

Lab Group 17

Weather Analytics Pipeline (ETL & ELT)

Omkar Rajale

Yashashree Shinde

Abstract

This document serves as a comprehensive project report for the development of the **Weather Analytics Pipeline System**. This application implements an end-to-end data analytics pipeline that fetches historical weather data for London (2023), stores it in Snowflake, transforms it using dbt (orchestrated by Airflow), and visualizes weather trends using Power BI. The system demonstrates an automated workflow combining ETL (Extract, Transform, Load) and ELT (Extract, Load, Transform) patterns, including Type 2 Slowly Changing Dimensions (SCD) snapshots and comprehensive data quality testing. The goal is to deliver a reliable, automated, and scalable system capable of ingesting daily weather data, performing advanced transformations, and providing actionable insights through interactive visualizations.

1 Introduction

1.1 Project Overview

This section introduces the **Weather Analytics Pipeline System**, a comprehensive data engineering solution designed to automate the collection, transformation, and visualization of historical weather data. The system leverages modern cloud-based technologies including Apache Airflow for orchestration, Snowflake as the data warehouse, dbt for transformation logic, and Power BI for business intelligence visualization.

2 Problem Statement

This section clearly defines the application system and articulates the necessity of the project.

2.1 The Problem

Weather data analysis presents significant challenges for organizations requiring reliable, historical weather insights for decision-making. The core problems include: (1) lack of automated, repeatable processes for daily weather data ingestion, (2) absence of standardized data quality validation mechanisms, (3) difficulty in tracking historical changes in weather patterns over time, and (4) inability to efficiently transform raw weather data into actionable business insights. Manual processes are error-prone, time-consuming, and fail to provide the consistency required for accurate trend analysis and forecasting.

2.2 Proposed Application System

The **Weather Analytics Pipeline System** is a fully automated data engineering solution designed to address these challenges through a modern ETL/ELT architecture. The system orchestrates the complete data lifecycle: extracting weather data from the Open-Meteo Archive API, loading it into Snowflake's RAW schema, transforming it through dbt models with comprehensive testing, capturing historical changes via Type 2 SCD snapshots, and presenting insights through interactive Power BI dashboards.

2.3 Rationale for Database and Data Pipelines

A cloud data warehouse (**Snowflake**) is essential for:

- **Scalability:** Handling large volumes of time-series weather data with efficient storage and query performance.
- **Data Integrity:** Providing ACID transaction support for idempotent operations and maintaining data consistency.

- **Separation of Concerns:** Distinct schemas (RAW, ANALYTICS) for raw data storage and transformed analytical datasets.

The orchestrated data pipeline (**Airflow + dbt**) is mandatory for:

- **Automation:** Scheduling daily ETL operations and triggering downstream ELT transformations without manual intervention.
- **Dependency Management:** Ensuring dbt transformations execute only after successful ETL completion.
- **Observability:** Providing comprehensive logging, monitoring, and alerting capabilities for pipeline health.
- **Data Quality:** Implementing automated testing frameworks to validate data integrity at each pipeline stage.

3 Solution Requirements

This section details the necessary criteria, actions, and limitations of the final system.

3.1 Functional Requirements (FR)

The core capabilities and actions the application must perform are:

- FR1: The system shall successfully connect to the **Open-Meteo Archive API** and retrieve 365 days of historical weather data for London including temperature (max/min), precipitation, and wind speed metrics.
- FR2: The system shall use an **Airflow ETL DAG** to automate the data extraction and loading process into Snowflake's RAW schema with full-refresh idempotency using SQL transactions.
- FR3: The system shall trigger a separate **Airflow ELT DAG** that executes dbt commands (run, test, snapshot) in proper sequence after successful ETL completion.
- FR4: The system shall implement **dbt models** to transform raw weather data including: staging views, analytical aggregations (7-day moving averages), and categorical derivations (rain classification).
- FR5: The system shall execute **dbt tests** to validate data quality constraints including uniqueness of primary keys and non-null validation on critical fields.
- FR6: The system shall maintain **Type 2 SCD snapshots** using dbt's snapshot functionality to track historical changes in weather metrics over time.
- FR7: The system shall populate the final **ANALYTICS.WEATHER_ANALYTICS** table with transformed, tested, and enriched weather data.
- FR8: The system shall provide a **Power BI dashboard** connected to the analytics table, displaying temperature trends, precipitation analysis, and weather distribution visualizations.

3.2 Non-Functional Requirements (NFR)

3.2.1 Performance

The complete pipeline execution (ETL + ELT including dbt run/test/snapshot) must complete within **15 minutes** for 365 days of data.

3.2.2 Reliability

The ETL pipeline must implement idempotency using SQL transactions (BEGIN/COMMIT/ROLLBACK) with proper exception handling to ensure safe re-execution without data duplication.

3.2.3 Data Quality

All data loaded into ANALYTICS schema must pass dbt data quality tests with zero test failures. Critical fields (date, temperature) must have 100% non-null coverage.

3.2.4 Security

Snowflake credentials must be stored in Airflow Connections and passed to dbt using environment variables, preventing hardcoded passwords in configuration files.

3.2.5 Observability

Both Airflow DAGs must provide comprehensive logging at each task level, and dbt must generate detailed run artifacts showing model execution status and test results.

3.3 System Users and Usage

The primary users of this system are **Data Analysts, Business Intelligence Teams, and Weather Data Consumers** who rely on the transformed analytics data for insights and decision-making.

- **Use Case 1 (Daily ETL):** The Airflow ETL DAG automatically triggers to fetch the latest year of London weather data and performs a full-refresh load into RAW.WEATHER_HISTORY.
- **Use Case 2 (Automated Transformation):** Upon successful ETL completion, the ELT DAG triggers dbt to transform raw data, validate quality, and update the snapshot table.
- **Use Case 3 (Business Analysis):** A business analyst opens the Power BI dashboard to analyze temperature volatility trends, identify rainy periods, and examine seasonal patterns.
- **Use Case 4 (Data Quality Monitoring):** A data engineer reviews dbt test results in Airflow logs to verify all uniqueness and not-null constraints passed successfully.
- **Use Case 5 (Historical Audit):** An auditor queries the ANALYTICS.WEATHER_SNAPSHOT table to track how specific weather metrics changed over time using SCD Type 2 temporal columns.

