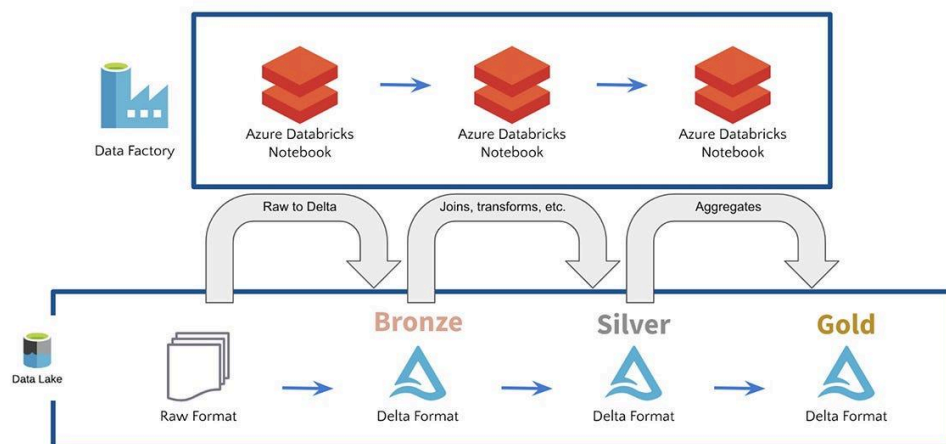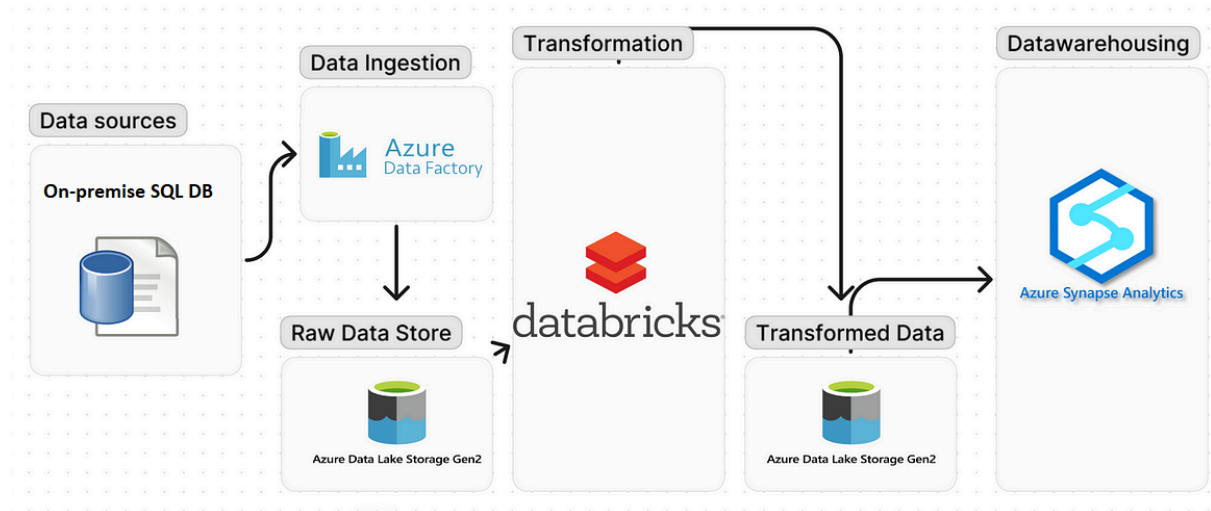# Project 1 and 2

## Project Architecture:





## Part 1 - Dataset Overview:

AdventureWorks database supports standard online transaction processing scenarios for a fictitious bicycle manufacturer - Adventure Works Cycles. Scenarios include Manufacturing, Sales, Purchasing, Product Management, Contact Management, and Human Resources.

```sql
SELECT TOP (1000) [ProductCategoryID]
      ,[ParentProductCategoryID]
      ,[Name]
      ,[rowguid]
      ,[ModifiedDate]
  FROM [AdventureWorksLT2017].[SalesLT].[ProductCategory]
```

## Create Database login

```sql
CREATE LOGIN mbr WITH PASSWORD = 'Vb1T3D1G3ROxi0s';
create user mbr for login mbr
```

## Give User access:



## Store secrets in Azure Key Vault:

**Part 2 - Data Ingestion:** From On-premise to Azure Data Lake Storage (Bronze container) (SQL DB to Parquet Format)

getschema.sql



Folder Structure in Data Lake:

```
bronze/Schema/Tablename/Tablename.parquet
Ex. bronze/SalesLT/Address/Address.parquet
```

Mahesh Raut - Project - 31 May 2024

Source dataset - sql db



Pipeline build for looping each table in sql db to copy all tables in bronze container

