CSA1526

# CLOUD COMPUTING

## Assignment - 5

Name : Sai Lokesh Nalabothu

Reg.No : 192365023

Branch : CSE - Cyber Security

Date : 24 - 02 - 2025

Serial.No : 19

# Hadoop and Real-time Data Processing

Hadoop is traditionally designed for batch Processing and is not inherently suited for Real-time data Processing.

## Hadoop and Real-time Data Processing

## Problem Statement

Hadoop is widely used for large-scale data Storage and batch Processing but is not designed for Real time data Processing.

1. Understanding the Components

* Hadoop distributed file system (HDFS): A Scalable Storage system for handling large datasts.

* Map Reduce: A batch Processing Model that is not inherently Real-time but can Process streamed data

in Micro - batches.

## Apache kafka

* kafka is a distributed messaging system designed for handling Real time data streams

* AcH as a buffer between data Procedurey and Real time Processing Systems.

## Apache Storm

* A Real time Stream Processing framework that Processes data as it arrives.

* Can Integrate with kafka to Consume and Process Real time data streams.

→ Latency

→ Scalability

→ Fault tolerance

→ Fraud detection.

1. How is Hadoop used in Real-time data Processing?

* Hadoop is traditionally designed for batch Processing Rather than Real-time data Processing.

It can be Integrated with other technologies to Enable near-Real time or Real time data Processing.

1. Integration with streaming technologies.

* APache katka

* APache storm

* APache spark

2. Real time data Processing workflow using Hadoop.

* Data Integration

* Stream Processing

* Storage in Hadoop

* Batch Processing for Deep Analysis.

3. Real world use cases

* Fraud detection

* Social Media Analytics

* IoT Data Processing

4. Challenges

* Latency

* Complex Architecture

* Scalability & Maintenance

5. Alternative technologies for Real-time

Processing

* APache flink

* APache Druid

* Google BigQuery & AWS kinesis

2. What role does Apache Kafka Play in Real time data Processing?

Apache kafka is a distributed Event Streaming Platform that Plays a Crucial Role in Real time data Processing by acting as a high - throughPut.

1. key Roles of kafka in Real time Data Processing

* Data Ingestion & streaming

* Message Buffering

* Data Storage for Replay & Fault tolerance

* Real time Processing with stream Processing Frameworks.

* Integration with Big data & Databases.

2. Use Cases of kafka

   * Fraud detection

   * IOT Data Processing

   * Social Media Analytics

   * Log Monitoring

3. Why use kafka ?

   * High ThroughPut

   * Scalability

   * Fault tolerance

   * Low Latency

4. Work Flow

   * Producers

   * topics

   * Consumers

   * Databases

3. Explain how Apache Storm integrates with Hadoop for Real-time Processing

Apache Storm is a Real time stream processing framework that integrates with Hadoop to Process and analyze data as it arrives.

1. How Apache Storm works in Real-time Processing

  * Spouts

  * Bolts

2. Integration of Apache Storm with Hadoop

1. Real time Data Ingestion

2. Stream Processing

3. Storing Processed Data in Hadoop.

4. Batch Processing & Analytics

3. Real World Use Cases of Storm With Hadoop

* Fraud detection

* IoT Data Analysis

* Real-time Log Processing

4. Advantages of using Apache Storm With Hadoop

* Real-time & Batch Hybrid

* scalability

* Fault tolerance

* Low Latency

4. Discuss the Challenges of using Hadoop for Real time analytics.

Hadoop is batch Processing framework designed for handling large-scale data storage and analytics.

1. High Latency in Processing

* Map Reduce

* large batches

* Low Latency

* Real time decision Making

2. Lack of Native streaming support

* Hadoop does not natively support Stream Processing.

* Live data streams.

3. Inefficiency in small file handling

* Hadoop's HDFS is optimized for large files but perform poorly with many small files.

* Each file's data is stored in NameNode.

WorkAround

* Use Apache HBase

Conclusion:

Hadoop alone is not well-suited for Real-time analytics due to its batch-oriented nature, high latency, and lack of native streaming support.