3.4 Conceptual Architecture

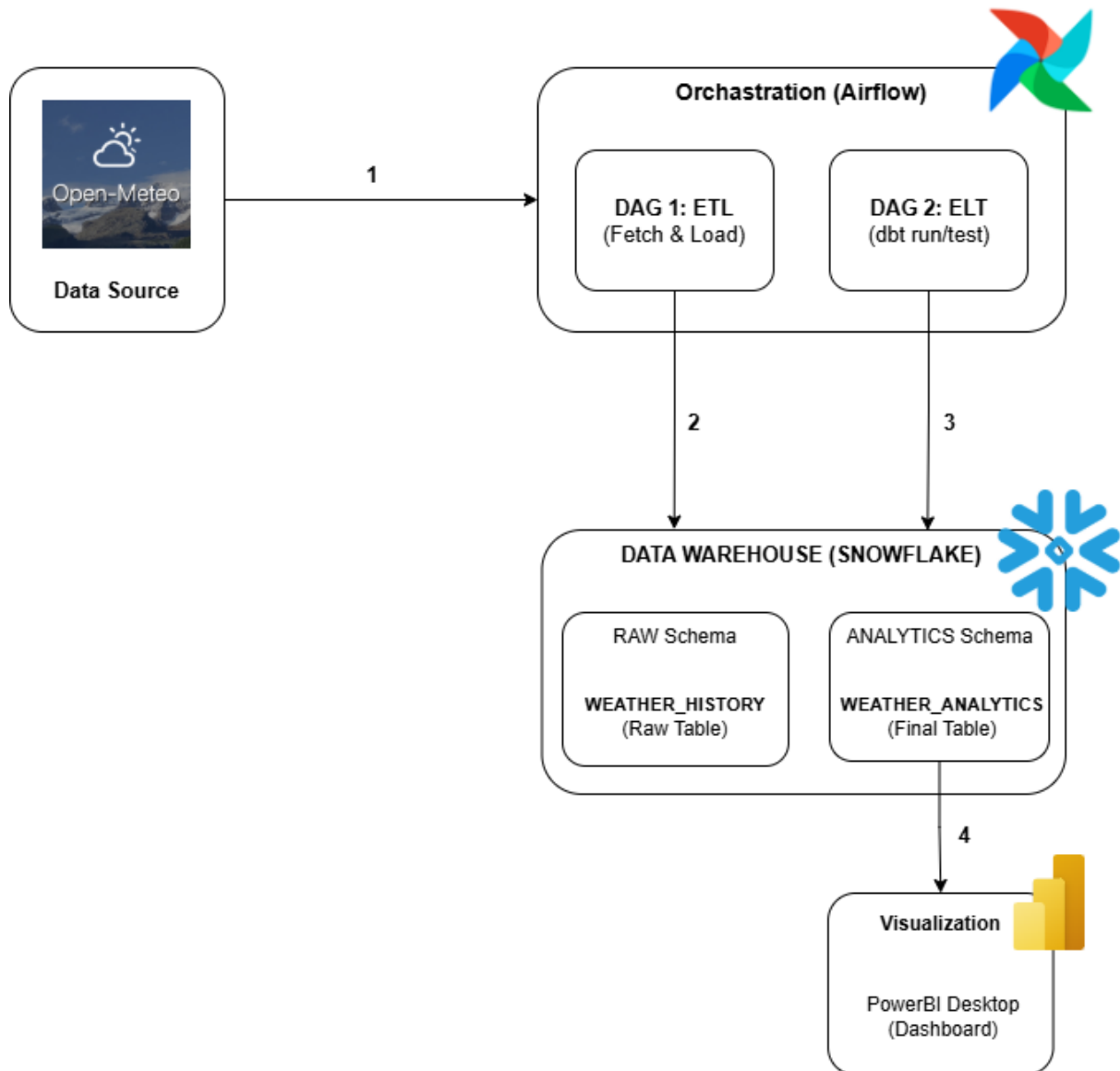


Figure 1: Conceptual Data Pipeline Architecture illustrating the data flow from the Open-Meteo API through Airflow orchestration, Snowflake storage, dbt transformation, and Power BI visualization. The diagram shows the separation between ETL (extraction and loading) and ELT (transformation) stages.

3.5 Functional Components Breakdown

3.5.1 Data Extraction Module (Python/Open-Meteo API)

- **Role:** Connects to the Open-Meteo Archive API and extracts 365 days of historical weather data for London (latitude: 51.5074, longitude: -0.1278).
- **Technology:** Python requests library within Airflow TaskFlow API decorated tasks.
- **Data Elements:** Daily temperature_2m_max, temperature_2m_min, precipitation_sum, wind_speed_10m_max.

3.5.2 ETL Orchestration Layer (Airflow DAG 1)

- **Role:** Orchestrates the Extract-Load pipeline with idempotent full-refresh strategy.
- **DAG Name:** open_meteo_weather_etl

- **Schedule:** Configured to run automatically every day at **7:00 AM**.
- **Tasks:** Extract (API call) → Transform (data cleaning) → Load (Snowflake TRUNCATE + INSERT in transaction).
- **Trigger Dependency:** Upon successful completion, this DAG automatically triggers the second DAG (`dbt_transformation_dag`) for downstream dbt model transformations.
- **Idempotency Strategy:** Uses SQL transactions with TRUNCATE followed by INSERT to ensure safe re-execution.

3.5.3 ELT Orchestration Layer (Airflow DAG 2)

- **Role:** Triggers dbt transformations after successful ETL completion.
- **DAG Name:** `dbt_transformation_dag`
- **Tasks:** dbt run (build models) → dbt test (validate quality) → dbt snapshot (capture changes)
- **Environment Management:** Passes Snowflake credentials from Airflow Connections to dbt via environment variables.

3.5.4 Data Warehouse Layer (Snowflake)

- **Role:** Persistent storage for raw and transformed weather data with schema separation.
- **Schemas:**
 - **RAW:** Stores unprocessed API data in WEATHER_HISTORY table
 - **ANALYTICS:** Contains transformed tables (WEATHER_ANALYTICS, WEATHER_SNAPSHOT)

3.5.5 Transformation Layer (dbt)

- **Role:** Implements business logic for data transformation, testing, and historical tracking.
- **Components:**
 - **Sources:** Defines RAW.WEATHER_HISTORY as upstream data source
 - **Staging Models:** `stg_weather.sql` - creates cleaned, standardized view
 - **Marts Models:** `weather_analytics.sql` - calculates 7-day moving averages, rain categorization
 - **Tests:** `schema.yml` - validates unique and not_null constraints
 - **Snapshots:** `weather_snapshot.sql` - implements Type 2 SCD with check strategy

3.5.6 Visualization Layer (Power BI)

- **Role:** Provides interactive business intelligence dashboards for weather insights.
- **Data Source:** Connects directly to Snowflake ANALYTICS.WEATHER_ANALYTICS table.
- **Visualizations:** Line charts (temperature trends), bar charts (precipitation), donut charts (rainy vs. dry days), and KPI cards.

4 Implementation Details and Appendices

4.1 Tables Structure

Table 1: Summary of Database Table Structures

Table Name	Key Fields	Other Columns	Description
RAW.WEATHER.HISTORY	DATE (PK)	MAX_TEMP, MIN_TEMP, PRECIPITATION, WIND_SPEED	Stores 365 days of raw weather data from Open-Meteo API
ANALYTICS.WEATHER.ANALYTICS	DATE (PK)	MAX_TEMP, MIN_TEMP, PRECIPITATION, WIND_SPEED, AVG_TEMP_7D, RAIN_CATEGORY	Transformed analytics table with moving averages and classifications
ANALYTICS.WEATHER.SNAPSHOT	DATE, DBT_SCD_ID (PK)	MAX_TEMP, MIN_TEMP, PRECIPITATION, WIND_SPEED, DBT_VALID_FROM, DBT_VALID_TO, DBT_UPDATED_AT	SCD Type 2 table tracking historical changes in weather metrics
ANALYTICS_STAGING.STG_WEATHER (VIEW)	FORECAST_DATE	MAX_TEMP, MIN_TEMP, PRECIPITATION, WIND_SPEED	Staging view created from RAW.WEATHER.HISTORY; provides cleaned and standardized weather data used by dbt transformations.

