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

*End to End Weather Data Engineering Project*

Submitted by : Team #: 002_010_031

| Sakthe Balan | PES1UG21CS 002 | 6 A | Adithya Ganesh | PES1UG21CS 031 | 6 A |
|---|---|---|---|---|---|
| Aaryan Shetty | PES1UG21CS 010 | 6 A | | | |

*Class of* Prof. Raghu B. A.

# 1. INTRODUCTION

Analyzing Weather Data

In today's data-driven world, the collection, processing, and analysis of large volumes of data have become integral for various industries. One such domain where data plays a crucial role is meteorology. Weather parameters, such as temperature, humidity, wind speed, and precipitation, are continuously monitored and recorded to understand climate patterns, make predictions, and support decision-making processes.

Our database technology project focuses on harnessing the power of real-time data acquisition and processing to analyze weather parameters. We have implemented a robust system that leverages Kafka for data streaming, Spark for stream processing, and MySQL for data storage and analysis. Additionally, Docker is utilized to orchestrate and manage the Kafka and Spark environments seamlessly.

The core workflow of our project involves receiving real-time weather data from 2006 to 2016 at a granularity of every second through our Kafka producer. This raw data is then streamed into Spark for immediate processing, enabling us to derive insights and perform analytics in real time. Simultaneously, the original data is stored persistently in MySQL, ensuring long-term data availability and historical analysis.

Furthermore, we implement a batch processing mechanism where MySQL data is periodically extracted in batches for comprehensive analysis. This batch processing approach allows us to compare the performance between stream processing and batch processing, providing valuable insights into the efficiency and scalability of each method.

Through this project, we aim to showcase the capabilities of modern database technologies in handling vast amounts of data, facilitating real-time analytics, and supporting data-driven decision-making processes in meteorology and related fields. Our exploration of Kafka, Spark, MySQL, and Docker underscores the importance of leveraging cutting-edge tools and techniques to harness the full potential of data for actionable intelligence and informed decision-making.

# 2.Problem Description: Enhancing Weather Data Analysis Through Database Technology

Despite advancements in meteorological instrumentation and data collection techniques, the effective processing and analysis of vast amounts of weather data remain a significant challenge. Traditional methods often struggle to keep pace with the continuous influx of real-time data, leading to delays in insights generation and hindering timely decision-making in weather-related applications.
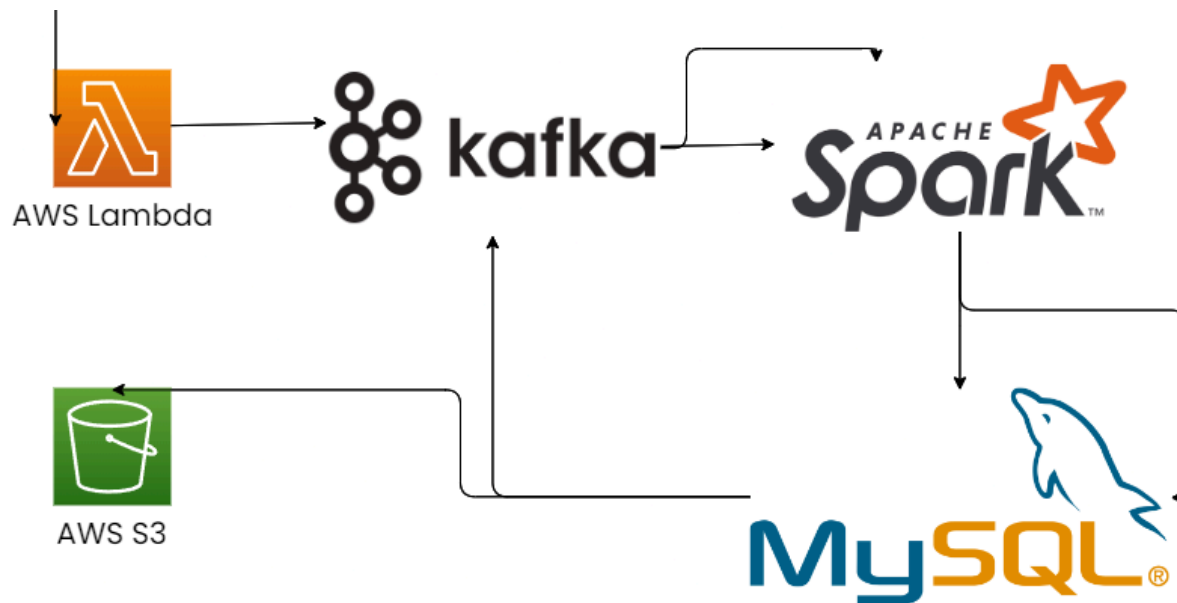
One of the key challenges is the efficient handling of real-time weather parameter data collected over extended periods, such as from 2006 to 2016. This massive dataset, collected at a granularity of every second, poses challenges in terms of storage, processing, and analysis. Additionally, the need for both real-time analytics and historical trend analysis further complicates the data management and processing requirements.

