

Assignment 3

Building ELT Data Pipelines with Airflow



Roy Hegde
Student ID: 24667610

Table of Contents

1. Project Overview	2
2. Dataset Overview	3
3. Data Ingestion and Staging (Bronze Layer)	4
a. Airflow DAG Setup for Data Ingestion	4
b. Schema and Table Design in the Bronze Layer	4
c. Loading Process and Data Quality Control	5
4. Data Cleaning and Transformation (Silver Layer)	6
a. Silver Layer Transformation Design	6
b. Model Execution and Data Quality Checks	7
5. Snapshot Implementation	8
a. Airflow DAG Setup for Data Ingestion	8
b. Benefits of Snapshotting	8
6. Star Schema Design (Gold Layer)	10
a. Fact and Dimension Tables	10
b. Schema Optimization and Data Integrity	11
7. Datamart Views	12
a. Dm_listing_neighbourhood	12
b. Dm_property_type	12
c. Dm_property_type	13
8. End-to-End Orchestration (Airflow DAG)	14
a. DAG Design and Sequential Data Processing	14
b. Triggering dbt Transformations	14
c. Monitoring and Error Handling	15
9. Challenges and Solutions	16
a. Sequential Data Loading and Order Integrity	16
b. Data Consistency and Type Validation	16
c. Integration of Airflow and dbt Cloud	16
d. Data Integrity in Slowly Changing Dimensions (SCD2)	17
10. Conclusion	18

1. Project Overview

The goal of this project is to create a reliable, automated data pipeline that can ingest, convert, and analyse census and Airbnb statistics to provide insightful information about neighbourhood demographics and Airbnb listing performance. This pipeline uses Dbeaver-PostgreSQL, Apache Airflow, dbt Cloud, and Google Cloud Storage (GCS) in conjunction with the Medallion Architecture (Bronze, Silver, Gold) to organise data for the best possible analysis.

The Medallion Architecture layers are designed as follows:

- **Bronze Layer:** Retains unprocessed raw data from a variety of sources, including local government entities, Airbnb, and census data.
- **Silver Layer:** Prepares organised tables for analysis by cleaning and standardising the data.
- **Gold Layer:** Uses certain datamart views to address business problems and arranges data into a star schema for effective querying.

Every stage of the project is described in this report, including data ingestion, transformation, Airflow orchestration, and analytical insights from dbt-generated views.



2. Dataset Overview

Both static and dynamic data were utilised in this project; they were all kept in Google Cloud Storage and accessible via Airflow for PostgreSQL ingestion.

- **Airbnb Listings:** A set of monthly CSV files with listing data, including listing_id, host_name, property_type, price, and availability metrics (e.g., 05_2020.csv to 04_2021.csv). To ensure chronological consistency, this data is fed into the pipeline month by month.
- **Census Data:** 2016Census_G01_NSW_LGA.csv and 2016Census_G02_NSW_LGA.csv primary census files include socioeconomic and demographic data on the Local Government Areas (LGAs) of New South Wales. While G02 covers measures like median age, income, and rent, G01 contains population statistics.
- **LGA Codes and Suburb Mapping:** Suburbs are mapped to LGAs using two more static datasets, NSW_LGA_CODE.csv and NSW_LGA_SUBURB.csv. Analysis of listings at the neighbourhood level is supported by this spatial reference.

By storing these datasets in GCS, scalability is guaranteed in the event that additional files are added or dataset sizes grow, and flexible access for ingestion is made possible.



3. Data Ingestion and Staging (Bronze Layer)

In order to ensure data integrity and facilitate traceability, the Bronze layer acts as the fundamental staging area for raw data intake, maintaining each dataset's original structure. With just few modifications, all datasets are fed into PostgreSQL tables inside the bronze schema.

a. Airflow DAG Setup for Data Ingestion

The data intake procedure is coordinated by a `load_to_bronze_schema_sequential` Apache Airflow Directed Acyclic Graph (DAG). In order to guarantee that data from all sources is present in the Bronze layer prior to downstream operations starting, the DAG comprises jobs for loading static files first, then monthly Airbnb datasets.

The Airflow DAG consists of:

- **Static File Ingestion:** Static files (census data and LGA mappings) are downloaded and loaded from GCS to PostgreSQL in separate operations. Each job loads data into PostgreSQL using a PostgresHook and downloads files using the GCSToLocalFilesystemOperator and a custom `load_to_postgres` function.
- **Sequential Monthly Airbnb Ingestion:** The DAG loads each Airbnb file one after the other in order to preserve chronological order. A loop that creates a task for every month and uses dependencies between jobs to guarantee the right sequence enforces this order.

The screenshot shows the Google Cloud Storage interface. On the left, there's a sidebar with 'Cloud Storage' selected under 'Buckets'. The main area is titled 'Bucket details' for 'australia-southeast1-bde-0c9d64c9-bucket'. It shows a 'Folder browser' with 'airbnb_part_1' expanded, revealing subfolders 'dags/' and 'logs/'. The 'data/' folder is also visible. Below this is a table listing 18 CSV files, each with a download icon, size (e.g., 4.5 MB), type (text/csv), creation date (Oct 26, 2024), and storage class (Standard). The table includes columns for Name, Size, Type, Created, and Storage class.

Name	Size	Type	Created	Storage class
01_2021.csv	4.5 MB	text/csv	Oct 26, 2024, 9:50:43 PM	Standard
02_2021.csv	4.5 MB	text/csv	Oct 26, 2024, 9:50:43 PM	Standard
03_2021.csv	4.4 MB	text/csv	Oct 26, 2024, 9:50:36 PM	Standard
04_2021.csv	4.3 MB	text/csv	Oct 26, 2024, 9:50:40 PM	Standard
05_2020.csv	4.7 MB	text/csv	Oct 21, 2024, 5:59:34 PM	Standard
06_2020.csv	4.6 MB	text/csv	Oct 26, 2024, 9:50:36 PM	Standard
07_2020.csv	4.5 MB	text/csv	Oct 26, 2024, 9:50:44 PM	Standard
08_2020.csv	4.2 MB	text/csv	Oct 26, 2024, 9:50:53 PM	Standard
09_2020.csv	4.7 MB	text/csv	Oct 26, 2024, 9:50:50 PM	Standard
10_2020.csv	4.6 MB	text/csv	Oct 26, 2024, 9:50:48 PM	Standard
11_2020.csv	4.5 MB	text/csv	Oct 26, 2024, 9:50:54 PM	Standard
12_2020.csv	4.5 MB	text/csv	Oct 26, 2024, 9:50:55 PM	Standard
2016Census_G01_NSW_LGA.csv	66.7 KB	text/csv	Oct 21, 2024, 5:59:20 PM	Standard
2016Census_G02_NSW_LGA.csv	5.8 KB	text/csv	Oct 21, 2024, 5:59:20 PM	Standard
NSW_LGA_CODE.csv	2.3 KB	text/csv	Oct 21, 2024, 5:59:08 PM	Standard
NSW_LGA_SUBURB.csv	215.2 KB	text/csv	Oct 21, 2024, 5:59:10 PM	Standard