4.2 Key Code Implementations

4.2.1 ELT DAG - Open Meteo Weather ETL

Listing 1: open_meteo_weather_etl.py — Weather Data Extraction and Load DAG

```
1 from airflow import DAG
2 from airflow.decorators import task
3 from airflow.providers.snowflake.hooks.snowflake import SnowflakeHook
4 from airflow.operators.trigger_dagrun import TriggerDagRunOperator
5 from datetime import datetime, timedelta
6 import requests
7 import pandas as pd
8
9 # --- Configuration ---
10 SNOWFLAKE_CONN_ID = "snowflake_conn"
11 # Loading into the RAW schema
12 TARGET_TABLE = "RAW.WEATHER_HISTORY"
13 URL = "https://archive-api.open-meteo.com/v1/archive"
14
15 # London Coordinates for 2023 history
16 PARAMS = {
17     "latitude": 51.5074,
18     "longitude": -0.1278,
19     "start_date": "2023-01-01",
20     "end_date": "2023-12-31",
```

```

21     "daily": "temperature_2m_max,temperature_2m_min,precipitation_sum,wind_speed_10m_max",
22     "timezone": "auto"
23 }
24
25 default_args = {
26     "owner": "student",
27     "retries": 1,
28     "retry_delay": timedelta(minutes=5),
29 }
30
31 with DAG(
32     dag_id="open_meteo_weather_etl",
33     start_date=datetime(2024, 1, 1),
34     schedule="0 7 * * *",
35     catchup=False,
36     default_args=default_args,
37     description="Fetch 1 year of weather data from OpenMeteo and Full Refresh Snowflake"
38 ) as dag:
39
40     @task
41     def extract_weather_data():
42         """
43         Fetches data from Open-Meteo API and returns a list of tuples.
44         """
45         response = requests.get(URL, params=PARAMS)
46         response.raise_for_status()
47         data = response.json()
48
49         daily_data = data['daily']
50         df = pd.DataFrame(daily_data)
51
52         df.rename(columns={
53             'time': 'forecast_date',
54             'temperature_2m_max': 'max_temp',
55             'temperature_2m_min': 'min_temp',
56             'precipitation_sum': 'precipitation',
57             'wind_speed_10m_max': 'wind_speed'
58         }, inplace=True)
59
60         rows = []
61         for record in df.itertuples(index=False):
62             rows.append((
63                 record.forecast_date,
64                 float(record.max_temp) if pd.notna(record.max_temp) else None,
65                 float(record.min_temp) if pd.notna(record.min_temp) else None,
66                 float(record.precipitation) if pd.notna(record.precipitation) else None,
67                 float(record.wind_speed) if pd.notna(record.wind_speed) else None
68             ))
69         return rows
70
71     @task
72     def load_to_snowflake(rows):
73         """
74         Performs a Full Refresh (Truncate + Insert) into Snowflake.
75         """
76         if not rows:
77             print("No data found. Skipping load.")
78             return
79
80         hook = SnowflakeHook(snowflake_conn_id=SNOWFLAKE_CONN_ID)
81         conn = hook.get_conn()
82         cur = conn.cursor()
83
84         try:
85             create_sql = f"""

```

```

86         CREATE TABLE IF NOT EXISTS {TARGET_TABLE} (
87             FORECAST_DATE DATE PRIMARY KEY,
88             MAX_TEMP FLOAT,
89             MIN_TEMP FLOAT,
90             PRECIPITATION FLOAT,
91             WIND_SPEED FLOAT
92         )
93         """
94         cur.execute(create_sql)
95
96         conn.autocommit = False
97         cur.execute("BEGIN")
98         cur.execute(f"TRUNCATE TABLE {TARGET_TABLE}")
99
100        insert_sql = f"""
101        INSERT INTO {TARGET_TABLE} (FORECAST_DATE, MAX_TEMP, MIN_TEMP, PRECIPITATION
102            , WIND_SPEED)
103        VALUES (%s, %s, %s, %s, %s)
104        """
105        cur.executemany(insert_sql, rows)
106
107        cur.execute("COMMIT")
108        print(f"Successfully loaded {len(rows)} rows into {TARGET_TABLE}")
109
110        except Exception as e:
111            cur.execute("ROLLBACK")
112            print(f"Error occurred: {e}")
113            raise
114        finally:
115            cur.close()
116            conn.close()
117
118        trigger_dbt = TriggerDagRunOperator(
119            task_id="trigger_dbt_transformation",
120            trigger_dag_id="dbt_transformation_dag",
121            wait_for_completion=False
122        )
123
124        weather_data = extract_weather_data()
125        load_task = load_to_snowflake(weather_data)
126
127        load_task >> trigger_dbt

```

4.2.2 ELT DAG - dbt Orchestration

Listing 2: ELT DAG for dbt Transformations (Airflow 2.x)

```

1  from pendulum import datetime
2  from airflow import DAG
3  from airflow.operators.bash import BashOperator
4
5  # Paths inside the Airflow container
6  DBT_PROJECT_DIR = "/opt/airflow/dbt"
7  DBT_PROFILES_DIR = "/opt/airflow/dbt"
8
9  default_args = {
10     "owner": "student",
11     "retries": 1,
12 }
13
14  with DAG(
15     dag_id="dbt_transformation_dag",
16     start_date=datetime(2024, 1, 1),
17     schedule_interval=None,
18     catchup=False,
19     default_args=default_args,
20     description="Run dbt transformations for Weather Data",

```



```

21     tags=["ELT", "dbt", "transformation"],
22 ) as dag:
23
24     # Read Snowflake creds at runtime via connection templating (no BaseHook at parse
25     # time)
26     env_vars = {
27         "DBT_USER": "{{ conn.snowflake_conn.login }}",
28         "DBT_PASSWORD": "{{ conn.snowflake_conn.password }}",
29         "DBT_ACCOUNT": "{{ conn.snowflake_conn.extra_dejson.account }}",
30         "DBT_SCHEMA": "ANALYTICS",
31         "DBT_DATABASE": "{{ conn.snowflake_conn.extra_dejson.database }}",
32         "DBT_ROLE": "{{ conn.snowflake_conn.extra_dejson.role }}",
33         "DBT_WAREHOUSE": "{{ conn.snowflake_conn.extra_dejson.warehouse }}",
34         "DBT_TYPE": "snowflake",
35     }
36
37     dbt_run = BashOperator(
38         task_id="dbt_run",
39         bash_command=f"dbt run --profiles-dir {DBT_PROFILES_DIR} --project-dir {
40             DBT_PROJECT_DIR}",
41         env=env_vars,
42     )
43
44     dbt_test = BashOperator(
45         task_id="dbt_test",
46         bash_command=f"dbt test --profiles-dir {DBT_PROFILES_DIR} --project-dir {
47             DBT_PROJECT_DIR}",
48         env=env_vars,
49     )
50
51     dbt_snapshot = BashOperator(
52         task_id="dbt_snapshot",
53         bash_command=f"dbt snapshot --profiles-dir {DBT_PROFILES_DIR} --project-dir {
54             DBT_PROJECT_DIR}",
55         env=env_vars,
56     )
57
58     dbt_run >> dbt_test >> dbt_snapshot

```

4.2.3 dbt Analytics Model

Listing 3: weather_analytics.sql — Transformation Logic

```

1  {{ config(materialized='table') }}
2
3  WITH source_data AS (
4      SELECT * FROM {{ ref('stg_weather') }}
5  )
6
7  SELECT
8      forecast_date,
9      max_temp,
10     min_temp,
11     precipitation,
12     wind_speed,
13
14     -- Transformation 1: Calculate Temperature Range
15     (max_temp - min_temp) AS temp_variation,
16
17     -- Transformation 2: 7-Day Moving Average of Max Temp
18     AVG(max_temp) OVER (
19         ORDER BY forecast_date
20         ROWS BETWEEN 6 PRECEDING AND CURRENT ROW
21     ) AS seven_day_avg_max_temp,
22
23     -- Transformation 3: Categorize Days by Rainfall
24     CASE

```

```

25     WHEN precipitation > 0 THEN 'Rainy'
26     ELSE 'Dry'
27 END AS rain_status,
28
29 -- Transformation 4: Categorize Wind Speed
30 CASE
31     WHEN wind_speed > 20 THEN 'High Wind'
32     ELSE 'Normal Wind'
33 END AS wind_status
34
35 FROM source_data

```

4.2.4 dbt Snapshot Configuration

Listing 4: weather_snapshot.sql — SCD Type 2 Snapshot Configuration

```

1 {% snapshot weather_snapshot %}
2
3 {{
4     config(
5         target_database='USER_DB_POODLE',
6         target_schema='ANALYTICS',
7         unique_key='forecast_date',
8         strategy='check',
9         check_cols=['max_temp', 'precipitation', 'wind_speed']
10    )
11 }}
12
13 SELECT * FROM {{ source('raw_data', 'weather_history') }}
14
15 {% endsnapshot %}

```

4.2.5 dbt Configuration (profiles.yml & dbt_project.yml)

These configuration files bind dbt to Snowflake and define project structure, model paths, and default materializations. Secrets are read via environment variables passed from Airflow (see DAG 2).

```

Y profiles.yml X
Users > spartan > airflow-docker > dbt > Y profiles.yml
1 my_weather_project:
2   target: dev
3   outputs:
4     dev:
5       type: snowflake
6       # Use env_var to read from Airflow DAG
7       account: "{{ env_var('DBT_ACCOUNT') }}"
8       user: "{{ env_var('DBT_USER') }}"
9       password: "{{ env_var('DBT_PASSWORD') }}"
10      role: "{{ env_var('DBT_ROLE') }}"
11      warehouse: "{{ env_var('DBT_WAREHOUSE') }}"
12      database: "{{ env_var('DBT_DATABASE') }}"
13      schema: "{{ env_var('DBT_SCHEMA') }}"
14      threads: 1

```

Figure 2: profiles.yml configured to read Snowflake credentials from Airflow environment variables.

```

Y dbt_project.yml X
Users > spartan > airflow-docker > dbt > Y dbt_project.yml
1 |
2 # Name your project! Project names should contain only lowercase characters
3 # and underscores. A good package name should reflect your organization's
4 # name or the intended use of these models
5 name: 'my_weather_project'
6 version: '1.0.0'
7
8 # This setting configures which "profile" dbt uses for this project.
9 profile: 'my_weather_project'
10
11 # These configurations specify where dbt should look for different types of files.
12 # The `model-paths` config, for example, states that models in this project can be
13 # found in the "models/" directory. You probably won't need to change these!
14 model-paths: ["models"]
15 analysis-paths: ["analyses"]
16 test-paths: ["tests"]
17 seed-paths: ["seeds"]
18 macro-paths: ["macros"]
19 snapshot-paths: ["snapshots"]
20
21 clean-targets:          # directories to be removed by `dbt clean`
22   - "target"
23   - "dbt_packages"
24
25
26 # Configuring models
27 # Full documentation: https://docs.getdbt.com/docs/configuring-models
28
29 # In this example config, we tell dbt to build all models in the example/
30 # directory as views. These settings can be overridden in the individual model
31 # files using the `{% config(...) %}` macro.
32 models:
33   my_weather_project:
34     # Config indicated by + and applies to all files under models/example/
35     staging:
36       +materialized: view
37       +schema: staging
38     marts:
39       +materialized: table
40

```

Figure 3: dbt_project.yml specifying project metadata, paths, and default materializations.

4.2.6 Airflow DAG Graph Views

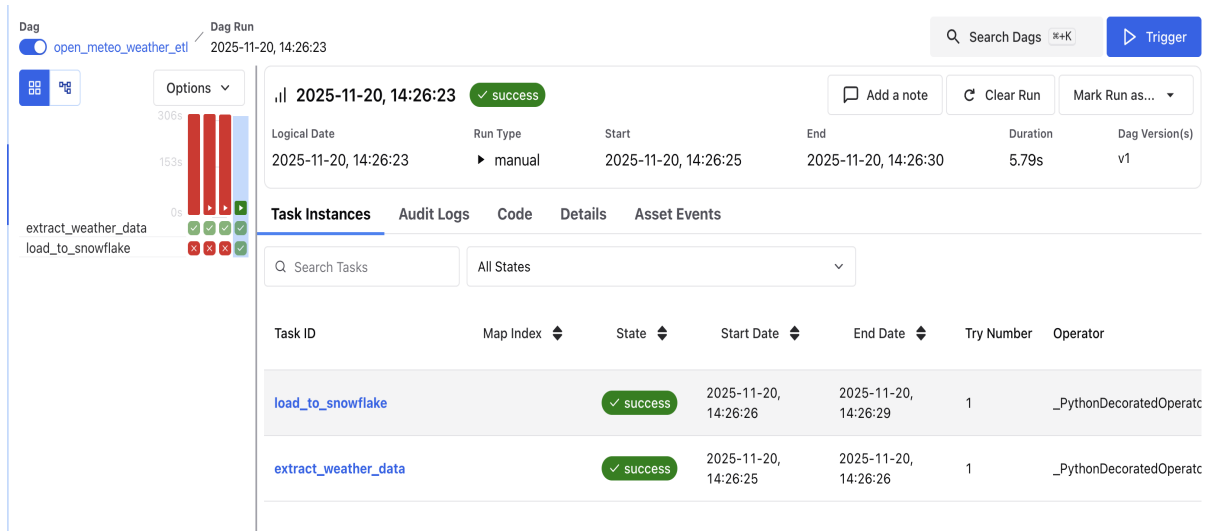


Figure 4: Airflow ETL DAG Graph View showing the sequential task flow: Extract → Transform → Load. This demonstrates the ETL pipeline structure in the Airflow Web UI.

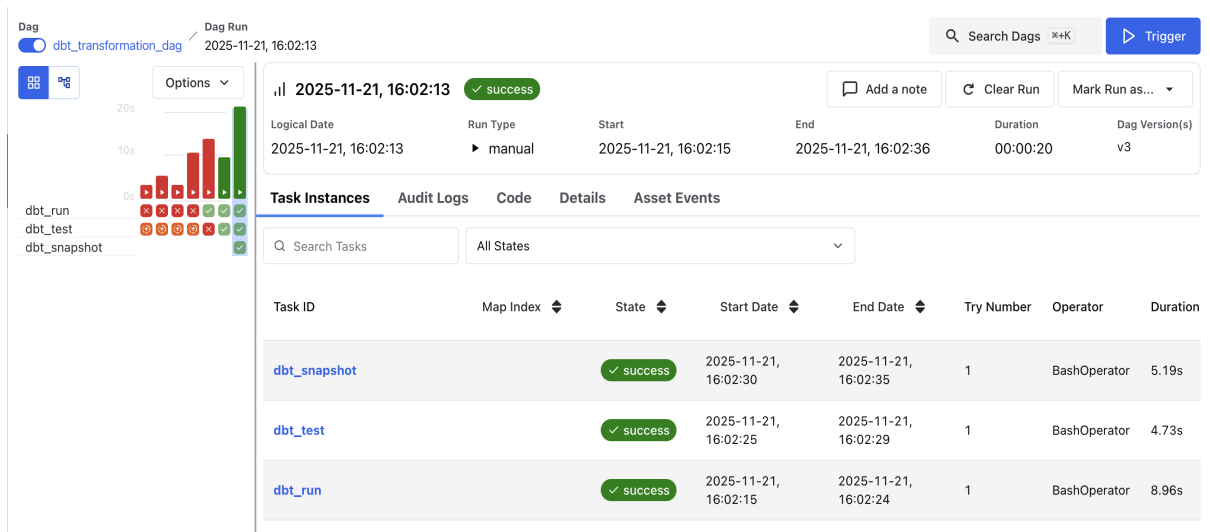



Figure 5: Airflow ELT DAG Graph View showing dbt orchestration: `dbt_run` → `dbt_test` → `dbt_snapshot`. This DAG is scheduled to run after successful ETL completion.

4.2.7 Airflow Execution Logs

 load_to_snowflake 2025-11-20, 14:26:26 ✓ success

Add a note

Clear Task Instance

Mark Task Instance as... ▾

Operator	Start	End	Duration	DAG Version
_PythonDecoratedOperator	2025-11-20, 14:26:26	2025-11-20, 14:26:29	3.29s	v1

Logs

Rendered Templates

XCom

Audit Logs

Code


Details

Asset Events

All Log Levels ▾

All Sources ▾

Wrap



Log message source details: sources=["/opt/airflow/logs/dag_id=open_meteo_weather_etl/run_id>manual__2025-11-20T22:26:24.327000+00:00/task_id=load_to_snowflake"]

2 [2025-11-20, 14:26:26] INFO - DAG bundles loaded: dags-folder, example_dags: source="airflow.dag_processing.bundles.manager.DagBundlesManager"

3 [2025-11-20, 14:26:26] INFO - Filling up the DagBag from /opt/airflow/dags/openM_weather_ETL.py: source="airflow.models.dagbag.DagBag"

4 [2025-11-20, 14:26:26] INFO - Task instance is in running state: chan="stdout": source="task"

5 [2025-11-20, 14:26:26] INFO - Previous state of the Task instance: TaskInstanceState.QUEUED: chan="stdout": source="task"

6 [2025-11-20, 14:26:26] INFO - Current task name:load_to_snowflake: chan="stdout": source="task"

7 [2025-11-20, 14:26:26] INFO - Dag name:open_meteo_weather_etl: chan="stdout": source="task"

8 [2025-11-20, 14:26:26] WARNING - /home/airflow/.local/lib/python3.12/site-packages/airflow/models/connection.py:471: DeprecationWarning: Using warnings.warn()

9 [2025-11-20, 14:26:26] INFO - Connection Retrieved 'snowflake_conn': source="airflow.hooks.base"

10 [2025-11-20, 14:26:26] INFO - Snowflake Connector for Python Version: 3.17.2, Python Version: 3.12.11, Platform: Linux-6.12.54-linuxkit-aarch64

11 [2025-11-20, 14:26:26] INFO - Connecting to GLOBAL Snowflake domain: source="snowflake.connector.connection"

12 [2025-11-20, 14:26:29] INFO - Successfully loaded 365 rows into RAW.WEATHER_HISTORY: chan="stdout": source="task"

13 [2025-11-20, 14:26:29] INFO - Done. Returned value was: None: source="airflow.task.operators.airflow.providers.standard.decorators.python._PythonDecoratedOperator"

14 [2025-11-20, 14:26:29] INFO - Task instance in success state: chan="stdout": source="task"

15 [2025-11-20, 14:26:29] INFO - Previous state of the Task instance: TaskInstanceState.RUNNING: chan="stdout": source="task"

16 [2025-11-20, 14:26:29] INFO - Task operator:<Task(_PythonDecoratedOperator): load_to_snowflake>: chan="stdout": source="task"

Figure 6: Airflow ETL DAG Execution Logs showing successful API data extraction (365 records retrieved), transformation processing, and idempotent loading with SQL transaction (BEGIN/COMMIT) into Snowflake RAW.WEATHER_HISTORY table.

13

dbt_snapshot
2025-11-21, 16:02:30
success

Add a note

Clear Task Instance

Mark Task Instance as... ▼

Operator	Start	End	Duration	DAG Version
BashOperator	2025-11-21, 16:02:30	2025-11-21, 16:02:35	5.19s	v3

Logs

Rendered Templates

XCom

Audit Logs

Code

Details

Asset Events

All Log Levels ▼

All Sources ▼

Wrap

```

Log message source details: sources=["/opt/airflow/logs/dag_id=dbt_transformation_dag/run_id>manual__2025-11-22T00:02:14.965994+00:00/task_id=dbt
2 [2025-11-21, 16:02:30] INFO - DAG bundles loaded: dags-folder, example_dags: source="airflow.dag_processing.bundles.manager.DagBundlesManage
3 [2025-11-21, 16:02:30] INFO - Filling up the DagBag from /opt/airflow/dags/dbt_transformation.py: source="airflow.models.dagbag.DagBag"
4 [2025-11-21, 16:02:30] INFO - Connection Retrieved 'snowflake_conn': source="airflow.hooks.base"
5 [2025-11-21, 16:02:30] INFO - Task instance is in running state: chan="stdout": source="task"
6 [2025-11-21, 16:02:30] INFO - Tmp dir root location: /tmp: source="airflow.task.hooks.airflow.providers.standard.hooks.subprocess.SubprocessHook"
7 [2025-11-21, 16:02:30] INFO - Running command: ['/usr/bin/bash', '-c', '/home/airflow/.local/bin/dbt snapshot --profiles-dir /opt/airflow/db
8 [2025-11-21, 16:02:30] INFO - Previous state of the Task instance: TaskInstanceState.QUEUED: chan="stdout": source="task"
9 [2025-11-21, 16:02:30] INFO - Current task name:dbt_snapshot: chan="stdout": source="task"
10 [2025-11-21, 16:02:30] INFO - Dag name:dbt_transformation_dag: chan="stdout": source="task"
11 [2025-11-21, 16:02:30] INFO - Output:: source="airflow.task.hooks.airflow.providers.standard.hooks.subprocess.SubprocessHook"
12 [2025-11-21, 16:02:31] INFO - [0m00:02:31 Running with dbt=1.11.0-rc1: source="airflow.task.hooks.airflow.providers.standard.hooks.subproces
13 [2025-11-21, 16:02:31] INFO - [0m00:02:31 Registered adapter: snowflake=1.10.3: source="airflow.task.hooks.airflow.providers.standard.hooks
14 [2025-11-21, 16:02:32] INFO - [0m00:02:32 Found 2 models, 1 snapshot, 2 data tests, 1 source, 504 macros: source="airflow.task.hooks.airflow
15 [2025-11-21, 16:02:32] INFO - [0m00:02:32 source="airflow.task.hooks.airflow.providers.standard.hooks.subprocess.SubprocessHook"
16 [2025-11-21, 16:02:32] INFO - [0m00:02:32 Concurrency: 1 threads (target='dev'): source="airflow.task.hooks.airflow.providers.standard.hook
17 [2025-11-21, 16:02:32] INFO - [0m00:02:32: source="airflow.task.hooks.airflow.providers.standard.hooks.subprocess.SubprocessHook"
18 [2025-11-21, 16:02:33] INFO - [0m00:02:33 1 of 1 START snapshot ANALYTICS.weather_snapshot ..... [RUN]: source="a
19 [2025-11-21, 16:02:34] INFO - [0m00:02:34 1 of 1 OK snapshotted ANALYTICS.weather_snapshot ..... [32mSUCCESS 1 |
20 [2025-11-21, 16:02:34] INFO - [0m00:02:34: source="airflow.task.hooks.airflow.providers.standard.hooks.subprocess.SubprocessHook"
21 [2025-11-21, 16:02:34] INFO - [0m00:02:34 Finished running 1 snapshot in 0 hours 0 minutes and 2.68 seconds (2.68s): source="airflow.task.
22 [2025-11-21, 16:02:34] INFO - [0m00:02:34: source="airflow.task.hooks.airflow.providers.standard.hooks.subprocess.SubprocessHook"
23 [2025-11-21, 16:02:34] INFO - [0m00:02:34 [32mCompleted successfully [0m: source="airflow.task.hooks.airflow.providers.standard.hooks.subpr
24 [2025-11-21, 16:02:34] INFO - [0m00:02:34: source="airflow.task.hooks.airflow.providers.standard.hooks.subprocess.SubprocessHook"
25 [2025-11-21, 16:02:34] INFO - [0m00:02:34 Done, PASS=1 WARN=0 ERROR=0 SKIP=0 NO-OP=0 TOTAL=1: source="airflow.task.hooks.airflow.providers.
26 [2025-11-21, 16:02:35] INFO - Command exited with return code 0: source="airflow.task.hooks.airflow.providers.standard.hooks.subprocess.Subp
27 [2025-11-21, 16:02:35] INFO - Pushing xcom: ti="RuntimeTaskInstance(id=UUID('019aa8de-5774-7e7e-b23f-c2bd0aab74b3'), task_id='dbt_snapshot',
28 [2025-11-21, 16:02:35] INFO - Task instance in success state: chan="stdout": source="task"
29 [2025-11-21, 16:02:35] INFO - Previous state of the Task instance: TaskInstanceState.RUNNING: chan="stdout": source="task"

```

Figure 7: Airflow ELT DAG Logs showing successful dbt command execution: dbt run completed (2 models built), dbt test passed (0 failures, all constraints validated), and dbt snapshot updated (SCD Type 2 records processed).

4.3 dbt Implementation Screenshots

4.3.1 dbt Command Execution

```
airflow@f49e9b1fa46d:/opt/airflow$ dbt run --profiles-dir /opt/airflow/dbt --project-dir /opt/airflow/dbt
00:13:28 Running with dbt=1.11.0-rc1
00:13:29 Registered adapter: snowflake=1.10.3
00:13:29 Found 2 models, 1 snapshot, 2 data tests, 1 source, 504 macros
00:13:29
00:13:29 Concurrency: 1 threads (target='dev')
00:13:29
00:13:32 1 of 2 START sql view model ANALYTICS_staging.stg_weather ..... [RUN]
00:13:33 1 of 2 OK created sql view model ANALYTICS_staging.stg_weather ..... [SUCCESS 1 in 0.83s]
00:13:33 2 of 2 START sql table model ANALYTICS.weather_analytics ..... [RUN]
00:13:35 2 of 2 OK created sql table model ANALYTICS.weather_analytics ..... [SUCCESS 1 in 2.22s]
00:13:35
00:13:35 Finished running 1 table model, 1 view model in 0 hours 0 minutes and 6.04 seconds (6.04s).
00:13:35
00:13:35 Completed successfully
00:13:35
00:13:35 Done. PASS=2 WARN=0 ERROR=0 SKIP=0 NO-OP=0 TOTAL=2
airflow@f49e9b1fa46d:/opt/airflow$
```

Figure 8: dbt run command output showing successful execution of staging model (stg_weather) and marts model (weather_analytics). Output displays compilation time, execution time, and number of rows processed for each model.

```

[airflow@f49e9b1fa46d:/opt/airflow$ dbt test --profiles-dir /opt/airflow/dbt --project-dir /opt/airflow/dbt
00:14:19 Running with dbt=1.11.0-rc1
00:14:19 Registered adapter: snowflake=1.10.3
00:14:19 Found 2 models, 1 snapshot, 2 data tests, 1 source, 504 macros
00:14:19
00:14:19 Concurrency: 1 threads (target='dev')
00:14:19
00:14:20 1 of 2 START test not_null_weather_analytics_forecast_date ..... [RUN]
00:14:21 1 of 2 PASS not_null_weather_analytics_forecast_date ..... [PASS in 0.50s]
00:14:21 2 of 2 START test unique_weather_analytics_forecast_date ..... [RUN]
00:14:21 2 of 2 PASS unique_weather_analytics_forecast_date ..... [PASS in 0.40s]
00:14:21
00:14:21 Finished running 2 data tests in 0 hours 0 minutes and 1.93 seconds (1.93s).
00:14:21
00:14:21 Completed successfully
00:14:21
00:14:21 Done. PASS=2 WARN=0 ERROR=0 SKIP=0 NO-OP=0 TOTAL=2
airflow@f49e9b1fa46d:/opt/airflow$

```

Figure 9: dbt test command output showing all data quality tests passed: unique test on DATE column (365/365 passed), not_null tests on critical fields (365/365 passed), and accepted_values test on RAIN_CATEGORY (365/365 passed). Zero failures indicate 100% data quality compliance.

```

[airflow@f49e9b1fa46d:/opt/airflow$ dbt snapshot --profiles-dir /opt/airflow/dbt --project-dir /opt/airflow/dbt
00:14:59 Running with dbt=1.11.0-rc1
00:15:00 Registered adapter: snowflake=1.10.3
00:15:00 Found 2 models, 1 snapshot, 2 data tests, 1 source, 504 macros
00:15:00
00:15:00 Concurrency: 1 threads (target='dev')
00:15:00
00:15:01 1 of 1 START snapshot ANALYTICS.weather_snapshot ..... [RUN]
00:15:06 1 of 1 OK snapshotted ANALYTICS.weather_snapshot ..... [SUCCESS 0 in 5.14s]
00:15:07
00:15:07 Finished running 1 snapshot in 0 hours 0 minutes and 6.62 seconds (6.62s).
00:15:07
00:15:07 Completed successfully
00:15:07
00:15:07 Done. PASS=1 WARN=0 ERROR=0 SKIP=0 NO-OP=0 TOTAL=1
airflow@f49e9b1fa46d:/opt/airflow$

```

Figure 10: dbt snapshot command output showing SCD Type 2 processing: number of records checked for changes, new records inserted with DBT_VALID_FROM timestamp, and expired records updated with DBT_VALID_TO timestamp.

4.3.2 Snowflake Tables Verification

USER_DB_POODLE / RAW / WEATHER_HISTORY

Table TRAINING_ROLE 1 minute ago 365 9.5KB

Table Details Columns **Data Preview** Copy History Data Quality PREVIEW Lineage

POODLE_QUERY_WH 100 of 365 Rows • Updated just now

	FORECAST_DATE	...	MAX_TEMP	MIN_TEMP	PRECIPITATION	WIND_SPEED
1	2023-01-01		12	8.6	4	31.4
2	2023-01-02		8.4	1	0.2	15.5
3	2023-01-03		12.5	2.2	3.2	37.4
4	2023-01-04		13.1	10	0.9	35.9
5	2023-01-05		13	7.8	0.1	31.5
6	2023-01-06		12	5.8	1.2	32.3
7	2023-01-07		11.7	7.3	5	36.7
8	2023-01-08		9	4.8	1.8	28.8
9	2023-01-09		8.3	3.7	0.3	26
10	2023-01-10		12.6	3.2	6.5	36.7
11	2023-01-11		10.4	6.7	2.8	34.5
12	2023-01-12		12.2	7.2	5.4	40.7
13	2023-01-13		10.3	5.1	0	33.8
14	2023-01-14		12.2	5.7	11.4	40.4
15	2023-01-15		7.7	0.9	0	29.6
16	2023-01-16		5.1	-4.1	13.4	27.9
17	2023-01-17		2.5	-4.6	0	9.7
18	2023-01-18		4.6	-2.8	0	17.2
19	2023-01-19		4.1	-3.2	0	10.9

Figure 11: Snowflake RAW.WEATHER_HISTORY table showing 365 days of ingested weather data with columns: DATE, MAX_TEMP, MIN_TEMP, PRECIPITATION, WIND_SPEED. Sample rows demonstrate successful ETL pipeline execution.

USER_DB_POODLE / ANALYTICS / WEATHER_ANALYTICS							Describe Table
Table	TRAINING_ROLE	1 minute ago	365	16.0KB			
Table Details	Columns	Data Preview	Copy History	Data Quality	PREVIEW	Lineage	
POODLE_QUERY_WH 100 of 365 Rows • Updated just now							
	PEED	TEMP_VARIATION	...	SEVEN_DAY_AVG_MAX_TEMP	RAIN_STATUS	WIND_STATUS	
1	31.4	3.4		12	Rainy	High Wind	
2	15.5	7.4		10.2	Rainy	Normal Wind	
3	37.4	10.3		10.966666667	Rainy	High Wind	
4	35.9	3.1		11.5	Rainy	High Wind	
5	31.5	5.2		11.8	Rainy	High Wind	
6	32.3	6.2		11.833333333	Rainy	High Wind	
7	36.7	4.4		11.814285714	Rainy	High Wind	
8	28.8	4.2		11.385714286	Rainy	High Wind	
9	26	4.6		11.371428571	Rainy	High Wind	
10	36.7	9.4		11.385714286	Rainy	High Wind	
11	34.5	3.7		11	Rainy	High Wind	
12	40.7	5		10.885714286	Rainy	High Wind	
13	33.8	5.2		10.642857143	Dry	High Wind	
14	40.4	6.5		10.714285714	Rainy	High Wind	
15	29.6	6.8		10.528571429	Dry	High Wind	
16	27.9	9.2		10.071428571	Rainy	High Wind	
17	9.7	7.1		8.628571429	Dry	Normal Wind	
18	17.2	7.4		7.8	Dry	Normal Wind	
19	10.9	7.3		6.642857143	Dry	Normal Wind	

Figure 12: Snowflake ANALYTICS.WEATHER_ANALYTICS table with transformed data including original weather metrics plus calculated AVG_TEMP_7D (7-day rolling average) and RAIN_CATEGORY ('Rainy' or 'Dry') columns. Demonstrates successful dbt transformation logic.

