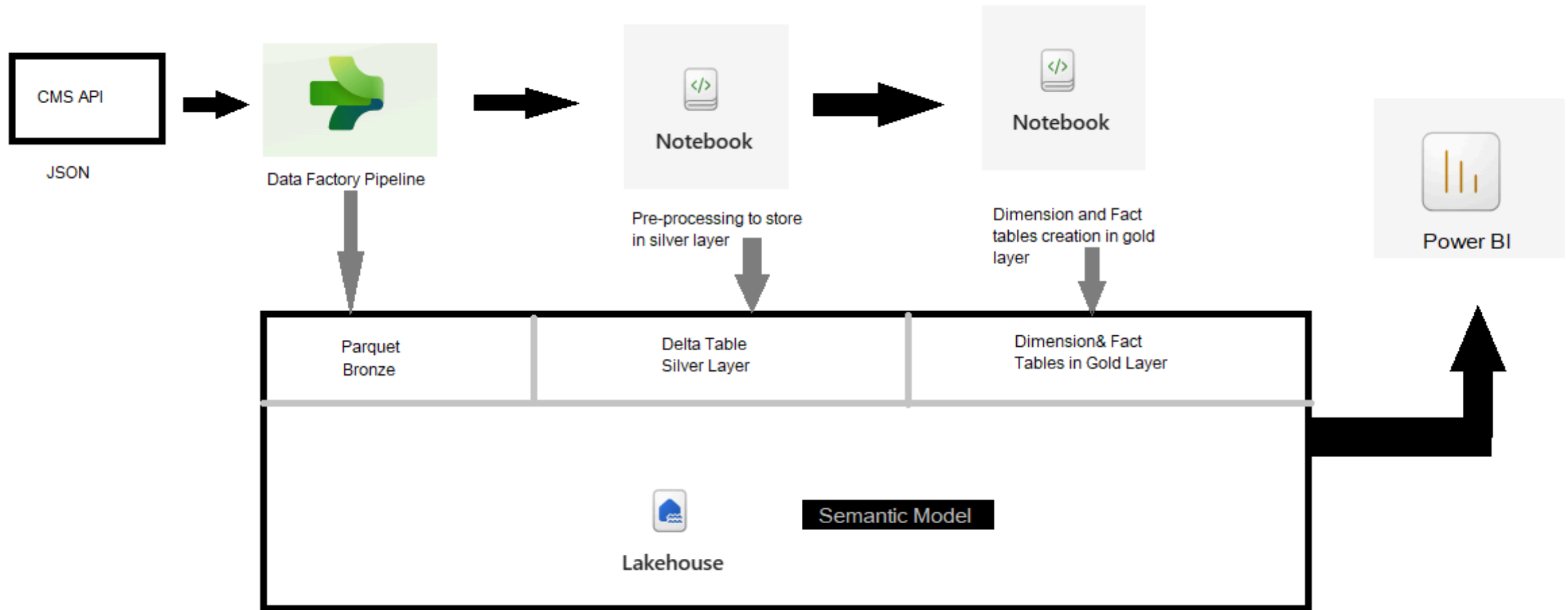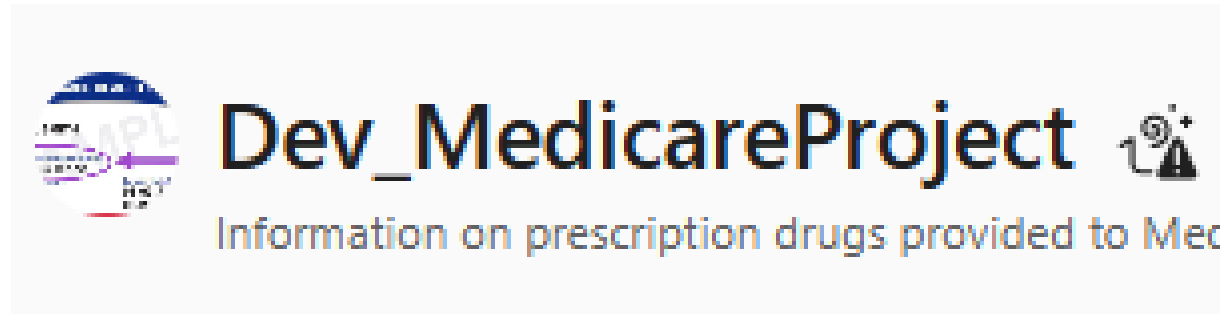# Objective

