Started on	Started on Wednesday, 10 November 2021, 2:32 PM
State	State Finished
Completed on	<b>Completed on</b> Wednesday, 10 November 2021, 4:00 PM
Time taken	<b>Fime taken</b> 1 hour 28 mins
Grade	Grade Not yet graded

Information

cards should contain 52 cards with four suits 1.Diamond 2.Clubs 3.Heart and 4.Spade. The details of the cards he is having with him are recorded in a file named deckdetails.csv with 10000 entries. The data of every card is mentioned in the following format [SUIT\_NAME RANK INTEGER]. For e.g.: SPADE 1. The suit name can be HEART or SPADE or CLUB or DIAMOND respectively and the rank is an integer value ranging from 1 to 13 for the values [A, 2, A second hand Playing Card seller James is having 10000 assorted cards. With which he would like to build valid decks of cards for reselling. A deck of 3.....J, Q, K]. A sample data set of cards in the file deckdetails csv is shown below in Fig.01



DIAMOND 3	CLUB 6	HEART 8	HEART 9	SPADE 6	HEART 7	DIAMOND 11	DIAMOND 13	CLUB 1	CLUB 2	HEART 4	CLUB 3	SPADE 11	
DIAMOND 1	CLUB 4	CLUB 5	HEART 12	CLUB 7	CLUB 8	HEART 2	CLUB 10	CLUB 11	CLUB 13	HEART 1	HEART 5	HEART 3	
SPADE 1	SPADE 2	SPADE 3	SPADE 4	DIAMOND 5	SPADE 5	<b>DIAMOND 12</b>	CLUB 12	CLUB 9	SPADE 7	<b>DIAMOND 8</b>	<b>DIAMOND 10</b>	SPADE 9	

Fig:01

Based on the above scenario answer the following.

Complete Marked out of 6 (6 Marks-[Ap/C,2])(CO:4; PO:1.4.1)

Assume, James process the sample data mentioned in fig:01 using MapReduce paradigm to generate a deck but he found few cards are missing. If the split size of the MapReduce job is 4, then depict the output of the Mapper phase, Shuffle and sort phase and Reduce Phase.

**DIAMOND 8** HEART 9 **CLUB 13** SPADE 9 FINAL **DIAMOND-8** REDUCE CLUB-13 SHUFFLE & SORT **DIAMOND 12** DIAMOND 13 **DIAMOND 10 DIAMOND 11** DIAMOND 3 **DIAMOND 5 DIAMOND 8 DIAMOND 1** HEART 1 HEART 12 CLUB 12 CLUB 13 CLUB 11 CLUB 10 CLUB 2 CLUB 5 CLUB 6 CLUB 8 CLUB 3 CLUB 4 CLUB 7 **DIAMOND 12 DIAMOND 10** DIAMOND 1 CLUB 4 DIAMOND 8 **DIAMOND 5** HEART 12 SPADE 5 SPADE 7 SPADE 9 MAPPER SPADE 2 SPADE 3 SPADE 4 CLUB 12 HEART 2 SPADE 1 CLUB 5 CLUB 8 CLUB 9 CLUB 7 **JIAMOND 12 DIAMOND 10 DIAMOND 5** DIAMOND 8 **DIAMOND 1 JEART 12** SPADE 5 SPADE 9 SPADE 7 SPADE 3 SPADE 4 CLUB 12 SPADE 2 **JEART 2** SPADE 1 CLUB 10 CLUB 11 **CLUB 13 JEART 1** CLUB 5 CLUB 8 CLUB 9 CLUB 7 CLUB 4 INPUT

		HEART-9										SPADE-9					
HEART 2	HEART 3	HEART 4	HEART 5	HEART 7	HEART 8	HEART 9		SPADE 1	SPADE 11	SPADE 2	SPADE 3	SPADE 4	SPADE 5	SPADE 6	SPADE 7	SPADE 9	
CLUB 10	CLUB 11	CLUB 13	HEART 1		HEART 5	HEART 3	DIAMOND 3	CLUB 6		HEART 8	HEART 9	SPADE 6	HEART 7		<b>DIAMOND 11</b>	DIAMOND 13	CLUB 1
HEART 3	DIAMOND 3	CLUB 6	HEART 8	HEART 9	SPADE 6	HEART 7	DIAMOND 11	DIAMOND 13	CLUB 1	CLUB 2	HEART 4	CLUB 3	SPADE 11				

