



BUDDHA INSTITUTE OF TECHNOLOGY

Gida Gorakhpur



Department- Computer Science and Allied (CSE AND DS)

Program & Semester- B.Tech 3rd Year (6th Semester)

Course and Code- BIG Data Analytics BCDS 061 & BCS 061

Course Outcome

CO No.	Course Outcome	Bloom's Knowledge Level (KL)
CO 1	Demonstrate knowledge of Big Data Analytics concepts and its applications in business.	K1, K2
CO 2	Demonstrate functions and components of Map Reduce Framework and HDFS.	K1, K2
CO 3	Discuss Data Management concepts in NoSQL environment.	K6
CO 4	Explain process of developing Map Reduce based distributed processing applications.	K2, K5
CO 5	Explain process of developing applications using HBASE, Hive, Pig etc.	K2, K5

UNIT-2

GRAPH AND MATRICES

Q1. Explain the characteristics of Big Data (5 V's) with suitable examples. (AKTU 2022, GATE – Concept Based)

Introduction:

Big Data refers to extremely large and complex datasets that cannot be efficiently processed using traditional data processing systems. With the rapid growth of digital technologies, enormous amounts of data are generated every second from sources such as social media, sensors, mobile devices, and business transactions. To understand the nature of Big Data, it is commonly described using a set of characteristics known as the “5 V’s”. These characteristics help in identifying the challenges and requirements of Big Data systems and are widely discussed in both academic examinations like AKTU and competitive exams such as GATE (concept-based).

Volume:

Volume refers to the massive amount of data generated and stored. In today’s digital world, organizations deal with data in terabytes, petabytes, and even exabytes. For example, companies like Facebook and Google generate enormous volumes of data daily in the form of posts, images, videos, and search queries. Traditional databases are not capable of handling such huge data efficiently, which necessitates the use of distributed storage systems like Hadoop Distributed File System (HDFS). The large volume of data requires scalable storage solutions and efficient data management techniques.

Velocity:

Velocity refers to the speed at which data is generated, collected, and processed. In many applications, data is produced continuously and must be processed in real-time or near real-time. For example, stock market trading systems require real-time analysis of data to make quick decisions. Similarly, online fraud detection systems must process transaction data instantly to identify suspicious activities. High velocity data requires advanced processing frameworks such as Apache Spark and stream processing tools to ensure timely analysis.

Variety:

Variety refers to the different types of data formats that exist in Big Data. Unlike traditional systems that handle only structured data, Big Data includes structured, semi-structured, and unstructured data. Structured data is organized in tables such as relational databases. Semi-structured data includes formats like XML and JSON, while unstructured data includes images, videos, emails, and social media content. Managing and analyzing such diverse data types is a major challenge and requires flexible storage systems like NoSQL databases and Hadoop.

Veracity:

Veracity refers to the quality, accuracy, and reliability of data. Since Big Data is collected from multiple sources, it often contains inconsistencies, noise, and incomplete information. For example, data collected from social media platforms may include fake news or misleading information. Poor data quality can lead to incorrect analysis and decision-making. Therefore, data cleaning and preprocessing techniques are essential to improve the reliability of

data before analysis.

Value:

Value is the most important characteristic of Big Data as it focuses on extracting meaningful insights from data. The ultimate goal of Big Data analytics is to generate value for organizations by improving decision-making, optimizing operations, and identifying new opportunities. For example, e-commerce companies analyze customer behavior to provide personalized recommendations, thereby increasing sales and customer satisfaction. Without extracting value, storing large amounts of data becomes useless.

Conclusion:

The 5 V's of Big Data—Volume, Velocity, Variety, Veracity, and Value—collectively define the complexity and significance of Big Data systems. These characteristics highlight the need for advanced technologies and frameworks to store, process, and analyze large datasets efficiently. Understanding these concepts is essential for students preparing for AKTU examinations as well as GATE, as they form the foundation of Big Data Analytics.

Q2. Differentiate between Traditional Data Processing and Big Data Processing. (AKTU 2021, GATE – Concept Based)**Introduction**

Traditional data processing systems were designed to handle structured data with limited size and complexity using centralized architectures. However with the exponential growth of data generated from modern applications such as social media IoT and e-commerce traditional systems became insufficient. This led to the development of Big Data processing which uses distributed computing and advanced frameworks to handle large scale diverse data efficiently.

Traditional Data Processing

Traditional data processing relies on relational database management systems (RDBMS) such as MySQL Oracle and SQL Server. These systems store data in structured format using tables with rows and columns. Data processing is performed using SQL queries and follows ACID properties to ensure consistency and reliability.

Characteristics

Centralized architecture where data is stored on a single system

Handles only structured data

Limited scalability mainly vertical scaling

High cost due to expensive hardware

Suitable for small to medium sized datasets

Processing is slower for very large datasets

Example

Banking systems maintaining customer records and transaction history in relational databases

Big Data Processing

Big Data processing is designed to handle massive volumes of structured semi-structured and unstructured data using distributed systems. It uses frameworks such as Hadoop and Spark which divide data into smaller chunks and process them in parallel across multiple machines.

Characteristics

- Distributed architecture using clusters of computers
- Handles structured semi-structured and unstructured data
- Highly scalable using horizontal scaling
- Cost effective using commodity hardware
- Supports real-time and batch processing
- Fault tolerant using data replication

Example

Social media platforms like Facebook processing user activity data including posts images and videos

Key Differences

Data Type

Traditional systems handle only structured data whereas Big Data systems support multiple data formats

Scalability

Traditional systems scale vertically by adding resources to a single machine whereas Big Data systems scale horizontally by adding more machines

Processing Speed

Traditional systems are slower for large datasets whereas Big Data systems provide faster processing using parallel computing

Cost

Traditional systems require expensive hardware whereas Big Data systems use low cost commodity machines

Fault Tolerance

Traditional systems have limited fault tolerance whereas Big Data systems ensure reliability through data replication

Technology

Traditional systems use SQL and RDBMS whereas Big Data systems use Hadoop Spark NoSQL databases

Conclusion

Traditional data processing is suitable for small structured datasets while Big Data processing is essential for handling large complex and diverse data in modern applications.

Q3. Explain Big Data Architecture in detail. (AKTU 2020, GATE – Concept Based)

Introduction

Big Data Architecture provides a structured framework for collecting storing processing and analyzing large volumes of data. It consists of multiple layers that work together to transform raw data into meaningful insights. This architecture is essential for handling the complexity and scale of modern data systems and is widely used in industries such as healthcare finance and e-commerce.

Layers of Big Data Architecture

1. Data Source Layer

This layer includes all sources from which data is generated. These sources may include social media platforms sensors IoT devices business transactions and web logs. Data generated at this stage is often large and diverse in nature.

Characteristics

Generates structured semi-structured and unstructured data

Continuous and high-speed data generation

Includes internal and external data sources

2. Data Ingestion Layer

This layer is responsible for collecting and transferring data from various sources into the storage system. Tools such as Apache Flume Kafka and Sqoop are commonly used.

Characteristics

Handles batch and real-time data ingestion

Ensures reliable data transfer

Supports data streaming

3. Storage Layer

This layer stores large volumes of data in distributed storage systems such as Hadoop Distributed File System (HDFS) and NoSQL databases.

Characteristics

Scalable and distributed storage

Fault tolerance using replication

Supports different data formats

4. Processing Layer

This layer processes data using frameworks such as MapReduce and Apache Spark. It can perform both batch processing and real-time processing.

Characteristics

Parallel data processing

Supports large scale computations

Provides high performance

5. Analytics Layer

In this layer data is analyzed using data mining machine learning and statistical techniques to extract insights.

Tools such as Hive Pig and Spark SQL are used.

Characteristics

Performs complex data analysis

Supports query processing

Generates insights and patterns

6. Visualization Layer

This layer presents the analyzed data in the form of dashboards charts and reports. Tools such as Tableau and Power BI are used.

Characteristics

User friendly representation

Supports decision making

Provides interactive visualization

Additional Components

Security ensures data privacy and protection

Data Governance ensures data quality and compliance

Conclusion

Big Data Architecture provides a comprehensive framework for handling large scale data efficiently. Each layer plays a crucial role in transforming raw data into valuable insights.

Q4. Discuss different types of data in Big Data. (AKTU 2019, GATE – Databases & Data Representation Concept)

Introduction

In Big Data Analytics data exists in multiple forms and structures which makes its processing and storage challenging. Unlike traditional systems that handle only structured data Big Data systems must deal with diverse data formats generated from various sources such as social media sensors and enterprise applications.

Understanding the types of data is essential for selecting appropriate storage models and processing techniques.