To get last 4 years Medicare's information on prescription drugs by drug and provider incrementally into Lakehouse so that reports can be created.

# Architecture

CMS API

JSON

Data Factory Pipeline

Notebook

Pre-processing to store in silver layer

Notebook

Dimension and Fact tables creation in gold layer

Power BI

| Parquet Bronze | Delta Table Silver Layer | Dimension& Fact Tables in Gold Layer |
|---|---|---|

Lakehouse

Semantic Model

# 1. Create new Workspace

**Dev_MedicareProject**

Information on prescription drugs provided to Med

# 2. Add Lakehouse

Add following folders

∨ **medicare_LH**

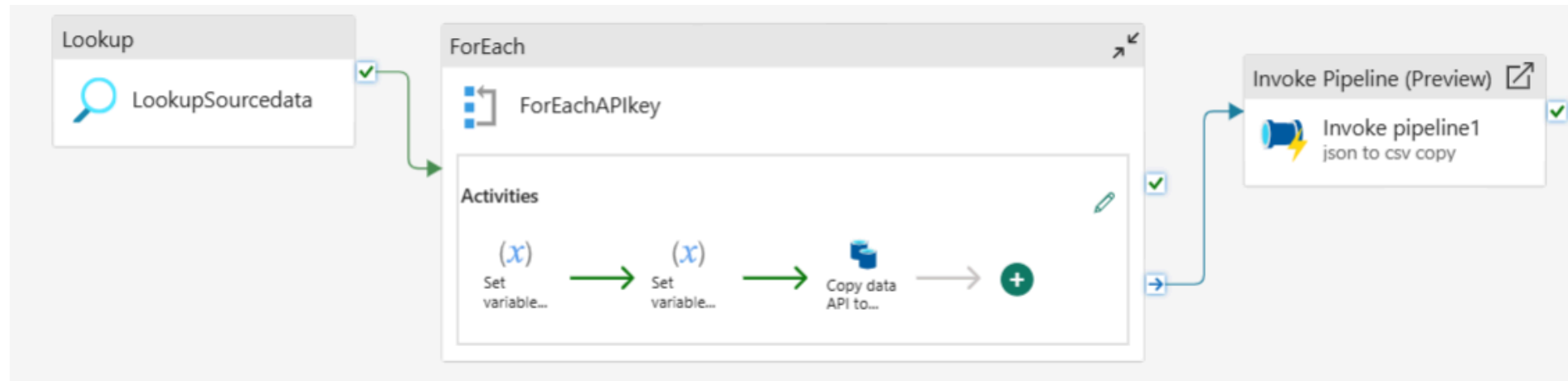> 📁 Tables

∨ 📁 Files

> 📁 bronze

> 📁 gold

> 📁 silver

> 📁 sourceData

# 3. Go to website data.cms.gov > Medicare part D prescribers - by provider and drug page
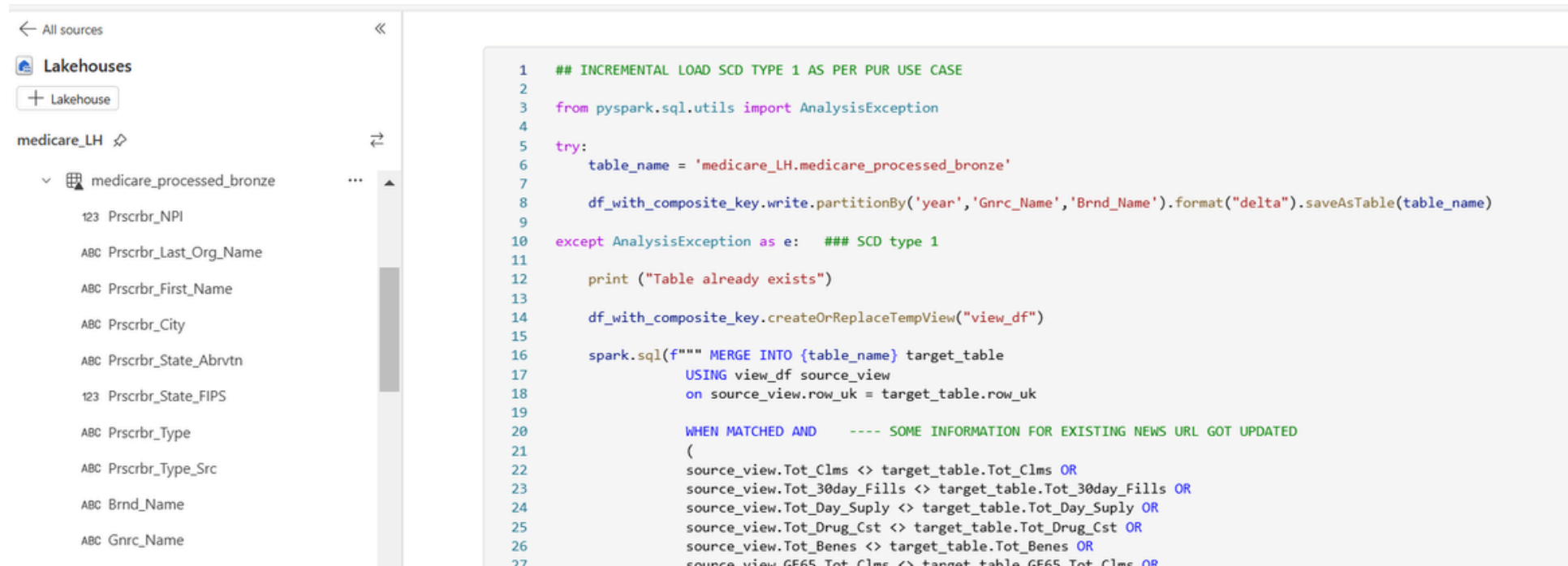


Get Access APIs year wise into a CSV file and store into sourceData directory in Lakehouse