## Trigger pipeline once - copy all tables



## Check bronze container



## Part 3 - Data Transformation

## Azure Databricks Overview

Create Compute Cluster

Clusters / New Compute   UI preview   Provide feedback

## data_transformation ✎

| Standard_DS3_v2 | 14 GB Memory, 4 Cores | ∨ |   ❓

☑ Terminate after [ 15 ] minutes of inactivity ❓

## Tags ❓

Add tags

| Key | Value | Add |

> Automatically added tags

▼ Advanced options

**Azure Data Lake Storage credential passthrough** ❓
☑ Enable credential passthrough for user-level data access

Spark    Logging    Init Scripts

**Spark config** ❓

```
spark.master local[*, 4]
spark.databricks.cluster.profile singleNode
```

[ Create Cluster ]   [ Cancel ]

**Storagemount notebook** - bronze, silver, gold



```python
configs = {
    "fs.azure.account.auth.type": "CustomAccessToken",
    "fs.azure.account.custom.token.provider.class":
spark.conf.get("spark.databricks.passthrough.adls.gen2.tokenProviderClas
sName")
}

dbutils.fs.mount(
    source = "abfss://bronze@mrkdatakagen2.dfs.core.windows.net/",
    mount_point = "/mnt/bronze",
    extra_configs = configs
)

dbutils.fs.ls("/mnt/bronze/SalesLT/")
```

**Silvermount**

**Goldmount**

```
1    configs = {
2      "fs.azure.account.auth.type": "CustomAccessToken",
3      "fs.azure.account.custom.token.provider.class": spark.conf.get("spark.databricks.passthrough.adls.gen2.
       tokenProviderClassName")
4    }
5
6    # Optionally, you can add <directory-name> to the source URI of your mount point.
7    dbutils.fs.mount(
8      source = "abfss://gold@mrkdatalakegen2.dfs.core.windows.net/",
9      mount_point = "/mnt/gold",
10     extra_configs = configs)

Out[5]: True
```

**Bronze to silver notebook**
Convert ModifiedDate column with Date Column

| ModifiedDate |
|---|
| 2002-06-01 00:00:00.000 |
| 2002-06-01 00:00:00.000 |
| 2002-06-01 00:00:00.000 |
| 2002-06-01 00:00:00.000 |
| 2002-06-01 00:00:00.000 |
| 2002-06-01 00:00:00.000 |
| 2002-06-01 00:00:00.000 |
| 2002-06-01 00:00:00.000 |
| 2002-06-01 00:00:00.000 |
| 2002-06-01 00:00:00.000 |
| 2002-06-01 00:00:00.000 |
| 2002-06-01 00:00:00.000 |

------------- to ---------------

| ModifiedDate |
|---|
| 2002-06-01 |
| 2002-06-01 |
| 2002-06-01 |
| 2002-06-01 |
| 2002-06-01 |
| 2002-06-01 |
| 2002-06-01 |
| 2002-06-01 |
| 2002-06-01 |
| 2002-06-01 |
| 2002-06-01 |
| 2002-06-01 |

```python
dbutils.fs.ls('mnt/bronze/SalesLT/')

dbutils.fs.ls('mnt/silver/')

input_path = '/mnt/bronze/SalesLT/Address/Address.parquet'

df = spark.read.format('parquet').load(input_path)

display(df)

from pyspark.sql.functions import from_utc_timestamp, date_format
from pyspark.sql.types import TimestampType

df = df.withColumn("ModifiedDate",
date_format(from_utc_timestamp(df["ModifiedDate"].cast(TimestampType()),
"UTC"), "yyyy-MM-dd"))
```

**Doing transformation for all tables**

```python
table_name = []
for i in dbutils.fs.ls('mnt/bronze/SalesLT/'):
        table_name.append(i.name.split('/')[0])
```

```python
from pyspark.sql.functions import from_utc_timestamp, date_format
from pyspark.sql.types import TimestampType

# Loop through table names
for i in table_name:
  # Construct path to the Parquet file
  path = '/mnt/bronze/SalesLT/' + i + '/' + i + '.parquet'

  # Read data from Parquet file
  df = spark.read.format('parquet').load(path)

  # Get list of columns in the DataFrame
  column = df.columns

  # Loop through each column
  for col in column:
    # Check if the column name contains "Date" or "date"
    if "Date" in col or "date" in col:
      # Convert the column to timestamp format (assuming it's not
already a timestamp)
      df = df.withColumn(col,
date_format(from_utc_timestamp(df[col].cast(TimestampType()), "UTC"),
"yyyy-MM-dd"))

  # Construct path to write the transformed data
  output_path = '/mnt/silver/SalesLT/' + i + '/'

  # Write the transformed DataFrame to Delta format in overwrite mode
  df.write.format('delta').mode("overwrite").save(output_path)

display(df)
```