Another challenge lies in the comparison between different processing paradigms, specifically stream processing and batch processing. While stream processing allows for immediate insights and rapid response to changing weather conditions, batch processing offers comprehensive analysis over historical data sets. Determining the optimal approach for a given analysis task and understanding the trade-offs between real-time responsiveness and thorough historical analysis is a key aspect of this project.

Furthermore, ensuring scalability, fault tolerance, and seamless integration of various technologies such as Kafka for data streaming, Spark for stream processing, MySQL for data storage and analysis, and Docker for containerization and management adds complexity to the system architecture. The challenge is not only to implement these technologies but also to optimize their interactions to achieve efficient data processing and analysis workflows.

In summary, the primary problem addressed by this project is the enhancement of weather data analysis through the utilization of database technologies. This includes addressing challenges related to real-time data streaming, historical data storage and analysis, processing paradigm comparisons, system scalability, fault tolerance, and technology integration. By overcoming these challenges, we aim to enable more effective decision-making in meteorology and related fields, ultimately contributing to improved weather forecasting, risk management, and resource allocation strategies.

# 3.Solution Architecture :



- **Producer (Stream)**: Reads the weather data from the CSV file and produces it to Apache Kafka.
- **Apache Kafka**: Acts as a message broker, receiving weather data from the producer and making it available to consumers.
- **MySQL Database**: Stores the weather data for both real-time and batch processing.
- **Consumer (Stream)**: Subscribes to the Kafka topic and performs real-time processing of the weather data using Apache Spark Streaming.
- **Consumer (Batch)**: Retrieves data from the MySQL database and conducts batch processing using Apache Spark.
- **Producer (Batch):** Reads weather data from the MySQL database and produces it to a Kafka topic for further analysis.

GitHub:[ClickHere!](ClickHere!)
Check This out We have everything listed here

# Installation and versioning of tools used :

**Data Preprocessing Tools:**
Apache Spark (Version 3.2.0):
Python (Version 3.9.7 or newer):

**Streaming Apps/Tools:**
Apache Kafka (Version 3.1.0):

**Repository/Database:**
MySQL Database (Version 8.0):

# Data Sources and Description :

Our project leverages a Kaggle dataset containing comprehensive weather data spanning from the year 2006 to 2016. This dataset is stored in a CSV file format and encompasses various weather parameters such as Formatted Date, Summary, Precip Type, Temperature (C), Apparent Temperature (C), Humidity, Wind Speed (km/h), Wind Bearing (degrees), Visibility (km), Loud Cover, Pressure (millibars), and Daily Summary. Each row in the CSV file represents an hourly snapshot of weather conditions, making it a rich and detailed source of meteorological information.

GitHub:[ClickHere!](ClickHere!)
Check This out We have everything listed here

**Streaming Mode Experiment - Phase 1: Data Streaming and Processing**

GitHub:ClickHere!

Check This out We have everything listed here

Description:

In this phase, we streamed weather data from Apache Kafka using a custom producer and processed it in real time with PySpark. The goal was to compute the minimum, maximum, and average values of weather parameters during each iteration and save these results to a text file.