Figure 3.1 Google Cloud Services (GCS) bucket

The screenshot shows the DBeaver Database Navigator interface. The left pane displays a tree view of the database structure. At the top level, there are 'Projects' and 'Database Navigator'. Under 'Database Navigator', it shows 'D Beamer Sample Database (SQLite)' and 'postgres 34.151.104.161:5432'. The 'postgres' entry is expanded, showing 'Databases', 'Schemas', and 'Tables'. The 'bronze' schema is selected, showing its sub-entities: 'Tables' (with entries like 'raw_airbnb_listings', 'raw_census_g01', etc.) and 'Foreign Tables', 'Views', 'Materialized Views', 'Indexes', 'Functions', 'Sequences', 'Data types', and 'Aggregate functions'. Other sections like 'public', 'Event Triggers', 'Extensions', 'Storage', 'System Info', 'Roles', 'Administer', and 'System Info' are also listed.

Figure 3.2 DBeaver Postgres Bronze Schema

The screenshot shows a database interface with a title bar 'fact_listings 1 X'. Below it is a search bar with the query 'SELECT DISTINCT month_year FROM fact_listings'. The main area is a grid table with the following data:

Grid	month_year
1	2020-05-01 10:00:00.000 +1000
2	2020-06-01 10:00:00.000 +1000
3	2020-07-01 10:00:00.000 +1000
4	2020-08-01 10:00:00.000 +1000
5	2020-09-01 10:00:00.000 +1000
6	2020-10-01 10:00:00.000 +1000
7	2020-11-01 11:00:00.000 +1100
8	2020-12-01 11:00:00.000 +1100
9	2021-01-01 11:00:00.000 +1100
10	2021-02-01 11:00:00.000 +1100
11	2021-03-01 11:00:00.000 +1100
12	2021-04-01 11:00:00.000 +1100

Below the grid, there are buttons for Refresh, Save, Cancel, and various database operations. A status message at the bottom indicates '12 row(s) fetched - 0.195s, on 2024-10-27 at 01:00:18'.

Figure 3.3 All months imported from Airflow

b. Schema and Table Design in the Bronze Layer

Tables in the Bronze schema reflect each CSV file's structure. For instance, Airbnb data with attributes like listing_id, host_name, and price is stored in raw_airbnb_listings. This table acts as a staging region, keeping any raw fields that could require cleaning or restructuring in the following levels while storing data precisely as it appears in the source.

Other tables in the Bronze schema include:

- **raw_census_g01** and **raw_census_g02**: These tables, which include variables like total_population_male, median_age_persons, and median_rent_weekly, hold demographic and socioeconomic data.
- **raw_lga_codes**: Includes LGA names and codes, making it possible to identify LGAs.
- **raw_lga_suburbs**: helps with neighborhood-level aggregations later on by mapping suburbs to their corresponding LGAs.

c. Loading Process and Data Quality Control

The `load_to_postgres` function manages row-wise insertion into PostgreSQL and eliminates unnecessary columns (such as unnamed columns from CSV outputs) to guarantee data quality. By allowing for precise control over the integrity of each row, this procedure guarantees that only structured data makes it to the Bronze tables.

```
from airflow import DAG
from airflow.providers.google.cloud.transfers.gcs_to_local import GCSToLocalFilesystemOperator
from airflow.providers.postgres.hooks.postgres import PostgresHook
from airflow.operators.python_operator import PythonOperator
from airflow.models import Variable
from datetime import datetime
import pandas as pd
import requests
import logging

# Default arguments for the DAG
default_args = {
    'start_date': datetime(2024, 1, 1),
    'catchup': False,
    'email_on_failure': True,
    'email_on_retry': False,
    'retries': 2,
}

# List of Airbnb files to be processed sequentially
airbnb_files = [
    '05_2020.csv', '06_2020.csv', '07_2020.csv', '08_2020.csv',
    '09_2020.csv', '10_2020.csv', '11_2020.csv', '12_2020.csv',
    '01_2021.csv', '02_2021.csv', '03_2021.csv', '04_2021.csv'
]

# Define the DAG
with DAG(
    dag_id='load_to_bronze_schema_sequential',
    schedule_interval=None, # No schedule interval, manually triggered
    default_args=default_args,
    description='Load raw data from GCS to Bronze schema in Postgres sequentially and trigger dbt job',
    catchup=False,
) as dag:

    # Task: Download and load static files (Census and LGA data)
    def load_to_postgres(file_path, table_name):
        postgres_hook = PostgresHook(postgres_conn_id='postgres')
        conn = postgres_hook.get_conn()
        cursor = conn.cursor()

        # Load CSV data into a Pandas DataFrame
        df = pd.read_csv(file_path)
        df = df.loc[:, ~df.columns.str.contains('^Unnamed')] # Drop unnamed columns if any

        # Insert DataFrame rows into Postgres table
        for _, row in df.iterrows():
            placeholders = ', '.join(['%s'] * len(row))
            columns = ', '.join(df.columns)
            sql = f'INSERT INTO bronze.{table_name} ({columns}) VALUES ({placeholders})'
            cursor.execute(sql, tuple(row))

        conn.commit()
        cursor.close()
        conn.close()

    # Task: Download and load static files (Census and LGA data)
    task_load_to_postgres = PythonOperator(
        task_id='load_to_postgres',
        python_callable=load_to_postgres,
        op_kwargs={'file_path': 'path/to/airbnb_files', 'table_name': 'bronze_airbnb'}
    )
```

Figure 3.4 Airflow DAG (`load_to_postgres`)

Airflow's logging features are used to keep an eye on job execution and debug any data ingestion problems, including missing files or corrupted rows, during the loading process.



4. Data Cleaning and Transformation (Silver Layer)

The Bronze layer's raw data is transformed into organised, analysis-ready tables by the Silver layer using data transformation and cleaning procedures. This layer makes ensuring that fields are correct, consistent, and structured for effective aggregation and joining in the Gold layer across many datasets.

