

Data Collection and Preparation - Final Project

Members:

1. Amreyeva Alina – 22B031240
2. Kemel Merey – 22B030615
3. Serik Dinmukhammed – 22B030439

1. Project Goal

The objective of this project was to build a robust, scalable data pipeline capable of handling real-world, frequently updating data. We implemented a streaming ingestion with processing to analyze the Carbon Intensity of the UK National Grid. The system automates the collection, cleaning, storage, and analysis of environmental data using Apache Airflow, Kafka, and SQLite.

2. API Selection

Chosen API: Carbon Intensity UK (National Grid ESO) - <https://api.carbonintensity.org.uk>

We selected this API because it is well suited for a streaming data project. First, the data is updated **every 30 minutes**, which fits the requirement of frequently updated data. Second, the API provides **real-world and meaningful values**, such as forecasted and actual carbon intensity, which allows us to perform proper analysis and comparisons. Third, the API returns **clean and well-structured JSON data**, which makes it easy to process in Python. Finally, the API is **stable, well-documented, and free**, and it does not require complex authentication, which simplifies development.

3. Architecture Overview

The system is containerized using Docker and consists of three main stages controlled by Airflow DAGs:

1. **Ingestion Layer:** A producer script fetches data from the API every 30 seconds and buffers the raw JSON into Apache Kafka.
2. **Processing Layer:** An hourly job consumes messages from Kafka, cleans the data, and stores it in a relational database (SQLite).
3. **Analytics Layer:** A daily job aggregates the data to generate insights, specifically comparing Forecast vs. Actual intensity across different parts of the day.

Technologies Used:

- **Orchestration:** Apache Airflow (LocalExecutor with Postgres backend).
- **Messaging:** Apache Kafka & Zookeeper.
- **Storage:** SQLite.
- **Language:** Python (Pandas, Requests, Kafka-Python).

4. Implementation Details

4.1 Job 1: Continuous Ingestion (DAG 1)

- **Schedule:** Runs every 15 minutes.
- **Logic:** The Python script runs an internal loop for 14.5 minutes. Inside this loop, it queries the API every 30 seconds.
- **Fault Tolerance:** Implemented retry logic (10 attempts) for Kafka connections and error handling for API network timeouts.
- **Output:** Sends a JSON message to the Kafka topic raw_events.

4.2 Job 2: Cleaning & Storage (DAG 2)

- **Schedule:** @hourly
- **Cleaning Rules:**
 1. **Flattening:** Extracts nested fields (data.intensity.forecast, data.intensity.actual, data.intensity.index) from the raw JSON.
 2. **Type Conversion:** Ensures numerical values are stored as Integers.
 3. **Handling Nulls:** Stores actual intensity if available; handles potential missing values gracefully.
- **Output:** Inserts cleaned rows into the SQLite events table.

4.3 Job 3: Daily Analytics (DAG 3)

- **Schedule:** @daily
- **Logic:** Reads historical data from SQLite using Pandas.
- **Feature Engineering:** Adds a part_of_day column based on the timestamp (Morning: 05-12, Afternoon: 12-17, Evening: 17-21, Night: 21-05).
- **Aggregation:** Computes the Mean, Min, and Max intensity for each part of the day, handling cases where actual data might be missing by falling back to forecast values.
- **Output:** Writes results to the daily_summary table.

5. Data Schemas

5.1 Kafka Topic Schema (raw_events)

The raw message sent to Kafka follows this JSON structure:

code JSON

```
{
  "ingestion_id": 123,
  "timestamp": "2025-12-17T10:00:00.123456",
  "source": "carbonintensity.org.uk",
  "payload": {
    "from": "2025-12-17T10:00Z",
    "to": "2025-12-17T10:30Z",
```

```
"intensity": {
  "forecast": 260,
  "actual": 255,
  "index": "moderate"
}
```

```
kemelmerey@Kemels-MacBook-Air DCP_FinalProject % docker-compose exec kafka kafka-console-consumer --bootstrap-server kafka:29092 --topic raw_events --from-beginning
{"ingestion_id": 1, "timestamp": "2025-12-17T16:39:30.967949", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 134, "index": "moderate"}}}
{"ingestion_id": 2, "timestamp": "2025-12-17T16:40:02.249700", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 134, "index": "moderate"}}}
{"ingestion_id": 3, "timestamp": "2025-12-17T16:40:33.380392", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 4, "timestamp": "2025-12-17T16:41:04.345638", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 5, "timestamp": "2025-12-17T16:41:35.327244", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 6, "timestamp": "2025-12-17T16:42:06.484277", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 7, "timestamp": "2025-12-17T16:42:37.583684", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 8, "timestamp": "2025-12-17T16:43:08.529812", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 9, "timestamp": "2025-12-17T16:43:39.360713", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 10, "timestamp": "2025-12-17T16:44:10.208540", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 11, "timestamp": "2025-12-17T16:44:41.198321", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 12, "timestamp": "2025-12-17T16:45:13.703512", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 13, "timestamp": "2025-12-17T16:45:44.636832", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 14, "timestamp": "2025-12-17T16:46:15.720661", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 15, "timestamp": "2025-12-17T16:46:46.732376", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 16, "timestamp": "2025-12-17T16:47:17.698617", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 17, "timestamp": "2025-12-17T16:47:48.714840", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 18, "timestamp": "2025-12-17T16:48:19.789368", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 19, "timestamp": "2025-12-17T16:48:50.864600", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 20, "timestamp": "2025-12-17T16:49:21.939832", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 21, "timestamp": "2025-12-17T16:49:53.015064", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 22, "timestamp": "2025-12-17T16:50:24.090296", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 23, "timestamp": "2025-12-17T16:50:55.165528", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 24, "timestamp": "2025-12-17T16:51:26.240760", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 25, "timestamp": "2025-12-17T16:51:57.316000", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 26, "timestamp": "2025-12-17T16:52:28.391232", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 27, "timestamp": "2025-12-17T16:52:59.466464", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
{"ingestion_id": 28, "timestamp": "2025-12-17T16:53:30.541696", "source": "carbonintensity.org.uk", "payload": {"from": "2025-12-17T16:00Z", "to": "2025-12-17T16:30Z", "intensity": {"forecast": 139, "actual": 135, "index": "moderate"}}}
```

5.2 SQLite Database Schema

Table 1: events (Cleaned Data)

Column	Type	Description
id	INTEGER PK	Auto-incrementing ID
ingestion_timestamp	TEXT	Timestamp of ingestion
forecast_intensity	INTEGER	Predicted carbon intensity
actual_intensity	INTEGER	Realized carbon intensity
index_intensity	TEXT	Category (e.g., 'moderate')
source	TEXT	API Source
created_at	TIMESTAMP	Record creation time

	cid	name	type	notnull	default_value	pk
1	0	id	INTEGER	0	<null>	1
2	1	ingestion_timestamp	TEXT	0	<null>	0
3	2	forecast_intensity	INTEGER	0	<null>	0
4	3	actual_intensity	INTEGER	0	<null>	0
5	4	index_intensity	TEXT	0	<null>	0
6	5	source	TEXT	0	<null>	0
7	6	created_at	TIMESTAMP	0	CURRENT_TIMESTAMP	0

Table 2: daily_summary (Aggregated Analytics)

Column	Type	Description
id	INTEGER PK	Auto-incrementing ID

summary_date	DATE	Date of the analysis
part_of_day	TEXT	Morning, Afternoon, Evening, Night
avg_forecast	REAL	Average predicted intensity
avg_actual	REAL	Average actual intensity
max_intensity	INTEGER	Max intensity recorded
min_intensity	INTEGER	Min intensity recorded

	cid	name	type	nonnull	dflt_value	pk
1	0	summary_date	DATE	0	<null>	0
2	1	part_of_day	TEXT	0	<null>	0
3	2	avg_forecast	REAL	0	<null>	0
4	3	avg_actual	REAL	0	<null>	0
5	4	max_intensity	INTEGER	0	<null>	0
6	5	min_intensity	INTEGER	0	<null>	0

6. Verification

6.1 Airflow DAGs

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks	Actions	Links
job1_ingestion_dag	dcp.team		/15 * * * *	2025-12-17, 16:45:00	2025-12-17, 16:45:00			...
job2_clean_store_dag	dcp.team		@hourly	2025-12-17, 17:22:11	2025-12-17, 17:00:00			...
job3_daily_summary_dag	dcp.team		@daily	2025-12-17, 17:22:13	2025-12-17, 00:00:00			...

