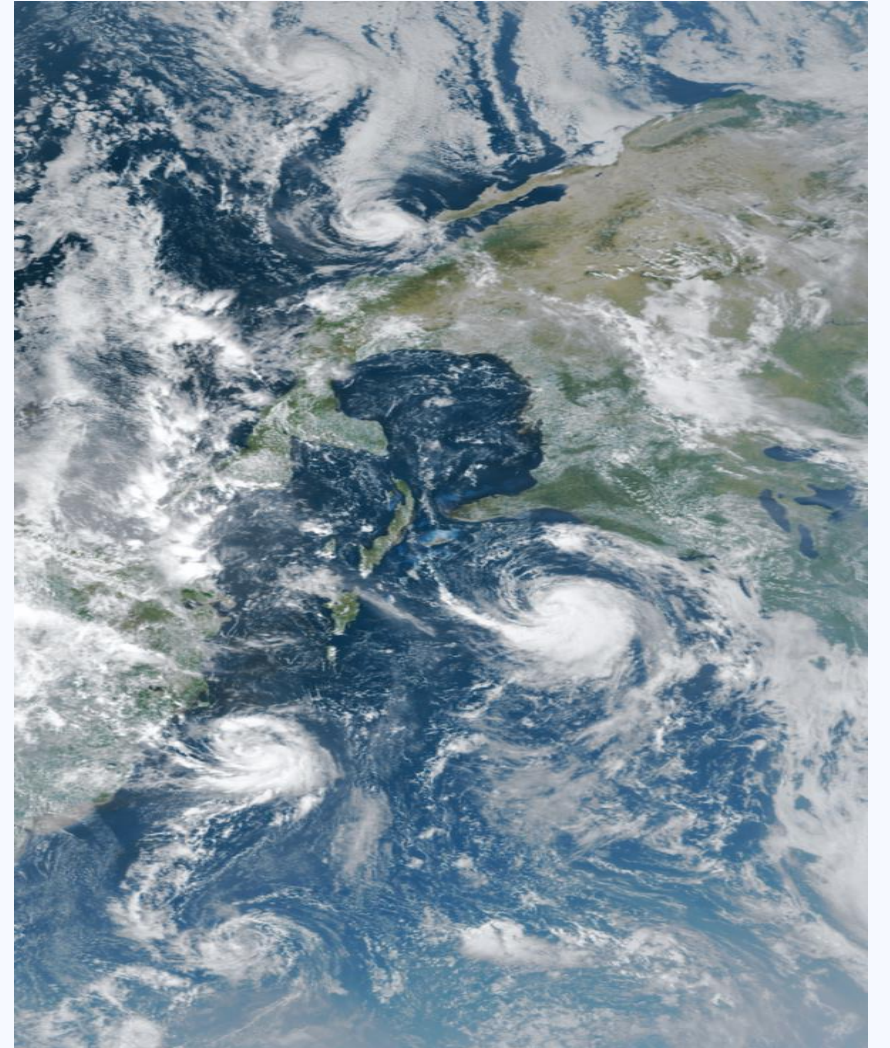# GLOBAL WEATHER REPOSITORY

GROUP 8:

- AISWARYA

- NEHA THAKUR

- SAUMYA VARSHNEY

- SHIKHA SINGH

- SREELEKHA

# INTRODUCTION

• Global Historical Climatology Network (GHCN) data from the National Oceanic and Atmospheric Administration (NOAA) website has been utilized for implementing the data pipeline.

• In this project we are working on Cloud Analytics and Data Warehouse Implementation by building a data pipeline to analyze archival and real time weather data and gather useful insights from this data.

• Temperature, Precipitation, Snow and Snow Depth patterns over the years will be analyzed and visualized for drawing useful insights which can be further utilized by industries.

Dataset link:

https://www.ncdc.noaa.gov/cdo-web/webservices/v2#gettingStarted

https://noaa-ghcn-pds.s3.amazonaws.com/index.html
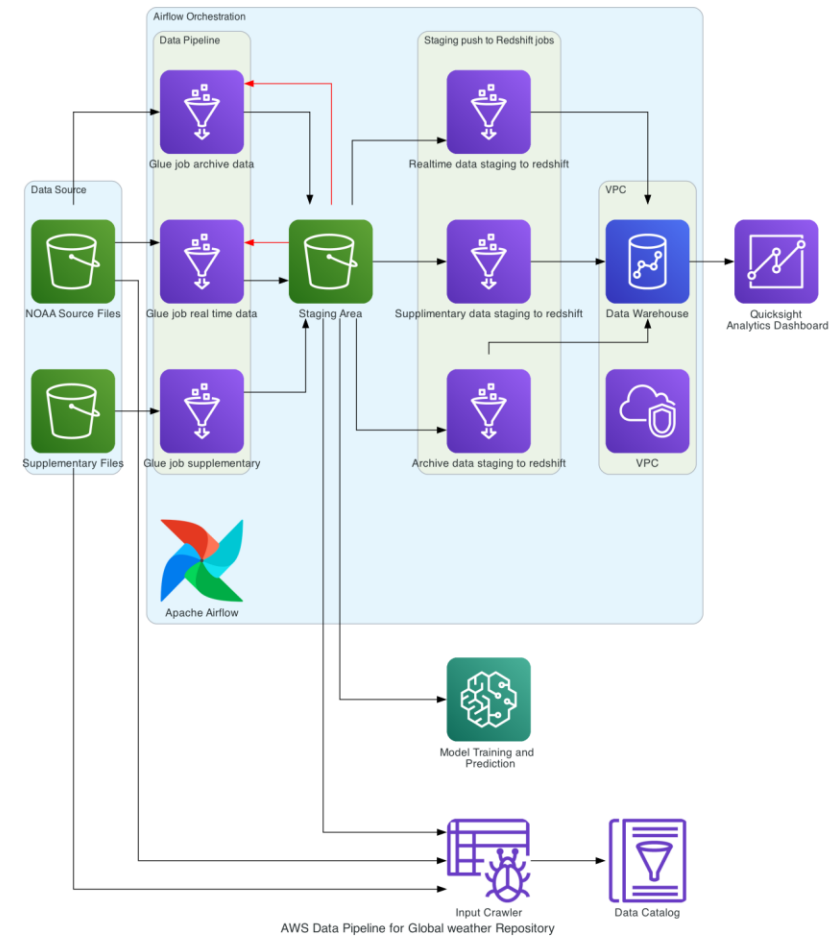
# Data Sources

NOAA GHCN-D AWS Bucket - This public bucket contains all the data collected by NOAA for the GHCN-D dataset from year 1750 to 2023. We are only considering the data from 2010 to 2023 for this project.

- Each year is stored as a CSV, and the headers comprise of
  - station id - id for the weather station where the data was collected from.
  - date id - date at which the data was collected in the format yyyymmdd (20231130).
  - element id / datatype id - what category of data was recorded (maximum temperature, minimum temperature, precipitation, snow fall and snow depth).
  - Data / value - the data recorded corresponding to element.

