



PES UNIVERSITY
100 feet Ring Road, BSK 3rd Stage
Bengaluru 560085 INDIA

Department of Computer Science and Engineering
B. Tech. CSE – 6th Semester
Jan – May 2024

UE21CS343BB3
DATABASE TECHNOLOGIES (DBT)

PROJECT REPORT
on

Stream and Batch Processing of Water Quality

Submitted by : Team #: 512_538_555_902

Sahana Parasuram	PES1UG21CS512	6I	Sharath M S	PES1UG21CS555	6I
Saransh Mehta	PES1UG21CS538	6I	Aditya Vishwanatha	PES1UG21CS902	6I

Class of Prof. Raghu B. A.

Stream and Batch Processing of Water Quality*Table of Contents*

Sl. No	Topic	Page No.
1.	Introduction Problem Description Solution Architecture	3
2.	Installation of Software Streaming Apps/Tools Database	4
3.	Input Data a. Source/s b. Description	7
4.	Streaming Mode Experiment a. Description b. Windows c. Results	8
5.	Batch Mode Experiment a. Description b. Data Size c. Results	11
6.	Comparison of Streaming & Batch Modes a. Results and Discussion	13
7.	Conclusion	14
8.	References	15

1. Introduction

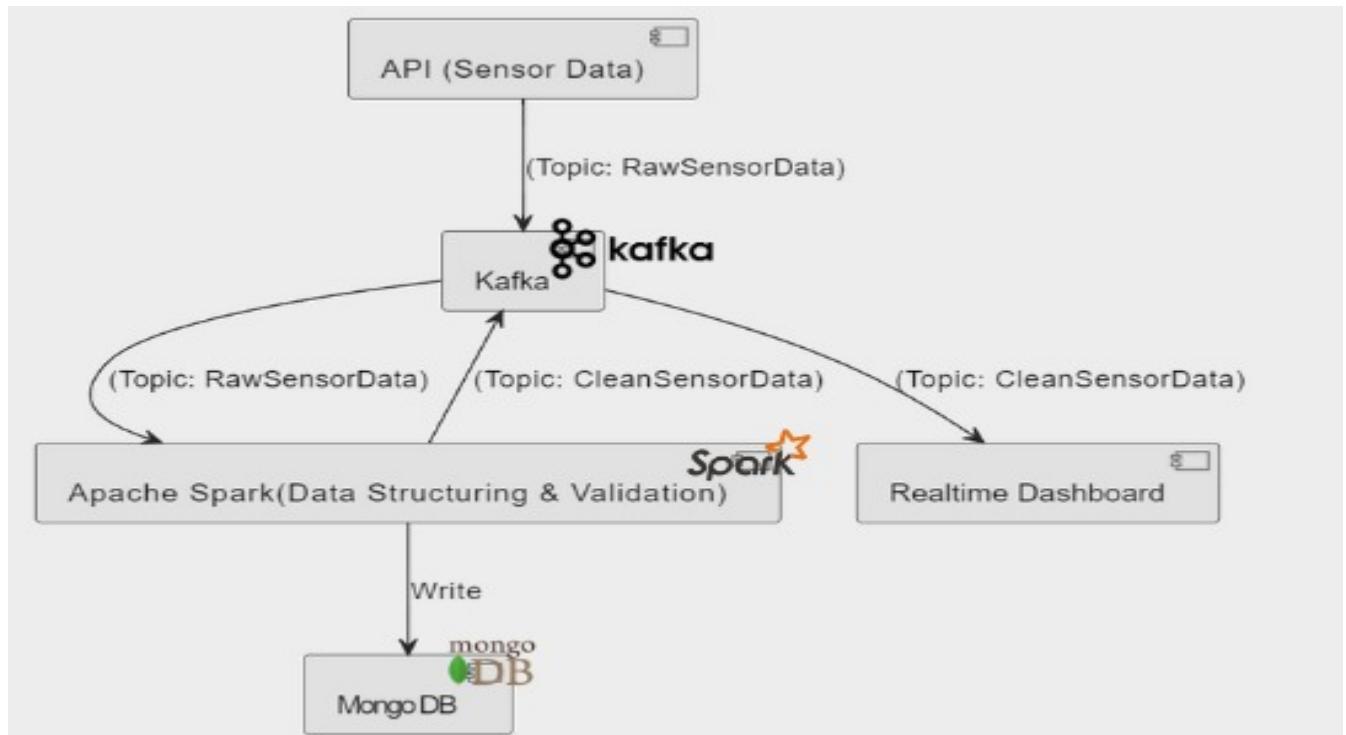
This project develops a real-time data pipeline utilizing Apache Kafka and Apache Spark to handle and analyze water quality data. The primary goal is to demonstrate effective real-time data processing compared to traditional batch processing, using simulated sensor data to model water quality metrics.

