

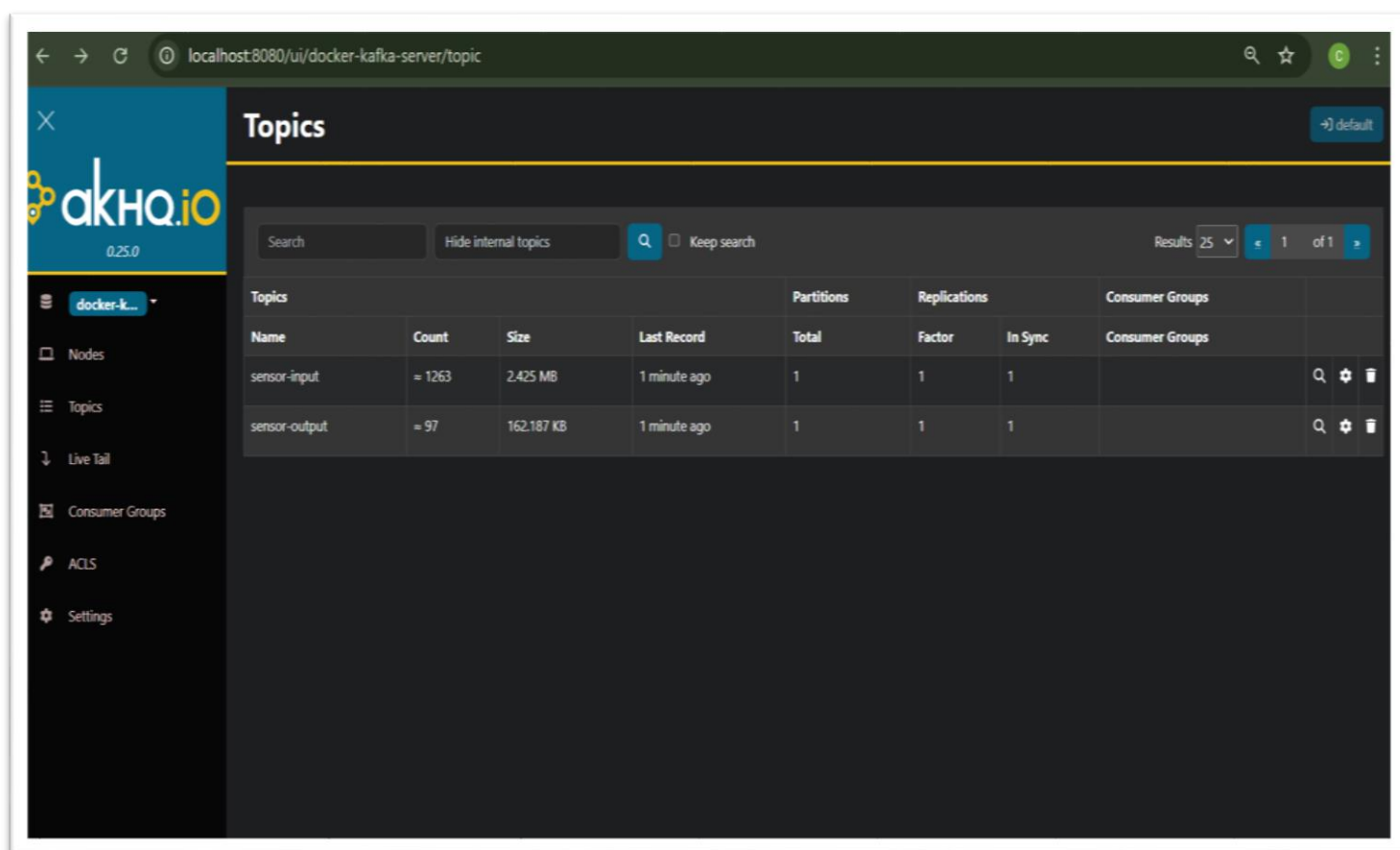
# Real-Time Sensor Data Processing

## Overview:

In this project, I developed a real-time data processing pipeline to handle sensor data streamed via Kafka. Each message includes a sensor ID, a value, and a timestamp in milliseconds. The goal was to compute the 1-minute average value for each sensor in real time. Using PySpark Structured Streaming, I read data from the Kafka input topic, applied window-based aggregation, and published the results to an output Kafka topic. Additionally, I persisted the results to MongoDB for further analysis. This solution demonstrates scalable, fault-tolerant stream processing suitable for real-world IoT applications.