# 4. Create Pipeline in the Lakehouse to copy API json format data for each file and store into bronze directory



I am first copying API json format data and then converting to .csv file, but we can directly copy from API json format and leave it at that in bronze directory
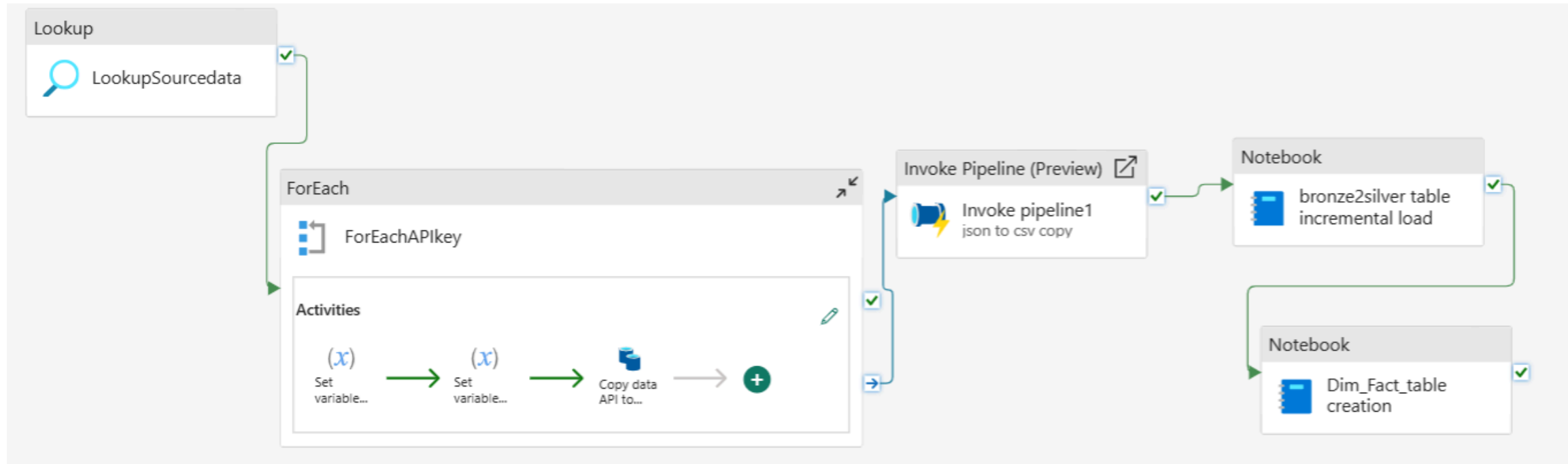
# 5. Perform some data cleansing and then save table in delta format enabling to be incrementally upserted when notebook runs
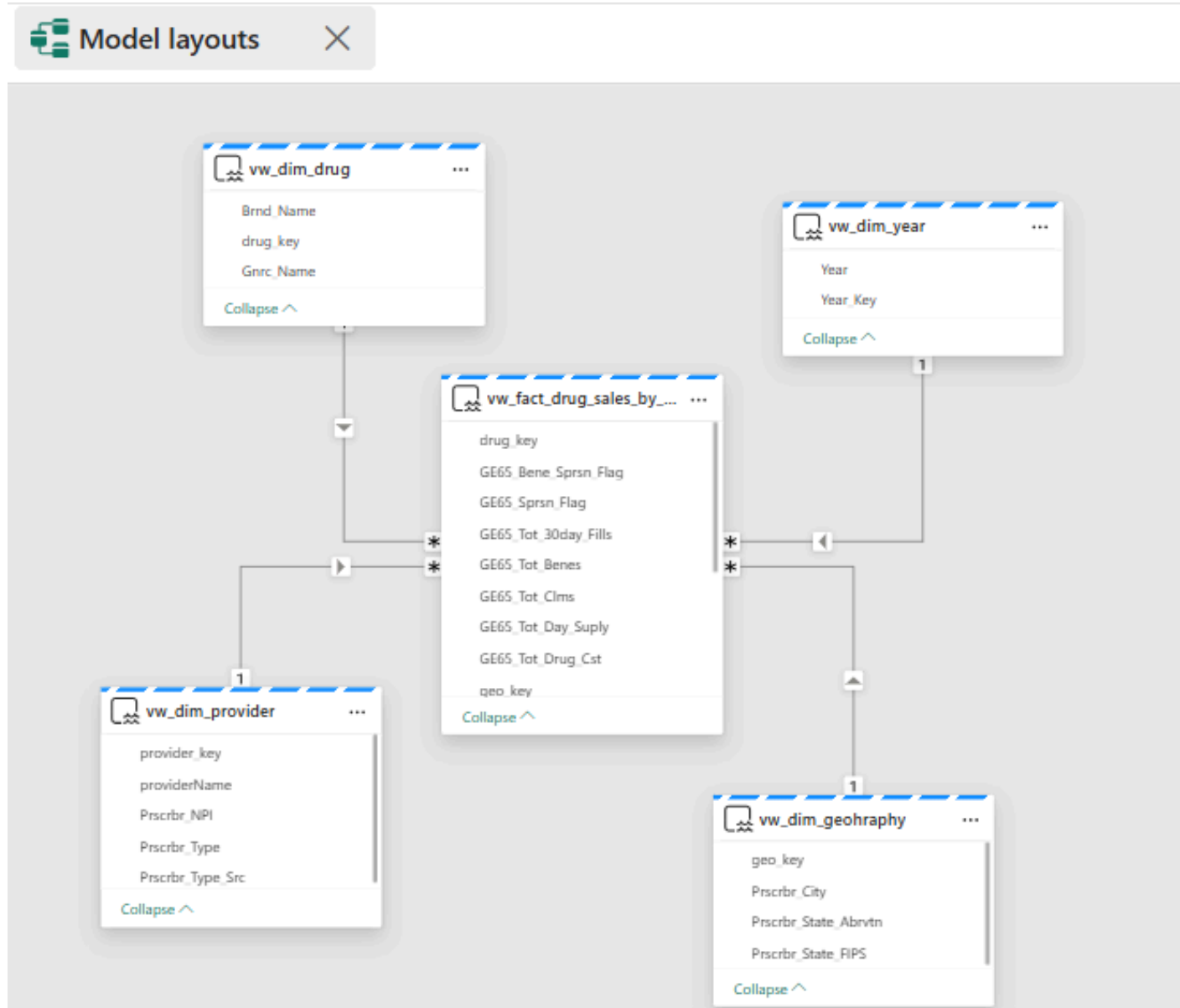
← All sources