Updated column:



**Silver container**

## Silver to gold notebook

Transform the column names for all tables (ex. ProductId to Product_Id)



```
for name in table_name:
    path = '/mnt/silver/SalesLT/' + name
    print(path)
    df = spark.read.format('delta').load(path)

    # Get the list of column names
    column_names = df.columns

    for old_col_name in column_names:
        # Convert column name from ColumnName to Column_Name format
        new_col_name = "".join(["_" + char if char.isupper() and not old_col_name[i - 1].isupper() else char for i, char in enumerate(old_col_name)]).lstrip("_")

        # Change the column name using withColumnRenamed and regexp_replace
        df = df.withColumnRenamed(old_col_name, new_col_name)

    output_path = '/mnt/gold/SalesLT/' +name +'/'
    df.write.format('delta').mode("overwrite").save(output_path)
```
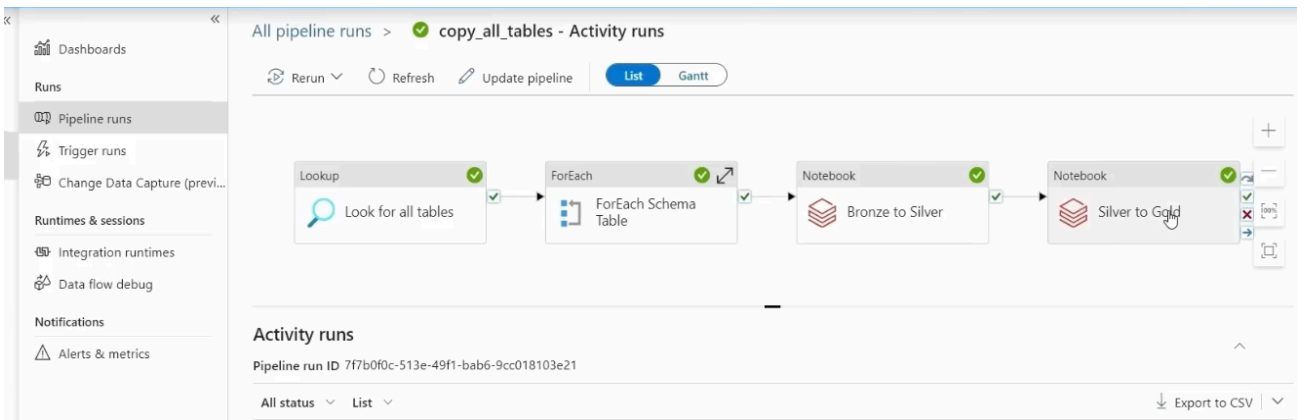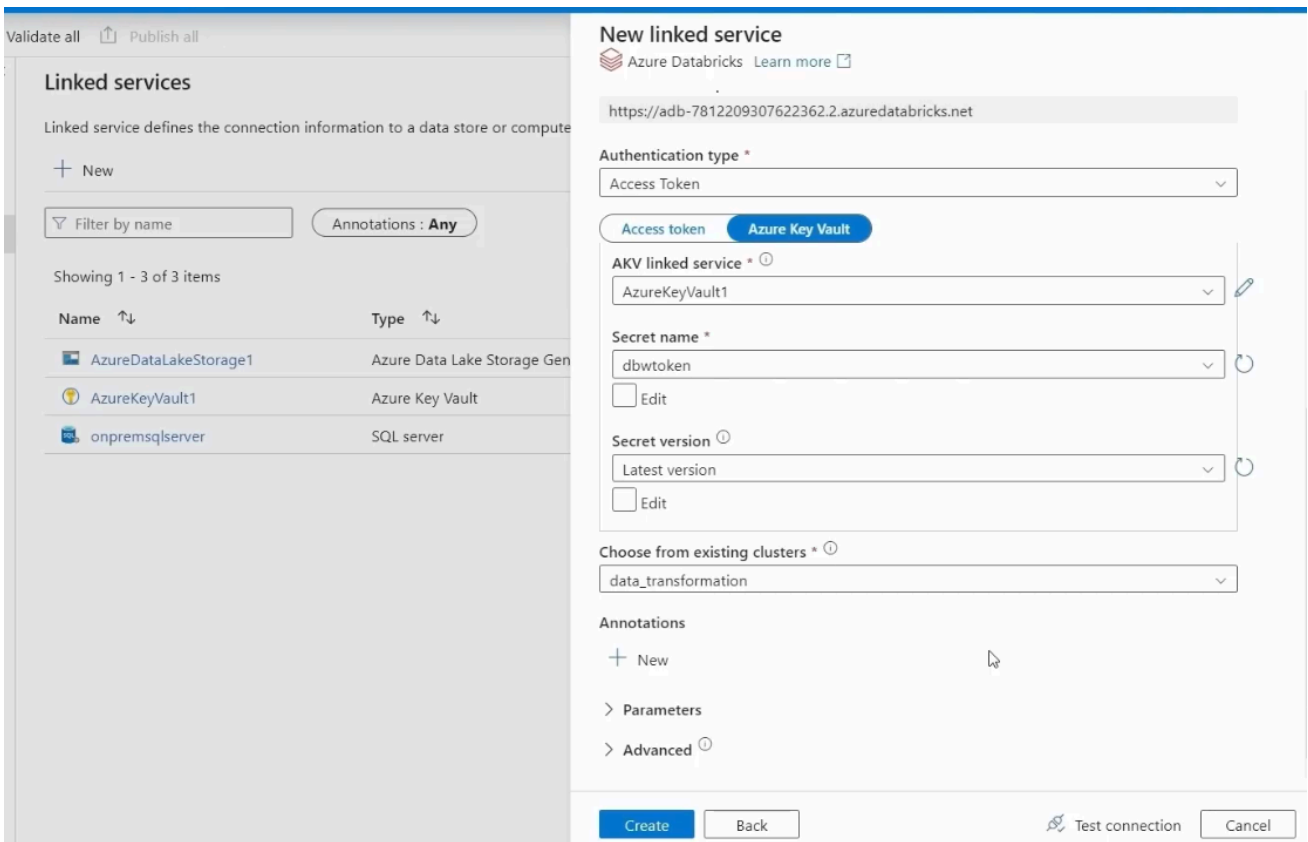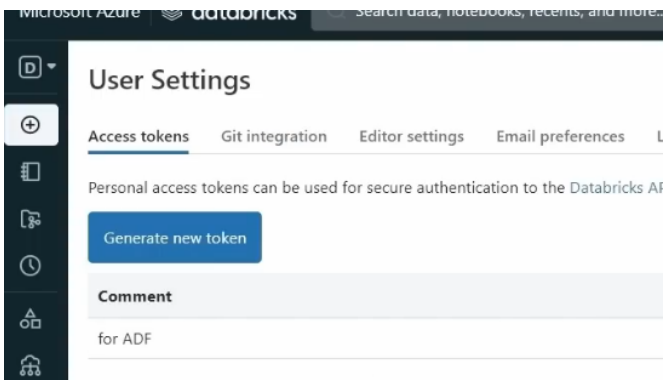
```python
table_name = []
for i in dbutils.fs.ls('mnt/silver/SalesLT/'):
  table_name.append(i.name.split('/')[0])
# Get list of column names in column_names
column_names = df.columns

# Loop through each column name
for old_col_name in column_names:
  # Convert column name from CamelCase to snake_case format
  new_col_name = "".join([word.lower() if word.isupper() and not i == 0
else word for i, word in enumerate(old_col_name)]).strip("_")

  # Change the column name using withColumnRenamed
  df = df.withColumnRenamed(old_col_name, new_col_name)

# Print transformed DataFrame
print(df)
```

