

**JADAVPUR UNIVERSITY**  
**DEPARTMENT OF INFORMATION TECHNOLOGY**  
**B.E. INFORMATION TECHNOLOGY, THIRD YEAR**  
**SECOND SEMESTER EXAM 2024**  
**SUBJECT: BIG DATA (HONS.)**

Full Marks: 100

Time: 3Hrs

<b>CO1 : (20)</b>	<p><b>Answer any five</b> <span style="float: right;"><b>5 X 4 =20</b></span></p> <p>Q.1</p> <ol style="list-style-type: none"> <li>i) What is the function of Job tracker and Task Tracker in HDFS?</li> <li>ii) What do you mean by a Hadoop Cluster and explain the various services hosted by the master node of a Hadoop Cluster.</li> <li>iii) What are the different modes in which Hadoop can be installed and what is the use of each mode from application and developer point of view?</li> <li>iv) Explain the uses of Name node, Data node and Secondary Name node in Hadoop Distributed File system.</li> <li>v) How does YARN allocate containers to the jobs?</li> <li>vi) Name the services running in slave part of Hadoop architecture.</li> </ol>
<b>CO2 : (20)</b>	<p><b>Answer any four</b> <span style="float: right;"><b>4 X 5 =20</b></span></p> <p>Q.2</p> <ol style="list-style-type: none"> <li>i) Consider a scenario where a company wants to analyze sales data alongside customer demographic data to understand buying patterns. Discuss how MapReduce can perform the relational join operation on these datasets.</li> <li>ii) Write a Map Reduce algorithm to get the Dot Product of two Large Vectors, assuming only non-zero elements of those vectors are given in input files and output file should show only non-zero entries ( assuming two vectors are same size) <math>v1=[ 5 4 0 1 2]</math> and <math>v2=[ 4 2 1 0 6]</math>.</li> <li>iii) Explain the role of driver code, mapper code and reducer code within a map reduce program model by a suitable example.</li> <li>iv) With a neat diagram describe the phases of Map Reduce Job processing.</li> <li>v) Explain how data is partitioned before it is sent to the reducer if no custom partitioner is defined in Hadoop.</li> </ol>
<b>CO3 : (30)</b>	<p><b>Answer any two</b> <span style="float: right;"><b>15 X 2 =30</b></span></p> <p>Q.3</p> <ol style="list-style-type: none"> <li>i) <ol style="list-style-type: none"> <li>a) State the procedural differences between sharding and replication.</li> <li>b) Explain the concept of quorum based write consistency in Cassandra.</li> <li>c) How the BASE properties of NoSQL databases do impact the design and scalability of distributed systems compared to the ACID properties of traditional relational databases?</li> <li>d) How can downtime be reduced in NoSQL databases?</li> </ol> </li> </ol> <p style="text-align: right;"><b>4+5+4+2=15</b></p>

	<p>ii)</p> <ul style="list-style-type: none"> <li>a) State how does MongoDB handles the issue of No schema, No DDL, No joins, and No transactions.</li> <li>b) Explain briefly the following features in MongoDB: Hierarchical Objects, MongoD, MongoS, Collections, BSON documents.</li> <li>c) Can you explain the concept of aggregation in MongoDB?</li> </ul> <p style="text-align: right;"><b>5+7+3 = 15</b></p> <p>iii)</p> <ul style="list-style-type: none"> <li>a) How can organizations effectively scale and manage HBase region Servers to accommodate the growth of their data stores, while maintaining fast data lookup and fault tolerance across the cluster?</li> <li>b) What is the Hierarchy of Tables in Apache HBase?</li> <li>c) What are the operational commands of HBase? Mention few.</li> <li>d) Define compaction in HBase?</li> <li>e) How does sorted maps work in HBase</li> </ul> <p style="text-align: right;"><b>5 X 3 =15</b></p>
<p><b>CO4 : (30)</b></p>	<p><b>Answer All</b> <span style="float: right;"><b>15 X 2 =30</b></span></p> <p>Q.4</p> <p>i)</p> <ul style="list-style-type: none"> <li>a) Consider a table named “website_traffic” in Hive, which contains records of website traffic with the following schema:  <i>session_id</i>: Unique identifier for each session  <i>browser_type</i>: Type of browser used by the visitor  <i>session_duration</i>: Duration of the session in seconds                      Formulate a Hive query to calculate the average session duration for each browser type used by visitors, filtering out sessions shorter than 30 seconds, and displaying the results along with the browser types.</li> <li>b) Can you explain the role and purpose of Thrift Server, Metastore and Managed tables in context of Apache Hive?</li> <li>c) Where does the data of a Hive table get stored?</li> </ul> <p style="text-align: right;"><b>6+6+3 =15</b></p> <p>ii)</p> <ul style="list-style-type: none"> <li>a) How does Spark leverage lineage information to recover from faults during computation?</li> <li>b) Mention few Spark Control operations with their functionalities.</li> <li>c) How does the DAG scheduler optimize task scheduling for performance improvement?</li> </ul> <p style="text-align: right;"><b>5+5+5 =15</b></p>