- Additionally, we scraped a lot of supplemental data from https://www.ncdc.noaa.gov/cdo-web/webservices/v2#gettingStarted using REST APIs.
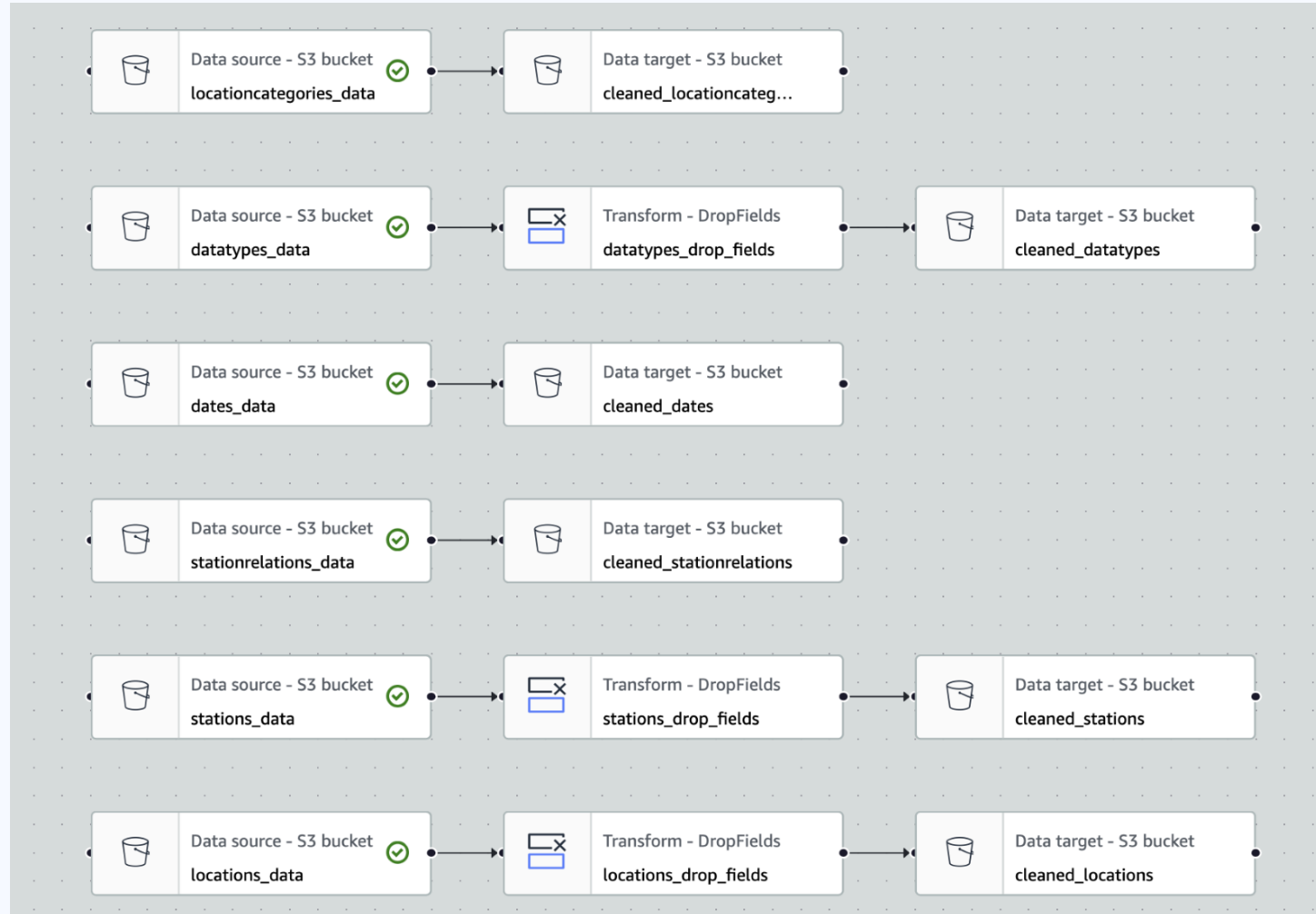
# CLOUD ARCHITECTURE

AWS Services used:

- VPC

- S3

- Glue

- Glue Crawler

- Glue Data Catalog

- Redshift

- Apache Airflow

- Quick Sight

- Sage Maker



AWS Data Pipeline for Global weather Repository

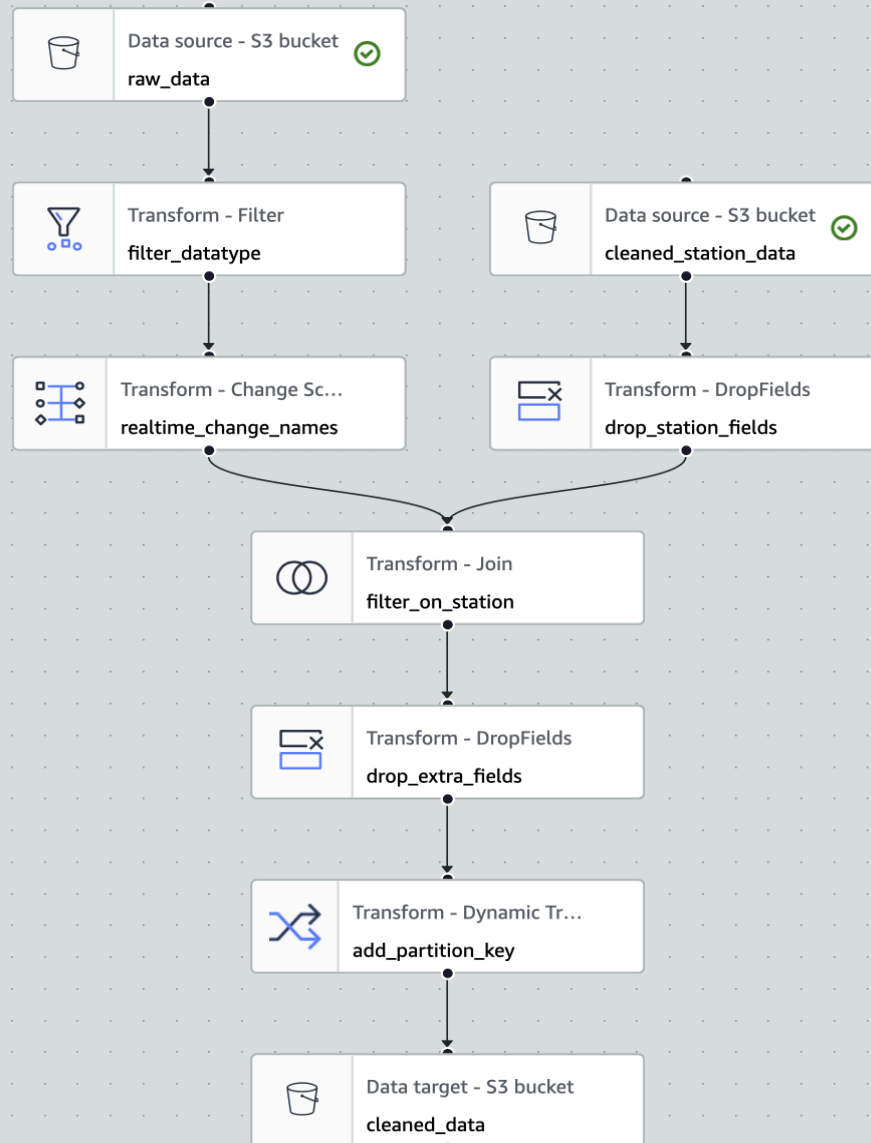# ETL (EXTRACT-TRANSFORM-LOAD)