# Mahesh Raut - Project - 31 May 2024

## Make connection with cluster in ADF

### Activity runs

Pipeline run ID 7f7b0f0c-513e-49f1-bab6-9cc018103e21

All status ⌄    List ⌄

Showing 1 - 14 items

| Activity name ↑↓ | Status ↑↓ | Activity type ↑↓ |
|---|---|---|
| Silver to Gold | ✅ Succeeded | Notebook |
| Bronze to Silver | ✅ Succeeded | Notebook |
| Copy Each Table | ✅ Succeeded | Copy data |
| Copy Each Table | ✅ Succeeded | Copy data |
| Copy Each Table | ✅ Succeeded | Copy data |
| Copy Each Table | ✅ Succeeded | Copy data |

**Part 4 - Data Loading - Gold container to Azure Synapse Serverless SQL Pool**

```sql
USE gold_db
GO
CREATE OR ALTER PROC CreateSQLServerlessView_gold @ViewName nvarchar(100)
AS
BEGIN
DECLARE @statement VARCHAR (MAX)
SET @statement = N'CREATE OR ALTER VIEW ' + @ViewName + N' AS
                            SELECT *
                            FROM                            OPENROWSET (
                                                                BULK
''https://mrkdatalakegen2.dfs.core.windows.net/gold/SalesLT/' + @ViewName + '/'',
                                                                FORMAT =
''DELTA"
                                                      ) as [result]'

EXEC (@statement)
END
GO
```

# Mahesh Raut - Project - 31 May 2024





## After trigger: