# Project Report: Weather Data Pipeline

1. Team

| Name | Role |
|---|---|
| Aubakirov Sanzhar | DAG 1 — Data Ingestion and Kafka Producer |
| Toremuratuly Abylay | DAG 2 — Data Cleaning and SQLite Writer |
| Toleu Bakhauddin | DAG 3 — Daily Analytics |

Roles description:

- DAG 1: Setup and run Airflow process for periodic data collection from WeatherAPI and sending to Kafka.

- DAG 2: Cleaning and processing data via Pandas, writing to SQLite.
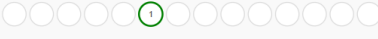
- DAG 3: Daily analytics — computing aggregates and saving results to summary table.



## 2. API

Selected API: WeatherAPI (http://api.weatherapi.com)

Selection criteria:

- Frequent data updates (every few minutes)

- Stability and well-documented

- JSON response format

- Provides real values for temperature, humidity, and weather conditions

- Not used in previous lab exercises

Example request:

GET http://api.weatherapi.com/v1/current.json?key=YOUR_API_KEY&q=Almaty&aqi=no

Example JSON response (key fields):

```json
{
  "location": {"name": "Almaty"},
  "current": {
    "temp_c": -1.9,
    "humidity": 92,
    "condition": {"text": "Freezing fog"},
    "wind_kph": 5.4,
    "pressure_mb": 1012.0,
    "feelslike_c": -5.0
  }
}
```

Comment: API is stable and provides real-time weather data, meeting project requirements.

---

## 3. Pipeline Architecture

Overall flow:

WeatherAPI → DAG 1 (Producer) → Kafka (raw_weather_events)

    → DAG 2 (Batch Processing) → SQLite (events)

    → DAG 3 (Daily Analytics) → SQLite (daily_weather_summary)

DAG descriptions:

- DAG 1: Periodic collection of weather data and sending to Kafka.
- DAG 2: Reading from Kafka, cleaning and validation, storing into SQLite.
- DAG 3: Reading from SQLite, aggregating data (min/max/avg), storing into summary table.

---

## 4. Kafka Topic Schema

| Field | Type | Description |
|---|---|---|
| timestamp | string | Request time |
| city | string | City |
| weather.current.temp_c | float | Temperature |
| weather.current.humidity | int | Humidity |
| weather.current.condition.text | string | Weather condition |
| metadata.source | string | Data source |

Topic: raw_weather_events

Comment: Each new API event is sent to Kafka with metadata.

---

## 5. Cleaning Rules (DAG 2)

Processing with Pandas:

1. Check mandatory fields: timestamp, city, weather.current.

2. Type conversion: temperature_c → float, humidity → int.

3. Allowed value ranges:

   o Temperature: -100 … 100 °C

   o Humidity: 0 … 100 %

4. Extract weather condition (condition.text).

5. Additional fields: wind_kph, pressure_mb, feelslike_c.
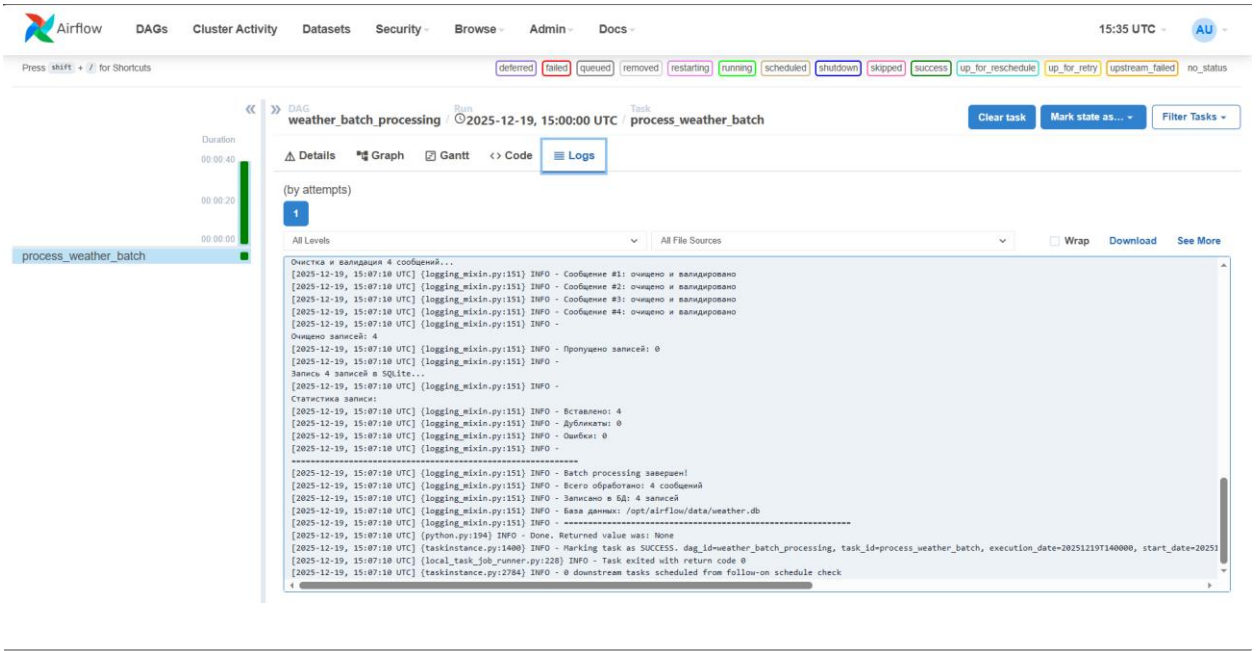
6. Skip invalid or empty records.

Example before/after cleaning:

*Before cleaning (JSON):*

| timestamp | city | weather.current.temp_c | weather.current.humidity | weather.current.condition.text |
|---|---|---|---|---|
| 2025-12-19T13:55:42 | Almaty | -1.9 | 92 | Freezing fog |

*After cleaning (DataFrame for SQLite):*

| timestamp | city | temperature_c | humidity | condition_text | wind_kph | pressure_mb | feelslike_c | source |
|---|---|---|---|---|---|---|---|---|
| 2025-12-19T13:55:42 | Almaty | -1.9 | 92 | Freezing fog | 5.4 | 1012.0 | -5.0 | weatherapi.com |

DAG
weather_batch_processing    Run    ©2025-12-19, 15:00:00 UTC    Task    process_weather_batch    Clear task    Mark state as... ▾    Filter Tasks ▾

⚠ Details    ᴥ Graph    ⧉ Gantt    <> Code    ☰ Logs

(by attempts)
1

All Levels        ▾        All File Sources        ▾        ☐ Wrap    Download    See More

```
Очистка и валидация 4 сообщений...
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - Сообщение #1: очищено и валидировано
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - Сообщение #2: очищено и валидировано
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - Сообщение #3: очищено и валидировано
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - Сообщение #4: очищено и валидировано
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO -
Очищено записей: 4
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - Пропущено записей: 0
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO -
Запись 4 записей в SQLite...
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO -
Статистика записи:
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - Вставлено: 4
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - Дубликаты: 0
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - Ошибки: 0
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO -
============================================================
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - Batch processing завершён!
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - Всего обработано: 4 сообщений
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - Записано в БД: 4 записей
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - База данных: /opt/airflow/data/weather.db
[2025-12-19, 15:07:10 UTC] {logging_mixin.py:151} INFO - ============================================================
[2025-12-19, 15:07:10 UTC] {python.py:194} INFO - Done. Returned value was: None
[2025-12-19, 15:07:10 UTC] {taskinstance.py:1400} INFO - Marking task as SUCCESS. dag_id=weather_batch_processing, task_id=process_weather_batch, execution_date=20251219T140000, start_date=20251...
[2025-12-19, 15:07:10 UTC] {local_task_job_runner.py:228} INFO - Task exited with return code 0
[2025-12-19, 15:07:10 UTC] {taskinstance.py:2784} INFO - 0 downstream tasks scheduled from follow-on schedule check
```

# 6. SQLite Schema

Table: events

| Field | Type | Description |
|---|---|---|
| timestamp | TEXT | Event time |
| city | TEXT | City |
| temperature_c | REAL | Temperature |
| humidity | INTEGER | Humidity |
| condition_text | TEXT | Weather condition |
| wind_kph | REAL | Wind speed |
| pressure_mb | REAL | Pressure |
| feelslike_c | REAL | Feels like |
| source | TEXT | Data source |

Table: daily_weather_summary

| Field | Type | Description |
|---|---|---|
| date | TEXT | Date |
| city | TEXT | City |
| min_temp | REAL | Minimum temperature |
| max_temp | REAL | Maximum temperature |
| avg_temp | REAL | Average temperature |
| avg_humidity | REAL | Average humidity |
| records_count | INTEGER | Number of records |
| created_at | TEXT | Inserted timestamp |
| UNIQUE(date, city) | - | Uniqueness per date and city |

Comment: DAG 3 generates daily aggregates for each location.

## 7. DAGs and Logs

DAG 1 — Weather Ingestion

- Periodic data collection from WeatherAPI, triggered every minute (one fetch per run).

- Sends data to Kafka topic raw_weather_events.

- Simulates streaming without blocking other DAGs.

- Log snippet:

Weather producer started (one iteration)

Sent: -1.9°C, Freezing fog

Task marked as SUCCESS

DAG 2 — Batch Processing

- Reads data from Kafka.

- Cleans and validates with Pandas.

- Stores into SQLite table events.

- Log snippet:

Total messages processed: 4

Inserted into DB: 4

Batch processing completed



DAG 3 — Daily Analytics

- Reads data from SQLite (events).

- Computes min/max/avg temperature and avg humidity.

- Writes to daily_weather_summary.

- Log snippet:

# Daily analytics completed

## 8. Analytics Example

Almaty, 19.12.2025

| City | min_temp | max_temp | avg_temp | avg_humidity | records_count |
|------|----------|----------|----------|--------------|---------------|
| Almaty | -2.0 | 0.5 | -1.2 | 90 | 24 |

## 9. Repository

Project structure:

```
project/
│ README.md
│ requirements.txt
├── src/
│   ├── job1_producer.py
│   ├── job2_processor.py
│   ├── job3_analytics.py
│   └── database.py
├── airflow/
│   └── dags/
│       ├── dag1_weather_ingestion.py
│       ├── dag2_weather_batch.py
│       └── dag3_daily_analytics.py
├── data/
│   └── weather.db
└── report/
    └── report.pdf
```