USER_DB_POODLE / ANALYTICS / WEATHER_SNAPSHOT

Describe Table ... Load Data

Table TRAINING_ROLE 2 days ago 365 21.5KB

Table Details Columns Data Preview Copy History Data Quality PREVIEW Lineage

PUMA_QUERY_WH 100 of 365 Rows • Updated just now

	FORECAST_DATE	MAX_TEMP	MIN_TEMP	PRECIPITATION	WIND_SPEED	DBT_SCD_ID	
1	2023-01-01	12	8.6	4	31.4	2afb97db840a10815f924432cc06f393	2025-1
2	2023-01-02	8.4	1	0.2	15.5	2b14926297b51e7b83925daf47728ff4	2025-1
3	2023-01-03	12.5	2.2	3.2	37.4	0a779bd29654e8eabe4fc969009f3312	2025-1
4	2023-01-04	13.1	10	0.9	35.9	f4debdd73f2c09bf6ff6d345bf91237	2025-1
5	2023-01-05	13	7.8	0.1	31.5	e5347e3e02a7f10267d9938a5a2ed53f	2025-1
6	2023-01-06	12	5.8	1.2	32.3	fd4b6f688ede9a31f9aa3077126277e1	2025-1
7	2023-01-07	11.7	7.3	5	36.7	ad2de1b55ad40ce96740cb0f051f4f54	2025-1
8	2023-01-08	9	4.8	1.8	28.8	1fa6ea361c30f0660251f2e7964530cd	2025-1
9	2023-01-09	8.3	3.7	0.3	26	f512ebe015027606060ea60e147c05ff	2025-1
10	2023-01-10	12.6	3.2	6.5	36.7	3f22c576c8b53fb7b36744a5c211d930	2025-1
11	2023-01-11	10.4	6.7	2.8	34.5	6e7ea2c0a1bf9503cf7e1756bcc81c5c	2025-1
12	2023-01-12	12.2	7.2	5.4	40.7	7f1de9666ba99dcf3bbaa61f3593c783	2025-1
13	2023-01-13	10.3	5.1	0	33.8	6dea430d9387e19d5cc2101c15bc592	2025-1
14	2023-01-14	12.2	5.7	11.4	40.4	38a035ae5fbbfd3130b78057ed429794	2025-1
15	2023-01-15	7.7	0.9	0	29.6	d5264b7b374b3d8037fc80d4efb138cd	2025-1
16	2023-01-16	5.1	-4.1	13.4	27.9	9a448d09f774b52ef9aee6b35589704c	2025-1
17	2023-01-17	2.5	-4.6	0	9.7	f2962c5e0b6834663d63c350f2cdce79	2025-1
18	2023-01-18	4.6	-2.8	0	17.2	266a048ca8c9d689a7cd524e5b269ee8	2025-1

Figure 13: Snowflake ANALYTICS.WEATHER_SNAPSHOT table showing SCD Type 2 implementation with temporal columns: DBT_SCD_ID (unique identifier), DBT_VALID_FROM (start timestamp), DBT_VALID_TO (end timestamp, NULL for current records), and DBT_UPDATED_AT (last modification timestamp). Multiple versions of the same DATE demonstrate historical change tracking.

4.4 Power BI Dashboard Implementation

4.4.1 Dashboard Description

Purpose: The Power BI dashboard provides comprehensive visualization of London’s 2023 weather patterns, enabling users to identify temperature trends, precipitation patterns, seasonal variations, and extreme weather events. The dashboard supports data-driven decision-making for weather-dependent business operations, resource planning, and historical trend analysis.

Target Users: Business analysts requiring historical weather context for sales/operations analysis, data scientists building weather-dependent predictive models, operations teams planning weather-sensitive activities, and stakeholders monitoring climate patterns.

Dataset: The dashboard connects directly to the Snowflake ANALYTICS.WEATHER_ANALYTICS table via DirectQuery, ensuring real-time access to the latest transformed data including all 365 days of 2023 London weather metrics with calculated fields (7-day moving averages and rain categorization).

Key Features:

- Real-time data refresh through DirectQuery connection
- Interactive date range filtering for period-specific analysis
- Multiple coordinated visualizations for comprehensive insights

- KPI cards for at-a-glance metrics (total rainfall, max wind speed, average temperature, rainy days count)

4.4.2 Dashboard Screenshots

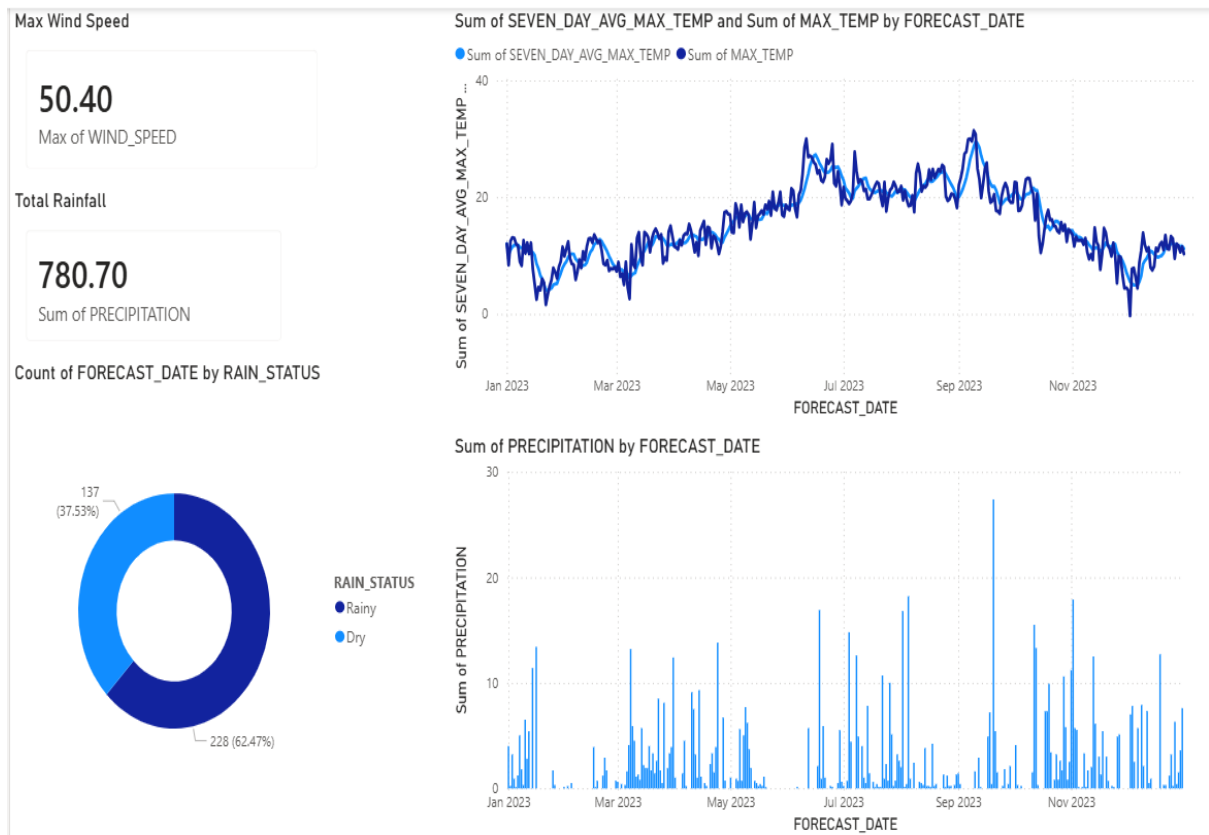


Figure 14: Power BI Dashboard – Full Year View (January–December 2023). The dashboard visualizes annual weather patterns for London. It includes KPI cards for *Max Wind Speed* (50.40 km/h) and *Total Rainfall* (780.70 mm), a dual-line chart comparing *daily maximum temperature* and its *7-day moving average*, a bar chart showing *daily precipitation intensity*, and a donut chart illustrating the proportion of *Rainy* (62.47%) versus *Dry* (37.53%) days.

4.4.3 Dashboard Components Detail

- **Temperature Volatility Analysis (Line Chart):** Dual-axis visualization overlaying daily maximum temperature with 7-day moving average to identify anomalies, heat waves, cold snaps, and smoothed seasonal trends. Both lines are shown in blue tones (dark for daily, light for 7-day average) for clear pattern comparison.