Problem Description

Water quality monitoring is vital for environmental health and public safety. Traditional batch processing methods can delay the action needed to address pollution or contamination events. This project explores a real-time processing solution to provide immediate insights into water quality data.

Solution Architecture

- API: Simulates water quality sensor data, mirroring actual sensor outputs.
- Streaming Apps/Tools: Uses Apache Kafka for real-time data ingestion and Apache Spark for on-the-fly data processing.
- Repository/Database: MongoDB for data persistence post-analysis, enabling further historical analysis.



Architecture Diagram

2. Installation of Software

- **Apache Kafka (3.7.0)** : a distributed streaming platform that makes it possible to create streaming apps and real-time data pipelines. We utilized this as a publisher and for data streaming.
https://downloads.apache.org/kafka/3.7.0/kafka_2.12-3.7.0.tgz
- **Apache Spark Streaming (3.5.1)** : Built on top of the Apache Spark platform, a real-time processing engine enables real-time data stream processing. permits batch processing as well.
<https://www.apache.org/dyn/closer.lua/spark/spark-3.5.1/spark-3.5.1-bin-hadoop3.tgz>
- **Pyspark (3.5.1)** : A streaming engine designed specifically for Python and its related features, similar to Apache Kafka.
TO INSTALL : pip install pyspark=="3.5.1"

Stream and Batch Processing of Water Quality

- **DBMS Used :** MongoDB - 7.0.8

<https://www.mongodb.org/static/pgp/server-7.0.asc>

Streaming Apps/Tools Used

Apache Kafka:

Apache Kafka facilitates the publishing and subscribing of streams of records, which enables real-time data processing. Kafka provides fault tolerance, scalability, and high throughput, making it suitable for building real-time data pipelines.

Apache Spark Streaming:

Apache Spark Streaming is used for processing real-time data streams from Kafka. It allows for scalable, high-throughput, fault-tolerant stream processing of live data streams. Spark Streaming enables the application of complex algorithms and analytics on real-time data.

MongoDB:

MongoDB is a NoSQL database used for storing real-time sensor data processed by Spark Streaming.

Zookeeper:

In this project, Zookeeper is used for managing configurations and coordination in the Kafka ecosystem.

The above figure represents starting the Kafka server using the specified configuration file.

