **'secrets/createScope'**, get redirected to the create secret scopes page.
3. Give a suitable name < give permission to all users < paste key vault URI under DNS name & key vault resource ID under Resource ID < hit create.



Create Cluster in Databricks:

1. Go to compute tab on the home page of your workspace and create cluster.
2. Give a suitable name and select runtime.
3. Choose worker type and driver type, we are selecting compute optimized standard F4.
4. Select the number of worker nodes that you need.
5. Create cluster.



# Phase 2 – Implementing Validation

Creating Notebook in Databricks:

Inside the notebook we need to follow a series of steps:
1. Getting the **input** file
2. **Fetching the schema** from database, so we make a connection to the database, with the help of a few database properties (usernames & passwords, URL, port number etc...)

We have used the scope we defined in databricks earlier to access the password from key vault to make the connection to the database, which is successful.

3. We **mount the storage account** as ADF will only fetch the file name and the real access to the files will come from the storage account (SAS token, storage info.)



4. Then at last we **implement the validation** logic to check for duplication and data format.

1. For implementing the **elimination of duplicate values**, we will take the total count of rows and distinct count of rows, if both values are equal, we can concur that there is no duplication of rows.

2. For finding **correct date format**, we will filter all records where filename is not what you require, in our case (product), we only need column name and date format, iterate row by row.

3. We use a **'todate'** function to pass the values in the rows and if a value is returned, we know that it's the correct format and value is <u>not null</u>.

```python
df1 = spark.read.csv('/mnt/landing/'+fileName, inferSchema=True, header=True)
#display(df1)

# Rule
errorFlag=False
errorMessage = ''
totalcount = df1.count()
print(totalcount)
distinctCount = df1.distinct().count()
print(distinctCount)
if distinctCount !=totalcount:
    errorFlag = True
    errorMessage = 'Duplication Found. Rule 1 Failed'
print(errorMessage)

# Rule 2
df2 =
df.filter(df.FileName==fileNameWithoutExt).select('ColumnName','ColumnDateFormat' )
rows = df2.collect()
for r in rows:
    colName = r[0]
    colFormat =r[1]
    print(colName, colFormat)
    #display(df1.filter(F.to_date(colName, colFormat).isNull() ==True))
    formatCount =df1.filter(F.to_date(colName, colFormat).isNotNull()
==True).count()
    if formatCount == totalcount:
        errorFlag = True
        errorMessage = errorMessage +' DateFormate is incorrect for {}
'.format(colName)
    else:
        print('All rows are good for ', colName)
print(errorMessage)

if errorFlag:
    dbutils.fs.mv('/mnt/landing/'+fileName,'/mnt/rejected/'+fileName )
    dbutils.notebook.exit('{"errorFlag": "true", "errorMessage":"'+errorMessage
+'"}')
else:
    dbutils.fs.mv('/mnt/landing/'+fileName,'/mnt/staging/'+fileName )
    dbutils.notebook.exit('{"errorFlag": "false", "errorMessage":"No error"}')
```
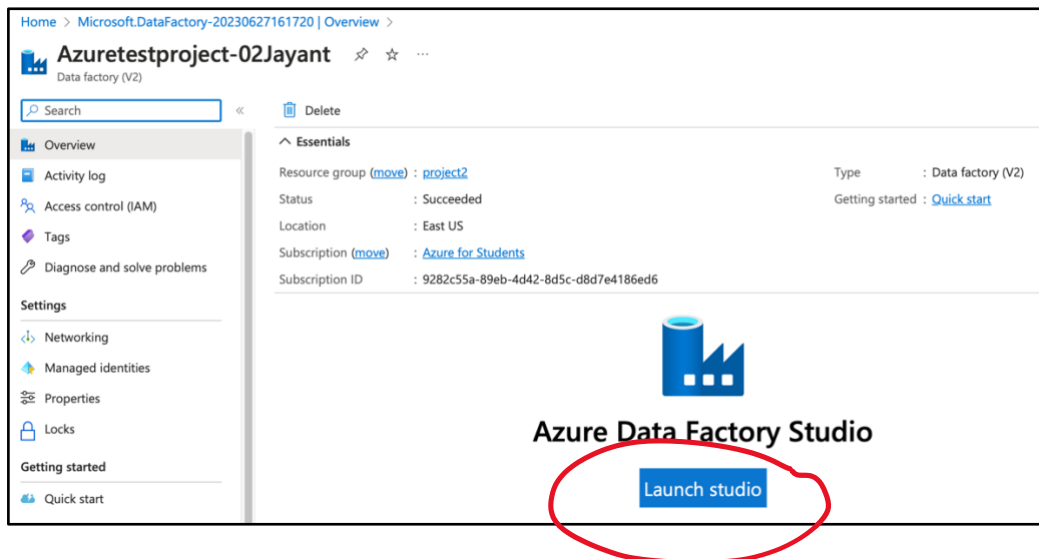
Phase 3 – Data Migration

