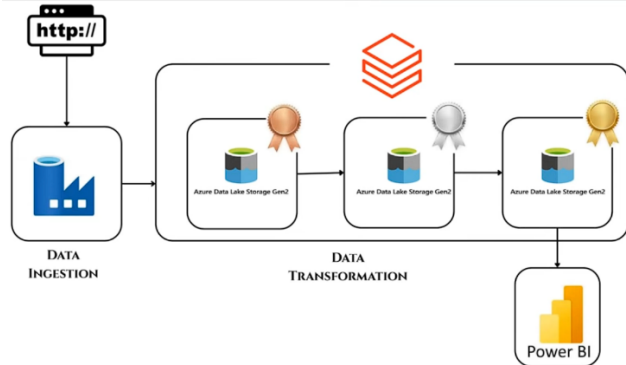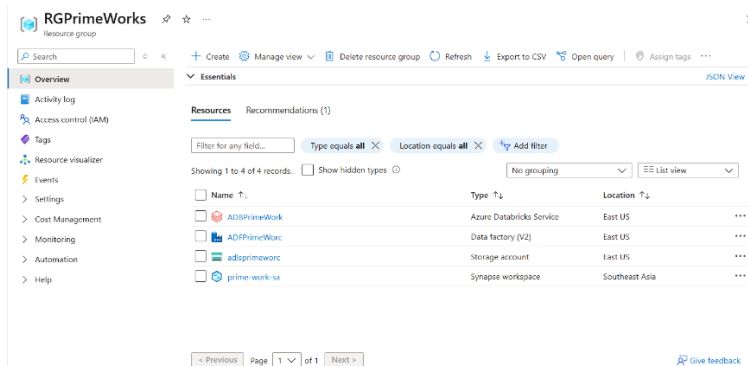# End-to-End Azure ETL Pipeline Project

Thilina Sasanka
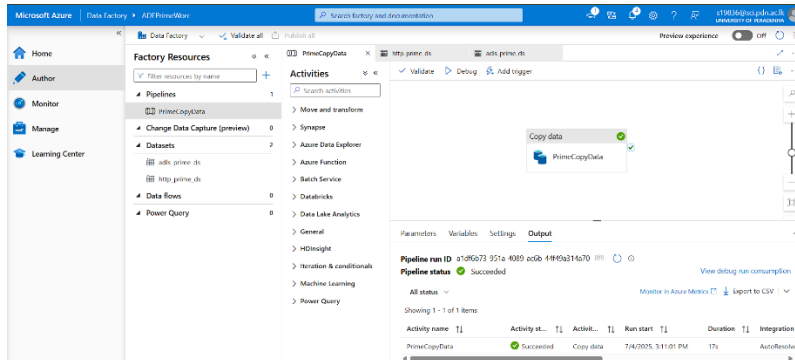UOP

## 1) ETL Pipeline Overview



## 2) Resource Group Creation
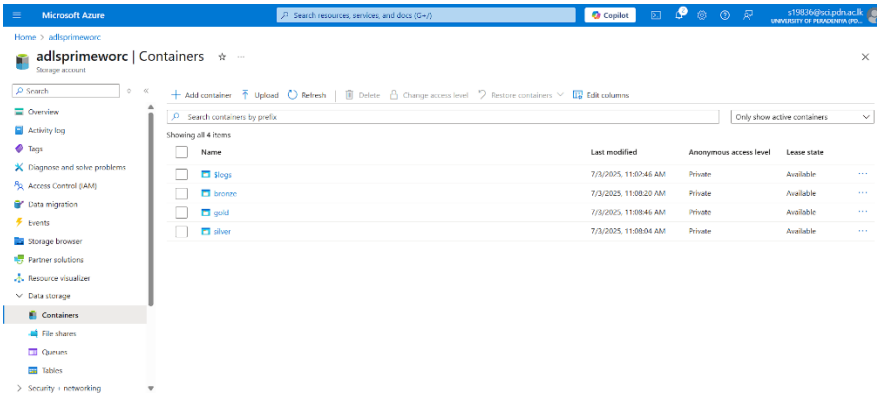


## 3) Azure Data Factory Creation

4) Azure Databricks Creation

5) Azure Databricks- Silver Layer - Cleaning & Transformations



```
✓  04:01 PM (<1s)                                          1

spark.conf.set("fs.azure.account.auth.type.adlsprimeworc.dfs.core.windows.net", "OAuth")
spark.conf.set("fs.azure.account.oauth.provider.type.adlsprimeworc.dfs.core.windows.net", "org.apache.
hadoop.fs.azurebfs.oauth2.ClientCredsTokenProvider")
spark.conf.set("fs.azure.account.oauth2.client.id.adlsprimeworc.dfs.core.windows.net",

spark.conf.set("fs.azure.account.oauth2.client.secret.adlsprimeworc.dfs.core.windows.net",

spark.conf.set("fs.azure.account.oauth2.client.endpoint.adlsprimeworc.dfs.core.windows.net", "https://
login.microsoftonline.com/                                                                   )
```

```
✓  04:01 PM (3s)                                           2

dbutils.fs.ls('abfss://bronze@adlsprimeworc.dfs.core.windows.net/')  # Check and update Azure
credentials if needed

[FileInfo(path='abfss://bronze@adlsprimeworc.dfs.core.windows.net/amazon_prime_titles.csv', name='amazon_pri
me_titles.csv', size=536676, modificationTime=1751525939000)]
```

```
✓  04:04 PM (<1s)                                          3

from pyspark.sql.functions import *
from pyspark.sql.types import *
```