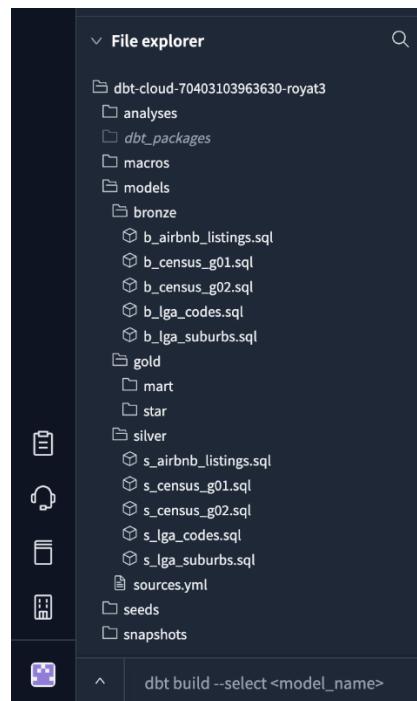


Figure 4.1 dbtcloud folder structure

a. Silver Layer Transformation Design

Every Silver layer model has a corresponding Bronze table and undergoes modifications that deal with:

- **Data Type Casting:** For computations and aggregations, fields such as IDs and prices are converted to the proper kinds (such as integer or numeric).
- **Handling Missing Values:** In fields like host_name and host_neighbourhood, null entries are swapped out for default values like "Unknown."

- **Standardizing Formats:** Ensuring that dates and category categories are formatted consistently across tables, including standardising text fields' case and converting dates to YYYY-MM-DD.
- **Validating Values:** Removing excessive or erroneous numbers, such as out-of-range review scores or negative pricing.

The dbt Cloud models created for this layer include:

- **s_airbnb_listings:** This table filters and converts data types for Airbnb listings. For example, pricing is verified to be within a suitable range, and listing_id and host_id are cast to integers. Review ratings are verified to make sure they fall within a certain range, and columns such as host_is_superhost are transformed to Boolean for precise analysis.
- **s_census_g01 and s_census_g02:** The process of cleaning census data involves converting population and socioeconomic indicators (such as total_population_male and median_rent_weekly) to integer or float types and eliminating non-numeric characters from lga_code_2016. To preserve data integrity, fields that are essential for joining, such as lga_code, must not be null.
- **s_lga_codes and s_lga_suburbs:** By removing whitespace from names and making sure that all lga_code and suburb_name values are non-null, these tables clean and format data pertaining to LGAs. When combining with other tables for geographic analysis, this ensures consistency.

b. Model Execution and Data Quality Checks

In DBT, every Silver layer model is set up to execute with data quality checks that confirm the existence and kind of important fields. Since transformations standardise the data and establish the groundwork for precise analysis in the Gold layer, data quality is essential in this tier. Validation checks on crucial fields (such as primary keys and foreign keys) guarantee that clean, structured data is accessible for additional processing while the models are run within dbt Cloud.

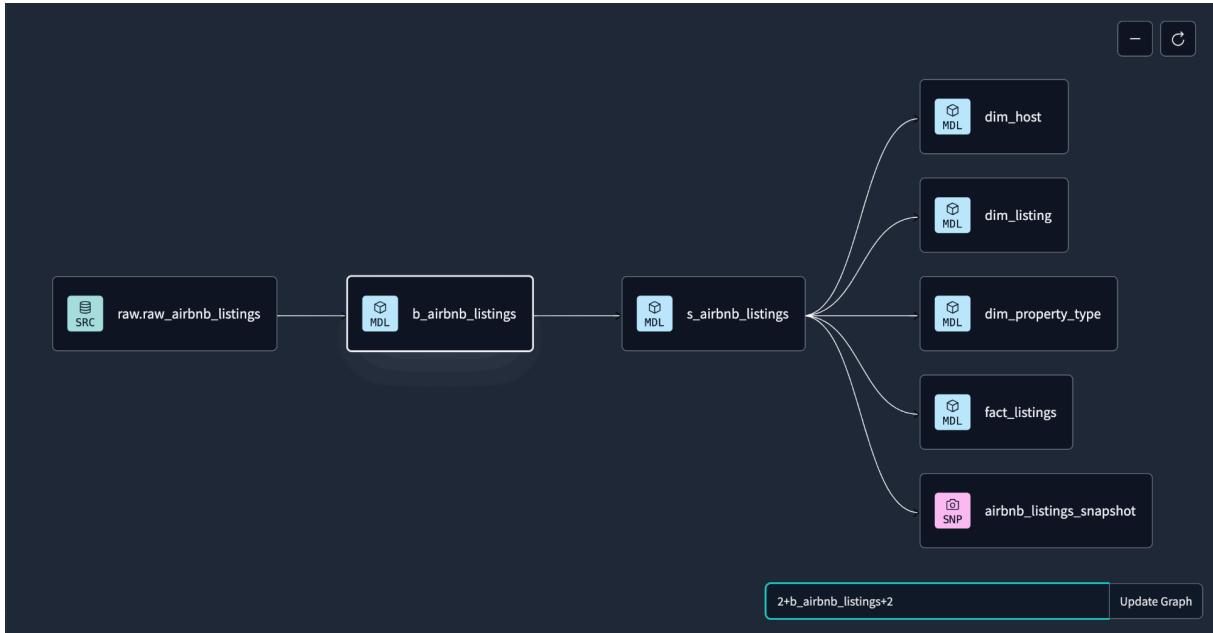


Figure 4.2 Lineage of `airbnb_listings`



Figure 4.3 Lineage of `census_g01`



Figure 4.4 Lineage of `census_g02`



Figure 4.5 Lineage of `lga_codes`

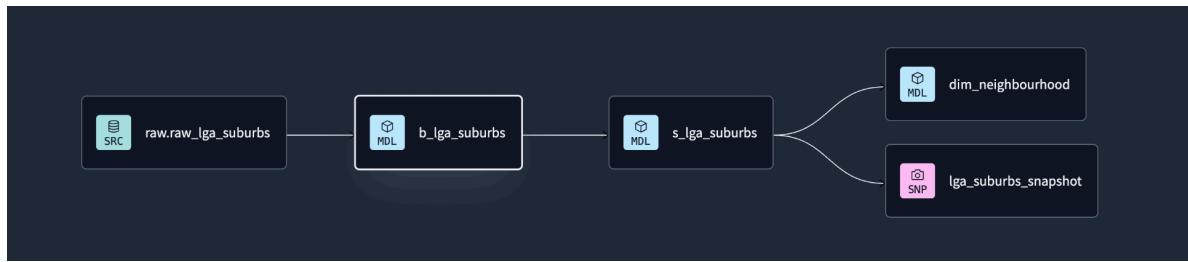


Figure 4.6 Lineage of `lga_suburbs`

■ ■ ■

5. Snapshot Implementation

In dbt, snapshots are crucial for monitoring changes over time, especially for variables that could vary over time, such as census demographics or listing pricing. Historical analysis and comparisons are made possible by this project's implementation of snapshots on Silver layer core tables.

