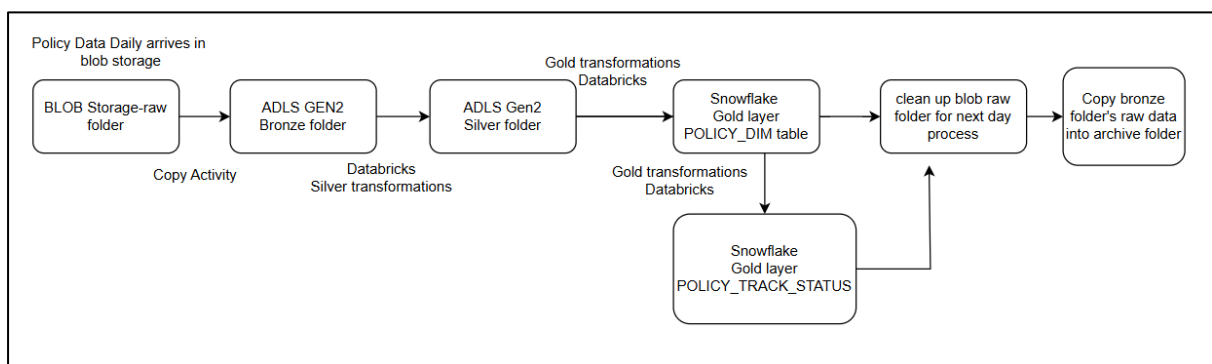


Policy Lifecycle Tracking Project | End-to-End Azure Pipeline with Databricks & Snowflake

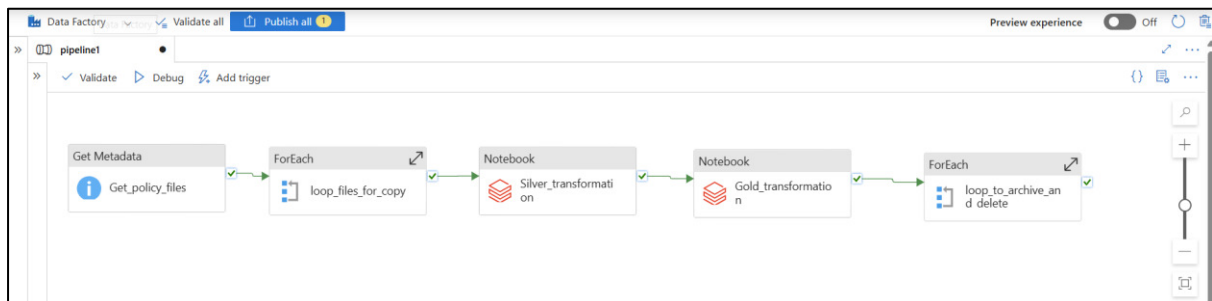
Project Summary

- **Objective:** To build a data pipeline for tracking policy status transitions over time by ingesting, transforming, and storing policy data from raw sources to a refined analytical state.
- **Source:** Daily policy data stored in a *Blob Storage (Raw Folder)*.
- **Destination:** A *Snowflake data warehouse*, with intermediate transformations and data quality layers handled in *Azure Data Lake Gen2* and *Databricks*.
- **Purpose:** Maintain a full lifecycle view of each policy, track changes in status, and compute the duration spent in each status phase.

Architecture diagram:



The end-to-End pipeline with dynamic parameters



Pipeline debug result for day-1 file

The screenshot shows the 'Output' tab of the pipeline debug results. The pipeline run ID is 774a623b-b615-4f49-82f5-dcb8ac04f06f. The pipeline status is 'Succeeded'. The table below shows the details of the activities and their execution results.