Supplementary Data Pipeline

Reads the raw CSV files from source S3 bucket, transforms and load the cleaned data to S3 staging area in parquet format.

- locationcategories.csv
- datatypes.csv
- dates.csv
- stationrelations.csv
- stations.csv
- locations.csv
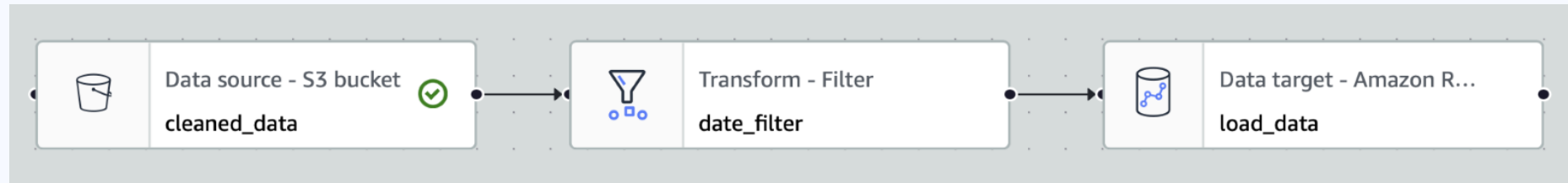
# ETL (EXTRACT-TRANSFORM-LOAD)

Archive (2010- 2022) and Real time (2023) Data Pipeline

Reads the raw CSV files from source S3 bucket, transforms and load the cleaned data to S3 staging area in parquet format.
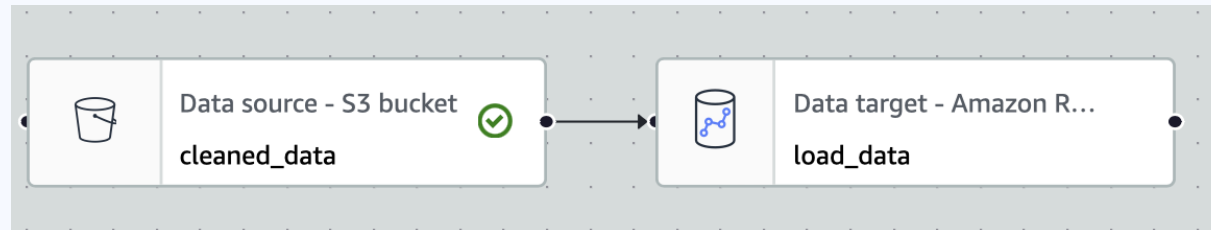
# LOAD TO DATA WAREHOUSE