6.2 Job 1 Logs (Ingestion)

Duration: 00:14:43

00:07:21

00:00:00

start_ingestion_loop

job1_ingestion_dag

2025-12-17, 17:00:00 UTC

start_ingestion_loop

Details

Graph

Gantt

Code

Logs

(by attempt)

1

All Levels

All File Sources

Wrap

Download

See More

[2025-12-17, 17:08:43 UTC] (job1_producer.py:24) INFO - Connected to Kafka at ['kafka:29092']

[2025-12-17, 17:08:44 UTC] (conn.py:400) INFO - <BrokerConnection client_id=kafka-python-producer-3, node_id=1 host=kafka:29092 <connecting> [IPv4 ('172.24.0.4', 29092)]>: connecting to kafka:29092 [IPv4 ('172.24.0.4', 29092)]>: Connection complete.

[2025-12-17, 17:08:44 UTC] (conn.py:461) INFO - <BrokerConnection client_id=kafka-python-producer-3, node_id=1 host=kafka:29092 <connected> [IPv4 ('172.24.0.4', 29092)]>: Connection complete.

[2025-12-17, 17:08:44 UTC] (conn.py:951) INFO - <BrokerConnection client_id=kafka-python-producer-3, node_id=bootstrap-0 host=kafka:29092 <connected> [IPv4 ('172.24.0.4', 29092)]>: Closing connection.

[2025-12-17, 17:08:44 UTC] (job1_producer.py:62) INFO - [1] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:09:15 UTC] (job1_producer.py:62) INFO - [2] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:09:49 UTC] (job1_producer.py:67) ERROR - Loop Error: HTTPConnectionPool(host='api.carbonintensity.org.uk', port=443): Max retries exceeded with url: /intensity (Caused by NewConnectionError('url [2025-12-17, 17:10:20 UTC] (job1_producer.py:62) INFO - [4] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:10:54 UTC] (job1_producer.py:67) ERROR - Loop Error: HTTPConnectionPool(host='api.carbonintensity.org.uk', port=443): Max retries exceeded with url: /intensity (Caused by NewConnectionError('url [2025-12-17, 17:11:25 UTC] (job1_producer.py:62) INFO - [6] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:11:56 UTC] (job1_producer.py:62) INFO - [7] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:12:26 UTC] (job1_producer.py:62) INFO - [8] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:12:57 UTC] (job1_producer.py:62) INFO - [9] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:13:28 UTC] (job1_producer.py:62) INFO - [10] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:13:59 UTC] (job1_producer.py:62) INFO - [11] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:14:30 UTC] (job1_producer.py:62) INFO - [12] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:15:01 UTC] (job1_producer.py:62) INFO - [13] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:15:32 UTC] (job1_producer.py:62) INFO - [14] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:16:03 UTC] (job1_producer.py:62) INFO - [15] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:16:34 UTC] (job1_producer.py:62) INFO - [16] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:17:05 UTC] (job1_producer.py:62) INFO - [17] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:17:36 UTC] (job1_producer.py:62) INFO - [18] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:18:07 UTC] (job1_producer.py:62) INFO - [19] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:18:38 UTC] (job1_producer.py:62) INFO - [20] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:19:09 UTC] (job1_producer.py:62) INFO - [21] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:19:40 UTC] (job1_producer.py:62) INFO - [22] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:20:11 UTC] (job1_producer.py:62) INFO - [23] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:20:42 UTC] (job1_producer.py:62) INFO - [24] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:21:13 UTC] (job1_producer.py:62) INFO - [25] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:21:44 UTC] (job1_producer.py:62) INFO - [26] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:22:15 UTC] (job1_producer.py:62) INFO - [27] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:22:46 UTC] (job1_producer.py:62) INFO - [28] SENT to 'raw_events': Forecast 136