## Creating Azure Data Factory:

1. Go to azure portal, search for data factories, and hit create, provide resource group that we created (project 1).
2. Give a suitable name, give Git configuration later and hit review + create.
3. Go to resource and launch ADF, where we will create pipelines to migrate data.



## Creating Linked Service to Databricks in Azure Data Factory:

*"We need to connect to databricks account so we will create a linked service"*

1. Go to ADF < launch studio, go to the manage tab and create new linked service.
2. In the compute tab search for Azure databricks, give a suitable name & description.
3. Choose 'AutoResolveIntegrationRuntime' < choose subscription < databricks workspace
4. Select 'New job cluster'

    4.1. Create access token:
        4.1.1. Go to user settings tab in the databricks homepage and click on < 'Generate Token' < give token description < create and copy token.

    4.2. Create secret and access policy in key vault:
        4.2.1. Following industry best practices, we store the databricks token as a secret in azure key vault
        4.2.2. Go to your key vault < secrets < generate < give name and paste the databricks token < create

*"Now we need to give Azure data Factory permission to access this key vault."*

    4.3. Create access Policy:
        4.3.1. In the key vault < access policy < Key, Secret & Certificate management.
        4.3.2. In the principal add the name of your Azure Data Factory account < add

*"Go back to your linked service and click now on key vault instead of access token."*

5. Create Linked service for key vault:
6. Create a new linked service for this key vault < select key vault name < Test connection    <create.

7. Select the name of your newly created linked service.
8. Select secret name "databrickstoken"
9. Select cluster version and cluster node type.

10. Select the number of workers.
11. Test connection and create.



## Creating Linked Service to Datalake storage in Azure Data Factory:

*"ADF need to connect to both landing folder and databricks, so we set up another linked service"*

Go to ADF < click on manage < create new.
Give a suitable name and description.
Select Azure data lake gen2 and give storage account name.
Test connection and hit create.

## Creating Pipeline to call the Databricks notebook:

1. Go to author tab in ADF < create new pipeline.
2. Search for 'notebook' < select databricks notebook and select the associated linked service.
3. Test the connection < under settings tab < provide notebook path

4. *"We want databricks to receive the name of the file from ADF, so we need to add a **base parameter**"*

5. Add a Base parameter:
   5.1. Create base parameter < under name write 'fileName' since that is how the code is written

6. Add a trigger:
   6.1.1. Since we want the pipeline execution to happen as soon as a file comes in landing    folder, we need add a storage trigger.
   6.1.2. Go to the add trigger button and add new trigger <give a suitable name.
   6.1.3. Type < storage events
   6.1.4. Select storage account name and container name (input data)

6.1.5. Blob path begins with 'landing/', so that any file in landing folder triggers pipeline.

6.1.6. Blob path ends with '.csv' to only execute for csv files.

6.1.7. Continue and create.

7. <u>Add a pipeline parameter:</u>

   7.1. Click outside notebook and click on parameters tab < create new.

   7.2. Under name write 'fileName'

   7.3. Go to trigger and edit value **'@triggerBody.fileName'** , this function will give the actual name of the file.

   7.4. Now go to the notebook parameter, we will pass the pipeline parameter in the value under the settings tab.

*"We are fetching the name of the file from trigger and feeding it to pipeline parameter → Notebook parameter, so that the code in our notebook is relevant for any file that is uploaded automatically."*

8. Publish.





*(Staging folder has been created automatically, moving the file from landing folder)*

# THE END