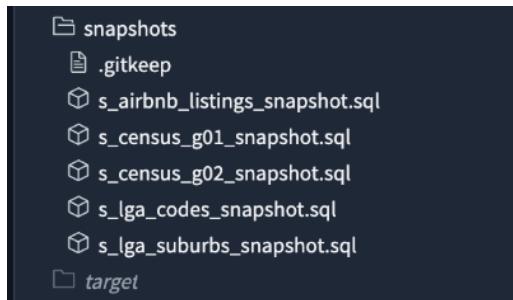


Figure 5.1 dbtcloud snapshots folder

a. Airflow DAG Setup for Data Ingestion

Strategy = 'check' is set up for each snapshot, comparing important fields (check columns) in every table to identify and document changes over time. Important snapshots consist of:

- **s_airbnb_listings_snapshot:** Monthly changes in listing parameters such as price, availability_30, number_of_reviews, and review_scores_rating are captured in this snapshot. Each record is guaranteed to belong to a specific listing by the unique key listing_id, and any modifications made to the check fields cause a new version of that listing to be recorded.
- **s_census_g01_snapshot and s_census_g02_snapshot:** Using distinct keys depending on lga_code, these snapshots document changes in the socioeconomic and demographic makeup of LGAs. To record any changes or new values over time, metrics like median_tot_hhd_inc_weekly and total_population_persons are monitored.
- **s_lga_codes_snapshot and s_lga_suburbs_snapshot:** In order to preserve any revisions to LGA names or suburb mappings for future use, these snapshots save geographic data for LGAs and suburbs. When administrative borders or naming standards change, this is crucial for longitudinal data analysis.

b. Benefits of Snapshotting

By preserving past data, snapshots enable in-depth trend analysis and comparisons across time. The project can address difficult business challenges by putting these pictures into practice, including:

- **Revenue and Price Trends:** Tracking market dynamics is made easier by monthly variations in price and expected revenue across several listings.
 - **Demographic Shifts:** As census data changes, socioeconomic assessments of neighbourhoods are made possible, which helps inform judgements about Airbnb activities in certain demographic regions.
- ■ ■

6. Star Schema Design (Gold Layer)

In order to facilitate analytics and optimise querying, the Gold layer arranges data into a star schema. This is especially useful for addressing business enquiries about demographic trends, revenue, and occupancy. With the help of dimension tables that offer descriptive qualities and fact tables that store metrics and foreign keys, this design facilitates easy aggregations and efficient joins.

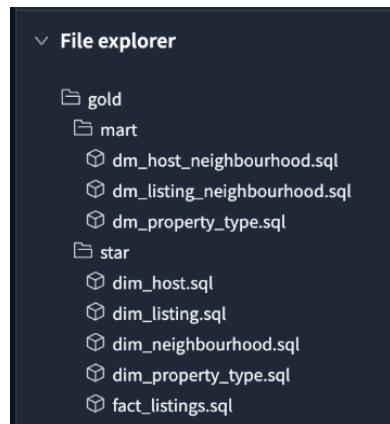


Figure 6.1 dbtcloud gold folder structure

a. Fact and Dimension Tables

The star schema in the Gold layer contains the following tables:

Fact Table: fact_listings:

- Key metrics and connections to other dimensions are included in this table via foreign keys like listing_id and host_id. Price, number of reviews, availability_30, number of stays, and expected income are among the metrics kept in this table.
- To enable effective analysis of monthly patterns, each row denotes a distinct listing for each month (month_year). For every listing, calculated data such as number_of_stays and estimated_revenue offer information on occupancy and revenue.

Dimension Tables:

- **dim_host:** host_id, host_name, host_neighbourhood, host_since, and host_is_superhost are among the host properties that are stored. Analysis based on host quality and involvement (e.g., superhost status) is supported by this table.
- **dim_listing:** Includes information on each listing, such as the listing's ID, neighbourhood, kind of property, type of room, and accommodations. This makes it possible to filter by listing attributes like neighbourhood and property type.
- **dim_neighbourhood:** Allows for mappings between suburb_name and lga_name, facilitating multi-level geographic analysis.
- **dim_property_type:** Includes room_type, accommodates, and property_type, which aid in classifying and evaluating postings according to the type of property.

b. Schema Optimization and Data Integrity

This layout reduces data duplication and creates a normalised structure by using foreign keys to connect the fact table to its dimensions. Consistent joins and trustworthy analytics are made possible by data types and non-null constraints, which guarantee that important values are appropriately referenced across tables.

The pipeline can handle large-scale queries with little processing time when using this star schema. This schema is perfect for the datamart views in the following section since it streamlines aggregations, filtering, and joins.



7. Datamart Views

The Gold layer's datamart views offer organised insights designed to address certain business queries. Each view provides high-level metrics for decision-making by combining information from fact and dimension tables. To facilitate temporal analysis and represent Slowly Changing Dimensions Type 2, these perspectives are represented in dbt as monthly snapshots (SCD2).

a. Dm_listing_neighbourhood

Data is compiled by listing neighbourhood and month in the dm_listing_neighbourhood view. Avg_estimated_revenue_per_active_listing, estimated_revenue, distinct_hosts, and percentage changes in active and inactive listings are among the metrics it offers.

Key metrics:

- Total Listings and Active Listings:** Monthly totals for all residential neighbourhoods' active and total listings.
- Revenue and Host Engagement:** Total and average revenue figures, superhost rates, and unique host counts are calculated.
- Trend Analysis:** The percentage change in active and inactive listings from month to month that enables the detection of growth and trends in the neighbourhood