[2025-12-17, 17:23:16 UTC] (job1_producer.py:71) INFO - Ingestion Job Finished

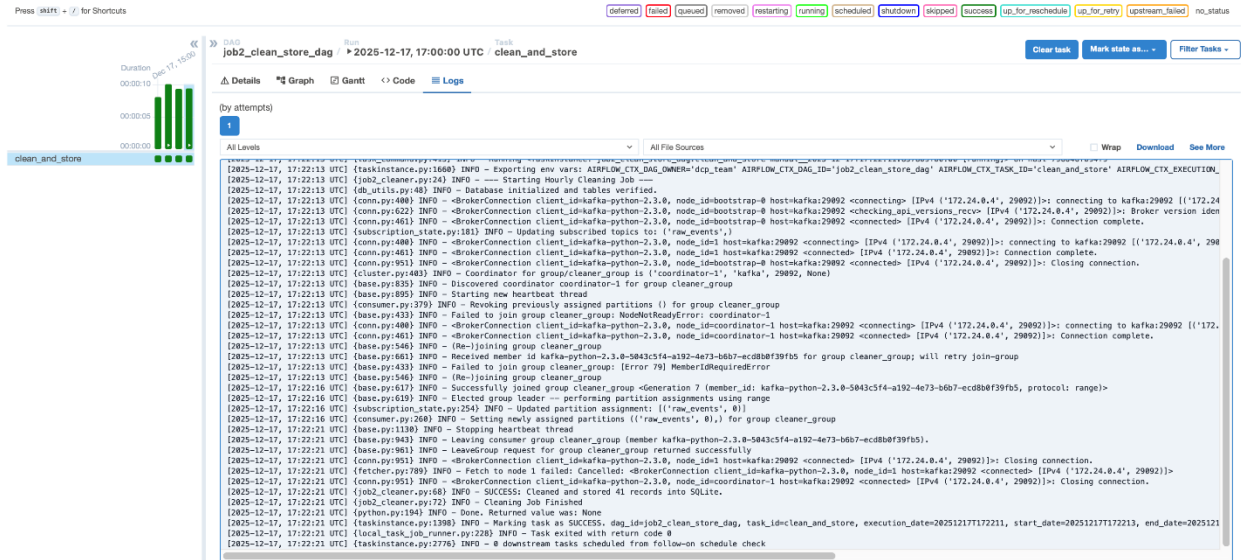
[2025-12-17, 17:23:16 UTC] (conn.py:951) INFO - <BrokerConnection client_id=kafka-python-producer-3, node_id=1 host=kafka:29092 <connected> [IPv4 ('172.24.0.4', 29092)]>: Closing connection.

[2025-12-17, 17:23:16 UTC] (python.py:194) INFO - Done. Returned value was: None

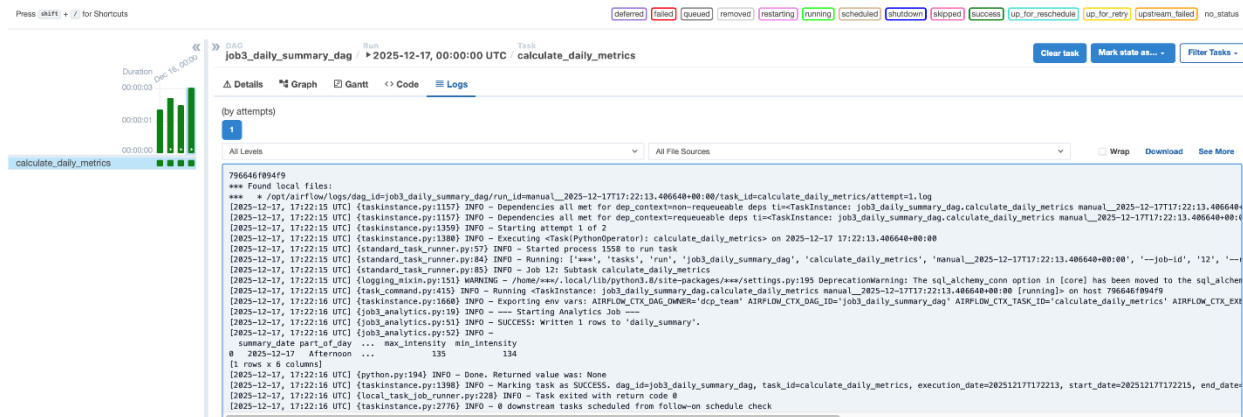
[2025-12-17, 17:23:16 UTC] (taskinstance.py:1398) INFO - Marking task as SUCCESS. dag_id=job1_ingestion_dag, task_id=start_ingestion_loop, execution_date=20251217T164500, start_date=20251217T170843, end_date=20251217, 17:23:16 UTC] (local_task_job_runner.py:228) INFO - Task exited with return code 0