CLUB 2

Marked out of 2 Complete

For finding the missing cards from the deck. The Map Reducer, first sorts and shuffles the input cards in various phases and then outputs the cards in all the four suits. MRUnitTesting is used to validate the execution of MapReduce separately in all its major phases to ensure the correctness of its working. Before testing any process the below given code annotated @Before is exeuted to initialize and set up the environment variables required for testing. Help James in completing the code so that he can set up his testing environment successfully in order to ensure that the execution is working as expected. Also justify the significance of the Code in Line 1 and Line 2

```
WordCountReducer reducer = new deck_countReducer();
reduceDriver = new ReduceDriver<Text, IntWritable, Text, IntWritable>();
                                                                                                               mapDriver = new MapDriver<LongWritable, Text, Text, IntWritable>();
                                                                                                                                                                                                                                                                                                                                                  mapReduceDriver = new MapReduceDriver<LongWritable,
                                                                           WordCountMapper = new deck_countMapper();
                                                                                                                                                                                                                                                                                                                                                                                        Text, Text,IntWritable, Text, IntWritable>();
                                                                                                                                                                                                                                                                                                                                                                                                                              .....\ Line 1
                                                                                                                                                                                                                                                                                                                                                                                                                                                              .....\ Line 2
                                                                                                                                                                                                                                                                         reduceDriver.setReducer(reducer);
                                                                                                                                                          mapDriver.setMapper(mapper);
public void setUp()
```

(2 Marks-[Ap/C,2])(CO:4; PO:1.4.1)

Line 1 and Line 2:

mapReduceDriver.setMapper(mapper);

mapReduceDriver.setReducer(reducer);

Significance of the Code in Line 1 and Line 2

when we want to test the MapRedcue job as a whole, we need to test both the mapper and the reducer. Therefore we need to set both the mapper and reducer class to the mapReduceDriver.

11/15/21, 8:39 PM

Complete
Marked out of 2

In local testing, among all possible job configuration parameters, setting the input and output file path to all phases is important as shown below.

```
JobConf conf = new JobConf(getConf(), getClass());
conf.setJobName("DeckCount");
FileInputFormat.addInputPath(conf, new
Path(args[0]));
FileOutputFormat.setOutputPath(conf, new
Path(args[1]));
conf.setOutputKeyClass(Text.class);
conf.setOutputValueClass(IntWritable.class);
conf.setMapperClass(DeckCountMapper.class);
conf.setCombinerClass(DeckCountReducer.class);
conf.setReducerClass(DeckCountReducer.class);
JobClient.runJob(conf);
```

Justify the statement in few lines.

(2 Marks-[An/C,2])(CO:4; PO:4.1.1)

Among the possible job configuration parameters, we set the input and output file paths, the mapper, reducer and combiner classes, and the output types (the input typesare determined by the input format, which defaults to TextInputFormat and has Long Writable keys and Text values).

Complete Marked out of 6

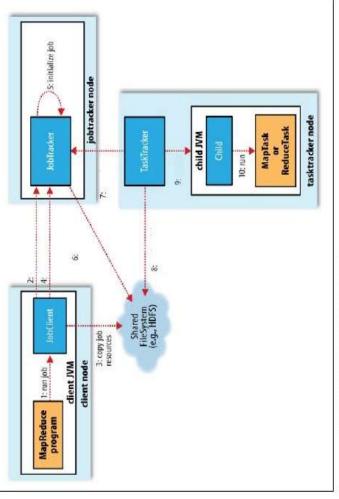
The anatomy of classic Hadoop job run can be listed as a sequence of steps as given below. Map the steps in the given diagram appropriately: Submit job

Retrieve job resources Get new job ID

Lannch

Retrieve the input files

Heartbeat (Returning task)



How Hadoop runs a MapReduce job

(6 Marks-[An/C,2])(CO:4; PO:4.1.1)

- 1. run job
- 2. Get new job ID
- 3. copy job resource
- 4. submit job

5. initialize job

6. Retrieve the input files

7.Heartbeat (Returning task)

8. Retrieve job resources

9.launch

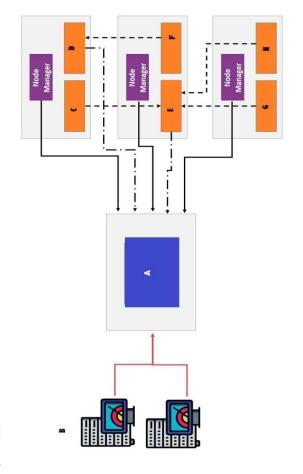
10. run

Complete

Marked out of 4

Apache Yarn is also called as Hadoop 2.0 which is mainly used for handling MapReduce, resource management and job scheduling in an efficient way. The following diagram depicts the architecture of Apache Yarn. Identify and Map the suitable components listed below in its appropriate blocks of the diagram: Question **5** 

- 1) Client
- 2) Resource Manager
- 3) Application Master
- 4) Container
- 5) Application Manager
- 6) Mapper Reducer



. (CO:4; PO: 1.4.1) (4 Marks -U/C,1])