**Lakehouses**

+ Lakehouse

medicare_LH ✎                                        ⇄

∨  ⊞ medicare_processed_bronze        ···

    123 Prscrbr_NPI

    ABC Prscrbr_Last_Org_Name

    ABC Prscrbr_First_Name

    ABC Prscrbr_City

    ABC Prscrbr_State_Abrvtn

    123 Prscrbr_State_FIPS

    ABC Prscrbr_Type

    ABC Prscrbr_Type_Src

    ABC Brnd_Name

    ABC Gnrc_Name

```python
1   ## INCREMENTAL LOAD SCD TYPE 1 AS PER PUR USE CASE
2
3   from pyspark.sql.utils import AnalysisException
4
5   try:
6       table_name = 'medicare_LH.medicare_processed_bronze'
7
8       df_with_composite_key.write.partitionBy('year','Gnrc_Name','Brnd_Name').format("delta").saveAsTable(table_name)
9
10  except AnalysisException as e:    ### SCD type 1
11
12      print ("Table already exists")
13
14      df_with_composite_key.createOrReplaceTempView("view_df")
15
16      spark.sql(f""" MERGE INTO {table_name} target_table
17                  USING view_df source_view
18                  on source_view.row_uk = target_table.row_uk
19
20                  WHEN MATCHED AND    ---- SOME INFORMATION FOR EXISTING NEWS URL GOT UPDATED
21                  (
22                  source_view.Tot_Clms <> target_table.Tot_Clms OR
23                  source_view.Tot_30day_Fills <> target_table.Tot_30day_Fills OR
24                  source_view.Tot_Day_Suply <> target_table.Tot_Day_Suply OR
25                  source_view.Tot_Drug_Cst <> target_table.Tot_Drug_Cst OR
26                  source_view.Tot_Benes <> target_table.Tot_Benes OR
27                  source_view.GE65_Tot_Clms <> target_table.GE65_Tot_Clms OR
```

# 6. Create a new notebook to transform table into Dimension and Fact in Gold layer

# 7. Adjust relationships in semantic data model layer



Create views on the dimension delta tables and the fact delta table using SQL analytics endpoint query

# 8. Add semantic model refresh activity to pipeline

# 9. Create a report in Power BI using views dimensions and measures

## Total Beneficiaries by State of Prescriber



**Generic Name of Drug**

- ☐ (Blank)
- ☐ Abacavir/Dolutegravir/Lamivudi...
- ☐ Acetaminophen With Codeine
- ☐ Acyclovir
- ☐ Adalimumab
- ☐ Albuterol Sulfate