[2025-12-17, 17:23:16 UTC] (taskinstance.py:2776) INFO - 0 downstream tasks scheduled from follow-on schedule check

6.3 Job 2 Logs (Cleaning)



6.4 Job 3 Logs (Analytics)



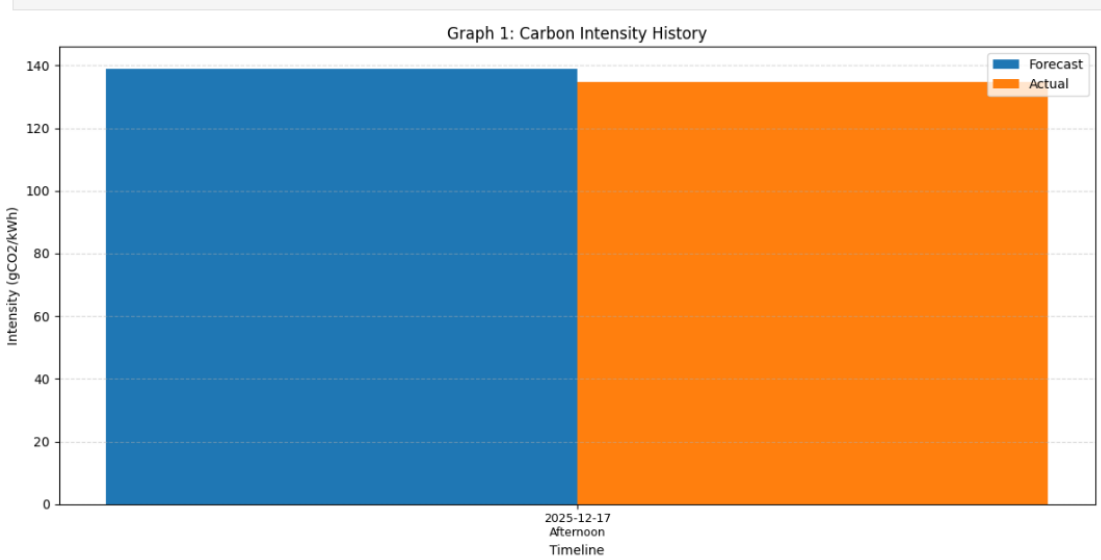
6.5 Database Content & Visualization

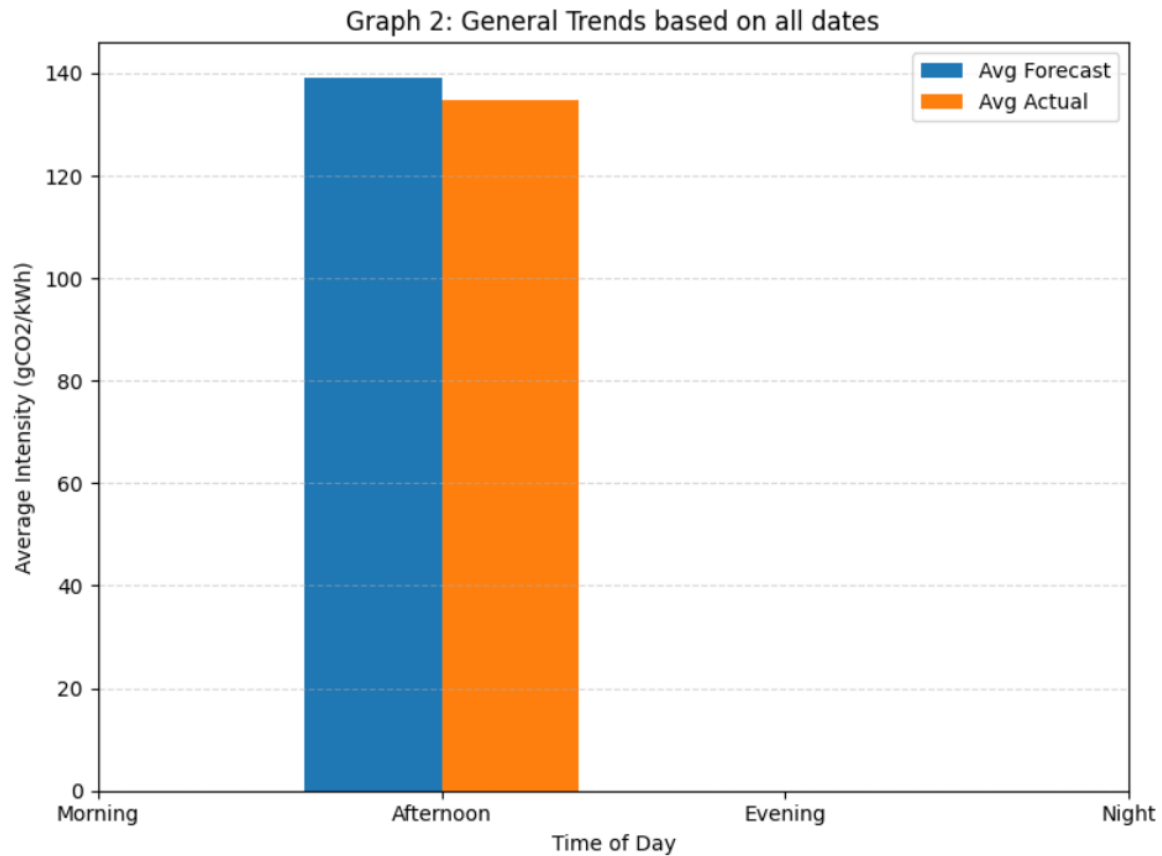
```
In [47]: query_events = "SELECT * FROM events ORDER BY id DESC LIMIT 10"
df_events = pd.read_sql_query(query_events, conn)

print("Last 10 Recorded Events (Cleaned Data)")
display(df_events)
```

Last 10 Recorded Events (Cleaned Data)

	id	ingestion_timestamp	forecast_intensity	actual_intensity	index_intensity	source	created_at
0	80	2025-12-17T17:22:15.201889	136	131	moderate	carbonintensity.org.uk	2025-12-17 17:22:21
1	79	2025-12-17T17:21:44.200970	136	131	moderate	carbonintensity.org.uk	2025-12-17 17:22:21
2	78	2025-12-17T17:21:13.305457	136	131	moderate	carbonintensity.org.uk	2025-12-17 17:22:21
3	77	2025-12-17T17:20:42.363494	136	131	moderate	carbonintensity.org.uk	2025-12-17 17:22:21
4	76	2025-12-17T17:20:11.319593	136	131	moderate	carbonintensity.org.uk	2025-12-17 17:22:21
5	75	2025-12-17T17:19:40.363466	136	131	moderate	carbonintensity.org.uk	2025-12-17 17:22:21
6	74	2025-12-17T17:19:09.158542	136	131	moderate	carbonintensity.org.uk	2025-12-17 17:22:21
7	73	2025-12-17T17:18:38.321542	136	131	moderate	carbonintensity.org.uk	2025-12-17 17:22:21
8	72	2025-12-17T17:18:07.356356	136	131	moderate	carbonintensity.org.uk	2025-12-17 17:22:21
9	71	2025-12-17T17:17:36.220188	136	131	moderate	carbonintensity.org.uk	2025-12-17 17:22:21





7. Conclusion

In this project, we built an end-to-end data pipeline that fully meets the project requirements. The system reliably collects, processes, and stores frequently updated data using Apache Airflow, Kafka, and SQLite. The pipeline is resilient to temporary failures and ensures stable data flow from ingestion to analytics. The final analytical results provide clear insights into carbon intensity patterns during different parts of the day, demonstrating a practical real-world data engineering solution