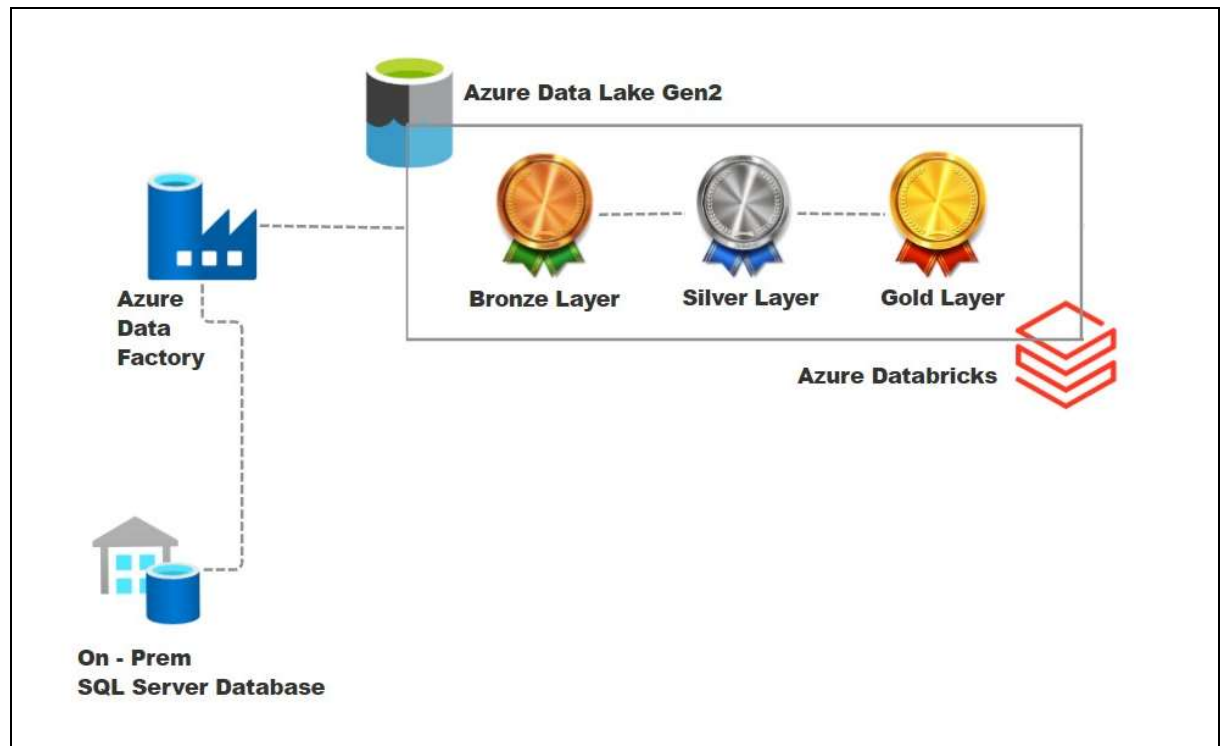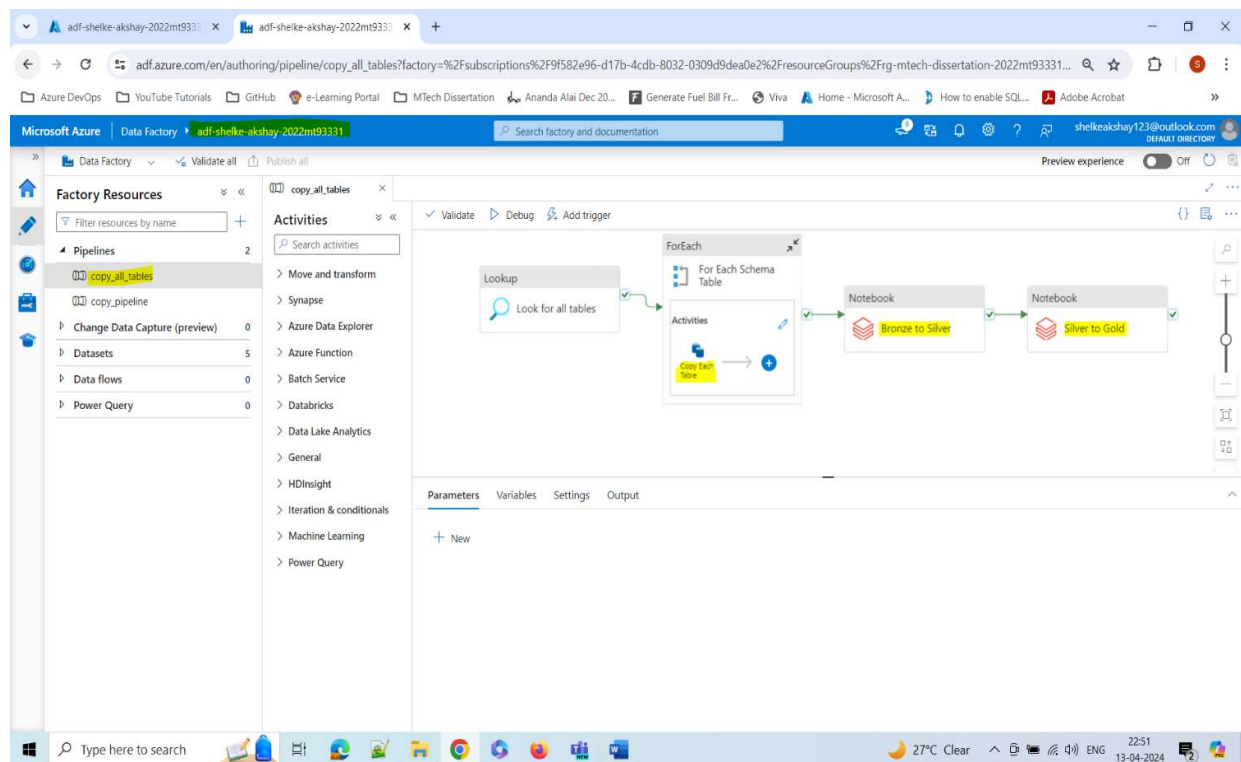