Stream and Batch Processing of Water Quality

Database

Mongodb is utilized in our project as a NoSQL database, which stores data in JSON-like documents with dynamic schemas (BSON), making the integration of data in certain types of applications easier and faster. In our project, MongoDB stores real-time sensor data, such as water temperature, in a collection within the 'RealTimeDB' database. This approach enables efficient data storage and fast retrieval capabilities, essential for performing real-time analytics and batch processing. The use of aggregation queries in MongoDB facilitates complex analytics, such as calculating average temperatures from accumulated data, directly within the database, thereby optimizing performance and scalability.

```
sharath@sharath-Dell-G15-5515:~/opt/spark
[RawData]: '2024-04-21T23:28:15.037687 6.80 109.57 5.28 Montrose_Beach 872939', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 15, 37687), 'WaterTemperature': 6.80, 'Turbidity': 109.57, 'BatteryLife': 5.28, 'Beach': 'Montrose_Beach', 'MeasurementID': 872939, '_id': ObjectId('662553b7c87cf0e6a5097ab1')}
[RawData]: '2024-04-21T23:28:16.041511 22.91 1559.22 6.94 Montrose_Beach 196615', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 16, 41511), 'WaterTemperature': 22.91, 'Turbidity': 1559.22, 'BatteryLife': 6.94, 'Beach': 'Montrose_Beach', 'MeasurementID': 196615, '_id': ObjectId('662553b7c87cf0e6a5097ab2')}
[RawData]: '2024-04-21T23:28:17.645512 4.39 852.70 11.45 Montrose_Beach 697042', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 17, 45512), 'WaterTemperature': 4.39, 'Turbidity': 852.70, 'BatteryLife': 11.45, 'Beach': 'Montrose_Beach', 'MeasurementID': 697042, '_id': ObjectId('662553b7c87cf0e6a5097ab3')}
[RawData]: '2024-04-21T23:28:19.03588 16.04 1409.51 16.04 Montrose_Beach 86099', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 19, 03588), '_id': ObjectId('662553b7c87cf0e6a5097ab4')}
[RawData]: '2024-04-21T23:28:19.652128 12.38 642.81 7.63 Montrose_Beach 989810', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 19, 52128), 'WaterTemperature': 12.38, 'Turbidity': 642.81, 'BatteryLife': 7.63, 'Beach': 'Montrose_Beach', 'MeasurementID': 989810, '_id': ObjectId('662553b7c87cf0e6a5097ab5')}
[RawData]: '2024-04-21T23:28:20.055707 28.95 1498.41 11.25 Montrose_Beach 483223', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 20, 55707), 'WaterTemperature': 28.95, 'Turbidity': 1498.41, 'BatteryLife': 11.25, 'Beach': 'Montrose_Beach', 'MeasurementID': 483223, '_id': ObjectId('662553b7c87cf0e6a5097ab6')}
[RawData]: '2024-04-21T23:28:21.058848 16.12 1945.15 11.36 Montrose_Beach 119813', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 21, 58848), 'WaterTemperature': 16.12, 'Turbidity': 1945.15, 'BatteryLife': 11.36, 'Beach': 'Montrose_Beach', 'MeasurementID': 119813, '_id': ObjectId('662553b7c87cf0e6a5097ab7')}
[RawData]: '2024-04-21T23:28:22.062941 21.52 1065.35 7.98 Montrose_Beach 522604', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 22, 62941), 'WaterTemperature': 21.52, 'Turbidity': 1065.35, 'BatteryLife': 7.98, 'Beach': 'Montrose_Beach', 'MeasurementID': 522604, '_id': ObjectId('662553b7c87cf0e6a5097ab8')}
[RawData]: '2024-04-21T23:28:23.066672 31.07 627.84 10.24 Montrose_Beach 934570', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 23, 66672), 'WaterTemperature': 31.07, 'Turbidity': 627.84, 'BatteryLife': 10.24, 'Beach': 'Montrose_Beach', 'MeasurementID': 934570, '_id': ObjectId('662553b7c87cf0e6a5097ab9')}
[RawData]: '2024-04-21T23:28:29.070658 30.29 917.77 8.64 Montrose_Beach 91732', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 24, 70658), 'WaterTemperature': 30.29, 'Turbidity': 917.77, 'BatteryLife': 8.64, 'Beach': 'Montrose_Beach', 'MeasurementID': 91732, '_id': ObjectId('662553b7c87cf0e6a5097ab0')}
[RawData]: '2024-04-21T23:28:25.074274 26.01 152.33 10.09 Montrose_Beach 641978', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 25, 74274), 'WaterTemperature': 26.01, 'Turbidity': 152.33, 'BatteryLife': 10.09, 'Beach': 'Montrose_Beach', 'MeasurementID': 641978, '_id': ObjectId('662553b7c87cf0e6a5097ab1')}
[RawData]: '2024-04-21T23:28:28.077589 20.94 1431.83 12.38 Montrose_Beach 769502', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 26, 77589), 'WaterTemperature': 20.94, 'Turbidity': 1431.83, 'BatteryLife': 12.38, 'Beach': 'Montrose_Beach', 'MeasurementID': 769502, '_id': ObjectId('662553b7c87cf0e6a5097ab2')}
[RawData]: '2024-04-21T23:28:28.082926 7.71 762.5 12.23 Montrose_Beach 887547', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 27, 82926), 'WaterTemperature': 7.71, 'Turbidity': 762.5, 'BatteryLife': 12.23, 'Beach': 'Montrose_Beach', 'MeasurementID': 887547, '_id': ObjectId('662553b7c87cf0e6a5097ab3')}
[RawData]: '2024-04-21T23:28:28.085670 11.27 836.99 11.04 Montrose_Beach 503087', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 28, 86670), 'WaterTemperature': 11.27, 'Turbidity': 826.99, 'BatteryLife': 11.04, 'Beach': 'Montrose_Beach', 'MeasurementID': 503087, '_id': ObjectId('662553b7c87cf0e6a5097ab4')}
[RawData]: '2024-04-21T23:28:29.090117 52.82 1168.82 12.26 Montrose_Beach 317770', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 29, 90117), 'WaterTemperature': 27.52, 'Turbidity': 1168.82, 'BatteryLife': 12.26, 'Beach': 'Montrose_Beach', 'MeasurementID': 317770, '_id': ObjectId('662553b7c87cf0e6a5097ab5')}
[RawData]: '2024-04-21T23:28:30.093617 15.89 766.31 7.28 Montrose_Beach 436856', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 30, 93617), 'WaterTemperature': 15.89, 'Turbidity': 766.31, 'BatteryLife': 7.28, 'Beach': 'Montrose_Beach', 'MeasurementID': 436856, '_id': ObjectId('662553b7c87cf0e6a5097ac1')}
[RawData]: '2024-04-21T23:28:31.097155 4.69 1562.38 10.71 Montrose_Beach 737766', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 31, 97155), 'WaterTemperature': 4.69, 'Turbidity': 1562.38, 'BatteryLife': 10.68, 'Beach': 'Montrose_Beach', 'MeasurementID': 737766, '_id': ObjectId('662553b7c87cf0e6a5097ac2')}
[RawData]: '2024-04-21T23:28:32.101322 21.56 341.75 9.48 Montrose_Beach 323422', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 32, 101322), 'WaterTemperature': 21.56, 'Turbidity': 341.75, 'BatteryLife': 9.48, 'Beach': 'Montrose_Beach', 'MeasurementID': 323422, '_id': ObjectId('662553b7c87cf0e6a5097ac3')}
[RawData]: '2024-04-21T23:28:33.105536 1.09 381.95 11.95 Montrose_Beach 747498', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 33, 105536), 'WaterTemperature': 0.19, 'Turbidity': 381.95, 'BatteryLife': 11.95, 'Beach': 'Montrose_Beach', 'MeasurementID': 747498, '_id': ObjectId('662553b7c87cf0e6a5097ac4')}
[RawData]: '2024-04-21T23:28:34.108661 20.29 119.62 12.27 Montrose_Beach 44289', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 34, 108661), 'WaterTemperature': 20.29, 'Turbidity': 119.62, 'BatteryLife': 12.27, 'Beach': 'Montrose_Beach', 'MeasurementID': 44289, '_id': ObjectId('662553b7c87cf0e6a5097ac5')}
[RawData]: '2024-04-21T23:28:35.113792 22.74 978.11 6.8 Montrose_Beach 945039', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 35, 113792), 'WaterTemperature': 22.74, 'Turbidity': 978.11, 'BatteryLife': 6.8, 'Beach': 'Montrose_Beach', 'MeasurementID': 945039, '_id': ObjectId('662553b7c87cf0e6a5097ac6')}
[RawData]: '2024-04-21T23:28:36.117933 26.07 1064.18 8.71 Montrose_Beach 269681', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 36, 117933), 'WaterTemperature': 26.07, 'Turbidity': 1064.13, 'BatteryLife': 7.1, 'Beach': 'Montrose_Beach', 'MeasurementID': 269681, '_id': ObjectId('662553b7c87cf0e6a5097ac7')}
[RawData]: '2024-04-21T23:28:37.120809 20.69 956.35 10.23 Montrose_Beach 90535', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 37, 120809), 'WaterTemperature': 20.69, 'Turbidity': 956.35, 'BatteryLife': 10.23, 'Beach': 'Montrose_Beach', 'MeasurementID': 90535, '_id': ObjectId('662553b7c87cf0e6a5097ac8')}
[RawData]: '2024-04-21T23:28:39.126163 29.63 341.85 8.04 Montrose_Beach 788486', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 39, 126163), 'WaterTemperature': 29.63, 'Turbidity': 341.85, 'BatteryLife': 8.04, 'Beach': 'Montrose_Beach', 'MeasurementID': 788486, '_id': ObjectId('662553b7c87cf0e6a5097ac9')}
[RawData]: '2024-04-21T23:28:39.130873 20.96 1166.39 9.45 Montrose_Beach 156693', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 39, 130873), 'WaterTemperature': 20.96, 'Turbidity': 1166.39, 'BatteryLife': 9.45, 'Beach': 'Montrose_Beach', 'MeasurementID': 156693, '_id': ObjectId('662553b7c87cf0e6a5097ac9')}
[RawData]: '2024-04-21T23:28:40.134188 1.82 652.08 6.74 Montrose_Beach 371682', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 40, 134188), 'WaterTemperature': 1.82, 'Turbidity': 652.08, 'BatteryLife': 6.74, 'Beach': 'Montrose_Beach', 'MeasurementID': 371682, '_id': ObjectId('662553b7c87cf0e6a5097ac9')}
[RawData]: '2024-04-21T23:28:41.138678 28.79 1469.65 11.46 Montrose_Beach 452497', 'TimeStamp': datetime.datetime(2024, 4, 21, 23, 28, 41, 138678), 'WaterTemperature': 28.79, 'Turbidity': 1469.65, 'BatteryLife': 11.46, 'Beach': 'Montrose_Beach', 'MeasurementID': 452497, '_id': ObjectId('662553b7c87cf0e6a5097ac9')}
```

