



**Continuous Assessment Test(CAT) – II APRIL 2024**

Programme	:	MCA	Semester	:	Winter 23-24
Course Code & Course Title	:	PMCA607L & Big Data Analytics	Class Number	:	CH2023240501412
Faculty	:	Dr.K.Madheswari	Slot	:	E2+TE2
Duration	:	90 minutes	Max. Mark	:	50

**General Instructions:**

- Write only your registration number on the question paper in the box provided and do not write other information.
- Use statistical tables supplied from the exam cell as necessary
- Use graph sheets supplied from the exam cell as necessary
- Only non-programmable calculator without storage is permitted

**Answer all questions**

Q. No	Sub Sec.	Description	Marks
1		<p>You are tasked with designing a comprehensive schema and queries for a multi-functional application in Cassandra, utilizing various data structures and concepts. The application aims to support user management, messaging, blogging, and task tracking functionalities. Construct a set of tables along with queries to perform operations as described below:</p> <p><b>User Management (5 marks):</b></p> <p>Design a schema for user management with a table named users. Each user should have a user_id generated as a UUID, a username, and a set of contact_ids representing the user's contacts. Implement the schema along with appropriate data types and constraints.</p> <p><b>Messaging (5 marks):</b></p> <p>Develop a schema for messaging functionality with a table named messages. Each message should be uniquely identified by a message_id generated as a UUID. The table should include columns for sender_id, receiver_id, timestamp, and a map of message metadata containing key-value pairs for message attributes. Construct the schema ensuring efficient storage and retrieval of messages.</p> <p><b>Blogging (5 marks):</b></p> <p>Create a schema for blogging capabilities with a table named posts to store blog posts. Each post should possess a post_id</p>	20

	<p>generated as a UUID, along with author_id, title, content, timestamp, and a set of tags associated with the post. Design the schema to support efficient querying of posts based on authors and tags.</p> <p><b>Task Tracking (5 marks):</b></p> <p>Enhance the schema to incorporate task tracking features with a table named tasks. Each task should be identified by a task_id generated as a UUID and include attributes such as title, description, due_date, and a list of assigned_users. Implement the schema to facilitate efficient management and retrieval of tasks.</p> <p>Ensure that your schema designs are optimized for scalability, data integrity, and query performance in a distributed Cassandra environment. Additionally, provide sample queries for each functionality to demonstrate the usage and effectiveness of your schema designs.</p>	
2	<p>Consider a scenario where you're analyzing a vast dataset containing user interactions on a social media platform. Each record in the dataset represents a user interaction event, including the user ID and the type of interaction (e.g., <u>likes, comments, shares</u>). Your goal is to process this dataset using MapReduce to derive meaningful insights</p> <p><b>Dataset:</b></p> <p>Record 1: (User1, Like)  Record 2: (User2, Comment)  Record 3: (User1, Share)  Record 4: (User3, Like)  Record 5: (User2, Like)  Record 6: (User1, Like)  Record 7: (User2, Comment)  Record 8: (User1, Share)  Record 9: (User3, Like)  Record 10: (User2, Like)</p> <ul style="list-style-type: none"> <li>● Apply the Map function to the given dataset, transforming each user interaction record into intermediate key-value pairs.</li> <li>● Illustrate how the intermediate key-value pairs generated in the Map phase are shuffled and sorted based on keys, preparing them for the Reduce phase.</li> <li>● Perform the Reduce operation on the shuffled and sorted key-value pairs to derive insights from the dataset.</li> <li>● Calculate relevant metrics or statistics based on the aggregated results obtained from the Reduce phase.</li> <li>● Extend your analysis by incorporating more complex operations into the MapReduce workflow. Consider scenarios such as filtering out specific user interactions, identifying user engagement patterns, or detecting anomalies in user behavior.</li> </ul> <p>Apply additional Map and Reduce functions as necessary to address these advanced analysis requirements.</p>	10

3	<p>You are tasked with writing a Hadoop MapReduce program to perform an inner join operation between two tables based on the "Student ID" field. The join must satisfy two filtering conditions simultaneously:</p> <p>a) The year of birth must be greater than (&gt;) 1990. b) The score of course 1 must be greater than (&gt;) 80, and the score of course 2 must be no more than (&lt;=) 95.</p> <p><b>Sample Data:</b></p> <p><b>Student Table:</b></p> <table><thead><tr><th>Student ID</th><th>Name</th><th>Year of Birth</th></tr></thead><tbody><tr><td>20170126453</td><td>Kristalee Copperwaite</td><td>2000</td></tr><tr><td>20170433596</td><td>Roeberta Naden</td><td>1997</td></tr></tbody></table> <p><b>Score Table:</b></p> <table><thead><tr><th>Student ID</th><th>score(course1)</th><th>Score (course2)</th><th>Score(course3)</th></tr></thead><tbody><tr><td>20170126453</td><td>93</td><td>97</td><td>80</td></tr><tr><td>20170140241</td><td>86</td><td>85</td><td>87</td></tr><tr><td>20170433596</td><td>82</td><td>60</td><td>80</td></tr></tbody></table> <ul style="list-style-type: none"><li>● Design a Hadoop MapReduce program to perform the inner join operation between the Student and Score tables.</li><li>● Filter records where the year of birth is greater than 1990.</li><li>● Filter records where the score of course 1 is greater than 80 and the score of course 2 is no more than 95.</li><li>● Utilize the provided sample data for both the Student and Score tables as input to the MapReduce program.</li><li>● Generate the join result, including the Student ID, Name, Year of Birth, and scores for courses 1, 2,</li></ul>	Student ID	Name	Year of Birth	20170126453	Kristalee Copperwaite	2000	20170433596	Roeberta Naden	1997	Student ID	score(course1)	Score (course2)	Score(course3)	20170126453	93	97	80	20170140241	86	85	87	20170433596	82	60	80	10
Student ID	Name	Year of Birth																									
20170126453	Kristalee Copperwaite	2000																									
20170433596	Roeberta Naden	1997																									
Student ID	score(course1)	Score (course2)	Score(course3)																								
20170126453	93	97	80																								
20170140241	86	85	87																								
20170433596	82	60	80																								
4	<p>Let's consider a sample table named "users" with the following schema:</p> <pre>CREATE TABLE users (   user_id UUID PRIMARY KEY,   username TEXT,   email TEXT,   age INT );</pre> <ul style="list-style-type: none"><li>● Create a sample keyspace and table in Cassandra to use for performing CRUD operations.</li><li>● Define the schema for the table, including primary key and any secondary indexes required for querying.</li><li>● Write functions in your Cassandra program to implement the following CRUD operations: Create: Insert new records into the table. Read: Retrieve records based on primary key or secondary Index. Update: Modify existing records with new data.</li></ul>	10																									



		<p>Delete: Remove records from the table.</p> <ul style="list-style-type: none"><li>● Test the Cassandra program with different datasets, covering various scenarios and edge cases.</li><li>● Verify the results of each CRUD operation to ensure that data is inserted, retrieved, updated, and deleted correctly.</li><li>● Validate the program's behavior when handling corner cases, such as missing records or invalid input parameters.</li></ul>	
--	--	---	--

\*\*\*\*\*All the best \*\*\*\*\*