Load Archive data to Datawarehouse



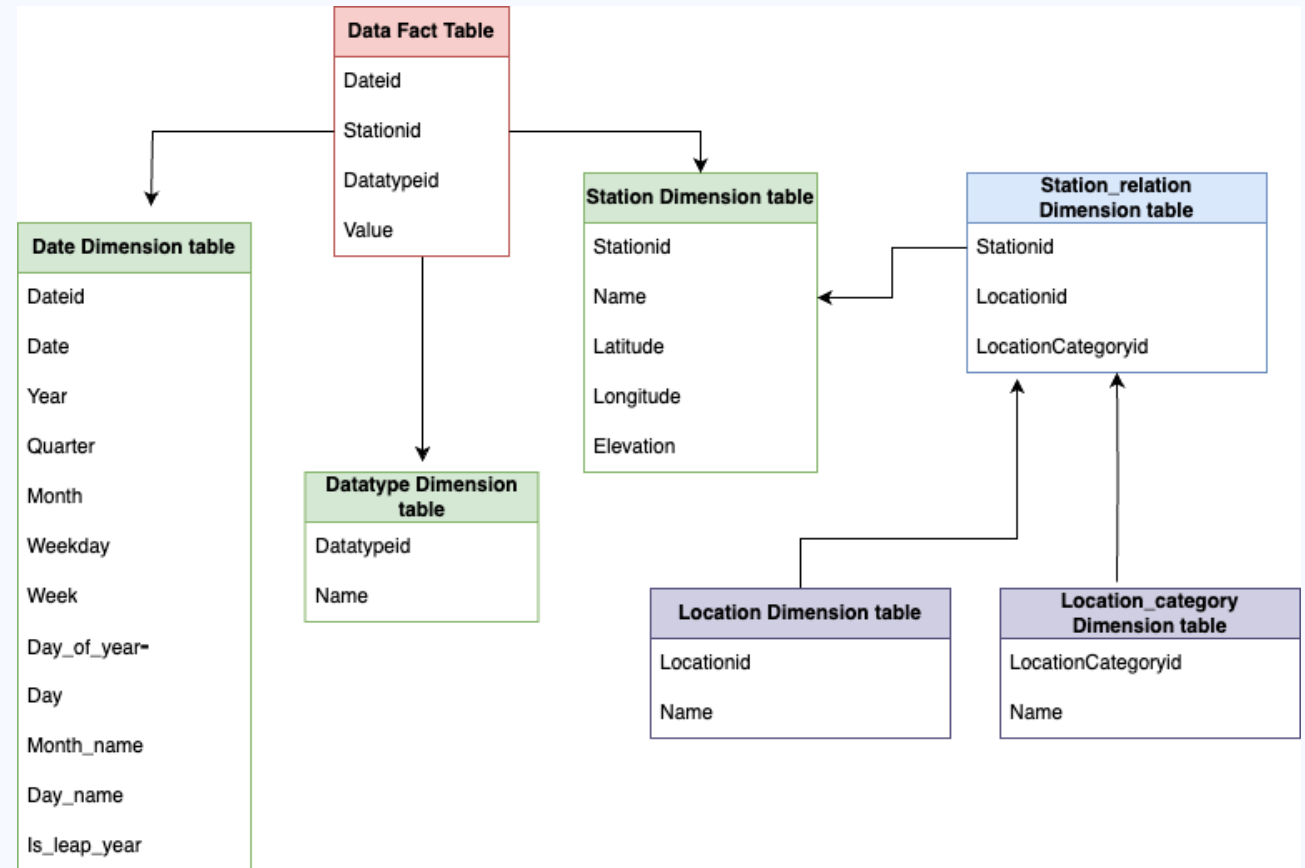Load Real time data to Data Warehouse



Similarly, we have loaded supplementary data to data Warehouse
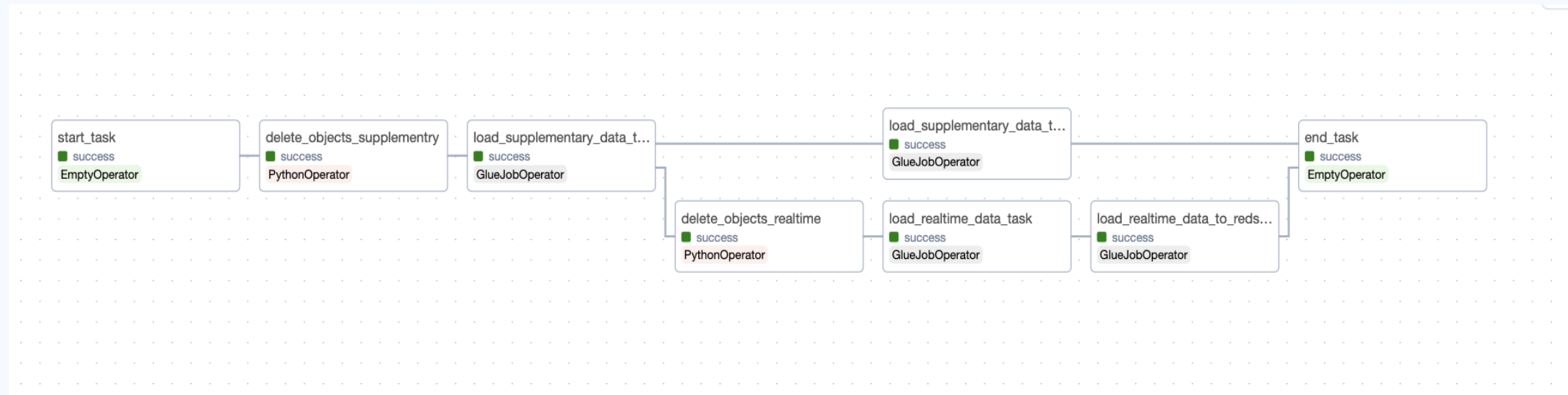
# Data Warehouse Schema

