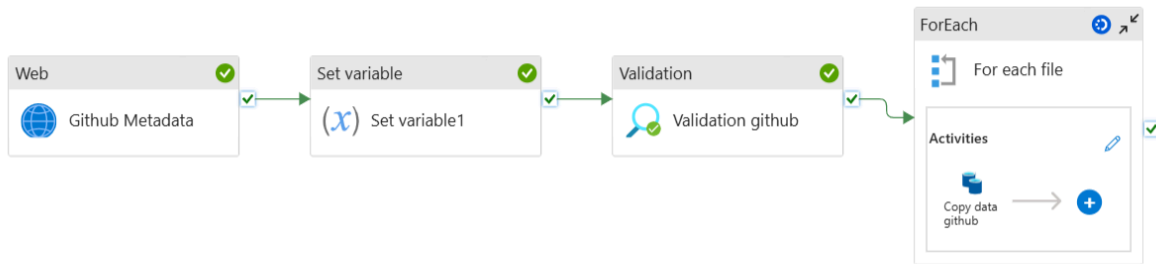


LLD

Low-Level Design (LLD)

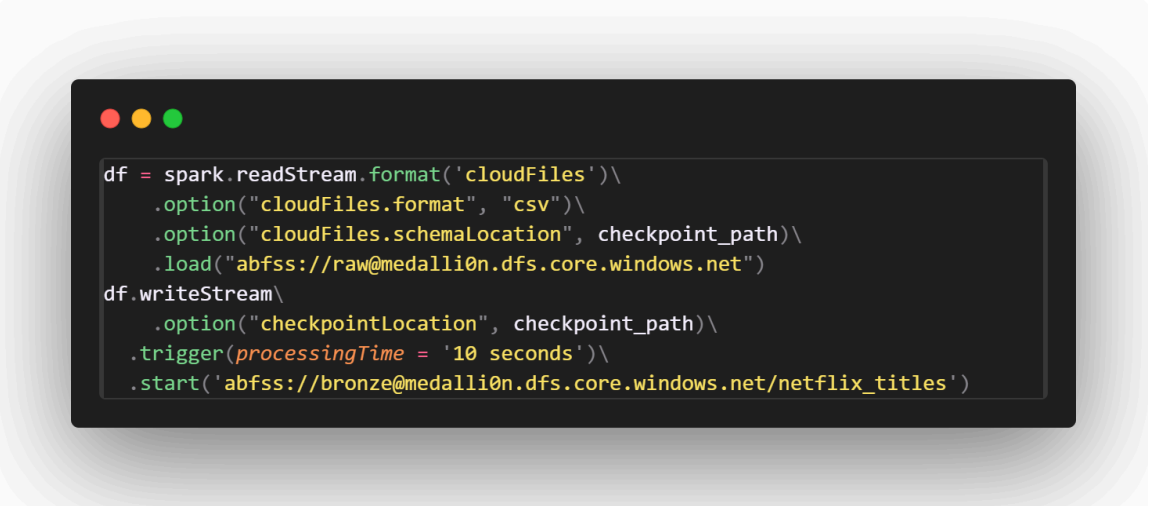
1. Ingestion Layer (Azure Data Factory)



- **Pipeline:** Netflix_Ingestion_Pipeline
 - **Activities:**
 - **Copy Data:** Pulls CSVs (e.g., netflix_titles.csv, netflix_cast.csv) from GitHub URLs.
 - **Parameters:**
 - **source_file:** Dynamic CSV path (e.g., `https://github.com/.../netflix_titles.csv, netflix_cast.csv`).
 - **sink_path:** Data Lake bronze folder (bronze@medallion.dfs.core.windows.net/).
 - **Validation:** Ensures file existence and schema consistency.
 - **Output:** CSVs in bronze layer. (e.g., bronze/netflix_titles/).
 - **Trigger:** Manual or scheduled (e.g., daily, weekly).