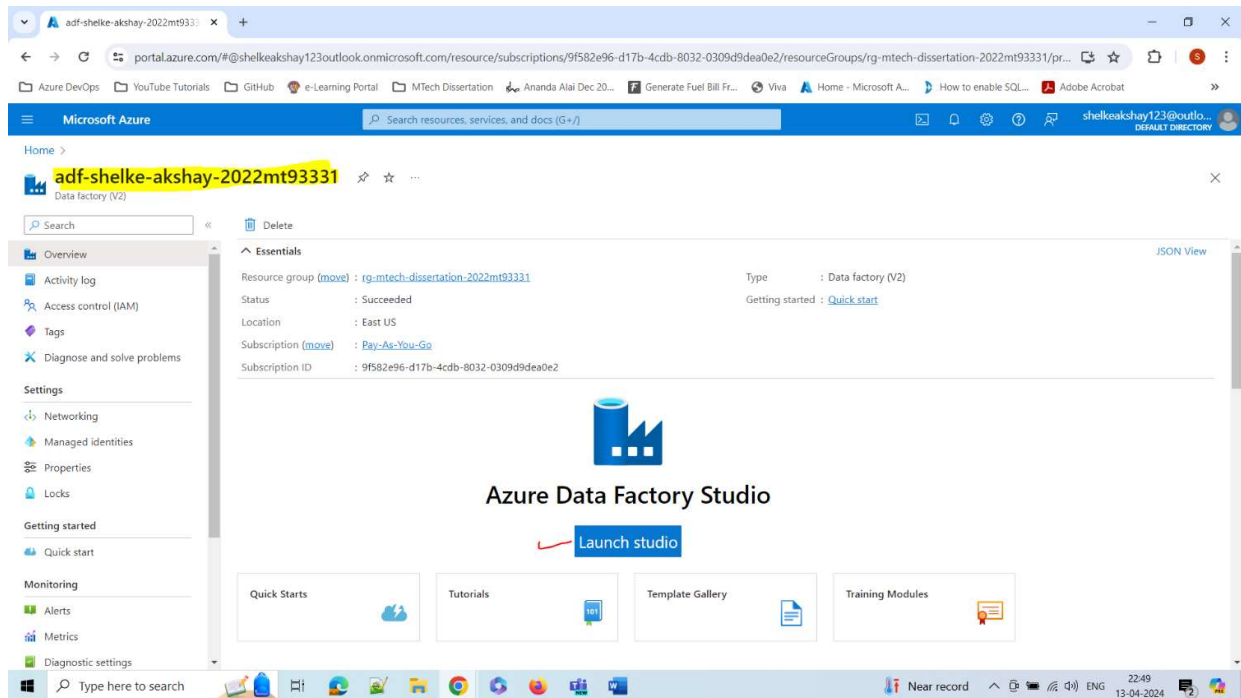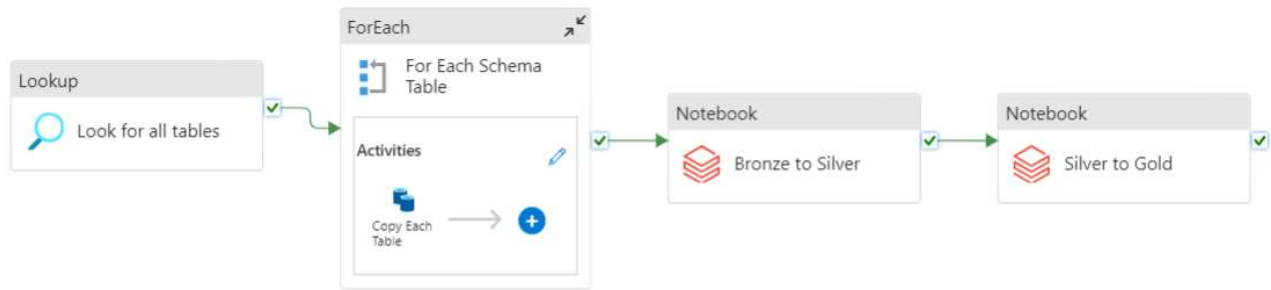# Data Transformation