The above figure represents the data stored in Mongodb after streaming.

Kafka Topics

RawSensorData:

This topic is used to store raw sensor data obtained from the API. Data published to this topic

includes information such as timestamp, water temperature, turbidity, battery life, beach, and measurement ID. Raw sensor data is streamed into this topic by the push_data_to_kafka.py script.

ValidatedSensorData:

This topic acts as an intermediate step for validated sensor data before it's persisted into MongoDB. After the sensor data is received from the RawSensorData topic, it undergoes validation to ensure its integrity and correctness. Validated sensor data is published to this topic by the structure_validate_store.py script.

CleanSensorData:

This topic stores the cleaned and validated sensor data after it's persisted into MongoDB. Once the data is validated and stored in MongoDB, it's published to this topic by the structure_validate_store.py script. Cleaned sensor data from this topic can be consumed by downstream systems, such as real-time dashboards, for visualization or further analysis.

3. Input Data

- a. Sources :** Beach water quality API from the sensors in the water at beaches along Chicago's lake Michigan

<https://data.world/cityofchicago/beach-water-quality-automated-sensors>

```
^Csharath@sharath-Dell-G15-5515:~/Desktop/Real-Time-Data-Pipeline-Using-Kafka-and-Spark-master$ python3 sensor.py
* Serving Flask app 'sensor'
* Debug mode: off
WARNING: This is a development server. Do not use it in a production deployment. Use a production WSGI server instead.
* Running on all addresses (0.0.0.0)
* Running on http://127.0.0.1:3030
* Running on http://192.168.29.30:3030
Press CTRL+C to quit
127.0.0.1 - - [21/Apr/2024 22:38:14] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:15] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:16] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:17] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:18] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:19] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:20] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:21] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:22] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:23] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:24] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:25] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:26] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:27] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:28] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:29] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:30] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:31] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:32] "GET /sensordata HTTP/1.1" 200 -
127.0.0.1 - - [21/Apr/2024 22:38:33] "GET /sensordata HTTP/1.1" 200 -
```

b. Description : The Chicago Park District maintains sensors in the water at beaches along Chicago's Lake Michigan lakefront. These sensors generally capture the indicated measurements hourly while the sensors are in operation during the summer. During other seasons and at some other times, information from the sensors may not be available. See <https://data.cityofchicago.org/d/k7hf-8y75> for a dataset with land-based weather measurements at selected beaches. The sensor locations are listed at <https://data.cityofchicago.org/d/g3ip-u8rb>. Please note that sensor locations change with the Park District's operational needs, primarily related to water quality.

4. Streaming Mode Experiment

a. Description : This experiment is designed to demonstrate the system's ability to handle and analyze continuous data streams using a combination of Apache Kafka and Apache Spark Streaming. Kafka serves as the initial ingestion point, reliably capturing live data streams that simulate environmental sensor outputs such as water temperature and turbidity. Once ingested, data is forwarded to Spark Streaming, which is configured to process this data in real time. The integration of these two technologies ensures seamless data flow and immediate processing, which is crucial for scenarios requiring rapid analytical feedback and decision-making.

b. Windows :

Time Windows: The streaming data is processed using a windowing technique, where data is grouped into fixed windows of one minute each. This approach allows the system to perform time-based aggregations and analyses, such as calculating the average temperature or detecting outliers and anomalies within each window. The use of windowing not only simplifies the processing of continuous streams by breaking down the infinite stream into manageable chunks but also enables the application of complex analytical models on a per-window basis.

Window Operations: Within each window, Spark Streaming executes several functions—aggregating temperature data to compute averages and applying statistical

Stream and Batch Processing of Water Quality

models to detect sudden changes that could indicate environmental anomalies. These operations are critical for applications like environmental monitoring, where timely detection of data trends and anomalies can lead to prompt action to mitigate potential hazards.

```
❸ dashboard.py
1  import random
2  import pandas as pd
3  from bokeh.driving import count
4  from bokeh.models import ColumnDataSource
5  from kafka import KafkaConsumer
6  from bokeh.plotting import curdoc, figure
7  from bokeh.models import DatetimeTickFormatter
8  from bokeh.models.widgets import Div
9  from bokeh.layouts import column, row
10 import ast
11 import time
12 import pytz
13 from datetime import datetime
14
15 tz = pytz.timezone('Asia/Calcutta')
16
17
18 UPDATE_INTERVAL = 1000
19 ROLLOVER = 10 # Number of displayed data points
20
21
22 source = ColumnDataSource({"x": [], "y": []})
23 consumer = KafkaConsumer('CleanSensorData', auto_offset_reset='earliest', bootstrap_servers=['localhost:9092'], consumer_timeout_ms=1000)
24 div = Div(
25     text='',
26     width=120,
27     height=35
28 )
29
30
31 @count()
32 def update(x):
33     for msg in consumer:
34         msg_value = msg
35         break
36
37 values = ast.literal_eval(msg_value.value.decode("utf-8"))
38 timestamp_iso = values["TimeStamp"]["$date"]
39
40 # Convert ISO 8601 timestamp to a datetime object
41 timestamp_datetime = datetime.fromisoformat(timestamp_iso.replace('Z', '+00:00'))
42
43 # Convert datetime to POSIX timestamp
44 timestamp_sec = timestamp_datetime.timestamp()
45
```