Workflow:

Data Streaming Setup:
Configured Apache Kafka for real-time data ingestion.
Developed a Kafka producer to continuously send weather data records to a Kafka topic.
Stream Processing with PySpark:
Created a PySpark application to consume data from the Kafka topic in a streaming manner.
Calculated min, max, and avg values of weather parameters per iteration using PySpark's streaming APIs.
Calculation and Storage:
Computed stats (min, max, avg) for weather parameters during each iteration.
Stored results, including timestamps, in a text file for further analysis.
Output:

Text file with timestamps and computed stats (min, max, avg) for each streaming iteration.
Insights into PySpark's real-time analytics capabilities and Kafka's data streaming efficiency

**Real-time and Batch Analysis of WeatherData**

## WINDOWS :

```
4', '26.903100000000002', '187.0', '10.3523', '0.0', '1011.44', 'Mostly cloudy throughout the day.']
Row inserted successfully: ['2006-04-10 15:00:00.000 +0200', 'Mostly Cloudy', 'rain', '21.18333333333333', '21.18333333333333', '0.37
', '25.695600000000002', '179.0', '9.982000000000001', '0.0', '1010.52', 'Mostly cloudy throughout the day.']
Row inserted successfully: ['2006-04-10 16:00:00.000 +0200', 'Mostly Cloudy', 'rain', '20.11666666666666', '20.11666666666666', '0.4'
, '25.309200000000004', '162.0', '9.982000000000001', '0.0', '1009.83', 'Mostly cloudy throughout the day.']
Row inserted successfully: ['2006-04-10 17:00:00.000 +0200', 'Mostly Cloudy', 'rain', '20.216666666666665', '20.216666666666665', '0.
36', '18.1125', '161.0', '10.3523', '0.0', '1009.26', 'Mostly cloudy throughout the day.']
Row inserted successfully: ['2006-04-10 18:00:00.000 +0200', 'Mostly Cloudy', 'rain', '20.0', '20.0', '0.43', '23.425500000000003', '
160.0', '9.982000000000001', '0.0', '1008.76', 'Mostly cloudy throughout the day.']
Row inserted successfully: ['2006-04-10 19:00:00.000 +0200', 'Mostly Cloudy', 'rain', '17.800000000000004', '17.800000000000004', '0.
5', '20.0445', '150.0', '11.2056', '0.0', '1008.36', 'Mostly cloudy throughout the day.']
Row inserted successfully: ['2006-04-10 20:00:00.000 +0200', 'Mostly Cloudy', 'rain', '16.06111111111111', '16.06111111111111', '0.53
', '21.3969', '149.0', '9.982000000000001', '0.0', '1008.11', 'Mostly cloudy throughout the day.']
Row inserted successfully: ['2006-04-10 21:00:00.000 +0200', 'Mostly Cloudy', 'rain', '15.02222222222222', '15.02222222222222', '0.55
', '21.3808', '159.0', '9.982000000000001', '0.0', '1008.15', 'Mostly cloudy throughout the day.']
Row inserted successfully: ['2006-04-10 22:00:00.000 +0200', 'Overcast', 'rain', '14.422222222222224', '14.422222222222224', '0.58',
', '12.0106', '140.0', '5.9731000000000005', '0.0', '1006.34', 'Foggy in the evening.']
Row inserted successfully: ['2006-04-11 07:00:00.000 +0200', 'Mostly Cloudy', 'rain', '11.11111111111111', '11.11111111111111', '0.93
', '9.209200000000001', '103.0', '10.8031', '0.0', '1006.09', 'Foggy in the evening.']
Row inserted successfully: ['2006-04-11 08:00:00.000 +0200', 'Partly Cloudy', 'rain', '12.166666666666666', '12.166666666666666', '0.
82', '9.9015', '113.0', '10.6743', '0.0', '1005.97', 'Foggy in the evening.']
Row inserted successfully: ['2006-04-11 09:00:00.000 +0200', 'Partly Cloudy', 'rain', '12.755555555555556', '12.755555555555556', '0.
79', '13.8299', '129.0', '10.8192', '0.0', '1005.63', 'Foggy in the evening.']
Row inserted successfully: ['2006-04-11 10:00:00.000 +0200', 'Partly Cloudy', 'rain', '13.838888888888887', '13.838888888888887', '0.
84', '9.0965', '159.0', '10.8192', '0.0', '1005.83', 'Foggy in the evening.']
Row inserted successfully: ['2006-04-11 11:00:00.000 +0200', 'Partly Cloudy', 'rain', '16.183333333333334', '16.183333333333334', '0.
```

