

CSA1526

CLOUD COMPUTING

Assignment - 4

Name: Sai Lokesh Alalabothu

Reg. No: 192365023

Branch: CSE - Cybersecurity

Date: 22-02-2025

Serial No: 19

# BIG DATA TOOLS AND FRAMEWORKS

Big data Processing involves managing vast volumes of structured, semi-structured, and unstructured data efficiently. Several tools and frameworks have been developed to store, process and analyze Big data.

## HADOOP

Apache Hadoop is an open-source framework that enables distributed storage and processing of large datasets.

### Key Concepts

- \* HDFS
- \* YARN
- \* MapReduce
- \* HBase

### Advantages

- \* Scalable
- \* Fault-tolerant
- \* Cost effective.

## use cases

- \* Log Analysis
- \* Data Warehousing
- \* Fraud detection

## Apache Spark

Apache Spark is a fast and general-purpose cluster computing system that provides in-memory data processing.

## Capability.

### Key Features

- \* RDD (Resilient Distributed Datasets)
- \* Spark SQL
- \* Spark streaming
- \* MLlib
- \* GraphX

### Advantages

- \* Faster than Hadoop
- \* Support multiple languages.



## NoSQL Databases

NoSQL Databases are designed to handle large scale, distributed data storage with high availability and flexibility.

### Types of NoSQL Databases

- \* Key-value stores
- \* Document oriented
- \* Column-family stores
- \* Graph databases

### Advantages

- \* Scalable
- \* Schema-less
- \* High Performance

### Use Cases

- \* Social media analytics
- \* Recommendation systems
- \* IOT

1. What is Apache Hadoop, and how does it support Big Data Processing?

Apache Hadoop is an open-source framework designed for storing and processing massive datasets in a distributed environment.

Hadoop is particularly useful for handling structured and unstructured data.

How Hadoop supports Big data Processing

Hadoop consists of several core components that enable efficient data storage and

Computation:

1. HDFS
2. MapReduce
3. YARN
4. HBase



## WHY HADOOP

- \* Scalability
- \* Fault tolerance
- \* Cost-effective
- \* Flexibility

## 2. Discuss the Role of Apache Spark in Big data Processing.

Apache Spark is an open-source, It supports Real-time analytics, machine learning and Stream Processing, making it a popular choice for modern Big data applications.

### Key Features

1. In-Memory Processing
2. Distributed and fault tolerance
3. Supports multiple workloads.

4. Compatible with Hadoop & No SQL Databases.
5. Multi-Language support.

### Components of Apache Spark

- \* Spark Core
- \* Spark SQL
- \* Spark Streaming
- \* MLlib
- \* GraphX

### Advantages

- \* Speed
- \* Flexibility
- \* Scalability
- \* Ease of Use

### Use cases

- \* Real time Analytics
- \* ETL
- \* Healthcare & Genomics



### 3. How do NOSQL databases support Big data storage and querying?

NoSQL databases are designed to handle large-scale, distributed, and unstructured data efficiently. NoSQL databases provide high scalability, flexible schema and faster performance.

#### Role of NOSQL in BigData Processing

1. Scalability
2. Flexible schema
3. High Performance
4. Real time Data Processing

#### Querying in NOSQL databases

\* Key value stores: GET key, set key value

```
db.users.find({"age": {"$gt": 25}})
```

```
select * from users where age > 25;
```



## Advantages

- \* Handles Large Volumes : scales efficiently with large data sets.
- \* High Availability : uses Replication for Fault tolerance.
- \* schema flexibility : No Need for predefined table structure.
- \* Optimized for Big data Analytics : works with Real time and Batch Processing frame works like Hadoop & Spark.

## Conclusion :

They are widely used in social media platforms, IOT systems, Recommendation engines, and analytics applications.

4. What are the advantages and challenges of using Big data tools?

Bigdata tools like Hadoop, Apache Spark, and NoSQL databases play a crucial role in processing, storing and analyzing vast amounts of data.

### Advantages of Big Data Tools

1. Scalability
2. Speed and Performance
3. Cost-effectiveness
4. flexibility in Data Processing
5. Fault tolerance and Reliability
6. Real time Data Processing
7. Advanced Analytics & AI/ML Support.

### Challenges of Big data tools

1. Complexity in implementation
2. High Infrastructure Costs



3. Data security and Privacy Risks

4. Data quality & Integration issues

5. Lack of standardization

6. Maintenance & troubleshooting

### Conclusion:

Big data tools Revolutionize data processing

by offering scalability, speed and

flexibility, but they come with challenges

like complexity, security risks and

high infrastructure costs.