# Project Report: Real-Time Air Quality Monitoring Pipeline

Team Members: Begali Dinmukhammed, Nazarov Dastan, Batyrbek Raiymbek

Date: December 19, 2025

## 1. API Justification

For this project, we selected the **Air Quality Index (AQI) API** ([https://aqicn.org/](https://aqicn.org/)).
- **Relevance:** Air quality is a critical environmental factor. Real-time monitoring allows for health alerts and urban planning.
- **Data Richness:** The API provides not only the overall AQI but also granular data like station coordinates, measurement time, and individual pollutant levels (though we focused on AQI for consistency).
- **Reliability:** It aggregates data from over 30,000 stations globally, providing a robust stream for our Kafka pipeline.

## 2. Kafka Topic Schema

The system uses a single Kafka topic named weather_raw_data. Data is produced in **JSON** format.
**Message structure:**
JSON

```json
{
 "city_raw": "String (e.g., 'Almaty, Kazakhstan')",
 "aqi": "Integer (0-500+)",
 "lat": "Float (Latitude)",
 "lon": "Float (Longitude)",
 "utime": "Integer (Unix Timestamp)",
 "timestamp": "String (ISO Format)"
}
```

- **Serialization:** UTF-8 encoded strings.
- **Partitions:** 1 (Default for development).

## 3. Data Cleaning Rules

Before storage, data undergoes a cleaning process in job2_cleaner.py:
1. **Normalization:** The city_raw field is split by commas to extract the primary city name. Leading/trailing spaces are removed, and the string is converted to **Title Case** (e.g., "ALMATY" -> "Almaty").
2. **Type Conversion:** AQI is explicitly cast to Integer. Coordinates are cast to Float.
3. **Handling Missing Values:** If AQI is missing, the record is discarded (filtering). If coordinates are missing, they are stored as NULL.
4. **Invalid Record Filtering:** Records with an AQI outside the logical range (0-999) are filtered out to remove sensor errors or "outlier" noise.

# 4. SQLite Storage Schema

The storage layer consists of two tables in app.db.
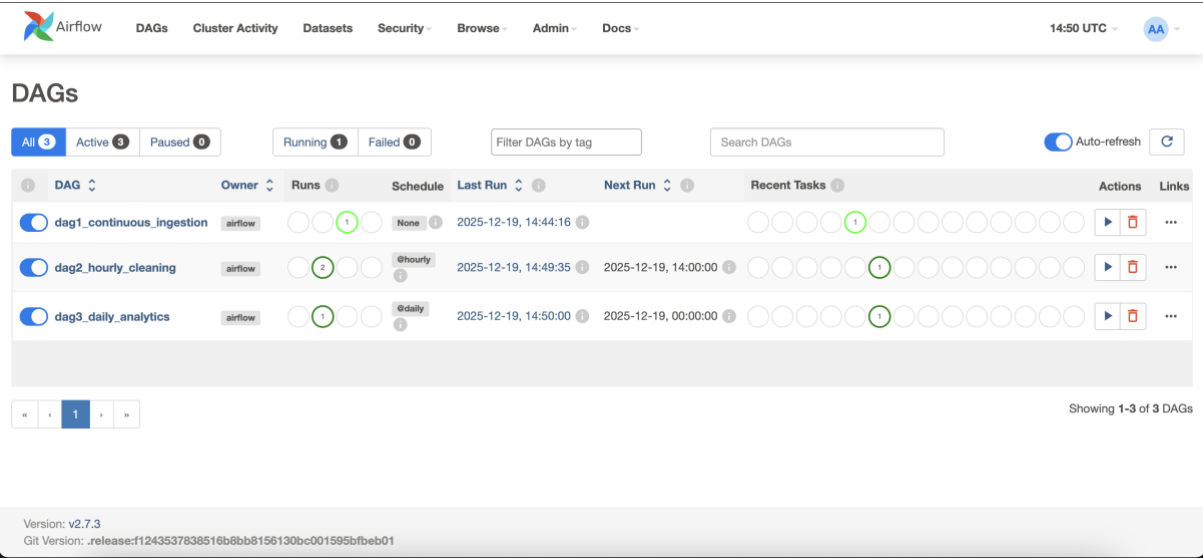
**Table: events (Cleaned Data)**