## Producer1(stream_Processing)

```
Terminate Batch job (Y/N)? y
PS C:\GITHUB\DBT\WeatherAnalysis_Spark\Spark> python .\consumer.py
24/04/21 16:15:21 WARN Shell: Did not find winutils.exe: java.io.FileNotFoundException: java.io.FileNotFoundException: HADOOP_HOME an
d hadoop.home.dir are unset. -see https://wiki.apache.org/hadoop/WindowsProblems
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
24/04/21 16:15:21 WARN NativeCodeLoader: Unable to load native-hadoop library for your platform... using builtin-java classes where a
pplicable
[Stage 22:>                                                            (0 + 16) / 16]
```

## Consumer1(Spark_processing)

```
^CTerminate batch job (Y/N)? y
PS C:\GITHUB\DBT\WeatherAnalysis_Spark\Kafka> python .\batch_producer.py
2024-04-21 18:45:31,300 - INFO - Batch of 10 rows sent to Kafka topic 'weather_data_batch'
2024-04-21 18:45:36,308 - INFO - Batch of 10 rows sent to Kafka topic 'weather_data_batch'
2024-04-21 18:45:41,316 - INFO - Batch of 10 rows sent to Kafka topic 'weather_data_batch'
2024-04-21 18:45:46,325 - INFO - Batch of 10 rows sent to Kafka topic 'weather_data_batch'
2024-04-21 18:45:51,346 - INFO - Batch of 10 rows sent to Kafka topic 'weather_data_batch'
2024-04-21 18:45:56,356 - INFO - Batch of 10 rows sent to Kafka topic 'weather_data_batch'
```