## Visual Output Overview:

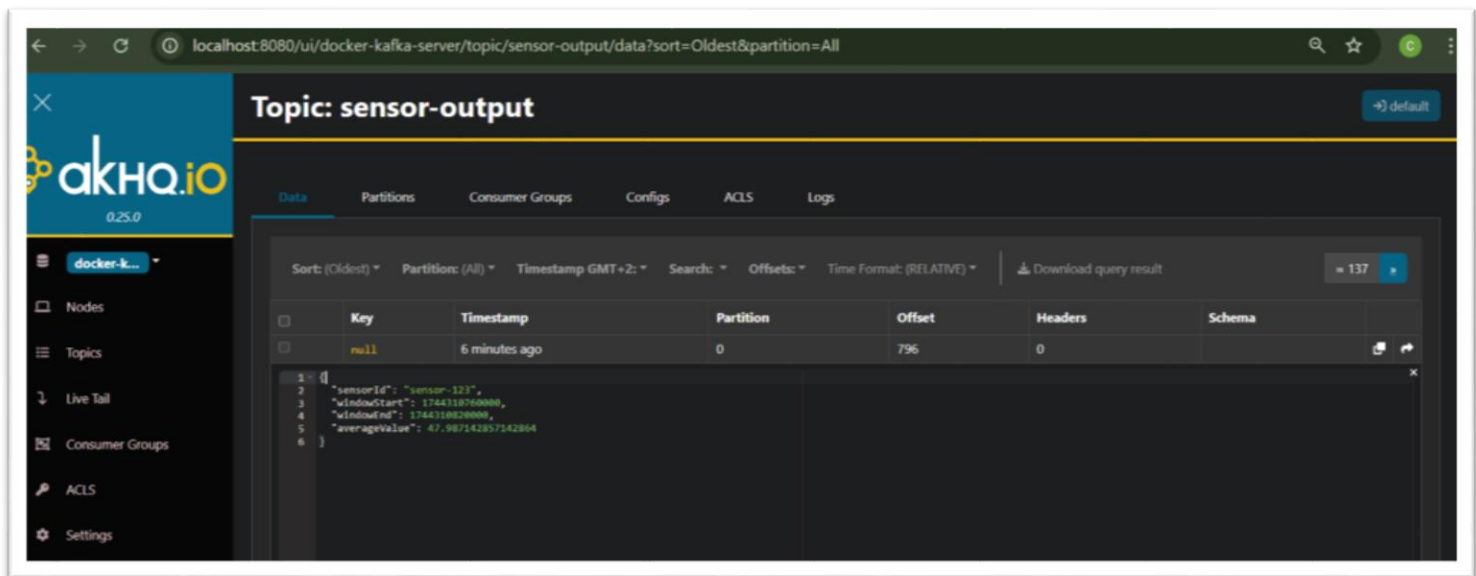
1. The sensor input data was read, processed according to the defined aggregation logic, and the results were sent to the sensor output.



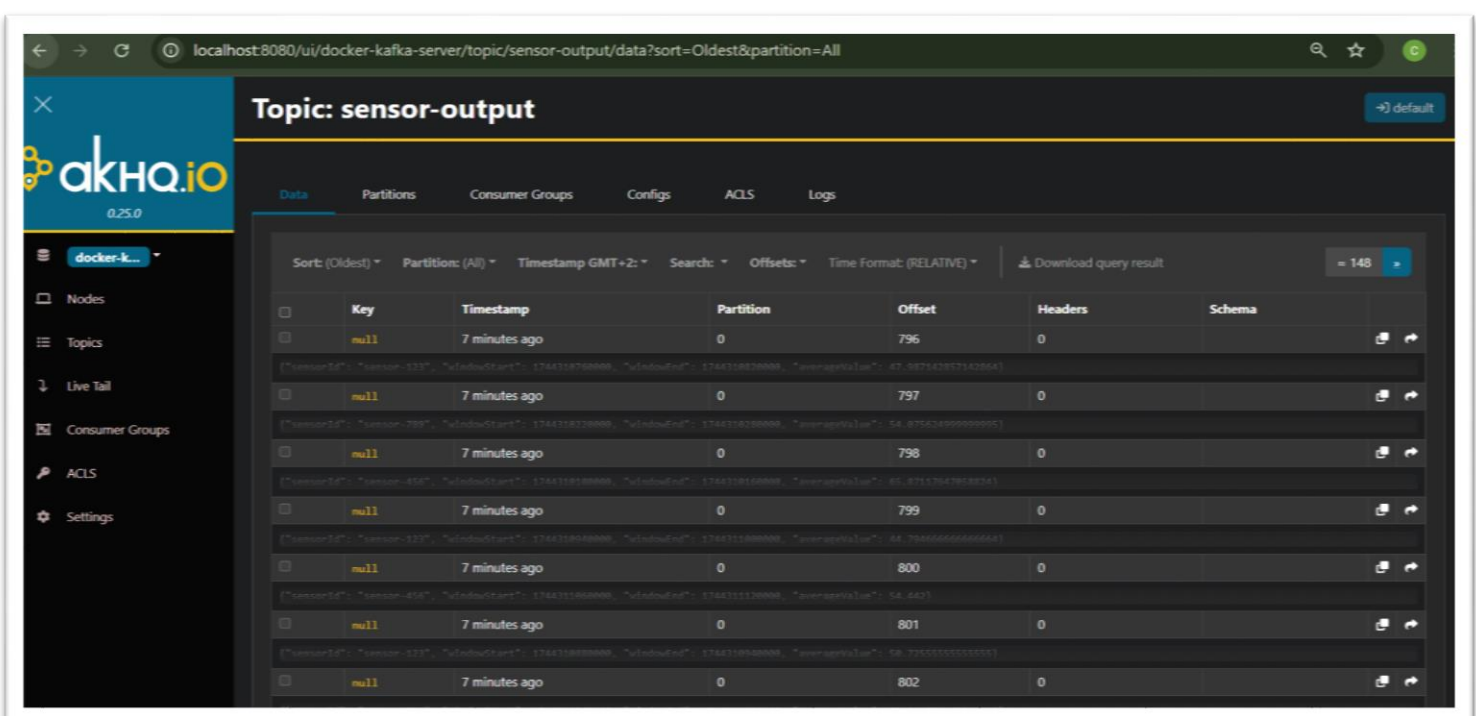
- The sensor-output topic contains the final processed data from the sensors. Each message is in JSON format and shows the average value of a sensor within a specific 1-minute time window.