## Total Beneficiaries by Generic Name of Drug



| Generic Name of Drug | Total Beneficiaries |
|---|---|
| Hydrocodone/Acetamin... | 2745 |
| Atorvastatin Calcium | 2613 |
| Gabapentin | 2238 |
| Amlodipine Besylate | 1955 |
| Oxycodone Hcl/Acetami... | 1849 |
| Tamsulosin Hcl | 1656 |
| Lisinopril | 1653 |
| Simvastatin | 1555 |
| Levothyroxine Sodium | 1519 |
| Metformin Hcl | 1303 |
| Omeprazole | 1271 |
| Losartan Potassium | 1252 |
| Baclofen | 1170 |
| Diclofenac Sodium | 1162 |
| Ciprofloxacin Hcl | 1143 |
| Fluticasone Propionate | 1073 |
| Tramadol Hcl | 1038 |
| Meloxicam | 1035 |
| Amoxicillin | 1032 |
| Azithromycin | 1022 |
| Oxycodone Hcl | 1009 |
| Metoprolol Tartrate | 956 |
| Hydrochlorothiazide | 947 |
| Albuterol Sulfate | 820 |
| Pantoprazole Sodium | 817 |
| Metoprolol Succinate | 801 |
| Sulfamethoxazole/Trime... | 792 |
| Morphine Sulfate | 757 |
| Methylprednisolone | 753 |
| Furosemide | 748 |
| Cephalexin | 708 |

## Total Claims by Generic Name of Drug

| Generic Name of Drug | Total Claims |
|---|---|
| Atorvastatin Calcium | 8786 |
| Gabapentin | 7755 |
| Levothyroxine Sodium | 6933 |
| Amlodipine Besylate | 6664 |
| Hydrocodone/Acetamino... | 6126 |
| Lisinopril | 6071 |
| Simvastatin | 5812 |
| Tamsulosin Hcl | 5755 |
| Metformin Hcl | 4849 |
| Losartan Potassium | 4324 |
| Oxycodone Hcl/Acetamin... | 4323 |
| Omeprazole | 4174 |
| Alprazolam | 3779 |
| Metoprolol Tartrate | 3595 |
| Hydrochlorothiazide | 3526 |
| Baclofen | 3101 |
| Pantoprazole Sodium | 3048 |
| Tramadol Hcl | 2974 |
| Metoprolol Succinate | 2886 |
| Apixaban | 2853 |
| Diclofenac Sodium | 2636 |
| Oxycodone Hcl | 2549 |
| Furosemide | 2517 |
| Albuterol Sulfate | 2455 |
| Meloxicam | 2413 |
| Temazepam | 2304 |
| Escitalopram Oxalate | 2172 |
| Fluticasone Propionate | 2115 |
| Finasteride | 1985 |
| Morphine Sulfate | 1961 |
| Pravastatin Sodium | 1871 |

## Total 30day Fills by Generic Name of Drug

| Generic Name of Drug | Total 30day Fills |
|---|---|
| Atorvastatin Calcium | 23294 |
| Amlodipine Besylate | 16935 |
| Levothyroxine Sodium | 16326 |
| Simvastatin | 16070 |
| Lisinopril | 15340 |
| Metformin Hcl | 12910 |
| Tamsulosin Hcl | 12831 |
| Losartan Potassium | 11479 |
| Gabapentin | 11261 |
| Omeprazole | 9760 |
| Hydrochlorothiazide | 9178 |
| Metoprolol Tartrate | 8444 |
| Metoprolol Succinate | 7334 |
| Pantoprazole Sodium | 6935 |
| Hydrocodone/Acetamino... | 6130 |
| Furosemide | 5357 |
| Pravastatin Sodium | 4768 |
| Glipizide | 4756 |
| Finasteride | 4628 |
| Clopidogrel Bisulfate | 4389 |
| Allopurinol | 4329 |
| Oxycodone Hcl/Acetamin... | 4328 |
| Rosuvastatin Calcium | 4326 |
| Alendronate Sodium | 4320 |
| Escitalopram Oxalate | 4190 |
| Carvedilol | 4117 |
| Alprazolam | 3949 |
| Montelukast Sodium | 3767 |
| Apixaban | 3570 |
| Atenolol | 3545 |
| Meloxicam | 3360 |