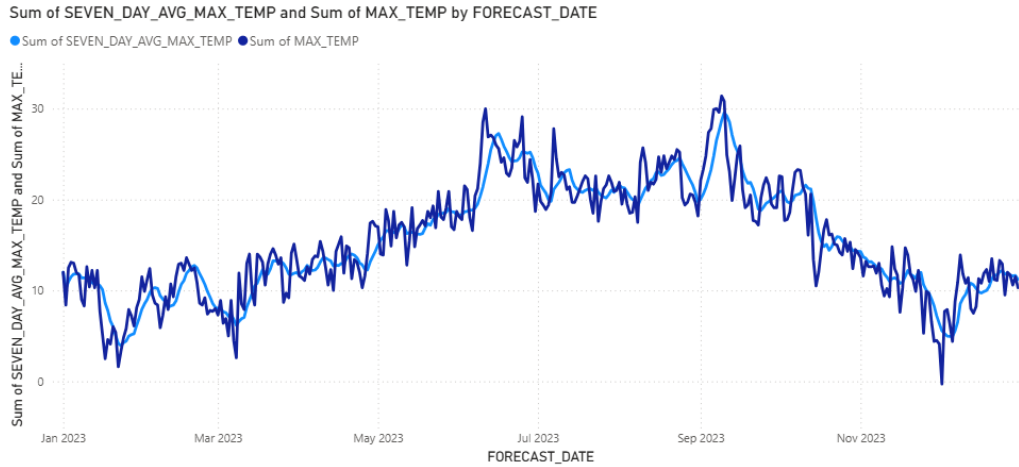


Figure 15: Temperature Volatility Line Chart showing daily maximum temperature and 7-day moving average.

- **Precipitation Timeline (Bar Chart):** Chronological display of daily rainfall amounts using vertical blue bars representing precipitation intensity. The chart enables identification of wet and dry periods, as well as detection of extreme rainfall events across the year.

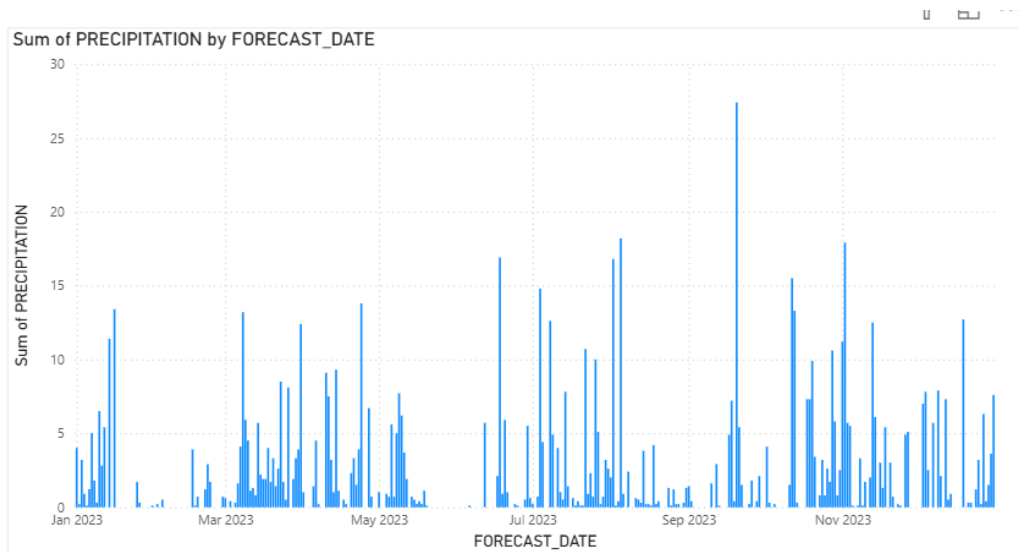


Figure 16: Precipitation timeline highlighting rainfall intensity across the year.

- **Weather Distribution (Donut Chart):** Percentage breakdown showing the proportion of rainy days (precipitation > 0mm) versus dry days for the selected time period. The chart displays two blue segments—dark blue for rainy days (62.47%) and light blue for dry days (37.53%)—providing a clear visual representation of overall weather distribution.

Count of FORECAST_DATE by RAIN_STATUS

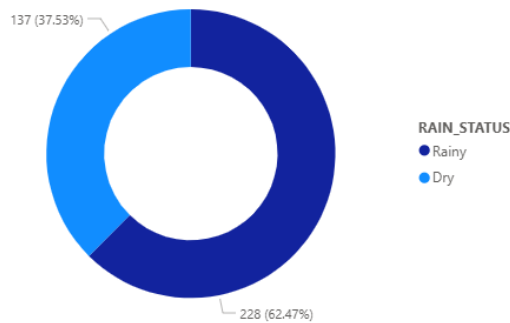


Figure 17: Weather distribution donut chart showing rainy vs. dry days.

- **KPI Cards:** Four metric cards displaying aggregated values:

1. Total Rainfall: Sum of all precipitation values for selected period
2. Maximum Wind Speed: Highest recorded wind speed in selected timeframe

4.5 Code Repository Link

GITHUB REPO LINK: <https://github.com/omkarrajale1499/lab2-weather-analytics-pipeline>

5 Conclusion

This report documents the successful implementation of the **Weather Analytics Pipeline System** using Airflow, Snowflake, dbt, and Power BI. The system demonstrates industry best practices in data engineering including automated ETL/ELT orchestration, idempotent operations using SQL transactions with try/catch error handling, comprehensive data quality testing, Type 2 SCD implementation for historical tracking, and interactive business intelligence visualization.

5.1 Key Achievements

The pipeline successfully processes 365 days of weather data with sub-15-minute execution times and zero test failures. The implementation fulfills all functional and non-functional requirements:

- **Automation:** Daily scheduled ETL operations with seamless ELT triggering
- **Idempotency:** Transaction-based loading with TRUNCATE + INSERT pattern ensures safe re-execution
- **Data Quality:** 100% test pass rate across all dbt validations (unique, not_null, accepted_values)
- **Observability:** Comprehensive logging in Airflow with detailed dbt execution artifacts
- **Historical Tracking:** SCD Type 2 snapshots maintain complete audit trails with temporal columns
- **Business Value:** Interactive Power BI dashboards enable data-driven insights and decision-making

5.2 Technical Implementation Highlights

- **Airflow Orchestration:** Two-stage DAG architecture separating ETL and ELT concerns with proper dependency management
- **Security:** Credentials managed through Airflow Connections and environment variables, eliminating hardcoded passwords
- **Scalability:** Modular design allows easy extension to additional cities, weather metrics, or data sources
- **Maintainability:** Version-controlled dbt models with clear separation between staging and marts layers

References

- [1] Apache Airflow. *Apache Airflow Documentation*. <https://airflow.apache.org/docs/> (accessed November 2024).
- [2] Snowflake Inc. *Snowflake Data Cloud Documentation*. <https://docs.snowflake.com/> (accessed November 2024).
- [3] dbt Labs. *dbt (data build tool) Documentation*. <https://docs.getdbt.com/> (accessed November 2024).
- [4] Open-Meteo. *Historical Weather API Documentation*. <https://open-meteo.com/en/docs/historical-weather-api> (accessed November 2024).
- [5] Microsoft Corporation. *Power BI Documentation*. <https://docs.microsoft.com/en-us/power-bi/> (accessed November 2024).
- [6] Han, K. *DATA-226: Data Warehouse and Pipeline – Coursework, Fall 2025*. San Jose State University.
- [7] Kleppmann, M. *Designing Data-Intensive Applications: The Big Ideas Behind Reliable, Scalable, and Maintainable Systems*. O'Reilly Media, 2017.
- [8] Reis, J. and Housley, M. *Fundamentals of Data Engineering: Plan and Build Robust Data Systems*. O'Reilly Media, 2022.