Activity name	Activity status	Activity type	Run start	Duration	Integration runtime	User properties	Activity run ID
Delete2	Succeeded	Delete	5/10/2025, 12:12:16 PM	5s	AutoResolveIntegrationRuntime (Canada Central)		294460a5-d597-40cd-b3e7-789c1c1fce58
Delete1	Succeeded	Delete	5/10/2025, 12:12:16 PM	5s	AutoResolveIntegrationRuntime (Canada Central)		6b254a02-e98f-47bc-a8b9-84c593c26ff0
Copy data2	Succeeded	Copy data	5/10/2025, 12:12:00 PM	15s	AutoResolveIntegrationRuntime (Canada Central)		f27ae7e3-0213-4e88-96fd-87aced72230b
loop_to_archive_and_delete	Succeeded	ForEach	5/10/2025, 12:11:59 PM	25s	AutoResolveIntegrationRuntime (Canada Central)		32b492f2-9f64-4894-b8fc-b97211e250b3
Gold_transformation	Succeeded	Notebook	5/10/2025, 12:11:02 PM	56s	AutoResolveIntegrationRuntime (Canada Central)		ea8811b5-ca25-438d-9581-56f23316dcf6
Silver_transformation	Succeeded	Notebook	5/10/2025, 12:10:22 PM	39s	AutoResolveIntegrationRuntime (Canada Central)		c364602-e106-4e18-adac-824be96bf5de
Copy data1	Succeeded	Copy data	5/10/2025, 12:10:07 PM	13s	AutoResolveIntegrationRuntime (Canada Central)		8e9f1e80-ec7d-4b1e-a058-13821a8c4b78
loop_files_for_copy	Succeeded	ForEach	5/10/2025, 12:10:06 PM	16s	AutoResolveIntegrationRuntime (Canada Central)		86eb5c4c-5a30-43ce-99ff-60e2488a5406
Get_policy_files	Succeeded	Get Metadata	5/10/2025, 12:09:55 PM	11s	AutoResolveIntegrationRuntime (Canada Central)		7795ebb2-233f-44d8-803a-4a7cb82d32ad

Data Transformation & Workflow

1. Bronze Layer (Raw Ingestion)

- Ingest daily policy data from Azure Blob Storage into ADLS Gen2 Bronze folder.
- Raw data is stored as-is, preserving original structure for audit and traceability.

2. Silver Layer (Cleansing & Standardization)

- Use Databricks to process Bronze layer data.
- Apply cleansing (null handling, schema enforcement) and standardization (field renaming, type casting).
- Store cleaned and standardized data in the Silver folder in ADLS Gen2.

Silver transformation logic is

Import the required functions and types

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

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

#Read policy file

```
policy_df = spark.read.format("csv").option("header", "true").option("inferSchema",  
"true").load("/mnt/bronze/policy_snapshot_*.csv")
```