Types of Data in Big Data

1. Structured Data

Structured data refers to data that is organized in a predefined format typically in rows and columns. It is easy to store process and analyze using relational database systems.

Characteristics

Highly organized format

Stored in tables with schema

Easy to query using SQL

Example

Banking records customer details transaction logs

2. Semi-Structured Data

Semi-structured data does not follow a strict tabular format but contains tags or markers to separate data elements.

It lies between structured and unstructured data.

Characteristics

Flexible schema

Self-describing structure

Requires special tools for processing

Example

XML JSON documents web data

3. Unstructured Data

Unstructured data does not have any predefined structure and is the most complex type of data in Big Data systems.

Characteristics

No fixed format

Difficult to store and analyze

Requires advanced techniques like NLP and image processing

Example

Images videos audio files emails social media posts

4. Quasi-Structured Data

Quasi-structured data contains irregular or inconsistent formats and cannot be easily categorized.

Characteristics

Irregular structure

Partially organized

Requires preprocessing

Example

Clickstream data log files

Importance in Big Data

Handling different types of data requires flexible storage systems such as NoSQL databases and distributed frameworks like Hadoop. Each type of data requires different processing techniques and tools.

Conclusion

The classification of data into structured semi-structured unstructured and quasi-structured forms is fundamental in

Big Data Analytics.

Q5. Explain Data Analytics Lifecycle in detail. (AKTU 2022, GATE – Data Analytics / Machine Learning Process Concept)

Introduction

The Data Analytics Lifecycle is a systematic approach used to extract meaningful insights from raw data. It consists of multiple phases that guide the process of data analysis from problem definition to deployment. This lifecycle ensures that data is processed efficiently and results are accurate and useful for decision making.

Phases of Data Analytics Lifecycle

1. Data Discovery

This phase involves understanding the problem identifying objectives and collecting relevant data from various sources.

Characteristics

Problem definition

Data identification

Initial exploration

2. Data Preparation

Raw data is cleaned transformed and organized into a suitable format for analysis.

Characteristics

Handling missing values

Removing duplicates

Data normalization

3. Model Planning

In this phase appropriate analytical techniques and models are selected based on the problem.

Characteristics

Selection of algorithms

Statistical methods

Feature selection

4. Model Building

The selected models are implemented using tools such as Python R or Hadoop frameworks.

Characteristics

Training models

Testing data

Performance tuning

5. Evaluation

The model is evaluated using metrics such as accuracy precision recall and error rate.

Characteristics

Model validation

Performance measurement

Refinement if required

6. Deployment

The final model is deployed in real-world systems to generate insights and support decision making.

Characteristics

Integration with systems

Monitoring performance

Continuous improvement

Applications

Fraud detection recommendation systems healthcare analytics

Conclusion

The Data Analytics Lifecycle provides a structured methodology for data analysis ensuring efficiency and accuracy.

Q6. Explain the challenges in Big Data processing. (AKTU 2021, GATE – Distributed Systems & Data Management Concept)

Introduction

Big Data processing involves handling extremely large volumes of data generated from diverse sources at high speed. Although Big Data technologies such as Hadoop and Spark provide powerful tools for processing data, several challenges arise due to the scale complexity and diversity of data. Understanding these challenges is essential for designing efficient Big Data systems and is an important topic in both AKTU and GATE examinations under distributed systems and data management.

Major Challenges in Big Data

1. Data Storage

One of the primary challenges is storing massive amounts of data efficiently. Traditional storage systems are not capable of handling petabytes or exabytes of data.

Issues

Limited storage capacity

High infrastructure cost

Need for distributed storage systems

Solution

Use of Hadoop Distributed File System (HDFS) and cloud storage

2. Data Processing

Processing large datasets requires high computational power and efficient algorithms.

Issues

Slow processing speed in traditional systems

Complex data transformations

Handling real-time data

Solution

Parallel processing using MapReduce and Apache Spark

3. Data Variety

Big Data includes structured semi-structured and unstructured data which makes integration difficult.

Issues

Different data formats

Complex data integration

Lack of uniform schema

Solution

Use of NoSQL databases and schema-less data models

4. Data Velocity

High speed data generation requires real-time processing capabilities.

Issues

Handling streaming data

Real-time analytics challenges

Latency issues

Solution

Stream processing tools like Apache Kafka and Spark Streaming

5. Data Quality (Veracity)

Ensuring accuracy and reliability of data is a major challenge.

