# Project: Adventure Works

# Foundation First !!! Four Key Questions

I. Where do we consolidate our data? > Storage

II. How will we get it there ? > Ingestion

III. How will we clean it up? > Transformation

IV. How will we analyze it? > Reporting

**BRIAN GWAYI** 

# Data Stack Popular Options

Storage > Snowflake, <u>BigQuery</u>, <u>s3</u>, Redshift Ingestion > Airbyte, <u>Airflow</u>, Fivetran Transformation > dbt Reporting > Tableau, Power BI, <u>Looker</u>, Superset

N/B This is not an exhaustive list.

**BRIAN GWAYI** 

# **Data Pipeline Architecture Design**

Source

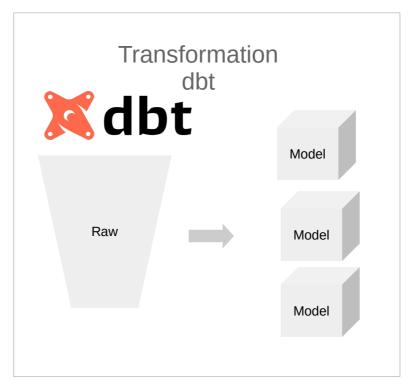


Ingestion



Storage
BigQuery

Amazon Redshift



Reporting



### Content

01 02 03

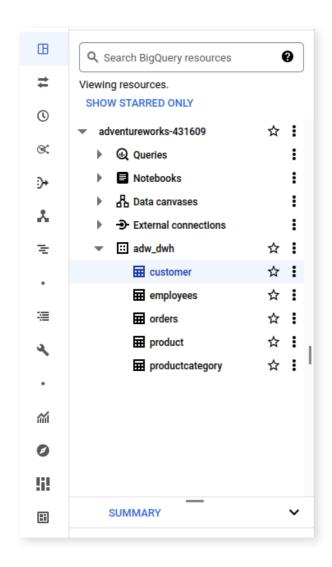
Storage/Database
Setting Google BigQuery

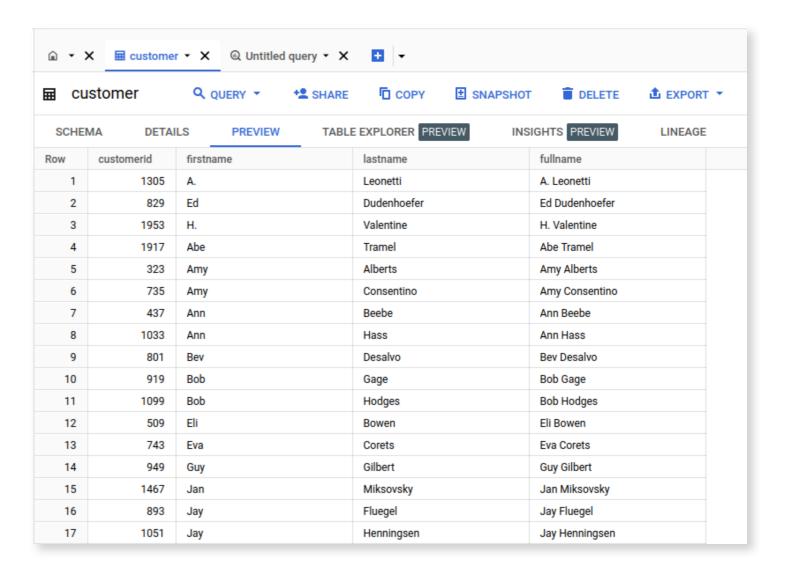
Ingestion
Setting up Apache Airflow
Writing elt Python script
Orchestrate data pipeline

Transformation
Setting up dbt
Transformation

Reporting Connecting Looker

# O1 Storage Set up BigQuery





# 02

# Ingestion Setting up Apache Airflow

- Airflow Documentation
- Production Deployment Documentation

### Writing ELT Python Script

- .py Code - Extract & Load

### # importing libraries

from airflow.decorators import dag, task
from datetime import datetime, timedelta
import requests
from google.cloud import bigquery
import pandas as pd
import psycopg2
from io import StringIO