## Total Day's Supply by Prescriber Type



| | |
|---|---|
| Family Practice | Internal Medicine |
| | Urology |
| | Rheumatol... |
| | Cardiology |
| | Endocrinology |
| | Anesthesiol... |
| | Inf... |
| | N... |
| | Obstetrics ... |
| | Pain Management |
| | Nephrology |
| | Neurology |
| | Phy... |
| | Ort... |

## Total Day's Supply by Prescriber Source



NPPES-Taxonomy
0M (0.09%)

NPPES-Specialty
4M (32.27%)

Claim-Specialty
8M (67.63%)

**Prescriber Source**
- Claim-Specialty
- NPPES-Specialty
- NPPES-Taxonomy

## Sum of Total Day's Supply by Year



Sum of Total Day's Supply

3.5M
3.0M
2.5M
2.0M
1.5M
1.0M
0.5M
0.0M

2019   2020   2021   2022

Year

# 10. Test Pipeline Run

Home   Activities   Run   View

💾  📝  ⚙️  ↶  ✓ Validate  ▷ Run  ▦ Schedule  ⚡ Add trigger (preview)  🕘 View run history  |  🔲 Copy data ⌄  ⌥ Dataflow  🔲 Notebook  🔍 Lookup  📢 Invoke Pipelin

LookupSourcedata

Invoke Pipeline (Previe ⬏ ✓
Invoke pipeline1
json to csv copy

Notebook ✓
Dim_Fact_table creation

ForEach ✓
ForEachAPIkey

Activities ✎
(x) Set variable... → (x) Set variable... → Copy data API to... → ⊕

Notebook ✓
bronze2silver table incremental load

Semantic model refresh...
Semantic model...

Parameters   Variables   Settings   **Output**

Pipeline run ID: 8b59ac4a-b5a7-4bb1-8a1a-c36e765ff95f ▥ ↻ ⓘ          Pipeline status ✓ Succeeded

🔍 Filter by keyword          Showing 17 items

| Activity name ↑↓ | Activity status ↑↓ | Duration ↑↓ |
|---|---|---|
| 🔍 LookupSourcedata | ✓ Succeeded | 12s |
| ▾ 🔲 ForEachAPIkey | ✓ Succeeded | 1m 49s |
| (x) Set variable year_id | ✓ Succeeded | Less than 1s |
| (x) Set variable year_name | ✓ Succeeded | Less than 1s |
| 🔲 Copy data API to JSON | ✓ Succeeded | 22s |
| (x) Set variable year_id | ✓ Succeeded | Less than 1s |
| (x) Set variable year_name | ✓ Succeeded | Less than 1s |
| 🔲 Copy data API to JSON | ✓ Succeeded | 29s |
| (x) Set variable year_id | ✓ Succeeded | Less than 1s |
| (x) Set variable year_name | ✓ Succeeded | Less than 1s |
| 🔲 Copy data API to JSON | ✓ Succeeded | 21s |
| (x) Set variable year_id | ✓ Succeeded | Less than 1s |
| (x) Set variable year_name | ✓ Succeeded | Less than 1s |
| 🔲 Copy data API to JSON | ✓ Succeeded | 24s |
| 📢 Invoke pipeline1 | ✓ Succeeded | 1m 14s |
| 🔲 bronze2silver table incremental load | ✓ Succeeded | 2m 40s |
| 🔲 Dim_Fact_table creation | ✓ Succeeded | 5m 47s |

# Thanks for reading