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	
dm_listing_neighbourhood_202410281714																	
1	listing_neighbourhood	month_year	active_listings	inactive_listings	total_listings	min_price	max_price	median_price	avg_price	distinct_hosts	superhost_rate	avg_review_scores_rating	total_stays	avg_estimated_revenue_per_active_listing	pct_change_active_listings	pct_change_inactive_listings	
2	Bayview	2020-05-01 10:00:00.000 +1000	3232	0	3232	0	2544	86.0	144.292980397020000	1200	0	92.05448513902110	61798	2487.7687490799500000			
3	Blacktown	2020-05-01 10:00:00.000 +1000	636	0	636	23	1099	63.0	95.301886792459000	239	0	89.92098400960230	7256	1037.92369729565070000			
4	Burwood	2020-05-01 10:00:00.000 +1000	542	0	542	15	2440	70.0	148.549815498159000	169	0	89.48096093032540	9230	2000.214021402200000			
5	Camden	2020-05-01 10:00:00.000 +1000	98	0	98	35	360	85.0	107.626552001294000	41	0	94.36666666666670	1212	1983.734990715500000			
6	Campbelltown	2020-05-01 10:00:00.000 +1000	238	0	238	29	2650	75.0	156.28071742857149000	75	0	91.92333233333320	3700	1291.7319931496700000			
7	Canada Bay	2020-05-01 10:00:00.000 +1000	928	0	928	24	2440	110.0	154.165948275960000	384	0	92.894915243730	17846	2838.644396517200000			
8	Canterbury-Bankstown	2020-05-01 10:00:00.000 +1000	1314	0	1314	12	1000	71.0	96.4736571747139000	502	0	90.63461038401540	18050	1338.185692141860000			
9	Cumberland	2020-05-01 10:00:00.000 +1000	1164	0	1164	15	2491	100.0	137.087779980189000	406	0	92.07471624468760	21448	2454.41622444960200000			
10	Fairfield	2020-05-01 10:00:00.000 +1000	102	0	102	20	561	76.0	102.351166669976000	56	0	89.6339031600140	2034	1237.552083333320000			
11	Georges River	2020-05-01 10:00:00.000 +1000	833	0	833	0	3500	80.0	130.912098030780000	312	0	92.16930175849100	14706	1854.209134015280000			
12	Hornsby	2020-05-01 10:00:00.000 +1000	930	0	930	20	5000	91.0	137.7643496900000	338	0	93.4901607843140	15162	2004.355373344100000			
13	Hunters Hill	2020-05-01 10:00:00.000 +1000	108	0	108	35	2440	185.0	206.57140140740740000	54	0	94.9722202232230	2182	7110.240740740740000			
14	Inner West	2020-05-01 10:00:00.000 +1000	4886	0	4886	15	9022	106.0	181.165451780598000	2010	0	93.43704121964900	108650	3745.191154111790000			
15	Lane Cove	2020-05-01 10:00:00.000 +1000	632	0	632	35	3000	110.0	230.072784610127000	249	0	90.0972202232230	10950	4419.930379146842000			
16	Liverpool	2020-05-01 10:00:00.000 +1000	308	0	308	15	586	98.0	125.467332407320000	113	0	91.23541094339600	4486	1786.337765331860000			
17	Mosman	2020-05-01 10:00:00.000 +1000	1020	0	1020	31	10000	190.0	448.041986497270000	370	0	94.80116894082000	21574	2809.145330194130000			
18	Northern Beaches	2020-05-01 10:00:00.000 +1000	10800	0	10800	12	7654	200.0	349.323539895040000	4152	0	95.153377844880	21294	6747.829917460200000			
19	North Sydney	2020-05-01 10:00:00.000 +1000	2700	0	2700	21	10000	180.0	221.323592598590000	1036	0	92.925602084210	53868	3852.508999996670000			
20	Penrith	2020-05-01 10:00:00.000 +1000	1210	0	1210	14	3000	97.0	106.874380165390000	427	0	90.673913040780	18530	2362.50279388400000			
21	Perth	2020-05-01 10:00:00.000 +1000	356	0	356	25	1500	111.5	140.286516850300000	116	0	95.5726841178470	4604	1964.44880247200000			
22	Ryde	2020-05-01 10:00:00.000 +1000	6566	0	6566	5	6000	119.0	191.262420409720000	2606	0	92.4700573042050	151280	1467.250457038900000			
23	Ryde	2020-05-01 10:00:00.000 +1000	1436	0	1436	20	2546	96.0	142.789986939510000	518	0	92.7154471547160	23862	2401.759029247900000			
24	Strathfield	2020-05-01 10:00:00.000 +1000	416	0	416	15	1000	98.5	123.422078292707000	157	0	88.542535291280	7194	2077.187500000000000			
25	Sutherland Shire	2020-05-01 10:00:00.000 +1000	1128	0	1128	35	2440	144.5	213.360629269500000	486	0	95.107855029240	20546	3856.437943264100000			
26	Sydney	2020-05-01 10:00:00.000 +1000	19186	0	19186	0	10000	135.0	200.0000000000000	6562	0	92.2573	383410	3559.600000000000000			
27	The Hills Shire	2020-05-01 10:00:00.000 +1000	650	0	650	24	6966	80.0	179.686461538462000	259	0	93.4227272727270	9596	1903.344615346200000			
28	Waverley	2020-05-01 10:00:00.000 +1000	10182	0	10182	5	10000	150.0	244.859860128680000	4210	0	93.8534528670460	240382	5386.759736794100000			
29	Willoughby	2020-05-01 10:00:00.000 +1000	1056	0	1056	26	2599	127.5	196.441287978800000	403	0	91.5734463278680	18410	3605.594969696700000			
30	Woolahra	2020-05-01 10:00:00.000 +1000	3098	0	3098	23	9998	150.0	314.229741802970000	1281	0	93.546926961270	68558	612.048644612400000			

Figure 7.1 dm_listing_neighbourhood table csv

b. Dm_property_type

The `dm_property_type` view provides information on listing performance based on property attributes by classifying data by `property_type`, `room_type`, `accommodates`, and `month_year`.

Key metrics:

- **Price and Occupancy Trends:** Contains the projected income for current listings, the total number of stays, and the minimum, maximum, median, and average prices.
 - **Host and Review Metrics:** Provides a quality indicator for every kind of property by displaying the average review ratings, distinct host count, and superhost rate.
 - **Trend Analysis:** Shows variations in market demand through monthly percentage changes in active and inactive listings by type of property.

Figure 7.2 dm_property_type table csv

c. Dm_host_neighbourhood

The datamart views answer several key business questions:

- **Revenue and Price Trends:** Stakeholders can monitor the most lucrative neighbourhoods or property types by looking at average price and average predicted revenue per active listing.
- **Host Engagement:** Particularly in areas with high visitor demand, the distinct_hosts counts and superhost_rate statistic provide evaluation of host quality and activity.
- **Occupancy and Demand:** Monitoring shifts in active_listings and inactive_listings reveals possible growth regions and offers insights into occupancy trends.