data cleansing

```
policy_clean_df=policy_df.withColumn("policy_status",trim(col("policy_status"))) \
```

```
.withColumn("policy_status",  
when(col("policy_status").isin("submitted","Submitted"),"Submitted").
```

```
when(col("policy_status").isin("Active","active"),"Active").
```

```
when(col("policy_status").isin("Canceled","canceled"),"Cancelled").
```

```
when(col("policy_status").isin("Mature","mature"),"Mature").
```

```
otherwise(initcap(col("policy_status")))) \
```

```
.withColumn("submission_date",col("submission_date").cast(DateType())) \
```

```
.withColumn("status_update_date",to_date(col("status_update_date"))) \
```

```
.withColumn("agent_id",trim(col("agent_id"))) \
```

```
.withColumn("agent_id", when(col("agent_id").isNull(),  
lit("UNKNOWN")).otherwise(col("agent_id"))) \
```

```
.fillna("UNKNOWN","policy_status") \
```

```
.dropna(subset=["submission_date", "status_update_date"])
```

```

#Read region file

region_df=spark.read.format("csv") \

.option("header", "true") \

.option("inferSchema", "true") \

.load("/mnt/region/regions.csv")

#join policy and region files

policy_region_df=policy_clean_df.join(region_df,policy_clean_df.region
==region_df.region_id,"inner").drop("region_id")

# validation rule

valid_statuses = ["Submitted","Active","Cancelled","Mature"]

policy_validated_df=policy_region_df.filter(col("policy_status").isin(valid_statuses))

# Add audit information

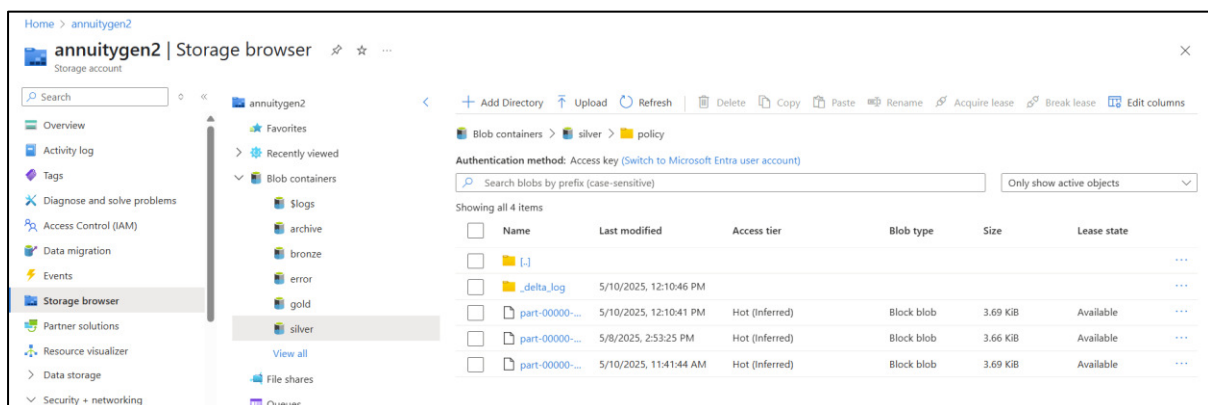
df_silver = policy_validated_df.withColumn("ingesttime", current_timestamp())

# Write to silver table

df_silver.write.mode("overwrite").format("delta").save("/mnt/silver/policy")

```

The transformed file is in silver folder as shown below



3. Gold Layer (Business Transformations)

- Read Silver data into Databricks.
- Apply business logic to generate the Policy Dimension (policy_dim) table:
 - Append daily data to maintain historical versions of policy records.
 - Capture all phase changes for each policy.
- Store policy_dim table into Snowflake as a dimensional table.

4. Policy Status Tracking

- Read policy_dim table from Snowflake.
- Compute previous status, current status, and date difference between transitions using window functions or lag operations.
- Store the output as policy_track_status table in Snowflake