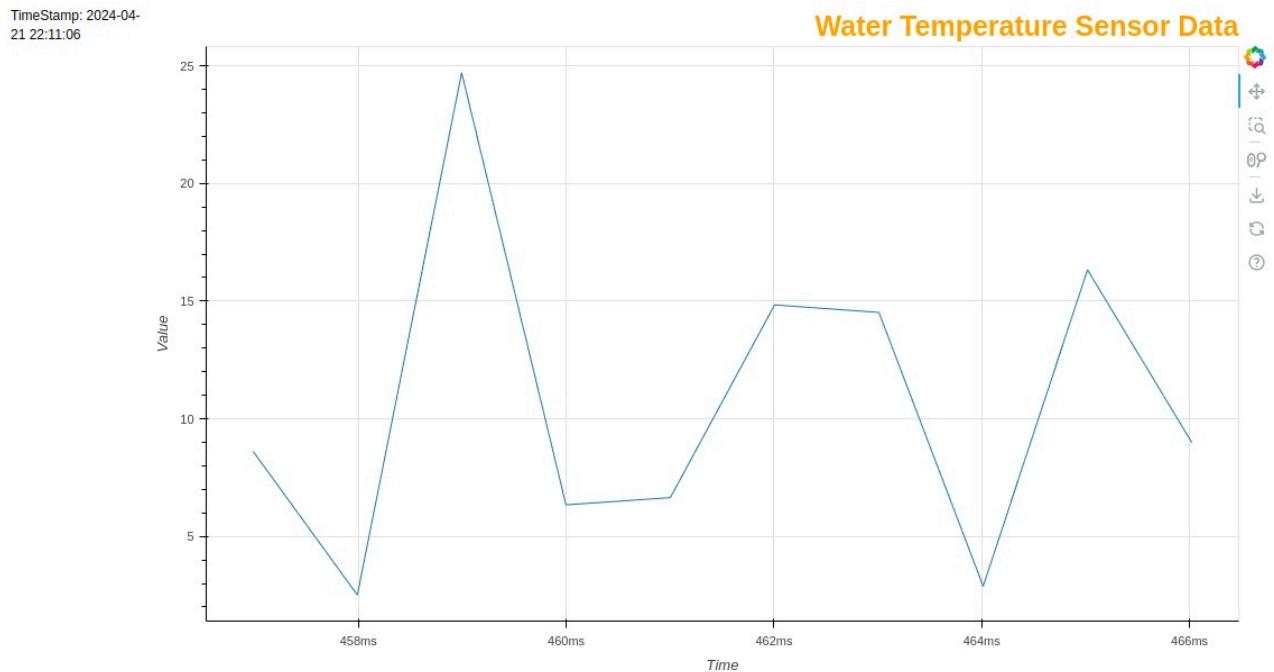


```
❹ dashboard.py
32     def update(x):
33
34         # Convert ISO 8601 timestamp to a datetime object
35         timestamp_datetime = datetime.fromisoformat(timestamp_iso.replace('Z', '+00:00'))
36
37         # Convert datetime to POSIX timestamp
38         timestamp_sec = timestamp_datetime.timestamp()
39
40         div.text = "TimeStamp: " + str(timestamp_datetime.strftime("%Y-%m-%d %H:%M:%S"))
41
42         y = values['WaterTemperature']
43         print(y)
44
45         source.stream({"x": [timestamp_sec], "y": [y]}, ROLLOVER)
46
47         p = figure(title="Water Temperature Sensor Data", x_axis_type="datetime", width=1000)
48
49         p.line("x", "y", source=source)
50
51         p.xaxis.formatter=DatetimeTickFormatter(hourmin = ['%H:%M'])
52         p.xaxis.axis_label = 'Time'
53         p.yaxis.axis_label = 'Value'
54         p.title.align = "right"
55         p.title.text_color = "orange"
56         p.title.text_font_size = "25px"
57
58         doc = curdoc()
59         #doc.add_root(p)
60
61         doc.add_root(
62             row(children=[div,p])
63         )
64         doc.add_periodic_callback(update, UPDATE_INTERVAL)
65
66
67
68
69
70
```



This is dashboard.py. It shows how to consume data from Kafka and prepare it for real-time visualization

c. Results :



The provided graph represents the real-time fluctuations in water temperature as captured and processed by the streaming system. The X-axis indicates the time, seemingly in milliseconds, while the Y-axis represents the water temperature values.

Real-Time Visualization: The outcomes of the streaming data analysis are dynamically visualized using a real-time dashboard. This visualization includes trend lines for water temperature over time, highlighted points where anomalies are detected, and alerts that are triggered based on specific threshold conditions.

System Responsiveness: The experiment highlights the system's responsiveness, as changes in the data stream are quickly reflected in the visualizations. For instance, a sudden spike in temperature due to an environmental factor is immediately captured, processed, and displayed, showcasing the potential of the system to support real-time environmental monitoring.