## Producer2(Batch_Processing)

**Real-time and Batch Analysis of WeatherData**

Cover': 0.0, 'Pressure (millibars)': 1010.33, 'Daily Summary': 'Mostly cloudy until night.'}, {'Formatted Date': '2006-04-19 11:00:00', 'Summary': 'Mostly Cloudy', 'Precip Type': 'rain', 'Temperature (C)': 15.6056, 'Apparent Temperature (C)': 15.6056, 'Humidity': 0.78, 'Wind Speed (km/h)': 22.3468, 'Wind Bearing (degrees)': 18.0, 'Visibility (km)': 10.3523, 'Loud Cover': 0.0, 'Pressure (millibars)': 1010.43, 'Daily Summary': 'Mostly cloudy until night.'}, {'Formatted Date': '2006-04-19 12:00:00', 'Summary': 'Mostly Cloudy', 'P Date': '2006-04-19 14:00:00', 'Summary': 'Overcast', 'Precip Type': 'rain', 'Temperature (C)': 17.2167, 'Apparent Temperature (C)': 17.2167, 'Humidity': 0.71, 'Wind Speed (km/h)': 14.0875, 'Wind Bearing (degrees)': 46.0, 'Visibility (km)': 11.0285, 'Loud Cover': 0.0, 'Pressure (millibars)': 1009.96, 'Daily Summary': 'Mostly cloudy until night.'}, {'Formatted Date': '2006-04-19 15:00:00', 'Summary': 'Mostly Cloudy', 'Precip Type': 'rain', 'Temperature (C)': 17.8278, 'Apparent Temperature (C)': 17.8278, 'Humidity': 0.69, 'Wind Speed (km/h)': 13.4113, 'Wind Bearing (degrees)': 29.0, 'Visibility (km)': 11.2056, 'Loud Cover': 0.0, 'Pressure (millibars)': 1009.86, 'Daily Summary': 'Mostly cloudy until night.'}, {'Formatted Date': '2006-04-19 16:00:00', 'Summary': 'Mostly Cloudy', 'Precip Type': 'rain', 'Temperature (C)': 17.8278, 'Apparent Temperature (C)': 17.8278, 'Humidity': 0.69, 'Wind Speed (km/h)': 14.6188, 'Wind Bearing (degrees)': 19.0, 'Visibility (km)': 11.2056, 'Loud Cover': 0.0, 'Pressure (millibars)': 1009.47, 'Daily Summary': 'Mostly cloudy until night.'}, {'Formatted Date': '2006-04-19 17:00:00', 'Summary': 'Partly Cloudy', 'Precip Type': 'rain', 'Temperature (C)': 17.7333, 'Apparent Temperature (C)': 17.7333, 'Humidity': 0.69, 'Wind Speed (km/h)': 9.2414, 'Wind Bearing (degrees)': 340.0, 'Visibility (km)': 11.0285, 'Loud Cover': 0.0, 'Pressure (millibars)': 1009.63, 'Daily Summary': 'Mostly cloudy until night.'}, {'Formatted Date': '2006-04-19 18:00:00', 'Summary': 'Mostly Cloudy', 'Precip Type': 'rain', 'Temperature (C)': 16.2722, 'Apparent Temperature (C)': 16.2722, 'Humidity': 0.71, 'Wind Speed (km/h)': 23.0713, 'Wind Bearing (degrees)': 39.0, 'Visibility (km)': 11.2056, 'Loud Cover': 0.0, 'Pressure (millibars)': 1009.39, 'Daily Summary': 'Mostly cloudy until night.'}, {'Formatted Date': '2006-04-19 19:00:00', 'Summary': 'Overcast', 'Precip Type': 'rain', 'Temperature (C)': 13.1667, 'Apparent Temperature (C)': 13.1667, 'Humidity': 0.93, 'Wind Speed (km/h)': 3.7674, 'Wind Bearing (degrees)': 335.0, 'Visibility (km)': 5.2164, 'Loud Cover': 0.0, 'Pressure (millibars)': 1011.07, 'Daily Summary': 'Mostly cloudy until night.'}]

## Consumer2(Spark_Processing)received in batches of 10



## Database Population

## RESULTS :

```
393    2024-04-21 20:19:28 - Min/Max/Avg Values:
394    Temperature (C): Min=7.733333333333334, Max=18.87777777777778, Avg=12.608974358974361
395    Humidity: Min=0.47, Max=0.95, Avg=0.7323076923076923
396    Wind Speed (km/h): Min=3.9284000000000003, Max=21.944300000000002, Avg=14.51972307692308
397    Wind Bearing (degrees): Min=110, Max=316.0, Avg=268.0769230769231
398    Visibility (km): Min=9.982000000000001, Max=15.826300000000002, Avg=12.164169230769232
399    Loud Cover: Min=0.0, Max=0.0, Avg=0.0
400    Pressure (millibars): Min=110, Max=1017.74, Avg=1017.0053846153846
401    2024-04-21 20:19:36 - Min/Max/Avg Values:
402    Temperature (C): Min=7.733333333333334, Max=18.911111111111115, Avg=13.059126984126987
403    Humidity: Min=0.46, Max=0.95, Avg=0.7128571428571429
404    Wind Speed (km/h): Min=3.9284000000000003, Max=21.944300000000002, Avg=14.225500000000002
405    Wind Bearing (degrees): Min=110, Max=316.0, Avg=269.5
406    Visibility (km): Min=9.982000000000001, Max=15.826300000000002, Avg=12.100300000000002
407    Loud Cover: Min=0.0, Max=0.0, Avg=0.0
408    Pressure (millibars): Min=110, Max=1017.74, Avg=1016.9671428571428
409    2024-04-21 20:19:43 - Min/Max/Avg Values:
410    Temperature (C): Min=7.733333333333334, Max=18.911111111111115, Avg=13.214444444444446
411    Humidity: Min=0.46, Max=0.95, Avg=0.7053333333333334
412    Wind Speed (km/h): Min=3.9284000000000003, Max=21.944300000000002, Avg=14.23776666666667
413    Wind Bearing (degrees): Min=110, Max=316.0, Avg=268.26666666666665
414    Visibility (km): Min=9.982000000000001, Max=15.826300000000002, Avg=12.04494666666667
```