Gold layer transformation logic is

```
# Import the required functions and types
from pyspark.sql.functions import *
from pyspark.sql.types import *
from pyspark.sql.window import Window

#read Policy file
policy_df = spark.read.format("delta").load("/mnt/silver/policy")

#Access snowflake account
sfOptions = {
    "sfURL": "https://lbnfhza-et12707.snowflakecomputing.com",
    "sfDatabase": "POLICY_DATA_DB",
    "sfSchema": "GOLD_LAYER",
    "sfWarehouse": "POLICY_WH",
    "sfRole": "ACCOUNTADMIN",
    "sfUser": "AAAAAAAAAAAA",
    "sfPassword": "XXXXXXXXXXXX",
}

#write dim current table
policy_df.write \
    .format("snowflake") \
    .options(**sfOptions) \
    .option("dbtable", "policy_dim") \
    .mode("append") \
    .save()
```

```

#Read the latest policy dim table

policy_dim_current = spark.read \

.format("snowflake") \

.options(**sfOptions) \

.option("dbtable", "policy_dim") \

.load()

#Lifecycle Status Tracking | get previous and next status

windowfunc=Window.partitionBy("policy_id").orderBy("status_update_date")

policy_status_df =

policy_dim_current.withColumn("prev_status",lag("policy_status").over(windowfunc)) \

.withColumn("prev_status_update_date",lag("status_update_date").over(windowfunc))

policy_diff_df =

policy_status_df.withColumn("No_of_day_in_status",datediff(col("status_update_date"),col("p

rev_status_update_date")))

policy_track_status_df = policy_diff_df \

.select("Policy_ID","Prev_status","policy_status","status_update_date","No_of_day_in_status")

#write dim current table

policy_track_status_df.write \

.format("snowflake") \

.options(**sfOptions) \

.option("dbtable", "policy_track_status") \

.mode("overwrite") \

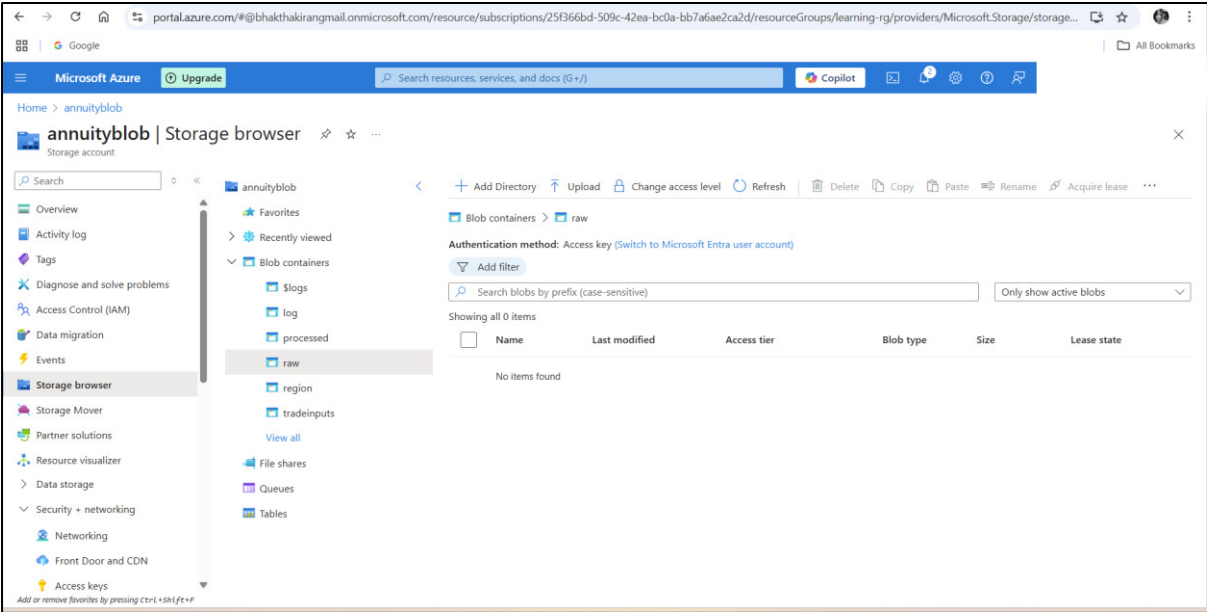
.save()

```

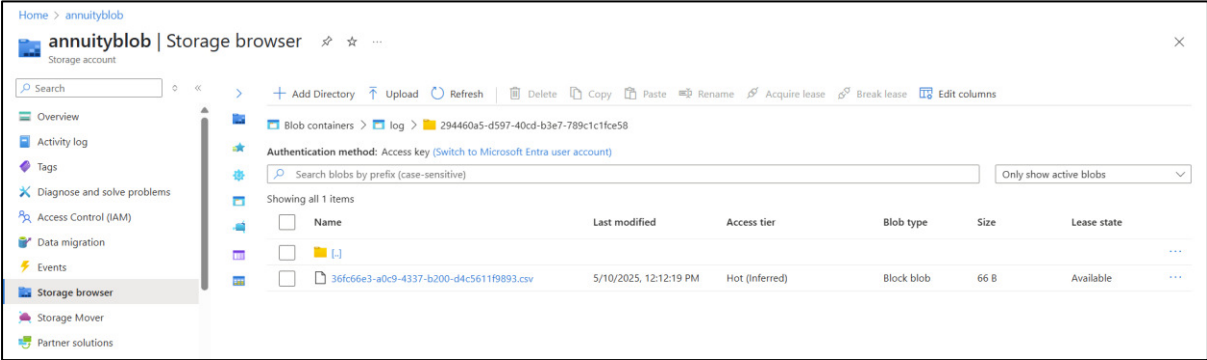
The transformed files are in snowflake as shown in below screenshots