Issues

Incomplete or inconsistent data

Noise and redundancy

Data cleaning complexity

Solution

Data preprocessing and validation techniques

6. Security and Privacy

Big Data often contains sensitive information which must be protected.

Issues

Data breaches

Unauthorized access

Privacy concerns

Solution

Encryption authentication and access control mechanisms

7. Scalability

Systems must be able to scale with increasing data volume.

Issues

Performance degradation

Resource management

Load balancing

Solution

Distributed computing and cloud-based infrastructure

Conclusion

Big Data processing faces multiple challenges including storage processing variety velocity security and scalability. Addressing these challenges requires advanced technologies and efficient system design.

Q7. Explain applications of Big Data in various domains. (AKTU 2020, GATE – Data Analytics**Applications Concept)****Introduction**

Big Data has transformed various industries by enabling organizations to analyze large datasets and extract valuable insights. Its applications span across multiple domains such as healthcare finance retail and social media. Understanding these applications helps in recognizing the practical importance of Big Data analytics.

Major Applications of Big Data**1. Healthcare**

Big Data is used to improve patient care and medical research.

Applications

Disease prediction

Personalized treatment

Medical imaging analysis

Example

Predicting disease outbreaks using patient data

2. Banking and Finance

Big Data helps in risk management fraud detection and customer analysis.

Applications

Credit scoring

Fraud detection

Algorithmic trading

Example

Detecting fraudulent transactions in real time

3. E-commerce

Big Data enhances customer experience and business operations.

Applications

Recommendation systems

Customer behavior analysis

Inventory management

Example

Amazon recommending products based on user preferences

4. Social Media

Big Data is widely used to analyze user interactions and trends.

Applications

Sentiment analysis

Trend detection

Influencer identification

Example

Analyzing Twitter data to identify trending topics

5. Transportation

Big Data improves traffic management and route optimization.

Applications

Traffic prediction

Smart transportation systems

Logistics optimization

Example

Google Maps suggesting fastest routes

6. Education

Big Data helps in improving learning outcomes and student performance.

Applications

Student performance analysis

Personalized learning

Dropout prediction

Example

Online platforms analyzing student progress

Conclusion

Big Data has wide-ranging applications across multiple industries and plays a vital role in decision making and optimization.

Q8. Explain Distributed Computing in Big Data. (AKTU 2019, GATE – Distributed Systems Concept)

Introduction

Distributed computing is a fundamental concept in Big Data where data processing is performed across multiple machines instead of a single system. This approach enables efficient handling of large datasets by dividing tasks into smaller sub-tasks and executing them in parallel.

Concept of Distributed Computing

In distributed computing multiple computers called nodes work together as a single system. Each node processes a portion of the data and results are combined to produce the final output.

Characteristics

Parallel processing

Scalability

Fault tolerance

Resource sharing

Working Mechanism

Data is divided into smaller chunks

Each chunk is processed on different nodes

Intermediate results are combined

Final output is generated

Advantages

High performance due to parallel execution

Scalability by adding more nodes

Reliability through fault tolerance

Challenges

Network communication overhead

Synchronization issues

Data consistency

Example in Big Data

Hadoop MapReduce framework where data is processed across clusters

Conclusion

Distributed computing is essential for Big Data processing as it enables efficient handling of large datasets.

Q9. Explain Parallel Processing in Big Data. (AKTU 2021, GATE – Parallel Computing Concept)

Introduction

Parallel processing is a computing technique in which multiple computations are carried out simultaneously to solve a problem faster. In the context of Big Data it plays a crucial role as large datasets require high computational power and time-efficient processing. Parallel processing divides a large task into smaller sub-tasks and executes them concurrently across multiple processors or machines. This concept is part of GATE syllabus under parallel computing and is frequently asked in AKTU examinations.

Concept of Parallel Processing

In parallel processing a large problem is broken down into independent smaller tasks. These tasks are executed simultaneously on different processors and the results are combined to produce the final output.