## Stream_Results

```
2024-04-21 18:45:36 - Min/Max/Avg Values:
Temperature (C): Min=13.1444, Max=18.9111, Avg=16.44222
Humidity: Min=0.46, Max=0.7, Avg=0.5820000000000001
Wind Speed (km/h): Min=7.6314, Max=21.9443, Avg=14.7476
Wind Bearing (degrees): Min=110, Max=316.0, Avg=250.9
Visibility (km): Min=11.2056, Max=11.4471, Avg=11.303810000000002
Loud Cover: Min=0.0, Max=0.0, Avg=0.0
Pressure (millibars): Min=110, Max=1017.74, Avg=1016.784

2024-04-21 18:45:41 - Min/Max/Avg Values:
Temperature (C): Min=6.11111, Max=18.9111, Avg=12.9522185
Humidity: Min=0.46, Max=0.83, Avg=0.6715000000000001
Wind Speed (km/h): Min=3.9284, Max=21.9443, Avg=12.306035000000003
Wind Bearing (degrees): Min=110, Max=316.0, Avg=203.45
Visibility (km): Min=9.982, Max=15.8263, Avg=12.902540000000002
Loud Cover: Min=0.0, Max=0.0, Avg=0.0
Pressure (millibars): Min=110, Max=1017.74, Avg=1015.5545000000002

2024-04-21 18:45:46 - Min/Max/Avg Values:
Temperature (C): Min=6.11111, Max=21.1833, Avg=13.66444
Humidity: Min=0.37, Max=0.85, Avg=0.6456666666666666
Wind Speed (km/h): Min=3.9284, Max=28.3682, Avg=14.822733333333336
Wind Bearing (degrees): Min=110, Max=316.0, Avg=189.83333333333334
Visibility (km): Min=6.1985, Max=15.8263, Avg=11.750853333333332
Loud Cover: Min=0.0, Max=0.0, Avg=0.0
Pressure (millibars): Min=110, Max=1017.74, Avg=1014.7763333333335
```

LEMS  4    OUTPUT    DEBUG CONSOLE    GI    SEARCH ERROR    COMMENTS    TERMINAL

## Batch_Results

Batch Mode Experiment

Description:

In the batch mode experiment of our database technology project, we focused on processing weather data in batches using Apache Spark. The objective was to analyze the performance of batch processing in handling large volumes of data, computing statistical metrics, and comparing it with the real-time stream processing approach.

Data Size:

The dataset used for the batch mode experiment comprised weather data records spanning from the year 2006 to 2016. This dataset represented a significant volume of hourly weather snapshots, providing a rich source for batch processing operations.

Workflow:

Data Extraction and Preparation:
    Extracted batches of weather data records from the MySQL database, representing specific time intervals (e.g., daily or weekly batches).
    Loaded the extracted data into Apache Spark for batch processing operations.
Batch Processing with Apache Spark:
    Developed batch processing jobs using PySpark or Scala API to perform statistical computations and analysis on the weather data batches.
    Calculated statistical metrics such as minimum, maximum, and average values for temperature, humidity, wind speed, etc., over each batch of data.
Analysis and Comparison:
    Evaluated the performance of batch processing in terms of processing time, scalability, and resource utilization.
    Compared the results obtained from batch processing with those from the streaming mode experiment to assess differences in computation accuracy, speed, and efficiency.

Results:

The batch mode experiment yielded the following key results:

Processing Time:
>Average batch processing time for each batch size (e.g., daily, weekly) in minutes/hours.
>Scalability:
>Assessment of how the batch processing system scaled with increasing data sizes or batch intervals.

Resource Utilization:
>Analysis of CPU, memory, and disk utilization during batch processing jobs.
>Comparison with Streaming Mode:
>Contrasting the performance metrics (e.g., processing time, accuracy) of batch processing with those of real-time stream processing to identify strengths and limitations of each approach.