| Column | Type | Description |
|---|---|---|
| id | INTEGER | Primary Key (Autoincrement) |
| city | TEXT | Cleaned city name |
| aqi | INTEGER | Air Quality Index |
| lat | REAL | Latitude |
| lon | REAL | Longitude |
| utime | INTEGER | Unix Time |
| timestamp | DATETIME | Formatted date/time |

**Table: daily_summary (Aggregated Analytics)**

| Column | Type | Description |
|---|---|---|
| city | TEXT | Primary Key (City name) |
| min_aqi | INTEGER | Lowest AQI recorded |
| max_aqi | INTEGER | Highest AQI recorded |
| avg_aqi | REAL | Average AQI |
| total_measurements | INTEGER | Count of observations |
| air_quality_category | TEXT | 'Good', 'Moderate', 'Unhealthy', etc. |
| last_sync | DATETIME | Time of last calculation |

# 5. Implementation Evidence (Screenshots)

## 5.1 Airflow DAGs



## 5.2 Successful Task Logs (Job 2 Cleaner)

## 5.3 Successful Task Logs (Job 3 Analytics)



## 5.4 Final Database Output

### Table events

| id | city | aqi | lat | lon | utime | timestamp |
|----|------|-----|-----|-----|-------|-----------|
| 1 | Almaty | 85 | 43.222 | 76.851 | 1766174400 | 2025-12-19 20:00:00 |
| 2 | Astana | 5 | 51.125286 | 71.46722 | 1766174400 | 2025-12-19 20:00:00 |
| 3 | Almaty | 85 | 43.222 | 76.851 | 1766174400 | 2025-12-19 20:00:00 |
| 4 | Astana | 5 | 51.125286 | 71.46722 | 1766174400 | 2025-12-19 20:00:00 |
| 5 | Almaty | 85 | 43.222 | 76.851 | 1766174400 | 2025-12-19 20:00:00 |
| 6 | Astana | 5 | 51.125286 | 71.46722 | 1766174400 | 2025-12-19 20:00:00 |
| 7 | Almaty | 85 | 43.222 | 76.851 | 1766174400 | 2025-12-19 20:00:00 |
| 8 | Astana | 5 | 51.125286 | 71.46722 | 1766174400 | 2025-12-19 20:00:00 |
| 9 | Almaty | 85 | 43.222 | 76.851 | 1766174400 | 2025-12-19 20:00:00 |
| 10 | Astana | 5 | 51.125286 | 71.46722 | 1766174400 | 2025-12-19 20:00:00 |

### Table daily_summary

| city | min_aqi | max_aqi | avg_aqi | total_measurements | air_quality_category | last_sync |
|------|---------|---------|---------|--------------------|-----------------------|-----------|
| Almaty | 85 | 85 | 85 | 10 | Moderate | 2025-12-19 14:50:01 |
| Astana | 5 | 5 | 5 | 10 | Good | 2025-12-19 14:50:01 |

# 6. Configuration & Extensibility

The pipeline is designed for high flexibility through two root-level configuration files:

- **.env (Security):** Stores the AQI_TOKEN privately. This ensures that sensitive API keys are never hard-coded, following security best practices and allowing for easy key rotation.
- **cities.txt (Scalability):** Contains a list of target cities (one per line). The system dynamically reads this file at runtime; to monitor new locations, a user simply updates this list without needing to modify the Python scripts.

# 7. Conclusion

The pipeline successfully demonstrates a full ETL/ELT flow:

1. **Ingestion:** Real-time data collection via API into Kafka.
2. **Processing:** Stream-like batch cleaning and normalization.
3. **Storage:** Reliable storage in a structured SQLite database.
4. **Analytics:** Automated metric calculation for end-user insights.