2. Bronze Layer (Databricks Autoloader)

- **Notebook:** Bronze_Ingestion



```
df = spark.readStream.format('cloudFiles')\
    .option("cloudFiles.format", "csv")\
    .option("cloudFiles.schemaLocation", checkpoint_path)\
    .load("abfss://raw@medallion.dfs.core.windows.net")
df.writeStream\
    .option("checkpointLocation", checkpoint_path)\
    .trigger(processingTime = '10 seconds')\
    .start('abfss://bronze@medallion.dfs.core.windows.net/netflix_titles')
```

- **Features:** Incremental loading, checkpointing, schema evolution (add new columns).
- **Output:** Delta tables (e.g., bronze_netflix_titles, bronze_netflix_cast).

3. Silver Layer (PySpark + Workflows)

- **Notebooks:**
 1. **Lookup Notebook:** Lookup_Silver

```

1 files = [
2     {
3         "sourcefolder" : "netflix_directors",
4         "targetfolder" : "netflix_directors"
5     },
6     {
7         "sourcefolder" : "netflix_cast",
8         "targetfolder" : "netflix_cast"
9     },
10    {
11        "sourcefolder" : "netflix_countries",
12        "targetfolder" : "netflix_countries"
13    },
14    {
15        "sourcefolder" : "netflix_category",
16        "targetfolder" : "netflix_category"
17    },
18 ]
19 dbutils.jobs.taskValues.set(key = "my_arr", value = files)

```

- **Purpose:** Defines source-target folder mappings.

2. Silver Transformation Notebook: Silver_Transformation

```

1 dbutils.widgets.text("sourcefolder", "netflix_directors")
2 dbutils.widgets.text("targetfolder", "netflix_directors")
3 var_src_folder = dbutils.widgets.get("sourcefolder")
4 var_tar_folder = dbutils.widgets.get("targetfolder")

```

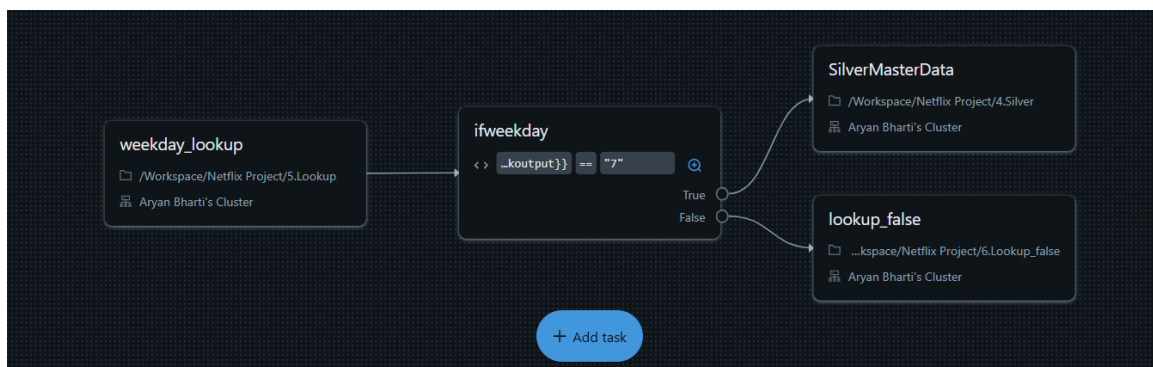
- **Transformations:**

```

1
2 df = df.fillna({'duration_minutes': 0, 'duration_seasons': 1}) # Null handling
3 df = df.withColumn('duration_minutes', col('duration_minutes').cast(IntegerType()))\
4   .withColumn('duration_seasons', col('duration_seasons').cast(IntegerType()))
5 # Type cast
6 df = df.withColumn('Short_Titles', split(col('title'), ':')[0]) # Split column
7 df = df.withColumn('rating', split(col('rating'), '-')[0]) # Split column
8 df = df.withColumn('flag_check', when(col('type') == 'Movie', 1)\
9   .when(col('type') == 'TV Show', 1)\
10  .otherwise(0)) # Conditional flag
11 df = df.withColumn("duration_rank", dense_rank().over(Window.orderBy(desc("
   duration_minutes")))) # Ranking