# Ingestion Setting up Apache Airflow

- Airflow Documentation
- Production Deployment Documentation

### Writing ELT Python Script

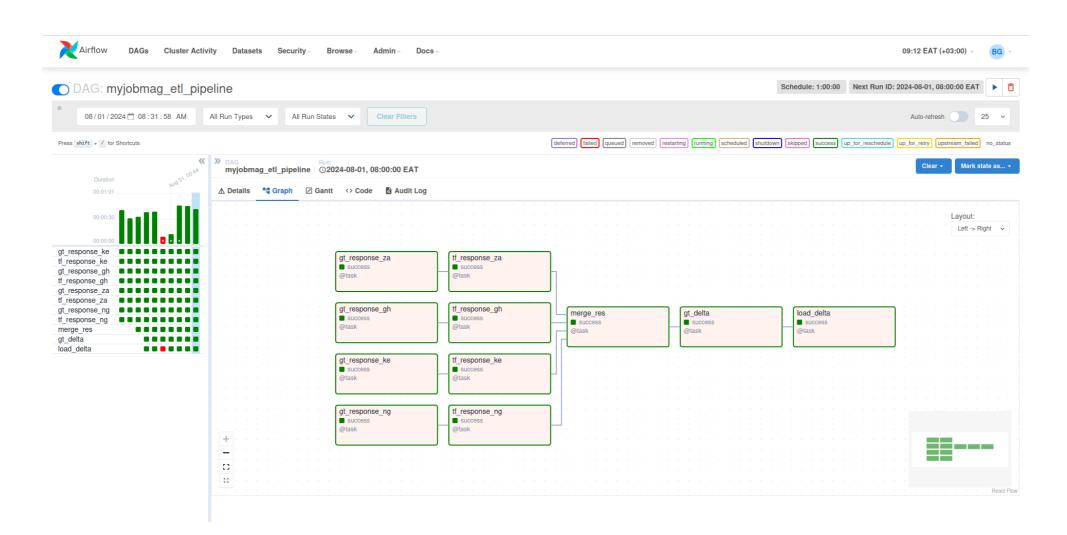
- .py Code - Extract & Load

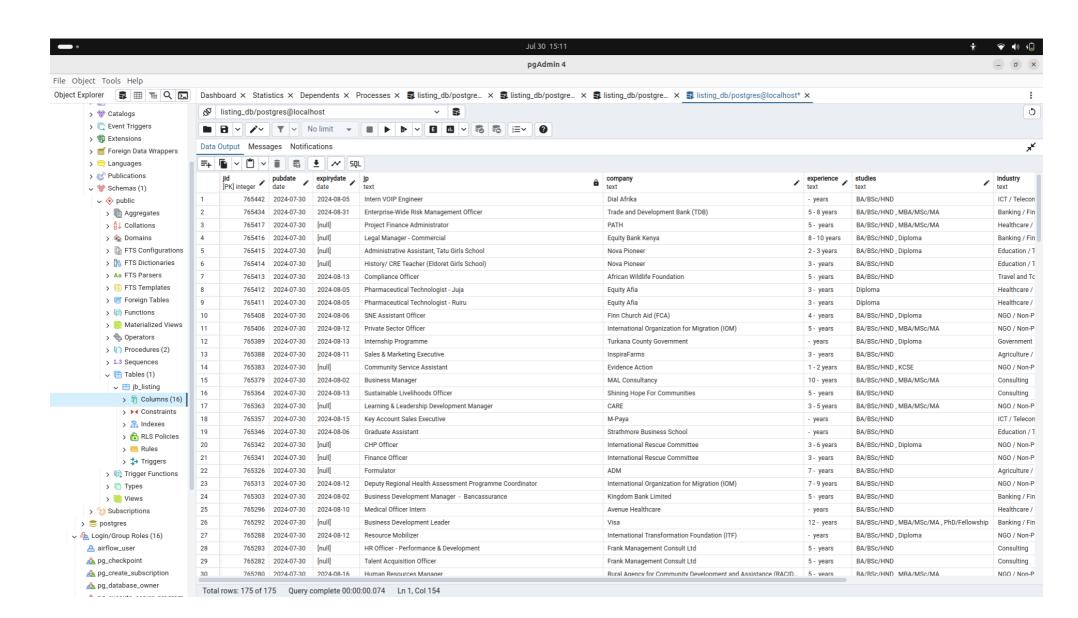
```
# instantiating DAG
@dag(
    schedule=timedelta(minutes=30),
    start_date=datetime(2024, 7, 29),
    catchup=False,
    tags=['Team B']
```