With these perspectives, the Gold layer provides useful information about Airbnb listings, allowing for data-driven choices to increase income and enhance visitor pleasure.

dm_host_neighbourhood_202410281712				
host_neighbourhood_lga	month_year	distinct_hosts	estimated_revenue	estimated_revenue_per_host
ARMIDALE REGIONAL	2020-05-01 10:00:00.000 +1000	1	3600	3600.000000000000000000
BATHURST REGIONAL	2020-05-01 10:00:00.000 +1000	10	376354	37635.400000000000000000
BAYSIDE	2020-05-01 10:00:00.000 +1000	43	237370	5520.2325581395300000
BAYSIDE	2020-05-01 10:00:00.000 +1000	24	150062	6252.5833333333300000
BAYSIDE	2020-05-01 10:00:00.000 +1000	9	40800	4533.3333333333300000
BAYSIDE	2020-05-01 10:00:00.000 +1000	108	1225390	11346.203703703700
BAYSIDE	2020-05-01 10:00:00.000 +1000	27	328704	12174.222222222200
BAYSIDE	2020-05-01 10:00:00.000 +1000	37	155584	4204.9729729729700000
BAYSIDE	2020-05-01 10:00:00.000 +1000	121	1031510	8524.8760330578500000
BAYSIDE	2020-05-01 10:00:00.000 +1000	55	343314	6242.0727272727300000
BAYSIDE	2020-05-01 10:00:00.000 +1000	44	315812	7177.5454545454500000
BAYSIDE	2020-05-01 10:00:00.000 +1000	4	14940	3735.0000000000000000
BAYSIDE	2020-05-01 10:00:00.000 +1000	7	48450	6921.4285714285700000
BAYSIDE	2020-05-01 10:00:00.000 +1000	23	89588	3895.1304347826100000
BAYSIDE	2020-05-01 10:00:00.000 +1000	2	9232	4616.0000000000000000
BAYSIDE	2020-05-01 10:00:00.000 +1000	4	18600	4650.0000000000000000
BAYSIDE	2020-05-01 10:00:00.000 +1000	11	74910	6810.0000000000000000
BAYSIDE	2020-05-01 10:00:00.000 +1000	41	309162	7540.5365853658500000
BAYSIDE	2020-05-01 10:00:00.000 +1000	92	516034	5609.0652173913000000
BAYSIDE	2020-05-01 10:00:00.000 +1000	15	916970	61131.33333333330000
BAYSIDE	2020-05-01 10:00:00.000 +1000	12	29566	2463.8333333333000000
BAYSIDE	2020-05-01 10:00:00.000 +1000	25	328680	13147.2000000000000000
BAYSIDE	2020-05-01 10:00:00.000 +1000	173	1360854	7866.2080924855500000
BAYSIDE	2020-05-01 10:00:00.000 +1000	3	10912	3637.3333333333300000
BAYSIDE	2020-05-01 10:00:00.000 +1000	2	39300	19650.0000000000000000
BURWOOD	2020-05-01 10:00:00.000 +1000	5	8580	1716.000000000000000000
BURWOOD	2020-05-01 10:00:00.000 +1000	13	71982	5537.0769230769200000
BURWOOD	2020-05-01 10:00:00.000 +1000	29	149574	5157.7241379310300000
BURWOOD	2020-05-01 10:00:00.000 +1000	57	418434	7340.0473484210500000

Figure 7.3 dm_host_neighbourhood table csv

8. End-to-End Orchestration (Airflow DAG)

The automated coordination of data intake, transformation, and loading using Apache Airflow is a crucial component of this project's design. From raw intake to analytics-ready views, data flows effectively thanks to our end-to-end orchestration, which guarantees a dependable and seamless data pipeline.

a. DAG Design and Sequential Data Processing

The following objectives guide the coordination of data input and transformation via the Airflow DAG, load_to_bronze_schema_sequential:

- **Chronological Data Loading:** To guarantee that each month's data is handled in the correct sequence, the DAG is set up to load Airbnb monthly data files sequentially. Maintaining chronological integrity is crucial, particularly for time-based metrics like number_of_reviews and monthly pricing.
- **Static Data Handling:** First, static datasets are loaded, such as census and LGA mappings. This information serves as the basis for joining and aggregations, which enable listings to be precisely linked to administrative areas and demographic data.

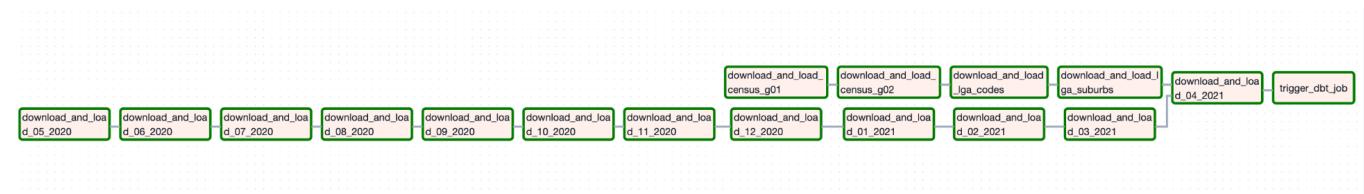


Figure 8.1 Airflow DAG execution of all monthly listings data

b. Triggering dbt Transformations

The DAG initiates a dbt Cloud task to carry out transformations when data loading is finished. By automating the data transfer from the Bronze layer (raw ingestion) to the Silver and Gold levels, this connection with dbt Cloud guarantees fast and reliable updates for all models.

- **API-Based Task:** The DAG pulls the required credentials from Airflow Variables for security and uses a PythonOperator to initiate the dbt Cloud task through an API request. Strong error handling is provided by this method, which logs any failures and raises exceptions to alert the team to problems.
- **Dependencies and Data Integrity:** The DAG guarantees that data is modified only when ingestion is finished, preserving data integrity throughout the pipeline, by imposing dependencies between jobs.

List Variable				
	Actions			Record Count: 5
	Key	Val	Description	Is Encrypted
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> api_key	*****		True
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> DBT_CLOUD_ACCOUNT_ID	70403103963630		True
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> DBT_CLOUD_API_TOKEN	*****		True
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> DBT_CLOUD_JOB_ID	70403104220169		True
<input type="checkbox"/>	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> DBT_CLOUD_URL	fc026.us1.dbt.com		True

Figure 8.2 Airflow Variables for API trigger (dbtcloud)

The screenshot shows the dbtcloud interface under the 'Jobs' section. At the top, there's a search bar labeled 'Search jobs...' and a dropdown for 'Environment' set to 'All'. A 'Create job' button is also visible. Below this, two job entries are listed:

- DBT Run**: Status is 'Last run succeeded 2m 38s ago', 'Not scheduled', and it's associated with 'Production' environment.
- DBT Run Test**: Status is 'Last run succeeded 3h 8m ago', 'Not scheduled', and it's associated with 'Production' environment.