```

- **Workflow: Silver_Workflow**



- **Tasks:**

- **Lookup_Locations:** Runs Lookup_Silver, outputs my_array.

```

1 dbutils.widgets.text('weekday', '7')
2 var = int(dbutils.widgets.get('weekday'))
3 dbutils.jobs.taskValues.set(key='weekoutput', value = var)

```

- **Silver_Notebook:** For-each loop over my_array, passes parameters to Silver_Transformation.
- **Conditional Task:** Weekday_Lookup checks weekday(); if 7 (Sunday), runs Silver_Transformation, else prints day.

- **Output:** Transformed Delta tables (e.g., silver/netflix_titles).

4. Gold Layer (Delta Live Tables)

- **Notebook:** DLT_Gold

```

1 import dlt
2 from pyspark.sql.functions import lit
3
4 # Lookup table rules for data quality
5 lookuptables_rules = {
6     "rule1": "show_id IS NOT NULL"
7 }
8
9 # Master data rules for the gold layer
10 masterdata_rule = {
11     "rule1": "newflag IS NOT NULL",
12     "rule2": "show_id IS NOT NULL"
13 }
14
15 # Netflix Directors table
16 @dlt.table(name="netflixdirectors")
17 @dlt.expect_all_or_drop(lookuptables_rules)
18 def netflix_directors():
19     df = spark.readStream.format("delta").load(
20         "abfss://silver@medallion.dfs.core.windows.net/netflix_directors"
21     )
22     return df
23
24 # Netflix Cast table
25 @dlt.table(name="netflixcast")
26 @dlt.expect_all_or_drop(lookuptables_rules)
27 def netflix_cast():
28     df = spark.readStream.format("delta").load(
29         "abfss://silver@medallion.dfs.core.windows.net/netflix_cast"
30     )
31     return df
32
33 # Netflix Countries table
34 @dlt.table(name="netflixcountries")
35 @dlt.expect_all_or_drop(lookuptables_rules)
36 def netflix_countries():
37     df = spark.readStream.format("delta").load(
38         "abfss://silver@medallion.dfs.core.windows.net/netflix_countries"
39     )
40     return df
41
42 # Netflix Category table
43 @dlt.table(name="netflixcategory")
44 @dlt.expect_all_or_drop({"rule1": "show_id IS NOT NULL"})
45 def netflix_category():
46     df = spark.readStream.format("delta").load(
47         "abfss://silver@medallion.dfs.core.windows.net/netflix_category"
48     )
49     return df
50
51 # Gold Staging Netflix Titles table
52 @dlt.table
53 def gold_stg_netflixtitles():
54     df = spark.readStream.format("delta").load(
55         "abfss://silver@medallion.dfs.core.windows.net/netflix_titles"
56     )
57     return df
58
59 # Gold Transformed Netflix Titles view
60 @dlt.view
61 def gold_trans_netflixtitles():
62     df = spark.readStream.table("LIVE.gold_stg_netflixtitles")
63     df = df.withColumn("newflag", lit(1))
64     return df
65
66 # Gold Netflix Titles table with master data rules
67 @dlt.table
68 @dlt.expect_all_or_drop(masterdata_rule)
69 # Corrected 'except_all_or_drop' to 'expect_all_or_drop'
70 def gold_netflixtitles():
71     df = spark.readStream.table("LIVE.gold_trans_netflixtitles")
72     return df

```

- **Pipeline:** DLT_Gold_Pipeline
 - **Output:** Gold Delta tables with dropped invalid records.
- **Features:** Streaming tables, views, and data quality enforcement.

Error Handling

- **ADF:** Retry on file fetch failure.
- **Databricks:** Debug via event logs (e.g., syntax errors in DLT).
- **DLT:** Drop invalid records; fail pipeline on critical errors.

Performance Considerations

- **Autoloader:** Checkpointing for efficient incremental loads.
- **Workflows:** Task parallelization via for-each loops.
- **DLT:** Minimizes recomputation with streaming updates.