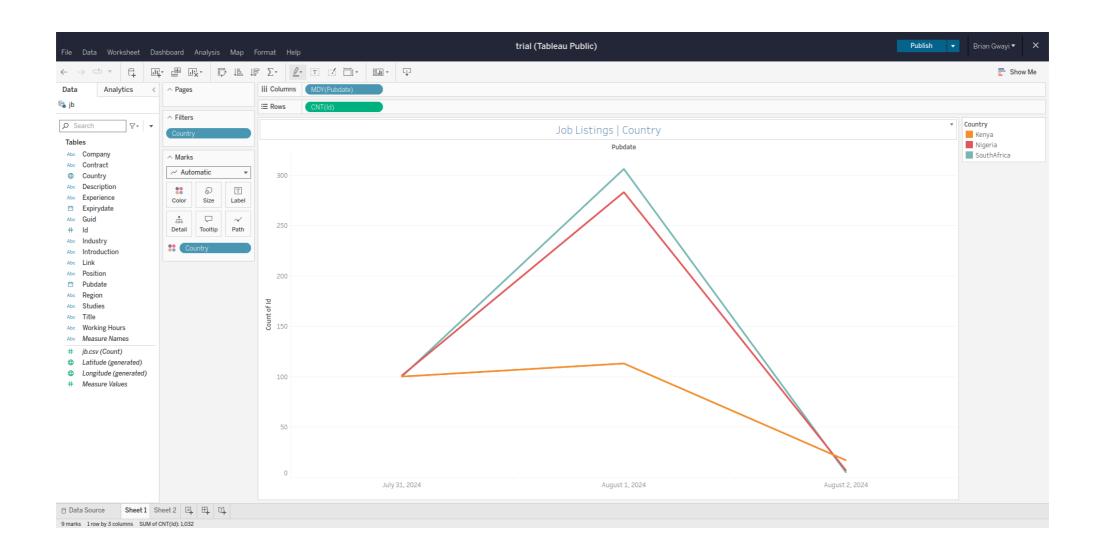
```
@task()
                                          def gt_tbls():
Task 1
                                                  conn = psycopg2.connect(
Get Table Lists
                                                     database = "adw_db",
                                                      user = "postgres",
from PostgreSQL
                                                      host= 'localhost',
                                                      password = "password",
                                                      Port = 5432)
                                                  cursor = conn.cursor()
sql = """SELECT table name
                                                  cursor.execute(sql)
     FROM information_schema.tables
                                                  tbls=cursor.fetchall()
     WHERE table_type = 'BASE TABLE'
     AND table_catalog = 'adventure_works'
                                                  conn.commit()
     AND table_schema NOT IN
                                                  conn.close()
     ('pg_catalog', 'information_schema');"""
                                                 tbls = [x[0]] for x in tbls]
                                                  Return thls
```

# Task 2 Extract tables

# Task 3 Upsert rowsBigQuery







# **Modern Data Stack Ecosystem 2024**

The right tools for building robust data stack architecture will be bases on Combination of budget, skillset, data sources and preferences.