#### Description of Fields:

- sensorId:** The ID of the sensor.
- windowStart:** Start time of the 1-minute window.
- windowEnd:** End time of the 1-minute window.
- averageValue:** The average of all readings from that sensor during the time window.



- This screenshot shows a list of JSON messages published to the sensor-output topic. Each message represents the average sensor reading within a specific 1-minute time window. Multiple such messages are generated in real time, one for each sensor and time window



4. This screenshot displays a list of JSON messages printed in the Visual Studio Code terminal. This output confirms that the data is being processed correctly and that the results are ready to be sent to the sensor-output Kafka topic and stored in a database (MongoDB).

The screenshot shows the Visual Studio Code interface with the file `kafka_streaming.py` open. The script is configured to read streaming data from Kafka and process it. The terminal output shows the following JSON messages being inserted into MongoDB and sent to Kafka:

```
48125}
Inserted into MongoDB: {"sensorId": "sensor-123", "windowStart": 1744311300000, "windowEnd": 1744311360000, "averageValue": 50.148125, "_id": ObjectId('67f8148544721a40c0c5b724')}
Sending to Kafka: {"sensorId": "sensor-321", "windowStart": 1744311360000, "windowEnd": 1744311420000, "averageValue": 73.08285714285715}
Inserted into MongoDB: {"sensorId": "sensor-321", "windowStart": 1744311360000, "windowEnd": 1744311420000, "averageValue": 73.08285714285715, "_id": ObjectId('67f8148544721a40c0c5b725')}
Sending to Kafka: {"sensorId": "sensor-123", "windowStart": 1744311360000, "windowEnd": 1744311420000, "averageValue": 62.087500000000006}
Inserted into MongoDB: {"sensorId": "sensor-123", "windowStart": 1744311360000, "windowEnd": 1744311420000, "averageValue": 62.087500000000006, "_id": ObjectId('67f8148544721a40c0c5b726')}
Sending to Kafka: {"sensorId": "sensor-789", "windowStart": 1744311360000, "windowEnd": 1744311420000, "averageValue": 55.194666666666666}
Inserted into MongoDB: {"sensorId": "sensor-789", "windowStart": 1744311360000, "windowEnd": 1744311420000, "averageValue": 55.194666666666666}
```

5. This screenshot shows a successful MongoDB connection established within Visual Studio Code. The data from the sensor-output Kafka topic has been written to the MongoDB database named `sensor_data_db`.

The screenshot shows the Visual Studio Code interface with the file `kafka_streaming.py` open. The terminal output shows the MongoDB connection and data retrieval process:

```
For mongosh info see: https://www.mongodb.com/docs/mongodb-shell/

-----
The server generated these startup warnings when booting
2025-04-10T11:50:10.683+02:00: Access control is not enabled for the database. Read and write access to data and configuration is unrestricted
-----

test> use sensor_data_db // switch to your DB
... db.sensor_readings.find().pretty() // view the inserted records
...
switched to db sensor_data_db
sensor_data_db> db.sensor_readings.find().pretty()
[
  {
    _id: ObjectId('67f79a4d901bf682f2a98a23'),
    sensorId: 'sensor-456',
    windowStart: Long('1744215120000'),
    windowEnd: Long('1744215180000'),
    averageValue: 50.024166666666666
  },
  {
    _id: ObjectId('67f79a4d901bf682f2a98a24'),
    sensorId: 'sensor-789',
    windowStart: Long('1744215120000'),
    windowEnd: Long('1744215180000'),
    averageValue: 54.20823529411765
  },
]
```

6. This screenshot shows the MongoDB Compass UI connected to the database `sensor_data_db`, specifically viewing the `sensor_readings` collection.