- Snowflake schema
- Fact table – Data
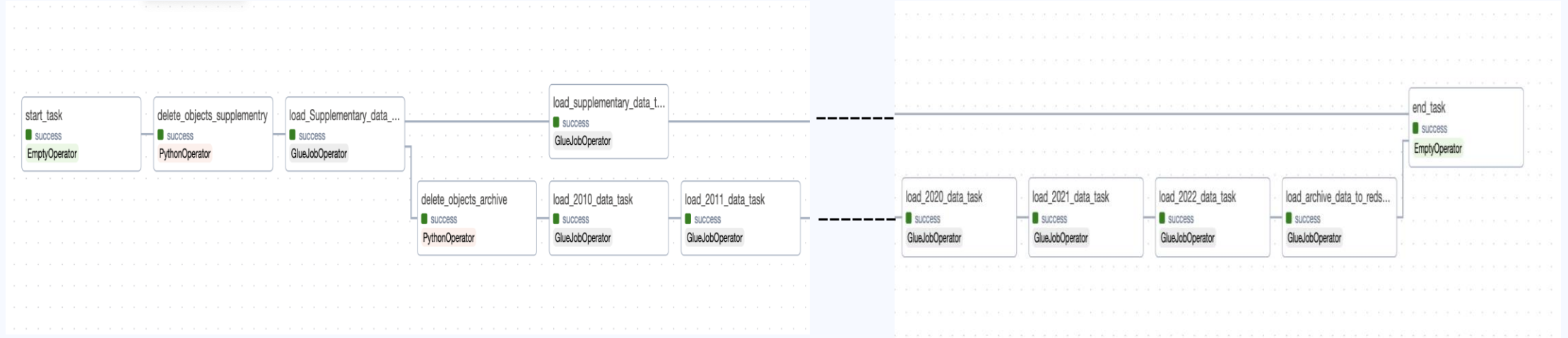- Dimension table – Date, Datatype and Station