Figure 8.3 dbtcloud jobs

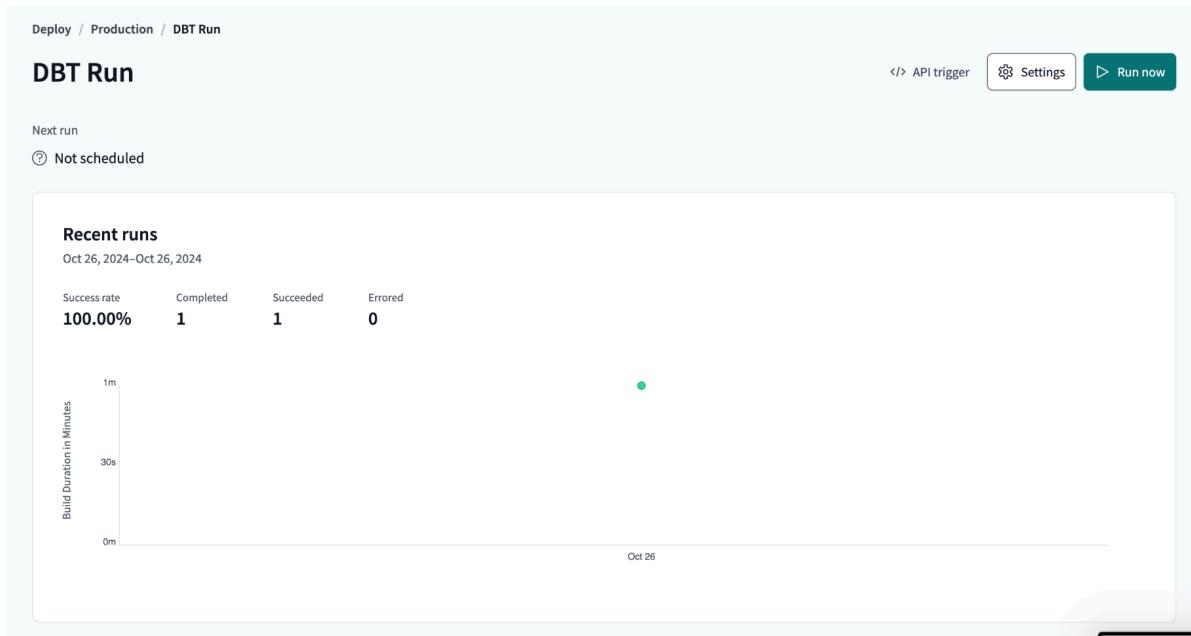


Figure 8.4 Successful dbtcloud job run

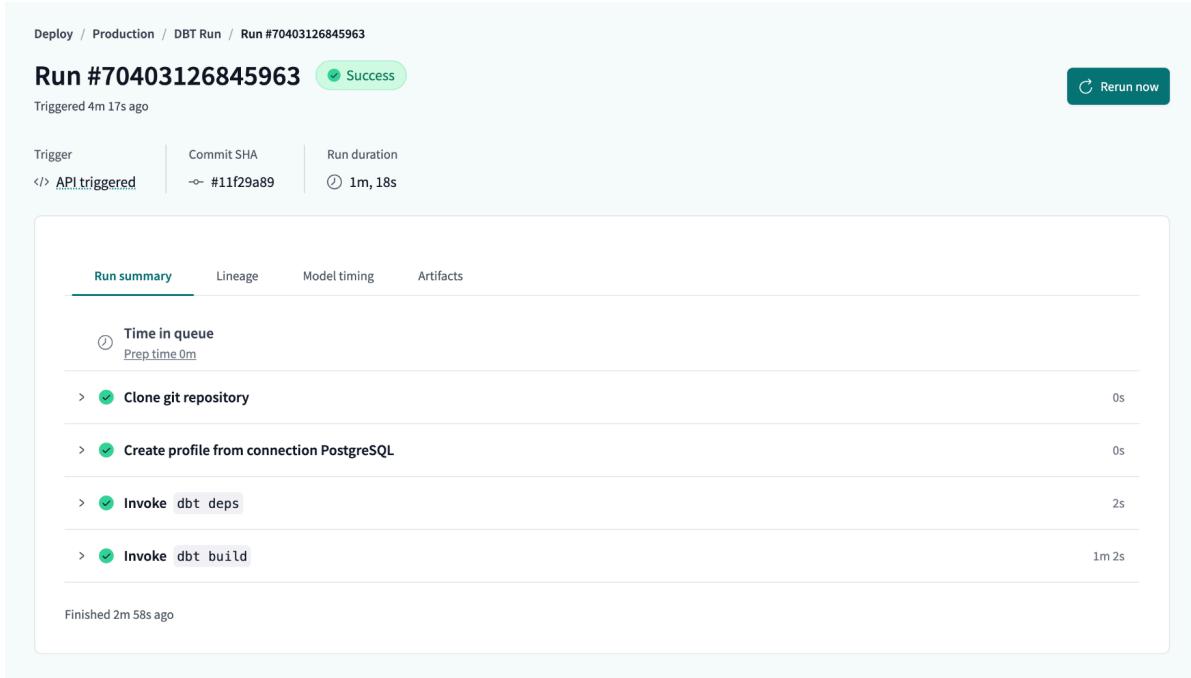


Figure 8.5 dbtcloud job run results

Invoke dbt build

1m 2s

Console Logs Debug Logs Search logs... Download logs

23 Successes

```

11:42:40 21 of 23 START sql table model gold.dim_property_type ..... [RUN]
11:42:40 22 of 23 START sql table model gold.fact_listings ..... [RUN]
11:42:46 19 of 23 OK created sql table model gold.dim_host ..... [SELECT 34193 in 6.22s]
11:42:46 23 of 23 START snapshot silver.airbnb_listings_snapshot ..... [RUN]
11:42:47 21 of 23 OK created sql table model gold.dim_property_type ..... [SELECT 81640 in 7.08s]
11:42:47 20 of 23 OK created sql table model gold.dim_listing ..... [SELECT 82515 in 7.42s]
11:42:48 22 of 23 OK created sql table model gold.fact_listings ..... [SELECT 899368 in 8.06s]
11:43:05 23 of 23 OK snapshotted silver.airbnb_listings_snapshot ..... [INSERT 0 1122574 in 19.35s]
11:43:07
11:43:07
11:43:07 Finished running 5 snapshots, 15 table models, 3 view models in 0 hours 0 minutes and 53.35 seconds (53.35s).
11:43:07
11:43:07
11:43:07 Completed successfully
11:43:07
11:43:07
11:43:07
11:43:07 Done. PASS=23 WARN=0 ERROR=0 SKIP=0 TOTAL=23