**Data Understanding**

```
✓  04:15 PM (20s)                                          5

df_silver = spark.read.format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("abfss://bronze@adlsprimeworc.dfs.core.windows.net/amazon_prime_titles.csv")
df_silver.display()
```

```
          ✓  04:16 PM (<1s)                                    6

    df_silver.printSchema()

root
 |-- show_id: string (nullable = true)
 |-- type: string (nullable = true)
 |-- title: string (nullable = true)
 |-- director: string (nullable = true)
 |-- cast: string (nullable = true)
 |-- country: string (nullable = true)
 |-- date_added: string (nullable = true)
 |-- release_year: string (nullable = true)
 |-- rating: string (nullable = true)
 |-- duration: string (nullable = true)
 |-- listed_in: string (nullable = true)
 |-- description: string (nullable = true)
```

## Data Cleaning

```
          ✓  04:18 PM (1s)                                     8

    df_silver.display()
```

```
          ✓  04:26 PM (4s)                                     10

    df_silver = df_silver.dropDuplicates()
    df_silver.display()
```

```
          ✓  04:34 PM (2s)                                     11

    df_silver = df_silver.na.fill({"rating": "Unrated", "country":"Unknown", "description":"Unknown",
    "date_added":"01/01/2025", "release_year":"2020", "duration":"Unknown", "cast":"Unknown",
    "listed_in":"Unknown"})
    df_silver.display()
```

```
          ✓  04:35 PM (1s)                                     12

    df_silver = df_silver.dropna('any')
    df_silver.display()
```

```
    ▶ ⌄   ✓  Just now (1s)                                     13

    df_silver = df_silver.withColumnRenamed('title', 'Content_title')
    df_silver.display()
  ▶ (2) Spark Jobs
  ▶ 🔳  df_silver:  pyspark.sql.dataframe.DataFrame = [show_id: string, type: string ... 10 more fields]
```

| Table ⌄ | + | | |
|---|---|---|---|
| ᴬᴮC show_id | ᴬᴮC type | ᴬᴮC Content_title | ᴬᴮC director |
| 1  s297 | Movie | Tomorrow When The War Began | Stuart Beattie |

```
▶ ✓ 05:21 PM (4s)                                              15
  df_silver.write.format("parquet")\
      .mode("append")\
      .option("path", "abfss://silver@adlsprimeworc.dfs.core.windows.net/amazon_prime_titles_silver.csv")\
      .save()
▶ (2) Spark Jobs
```

6) Azure Databricks - Gold Layer – Transformations

```
Gold_Layer_Transformation    Python ∨   Tabs: OFF ∨   ☆                        ⊞    ▶ Run all    ● Thilina Pathirana's Clus... ∨
File  Edit  View  Run  Help   Last edit was 6 minutes ago

▶ ✓ 05:46 PM (1s)                                              1
  spark.conf.set("fs.azure.account.auth.type.adlsprimeworc.dfs.core.windows.net", "OAuth")
  spark.conf.set("fs.azure.account.oauth.provider.type.adlsprimeworc.dfs.core.windows.net", "org.apache.hadoop.fs.azurebfs.oauth2.
  ClientCredsTokenProvider")
  spark.conf.set("fs.azure.account.oauth2.client.id.adlsprimeworc.dfs.core.windows.net", "")
  spark.conf.set("fs.azure.account.oauth2.client.secret.adlsprimeworc.dfs.core.windows.net", "")
  spark.conf.set("fs.azure.account.oauth2.client.endpoint.adlsprimeworc.dfs.core.windows.net", "https://login.microsoftonline.com/
  ")
```

```
▶ ✓ 05:47 PM (4s)                                              2
  dbutils.fs.ls('abfss://silver@adlsprimeworc.dfs.core.windows.net/')

[FileInfo(path='abfss://silver@adlsprimeworc.dfs.core.windows.net/amazon_prime_titles_silver.csv/', name='amazon_prime_titles
odificationTime=1751543484000)]
```

```
▶ ✓ 05:47 PM (<1s)                                             3
  from pyspark.sql.functions import *
  from pyspark.sql.types import *
```

```
▶ ✓ 05:47 PM (19s)                                             4
  df_gold = spark.read.format("parquet")\
      .option("header", "true")\
      .option("inferSchema", "true")\
      .load("abfss://silver@adlsprimeworc.dfs.core.windows.net/amazon_prime_titles_silver.csv")
  df_gold.display()
```

```
▶ ✓ 4 minutes ago (2s)                                         5
  df_gold = df_gold.withColumn("date_added", to_date(df_gold["date_added"], "mm/dd/yyyy"))
  df_gold = df_gold.withColumn("year_added", year(df_gold["date_added"]))
  df_gold.display()
▶ (1) Spark Jobs
```

```
▶ ✓ Just now (1s)                                              6
  df_gold = df_gold.withColumn("category_1",split(df_gold["listed_in"], ",")[0])
  df_gold = df_gold.withColumn("category_2",split(df_gold["listed_in"], ",")[1])

  df_gold.display()
```

```
▶ ✓ Just now (1s)                                              7
  df_gold = df_gold.withColumn("category_2", when(df_gold["category_2"].isNull(), "Unknown").otherwise(df_gold["category_2"]))
  df_gold.display()
```

```
▶  ∨  ✓  08:57 PM (24s)                                    8

  df_gold.write.format("delta")\
      .mode("append")\
      .option("path", "abfss://gold@adlsprimeworc.dfs.core.windows.net/amazon_prime_titles_gold.csv")\
      .save()
```

7) Load data to the Synapse Analytics





8) Connecting to Power BI