Key Features

- Simultaneous execution of tasks
- Reduction in processing time
- Efficient utilization of resources
- Improved system performance

Types of Parallel Processing

1. Data Parallelism

Same operation is performed on different pieces of distributed data

Example

Processing chunks of a large dataset across multiple nodes

2. Task Parallelism

Different tasks are executed simultaneously

Example

Running different algorithms on the same dataset

Working in Big Data

Big Data frameworks like Hadoop and Spark use parallel processing to handle massive datasets

Data is divided into blocks

Each block is processed in parallel

Intermediate results are merged

Advantages

Faster execution of large tasks

Scalability

Efficient resource utilization

Challenges

Synchronization issues

Communication overhead

Complex programming

Conclusion

Parallel processing significantly improves the efficiency of Big Data systems by reducing execution time and enabling scalability.

Q10. Explain the role of Cloud Computing in Big Data. (AKTU 2022, GATE – Cloud Computing Concept)

Introduction

Cloud computing provides on-demand access to computing resources such as storage processing power and networking over the internet. In Big Data it plays a significant role by offering scalable and cost-effective infrastructure to store and process large datasets. The integration of cloud computing with Big Data enables organizations to handle data efficiently without investing in expensive hardware.

Role of Cloud Computing in Big Data

1. Scalable Storage

Cloud platforms provide virtually unlimited storage capacity

Example

Amazon S3 Google Cloud Storage

2. Distributed Processing

Cloud enables distributed processing using clusters of virtual machines

Example

Running Hadoop or Spark on cloud platforms

3. Cost Efficiency

Pay-as-you-go model reduces infrastructure cost

Benefits

No need for physical hardware

Reduced maintenance cost

4. Flexibility

Resources can be scaled up or down based on demand

Example

Increasing storage during peak data generation

5. Accessibility

Data and applications can be accessed from anywhere

Example

Cloud-based analytics platforms

6. Reliability and Backup

Cloud provides data replication and backup services

Benefits

High availability

Disaster recovery

Applications

Big Data analytics machine learning real-time processing

Conclusion

Cloud computing enhances Big Data processing by providing scalable flexible and cost-effective solutions

Q11. Explain Big Data Tools and Technologies. (AKTU 2020, GATE – Distributed Systems & Data Tools

Concept)

Introduction

Big Data tools and technologies are designed to handle the storage processing and analysis of large datasets. These tools form the backbone of Big Data ecosystems and enable efficient data management and analytics.

Understanding these tools is essential for both academic and competitive exams.

Major Big Data Tools

1. Hadoop

An open-source framework for distributed storage and processing

Components

HDFS for storage

MapReduce for processing

2. Apache Spark

Fast in-memory data processing framework

Features

Real-time processing

Supports machine learning

3. NoSQL Databases

Non-relational databases designed for flexible data storage

Examples

MongoDB Cassandra

4. Apache Hive

Data warehouse tool for querying large datasets using SQL-like language

5. Apache Pig

High-level platform for data processing using scripting language

6. Apache Kafka

Distributed streaming platform for real-time data processing

Characteristics of Big Data Tools

Scalability

Fault tolerance

High performance

Flexibility

Applications

Data analytics machine learning real-time processing

Conclusion

Big Data tools and technologies provide efficient solutions for handling large datasets.

Q12. Explain the difference between Data Mining and Big Data Analytics. (AKTU 2019, GATE – Data Mining & Databases Concept)

Introduction

Data Mining and Big Data Analytics are closely related concepts but differ in their scope techniques and applications. Both are used to extract useful information from data but operate at different scales and complexities. Understanding their differences is important for designing data-driven systems.

Data Mining

Data Mining refers to the process of discovering patterns and knowledge from structured datasets using statistical and machine learning techniques.

Characteristics

Works on structured data

Uses algorithms like classification clustering association

Focuses on pattern discovery

Example

Finding customer purchasing patterns in a retail database

Big Data Analytics

Big Data Analytics refers to analyzing large volumes of structured semi-structured and unstructured data using advanced tools and technologies.

Characteristics

Handles massive datasets

Uses distributed computing frameworks

Supports real-time and batch processing

Example

Analyzing social media data for sentiment analysis

Key Differences

Data Size

Data Mining works on smaller datasets while Big Data Analytics handles massive datasets

Data Type

Data Mining mainly deals with structured data while Big Data Analytics handles all types of data

Tools

Data Mining uses traditional tools while Big Data Analytics uses Hadoop Spark etc

Processing

Data Mining is mostly sequential while Big Data Analytics uses parallel processing

Conclusion

Data Mining and Big Data Analytics are complementary techniques where Data Mining focuses on pattern extraction and Big Data Analytics provides scalable solutions for handling large datasets.