# APACHE AIRFLOW

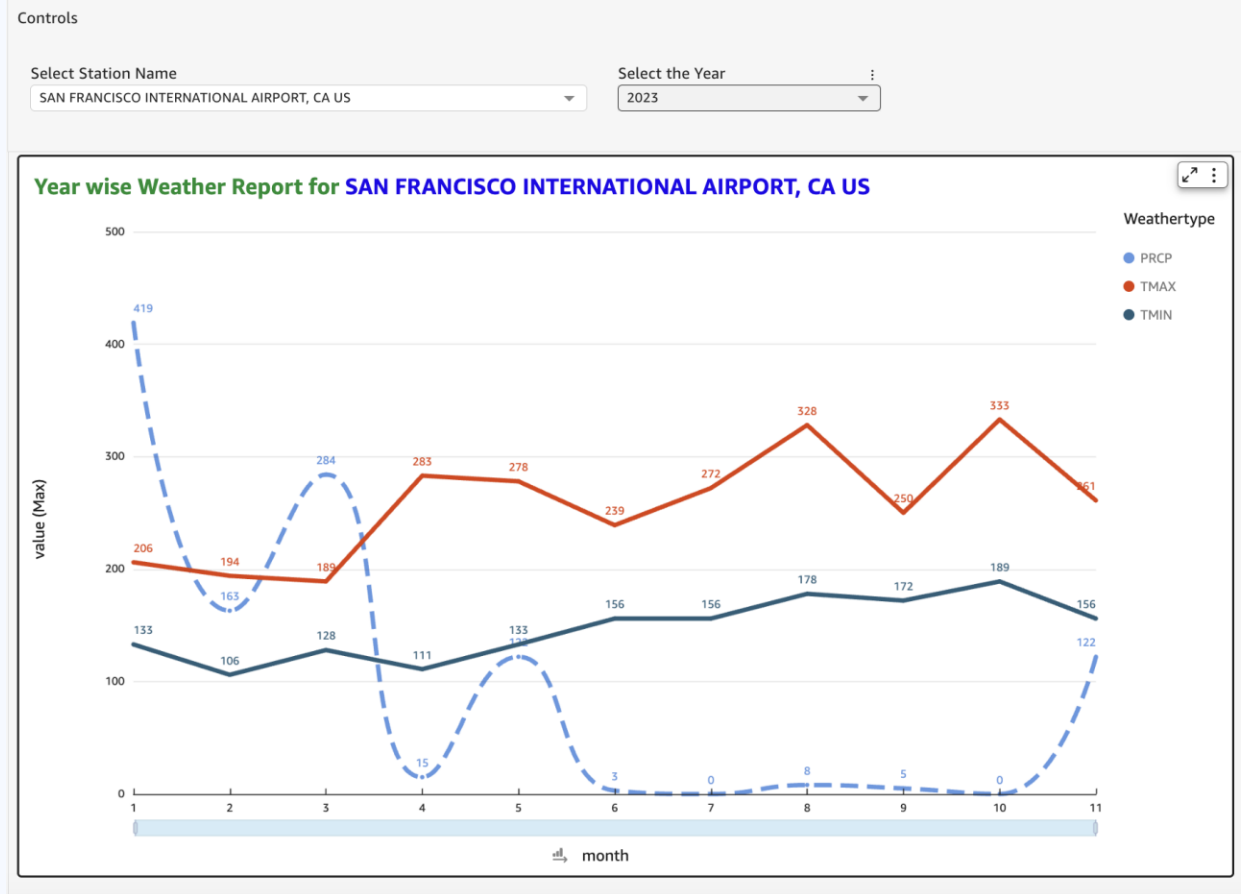Real-time data pipeline

# APACHE AIRFLOW

Archive data pipeline

# VISUALIZATION



Analysis and Visualizations on Weather Parameters

Year wise Weather Analysis

# VISUAL ANALYSIS

# PREDICTIVE ANALYTICS

Prediction Model

We have used the XGBoost Algorithm to train regression model to predict maximum and minimum temperature, precipitation, snow fall and snow depth. We used this model as this is known to work best for tabular data.

Final Schema for Training Data Features

| year | quarter | month | week | day_of_year | leap_year | latitude | longitude | elevation |
|------|---------|-------|------|-------------|-----------|----------|-----------|-----------|