Potential for Immediate Environmental Monitoring: The ability to monitor environmental conditions in real time has significant implications. For industries such as

water quality management, having immediate access to data about changes in water conditions can be crucial for pollution control and ecological conservation. The system's capability to provide continuous and instant analysis makes it an invaluable tool in such contexts.

5. Batch Mode Experiment

a. Description :

The batch mode experiment in this project is designed to process accumulated data stored in MongoDB using Apache Spark, one of the leading platforms for large-scale data processing. This setup allows for periodic processing of large data batches, which are structured similarly to the real-time streams but are processed in intervals (e.g., daily or hourly). By performing the same types of aggregations as those done in the streaming experiment, such as calculating average temperatures or detecting outliers, the batch mode experiment provides a basis for directly comparing performance and outcomes between the two processing modes.

Spark Integration with MongoDB: Apache Spark is configured to read data directly from MongoDB using the MongoDB Connector for Spark. This integration facilitates seamless data transfer into Spark's dataframes, which are then used for further analysis.

Data Aggregation and Analysis: Spark employs sophisticated data processing algorithms to efficiently aggregate and analyze data. The experiment focuses on calculating average temperatures and other statistical metrics from sensor data, leveraging Spark's ability to handle complex computations over large datasets.

b. Data Size :

Handling Large Datasets: The batch processing experiment is specifically designed to manage and process batches containing large records. This scale is chosen to demonstrate Spark's capabilities in handling substantial volumes of data without compromising on processing speed or accuracy.

Data Batching Strategy: The data is grouped into batches based on either time intervals or record count, ensuring consistent and manageable batch sizes for processing. This strategy helps in optimizing the processing workload and aligning it with the system's

Stream and Batch Processing of Water Quality

capacity.

```

❸ structure_validate_store.py
1  import json
2  from bson import json_util
3  from dateutil import parser
4  from pyspark import SparkContext
5  from kafka import KafkaConsumer, KafkaProducer
6
7  #Mongo DB
8  from pymongo import MongoClient
9
10 client = MongoClient('localhost', 27017)
11 db = client['RealTimeDB']
12 collection = db['RealTimeCollection']
13
14 def timestamp_exist(TimeStamp):
15     if collection.count_documents({'TimeStamp': {'$eq': TimeStamp}}) > 0:
16         return True
17     else:
18         return False
19
20 def structure_validate_data(msg):
21     data_dict = {}
22
23     # Create RDD
24     rdd = sc.parallelize([msg.value.decode("utf-8").split()])
25     data_dict["RawData"] = str(msg.value.decode("utf-8"))
26
27     # Data validation and create json data dict
28     try:
29         data_dict["TimeStamp"] = parser.isoparse(rdd.collect()[0])
30     except Exception as error:
31         data_dict["TimeStamp"] = "Error"
32
33     try:
34         data_dict["WaterTemperature"] = float(rdd.collect()[1])
35         if data_dict["WaterTemperature"] > 99 or data_dict["WaterTemperature"] < -10:
36             data_dict["WaterTemperature"] = "Sensor Malfunctions"
37     except Exception as error:
38         data_dict["WaterTemperature"] = "Error"
39
40     try:
41         data_dict["Turbidity"] = float(rdd.collect()[2])
42         if data_dict["Turbidity"] > 5000:
43             data_dict["Turbidity"] = "Sensor Malfunctions"
44     except Exception as error:
45
❹ structure_validate_store.py
20     def structure_validate_data(msg):
21         try:
22             if data_dict["Turbidity"] > 5000:
23                 data_dict["Turbidity"] = "Sensor Malfunctions"
24         except Exception as error:
25             data_dict["Turbidity"] = "Error"
26
27         try:
28             data_dict["BatteryLife"] = float(rdd.collect()[3])
29         except Exception as error:
30             data_dict["BatteryLife"] = "Error"
31
32         try:
33             data_dict["Beach"] = str(rdd.collect()[4])
34         except Exception as error:
35             data_dict["Beach"] = "Error"
36
37         try:
38             data_dict["MeasurementID"] = int(str(rdd.collect()[5]).replace("Beach",""))
39         except Exception as error:
40             data_dict["MeasurementID"] = "Error"
41
42         return data_dict
43
44     sc = SparkContext.getOrCreate()
45     sc.setLogLevel("WARN")
46
47     consumer = KafkaConsumer('RawSensorData', auto_offset_reset='earliest', bootstrap_servers=['localhost:9092'], consumer_timeout_ms=1000)
48     producer = KafkaProducer(bootstrap_servers=['localhost:9092'])
49
50     for msg in consumer:
51         if msg.value.decode("utf-8") != "Error in Connection":
52             data = structure_validate_data(msg)
53
54             if timestamp_exist(data['TimeStamp']) == False:
55                 # Push data to MongoDB
56                 collection.insert_one(data)
57                 # Push data to ValidatedSensorData Kafka topic
58                 producer.send("ValidatedSensorData", json.dumps(data, default=json_util.default).encode('utf-8'))
59                 # Push data to CleanSensorData Kafka topic
60                 producer.send("CleanSensorData", json.dumps(data, default=json_util.default).encode('utf-8'))
61
62             print(data)
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83