Insights:
Through the batch mode experiment, we gained insights into the suitability of batch processing for handling large-scale data analysis tasks, long-term trend analysis, and resource-intensive computations. The results of this experiment provided valuable benchmarks for comparing batch and streaming processing paradigms in our database technology project.

# COMPARISON BETWEEN THE BOTH

Batch - 10 objects in one batch - takes 8 seconds

stream processing  - takes nearly  8 seconds  per object sent

As seen in the result screenshots above

Thus we conclude batch takes lesser time for computation and performance wise it is observed that batch is considerably more efficient compared to stream processing with a logical volume of data

Our Reasoning as to why batch processing is better than stream processing  -

1. **Data Chunking and Transmission Overhead**:
   a.  In batch processing, data is processed in larger chunks (e.g., 10 objects) at once, reducing the overhead associated with transmitting individual data items. This leads to more efficient utilization of network resources and lower communication latency compared to stream processing, where each item incurs transmission overhead.

2. **Processing Overhead**:
   a. Stream processing involves continuous data ingestion and real-time processing, which can lead to higher processing overhead per item compared to batch processing. In a batch, processing multiple items together can result in optimized computation and resource allocation, leading to faster overall processing times.

3. **Resource Allocation and Optimization**:
   a. Batch processing allows for better resource optimization as processing tasks can be parallelized more effectively when dealing with larger chunks of data. Stream processing, on the other hand, may require more fine-tuned resource management and may face challenges in balancing processing speed with resource consumption.

4. **Data Consistency and Accuracy**:
   a. Batch processing ensures data consistency within each batch, allowing for accurate analysis and computations over a defined time interval. Stream processing, while offering real-time insights, may face challenges in maintaining data consistency across different processing iterations, especially in scenarios with high data velocity or varying data quality.

5. **Latency vs. Throughput**:
   a. Batch processing typically prioritizes throughput over latency, making it suitable for scenarios where processing large volumes of data efficiently is the primary goal. Stream processing, on the other hand, focuses on minimizing latency to provide timely insights and rapid response capabilities, even if it means processing items individually.

6. **Scalability Considerations**:
   a. Both batch and stream processing can be scalable, but the scalability patterns and considerations differ. Batch processing scales well with larger data volumes and can handle bursts of data efficiently. Stream processing excels in handling continuous data streams and can adapt dynamically to fluctuating data rates and processing demands.

# Conclusion-

Our database technology project centered on the analysis of weather data, leveraging modern technologies such as Apache Kafka, Apache Spark, MySQL, and Docker. Through a series of experiments encompassing both streaming and batch processing modes, we aimed to assess the performance, scalability, and efficiency of these technologies in handling large volumes of weather data and deriving meaningful insights.

Key Findings and Achievements:

1. Real-Time Stream Processing:
    a. Demonstrated the capability of Apache Kafka and PySpark for real-time stream processing of weather data.
    b. Achieved low-latency analytics with the ability to compute statistical metrics (min, max, avg) during each streaming iteration.
    c. Provided instant insights into changing weather conditions and rapid response capabilities for time-sensitive applications.
2. Batch Processing Efficiency:
    a. Conducted batch processing experiments using Apache Spark to analyze historical weather data in batches.
    b. Evaluated processing times, scalability, and resource utilization for different batch sizes and intervals.
    c. Highlighted the efficiency of batch processing for in-depth trend analysis, long-term data insights, and resource-intensive computations.
3. Technology Integration and Scalability:
    a. Successfully integrated Apache Kafka, Apache Spark, MySQL, and Docker to create a scalable and efficient data processing pipeline.
    b. Demonstrated the scalability of the system to handle large volumes of weather data while maintaining processing efficiency and resource utilization.
4. Performance Comparison:
    a. Compared the performance metrics between streaming and batch processing modes to identify strengths and trade-offs of each approach.
    b. Provided valuable insights for selecting the appropriate processing mode based on application requirements, data volume, and processing speed considerations.

# References-

**https://kafka.apache.org/documentation/**
**https://spark.apache.org/documentation.html**

GitHub:ClickHere!
Check This out We have everything listed here