```

Finished 2m 58s ago

Figure 8.6 Logs of dbtcloud trigger

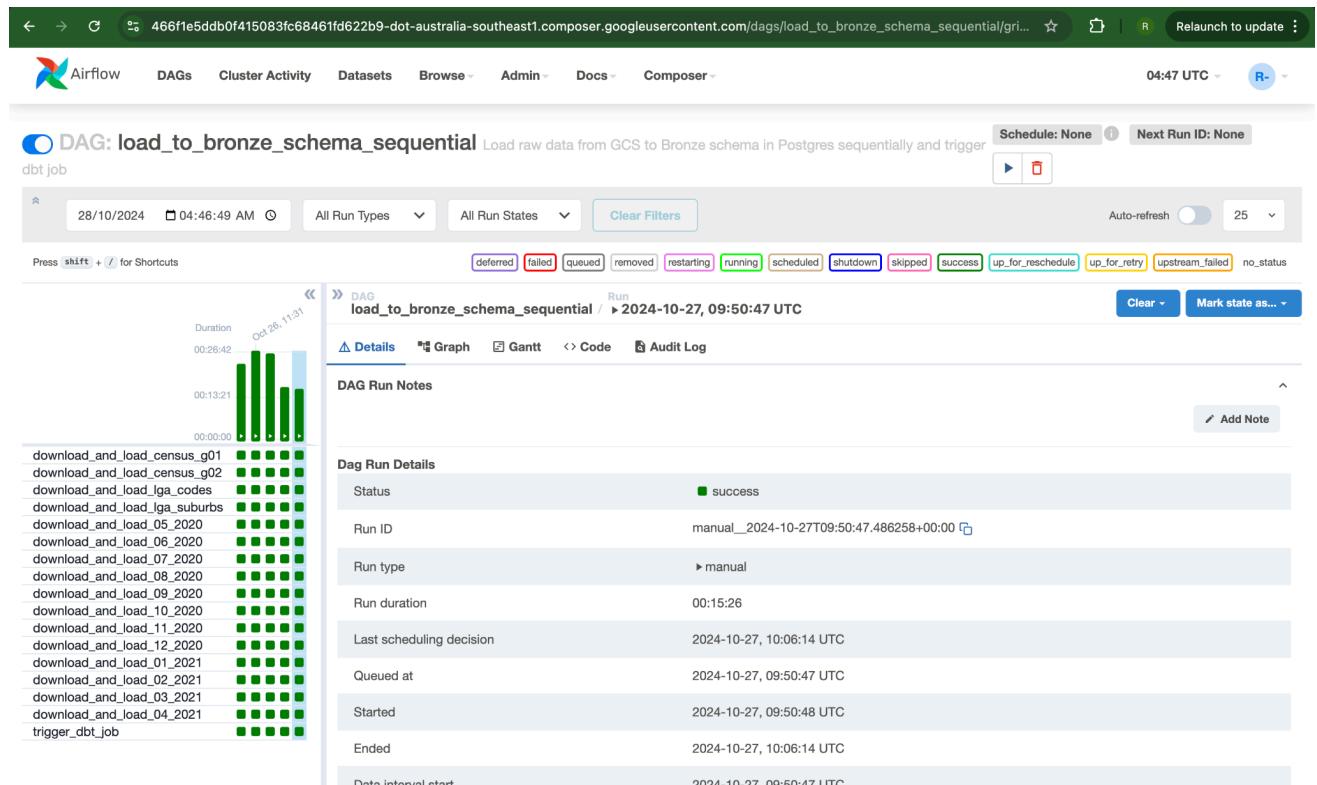


Figure 8.7 Airflow DAG status

c. Monitoring and Error Handling

The way the Airflow DAG is set up to deal with faults is by:

- **Retries and Alerts:** To make sure that team members are informed in the event of problems, the DAG's default_args setting incorporates email notifications upon failure as well as retries for temporary failures.
- **Logging:** Every job keeps track of its progress, giving insight into the transformation and loading procedures. Especially during data import, this recording is useful for pipeline monitoring and problem diagnosis.

By minimising manual involvement and guaranteeing that all transformations and views are promptly updated for precise analysis, this automated orchestration streamlines the data flow.



9. Challenges and Solutions

This project encountered several technical challenges, each addressed with strategic solutions to ensure data integrity, automation, and performance.

a. Sequential Data Loading and Order Integrity

- **Challenge:** Maintaining chronological correctness in the dataset, particularly for time-dependent studies, required that Airbnb data be fed sequentially per month.
- **Solution:** To enforce a rigorous sequence, the Airflow DAG was constructed with dependencies between monthly data loading operations. Processing of the data file for each month didn't start until the preceding file was completely loaded. The integrity of time-based measurements and trends was preserved by this sequential processing.

b. Data Consistency and Type Validation

- **Challenge:** A number of data types and certain inconsistent or missing values (such as wrong dates or negative prices) were included in the raw Airbnb dataset, which may have affected analysis and transformations.
- **Solution:** For fields like price, review_scores_rating, and availability_30, the Silver layer changes used range checks, null handling, and casting. In order to maintain a clean and consistent dataset for analysis, dbt tests were also included to check data types and make sure needed fields weren't null.

c. Integration of Airflow and dbt Cloud

- **Challenge:** To integrate Airflow with dbt Cloud in a secure and dependable manner so that dbt transformations would start instantly after data loading.
- **Solution:** via a PythonOperator and dbt credentials safely kept in Airflow Variables, Airflow initiated dbt Cloud tasks via an API-based methodology. This method made sure that the dbt transformations only executed once the data intake process was successfully finished.

d. Data Integrity in Slowly Changing Dimensions (SCD2)

- **Challenge:** Some tables, such Airbnb listings and census demographics, required historical monitoring in order to be analysed over time.
- **Solution:** SCD Type 2 was set up to record historical changes in particular fields in dbt snapshots. As a result, precise longitudinal analysis was made possible, supporting time-based patterns and demographic shifts in the data with historical context.



10. Conclusion

Using Google Cloud Storage, PostgreSQL, Apache Airflow, and dbt Cloud, this project effectively developed a fully coordinated data pipeline that was organised in accordance with the Medallion Architecture (Bronze, Silver, Gold). The pipeline ensures a stable flow of data from raw input to analytical insights by utilising an effective, tiered design that facilitates automated data ingestion, transformation, and aggregation.

Datamart views and the star schema in the Gold layer offer useful information about neighbourhood demographics and Airbnb listings. Targeted measurements and optimised data formats efficiently address business issues about occupancy trends, revenue trends, and property performance.

The project is positioned as a scalable solution for continuous Airbnb and demographic data analysis because of its automation, data validation, and historical recording, which allow for accurate, repeatable analysis. Future improvements, like adding more years to the dataset or including other pertinent data sources, may be built upon the strong foundation this pipeline offers.