```

This is structure_validate_store.py.

c. Result :

```
sharath@sharath-Dell-G15-5515:~/Desktop/Real-Time-Data-Pipeline-Using-Kafka-and-Spark-master$ python3 batch_mode.py
Average Water Temperature: 15.85
Batch Query Execution Time: 0.0046062 seconds
○ sharath@sharath-Dell-G15-5515:~/Desktop/Real-Time-Data-Pipeline-Using-Kafka-and-Spark-master$ █
```

The batch mode processing was performed on a dataset using Apache Spark to calculate aggregated metrics, such as the average water temperature. The results of the execution yielded an average water temperature of 15.85°C. This figure is representative of the environmental conditions within the timeframe of the accumulated data. The calculation was executed in a notably rapid execution time of 0.0046062 seconds, demonstrating the batch processing's capability for high-speed computation on stored datasets.

The quick execution time is particularly noteworthy, suggesting that the Spark-based batch processing system is well-optimized for such aggregative computations. Furthermore, this fast processing capability of batch operations, when applied to large volumes of data that do not require real-time analysis, could offer organizations a time-efficient and cost-effective solution for periodic data analysis tasks.

6. Comparison of Streaming & Batch Modes

```
● sharath@sharath-Dell-G15-5515:~/Desktop/Real-Time-Data-Pipeline-Using-Kafka-and-Spark-master$ python3 comparison.py
Streaming Temperature: 16.473300000000005, Streaming Execution Time: 0.13421154022216797 seconds
Batch Temperature: 15.850992662926197, Batch Execution Time: 0.0046274662017822266 seconds
○ sharath@sharath-Dell-G15-5515:~/Desktop/Real-Time-Data-Pipeline-Using-Kafka-and-Spark-master$ █
```

Temperature Results Comparison:

The average temperatures obtained from the streaming and batch modes are relatively close, with a minor difference that could be attributed to the datasets' time frame or variations in data sampled during the streaming window versus the entire batch dataset.

Streaming Mode: Produced an average temperature of 16.47°C, reflecting the near-instantaneous analysis of the data as it streams in real time.

Batch Mode: Yielded an average temperature of 15.86°C, indicating the processing of accumulated data over a given period.

Execution Time Comparison:

There is a significant difference in execution times between the two modes, which highlights the trade-offs between real-time processing and batch data analysis.

Streaming Mode: Demonstrates a processing time of approximately 0.13 seconds. This duration is reflective of the near-real-time requirement for streaming data processing, where data is continually ingested and processed in small windows, necessitating quick turnarounds.

Batch Mode: Showcases a faster execution time of 0.005 seconds. This efficiency can be largely attributed to the processing of a predetermined dataset in one go, allowing for more optimized computation and resource allocation.

Discussion:

Timeliness: Streaming processing is designed to provide ongoing insights, which is crucial for time-sensitive decisions. The slightly higher average temperature in streaming could indicate a real-time response to a transient environmental condition.

Resource Efficiency: Batch processing appears to be more resource-efficient in this instance, handling a large dataset with a quicker execution time than streaming. This suggests that for tasks where real-time analysis is not critical, batch processing could be more favorable in terms of both speed and computational load.

Use Case Suitability: The choice between streaming and batch processing should be informed by the specific use case requirements. For operations that rely on the most up-to-date information, such as fraud detection or emergency response systems, streaming offers a distinct advantage. Conversely, for analytical tasks that can be scheduled and do not require immediate action, batch processing is advantageous due to its efficiency and speed.

7. Conclusion:

The project demonstrates the benefits of real-time data processing for quick water quality monitoring, while also showing how batch processing is useful for detailed analysis. This

helps inform decisions on setting up data processing systems for water quality monitoring.

8. References

a. URLs

- [1] https://downloads.apache.org/kafka/3.7.0/kafka_2.12-3.7.0.tgz
- [2] <https://www.apache.org/dyn/closer.lua/spark/spark-3.5.1/spark-3.5.1-bin-hadoop3.tgz>
- [3] <https://www.mongodb.org/static/pgp/server-7.0.asc>
- [4] <https://data.world/cityofchicago/beach-water-quality-automated-sensors>
- [5] <https://data.cityofchicago.org/d/k7hf-8y75>
- [6] <https://data.cityofchicago.org/d/g3ip-u8rb>.
- [7] https://www.youtube.com/watch?v=udnX21_SuU&ab_channel=LearningJournal
- [8] https://www.youtube.com/watch?v=QaoJNXW6SQo&ab_channel=Simplilearn