1) Azure Data Factory(ADF) Pipeline

a. Look for all tables.
b. Copy Each table into BRONZE container.
c. Data Transformation – BRONZE to SILVER
d. Data Transformation – SILVER to GOLD

2)  Azure Databricks(ADB) – Code Notebooks for Data Transformation

| Azure Databricks Instance |
| --- |
|  |

| Azure Databricks Compute Cluster |
| --- |
|  |

| Python Code To create Mount point for Azure Blob Storage Container – For Bronze Container |
| --- |

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


try:
  dbutils.fs.mount(
    source = "abfss://bronze@storageaccmtechdemo.dfs.core.windows.net/",
    mount_point = "/mnt/bronze",
    extra_configs = configs)
  print("Mount Point created successfully")
except:
  print("Mount Point already exists")
```

| Python Code To create Mount point for Azure Blob Storage Container – For Silver Container |
|---|

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

try:
  dbutils.fs.mount(
    source = "abfss://silver@storageaccmtechdemo.dfs.core.windows.net/",
    mount_point = "/mnt/silver",
    extra_configs = configs)
  print("Mount Point created successfully")
except:
  print("Mount Point already exists")
```

| Python Code To create Mount point for Azure Blob Storage Container – For Gold Container |
|---|

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

try:
  dbutils.fs.mount(
    source = "abfss://gold@storageaccmtechdemo.dfs.core.windows.net/",
    mount_point = "/mnt/gold",
    extra_configs = configs)
  print("Mount Point created successfully")
except:
  print("Mount Point already exists")
```

**Python Code for Data Transformation – BRONZE to SILVER**

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


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


table_name = []

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



for i in table_name:
    path = "/mnt/bronze/SalesLT/" + i + "/" + i + ".parquet"
    df = spark.read.format("parquet").load(path)
    column = df.columns

    for col in column:
        if "Date" in col or "date" in col:
            df = df.withColumn(col,
date_format(from_utc_timestamp(df[col].cast(TimestampType()), "UTC"), "yyyy-MM-dd"))

    output_path = "/mnt/silver/SalesLT/" + i + "/"
    df.write.format("delta").mode("overwrite").save(output_path)

display(df)
```

**Python Code for Data Transformation – SILVER to GOLD**

```python
from pyspark.sql import SparkSession
from pyspark.sql.functions import col, regexp_replace


dbutils.fs.ls("mnt/silver/SalesLT/")

table_name = []

for i in dbutils.fs.ls("mnt/silver/SalesLT/"):
    table_name.append(i.name.split('/')[0])


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)

display(df)
```

3) Azure Data Factory pipeline RUNS

## 4) Delta tables stored in Silver Container



## 5) Delta tables stored in Gold Container – Data Transformation Completed