

# Nirma University

## Institute of Technology

Semester End Examination (IR), December - 2022

B. Tech. in Computer Science and Engineering, Semester-VII

2CS702 Big Data Analytics

Roll /  
Exam No.  
Time: 3 Hours

Supervisor's Initial  
with Date

Max Marks: 100

- Instructions:
1. Attempt all questions
  2. Figures to right indicate full marks
  3. Assume necessary data.
  4. Use section-wise separate answer book.
  5. Draw neat sketches wherever necessary.

### SECTION - I

**Q-1. Do as directed. [16]**

- (A) Explain the terms, (8)  
CO2 BL2
- 1) HDFS block
  - 2) Input split
  - 3) RDD
  - 4) Record reader
- (B) Explain the sort and shuffle phase in the MapReduce programming model. (8)  
CO2 BL3
- Discuss the problems that may be occurred. How to optimize the function of the sort and shuffle process?

**Q-2. Do as directed. [18]**

- (A) How CAP theorem is supported in NoSQL databases? List out the pros, (6)  
CO2 BL3
- cons and database names of NoSQL databases.
- (B) Explain the types of operations performed on RDD in Spark architecture. (6)  
CO2 BL4

**OR**

- (B) Describe any three types of input format in map-reduce programming with (6)  
CO2BL4
- its appropriate application.
- (C) Explain the anatomy of file write operation in HDFS architecture. (6)  
CO2 BL2

**OR**

- (C) Describe the major limitations of HDFS architecture. (6)  
CO2 BL2

**Q-3. Do as directed. [16]**

- (A) Briefly discuss scalable platforms with their features and drawbacks. (8)  
CO1BL3
- (B) What are the major drawbacks of Hadoop v1.0? Describe the tool used to (8)  
CO3BL3
- overcome the listed drawbacks.

**SECTION - II****Q-4. Do as directed. [18]**

- (A) Consider the following dataset, (8)  
CO3BL5 ***Employee (Emp\_id, Emp\_name, Qualification, Experience).***  
Write a MapReduce pseudo code to filter employee names who are Ph.D. with more than 10 years of experience.
- (B) Define Fault Tolerance. At what extent fault tolerance is supported in (5)  
CO3 BL4 horizontal scaling platforms and vertical scaling platforms? Also justify your answer.
- (C) Explain the types of data analytics with suitable case studies. (5)  
CO1BL4

**Q-5. Do as directed. [16]**

- (A) Discuss any one application in detail which present the limitation of (6)  
CO1BL3 MapReduce programming. How Apache Spark overcome this limitation? Explain the features and advantages of using Spark.
- (B) Differentiate between MapReduce, Hive, and PIG architectures. (5)  
CO2BL2

**OR**

- (B) Differentiate between MapReduce, Cassandra, and RDBMS. (5)  
CO2BL2
- (C) What is the default size of the HDFS block? How to configure HDFS block (5)  
CO2BL5 size? If the input file size is of 1TB and the block size is 0.5 GB, then find the number of mapper tasks. Assume MapReduce input split size is equal to HDFS block size.

**OR**

- (C) What is the difference between MapReduce Input Split and HDFS block? If (5)  
CO2BL5 the number of mappers tasks = 100, then how many MapReduce Input Splits are there?

**Q-6. Do as directed. [16]**

- (A) Consider a blog database in MongoDB. Create collection posts in the blog (5)  
CO4BL 4 database. Assume that each document can be described by title, description, category, likes, tags, and date attributes. Assume data as per the requirement. Write a query for the followings,
- 1) Print all posts posted in the month of November 2022.
  - 2) Update the category of the posts with 'Recent' which were posted one day before the current date.
  - 3) Delete all posts which have no tags.
  - 4) Sort all the posts based on likes.
  - 5) Add one post under the category News.
  - 6) Filter all posts which have likes more than 15.
  - 7) Display all posts of the collection posts.
  - 8) Filter all posts which contain pandemic words in the description.