Post-Processing Logic:

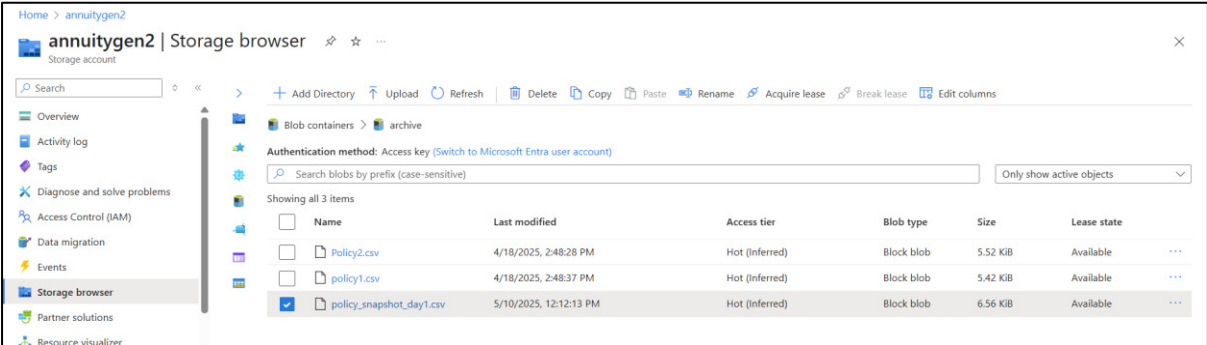
- After processing each day's data: ○ Files are deleted from the raw folder in Blob Storage.
- Metadata/logs of processed files are stored in a log folder.
- Bronze-level files are moved to an archive folder for backup and auditing.



Logged in log file regarding file deletion from raw folder



File was moved to archive folder



Processing day-2 file:Day-2 file was processed and debug result as shown below

</

Day-2 data is appended in snowflake Policy_dim table

POLICY_DATA_DB.GOLD_LAYER

Settings

Code Versions

1

select * FROM POLICY_DIM;

Results

Chart

</

The below table has the previous status and no of days that Policy was in previous status

ACCOUNTADMIN

COMPUTE_WH (X-Small)

Share

POLICY_DATA_DB.GOLD_LAYER

Settings

Code Versions

1

select * FROM POLICY_TRACK_STATUS;

Results

Chart

	POLICY_ID	PREV_STATUS	POLICY_STATUS	STATUS_UPDATE_DATE	# NO_OF_DAY_IN_STATUS	POLICY_ID
15	POL1016	null	Submitted	2024-08-30	null	POL1079
16	POL1017	null	Cancelled	2022-10-21	null	
17	POL1017	Cancelled	Submitted	2024-02-21	488	
18	POL1018	null	Active	2023-02-02	null	
19	POL1018	Active	Submitted	2024-09-23	599	
20	POL1021	null	Active	2025-04-25	null	
21	POL1024	null	Submitted	2023-01-02	null	
22	POL1024	Submitted	Active	2024-10-30	667	
23	POL1025	null	Submitted	2024-01-03	null	
24	POL1026	null	Submitted	2023-08-26	null	
25	POL1026	Submitted	Active	2024-12-08	470	

Conclusion

- End-to-end pipeline ensures clean, historical, and traceable policy data.
- Intelligent logging, archival, and deletion mechanisms ensure efficient storage management and audit readiness.
- Gold layer delivers value-added insights into policy lifecycle stages, enabling strategic decision-making.