Final Schema for Training Data Targets

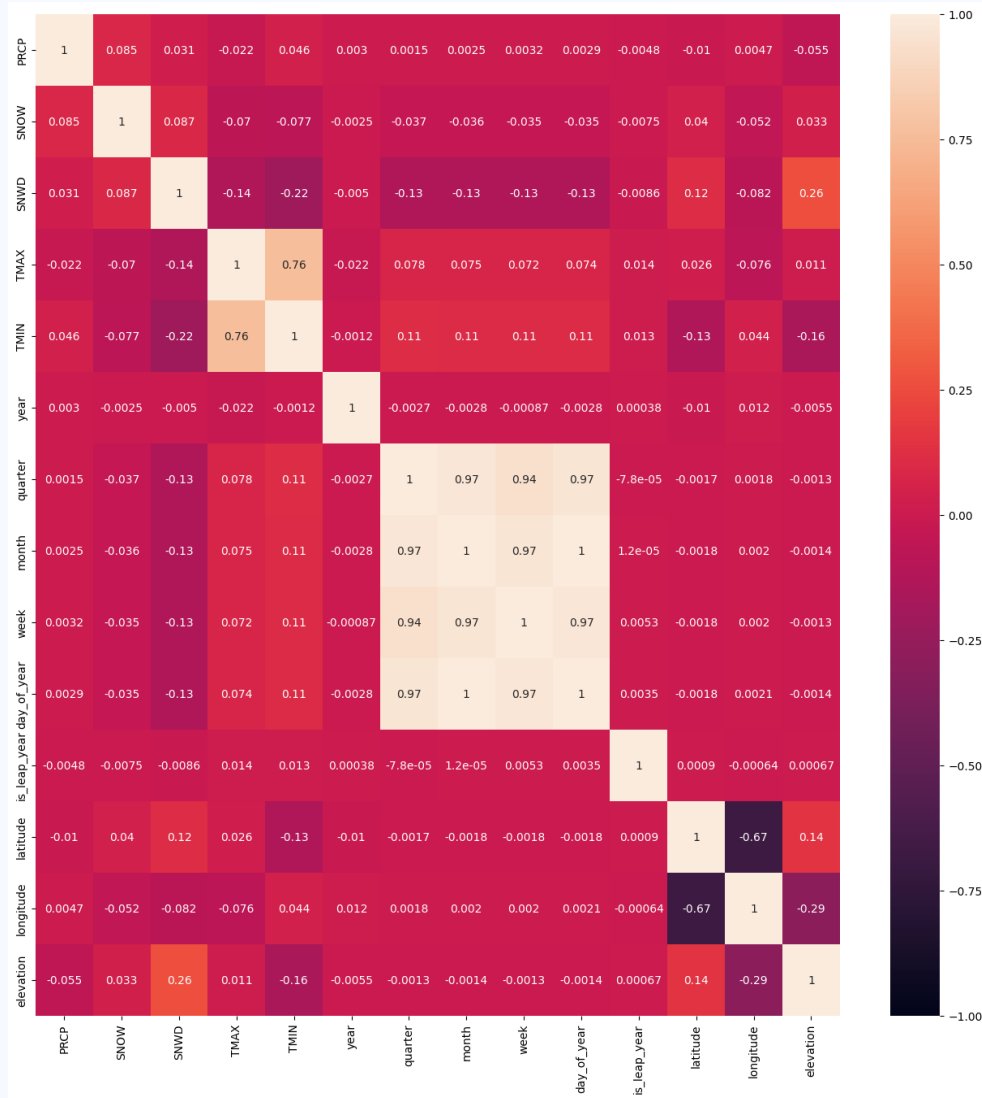| TMAX | TMIN | PRCP | SNOW | SNWD |
|------|------|------|------|------|

**R2 score for your model is 0.398**

# PREDICTIVE ANALYTICS

```
locationcategory_name = "CITY"
location_name = "San Francisco"
start_date = "2023-12-30"
duration = "7"
predict(locationcategory_name, location_name, start_date, duration, station_relations)
```

| | Location | Date | TMAX (C) | TMIN (C) | PRCP (cm) | SNOW (cm) | SNWD (cm) |
|---|---|---|---|---|---|---|---|
| 0 | San Francisco, CA US | 2023-12-30 | 7.196602 | 4.007489 | 16.250647 | 0.009818 | -0.681265 |
| 1 | San Francisco, CA US | 2023-12-31 | 7.196602 | 4.007489 | 16.250647 | 0.009818 | -0.681265 |
| 2 | San Francisco, CA US | 2024-01-01 | 8.431298 | 2.279749 | 3.353752 | -0.052752 | 1.545340 |
| 3 | San Francisco, CA US | 2024-01-02 | 8.431298 | 2.279749 | 3.353752 | -0.052752 | 1.545340 |
| 4 | San Francisco, CA US | 2024-01-03 | 8.431298 | 2.279749 | 4.143710 | -0.009639 | 1.545340 |
| 5 | San Francisco, CA US | 2024-01-04 | 9.103164 | 3.137930 | 3.751980 | -0.003760 | 1.545340 |
| 6 | San Francisco, CA US | 2024-01-05 | 9.103164 | 3.137930 | 3.751980 | -0.003760 | 1.545340 |
| 7 | San Francisco, CA US | 2024-01-06 | 9.103164 | 3.137930 | 3.751980 | 0.386499 | 1.460063 |

# Correlation Analysis

- We can see that there is a white patch in the center.

- All those fields are highly positively correlated.

- Latitude and Longitude are negatively correlated.

# CONCLUSION

AWS powered system efficiently manages data from various sources, making it easy for businesses . We used tools like AWS Glue for pipelining, Redshift as data warehouse , Airflow to Orchestrate the pipeline and many more. Our data system on AWS helps companies organize information easily, make smarter decisions using advanced analysis, future predictions and stay flexible as their needs evolve. It's like having a reliable assistant that keeps everything in order and helps the business grow smarter.

# THANK YOU