- A Resource Manager
- B Client
- C Container
- D Application Master

E - Application Master

F - Container

G - Container

H - Container

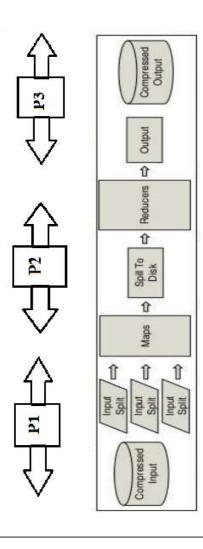
Information

A climate model is a mathematical representation of climate systems based on various factors that impacts the climate of the Earth. Basically, it describes project climate conditions by simulating the climate changes based on factors that affect climate. Regional Climate Model Evaluation System (RCMES) is the interaction of various drivers of climate like ocean, sun, atmosphere, etc. to provide an insight into the dynamics of the climate system. It is used to required for analysis and evaluation of the climate output model against billions of remote sensing data present in various external repositories.

Based on the above scenario answer the following questions

Marked out of 3 Complete

RCMET (Regional Climate Model Evaluation Toolkit). They used Hadoop for loading billions of remote sensor data and evaluating them based on different parameters required for modelling climate. While storing the data they have applied a suitable compression technique. Identify the different phases PI, P2, NASA's Jet Propulsion Laboratory has developed RCMES with two components such as RCMED (Regional Climate Model Evaluation Database) and P3 given in the MapReduce pipeline while applying the process of compression.



(3 Marks -[Ap/C,2])(CO-3;PO;1.4.1)

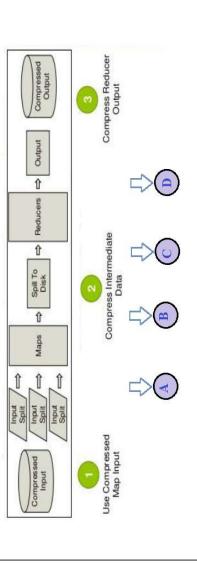
Р1 - Мар

P2 - Shuffle and sort

P3 - Reduce

Marked out of 2 Complete

Derive the type of Input / Output obtained in A, B, C and D during the process of compressing intermediate data given in the workflow..



(2 Marks -[Ap/C,2]) (CO:3; PO: 1.4.1)

- B Compress output
- C Decompress input
- D Compress output

Marked out of 4 Complete

Compression reduces the space needed to store files in the file system and speeds-up the data transfer from one node to another in the network. Match the following compression techniques with its algorithm, tool and its feature.

S S	Compression Format	Algorithm	Tool	Feature
-	Gzip	124	Lzop	Non Splittable
2	077	DEFLATE	NA	Best Compression ratio
33	LZ4	Bzip2	Bzip2	Slower than LZO and LZ4
4	Bzip2	027	Gzip	Balanced Compression/Decompression

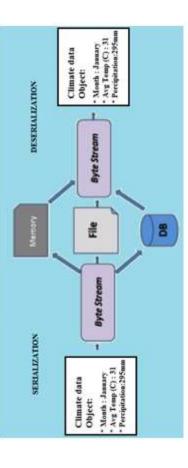
(4 Marks -[An/C,2]) (CO:3; PO: 2.3.1)

## S.NO Compression Format Algorithm Tool Feature

slower than Izo and Iz4	non spilttable	balanced compression/ decompression	best compression ratio
gzip	doz	na	bzip2
deflate	CZO	LZ4	bzip2
Gzip	PZO	LZ4	Bzip2
_	7	က	4

Question 9 Complete Marked out of 1

Serialization and deserialization is applied in inter process communication between objects while storing or transmitting them from one file to another.

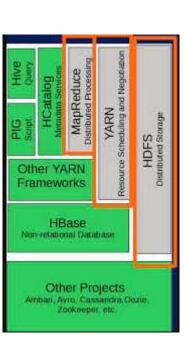


Determine the format to which the data objects are converted while saving the state of the object. (1 Marks-[An/C,2])(CO:3;PO:2.3.1)

**JSON FILE FORMAT** IS USED WHILE SAVING

Information

Hadoop is an open-source framework to store and process Big Data in a distributed environment. The Hadoop ecosystem contains different sub-projects (tools) that are used to help Hadoop modules.



Marked out of 4 Question 10 Complete

Match the following Hadoop modules based on their data model, architecture Rad-Write modes, Data Access etc:

Hadoop	Data Model	Architecture	R/W	Data Access
HBASE	Data is partitioned	Metastore	Better WRITE Cell -level	Cell -level
CASSANDR/	CASSANDRA Nested relational	Master based	Query	Schema level
HIVE	Datawarehouse Tool	Multi-query approach Scripts	n Scripts	Row level
PIG	Columns are grouped Masterless	d Masterless	Better Read	Better Read Schema level

(CO:5; PO:4.1.1)(4 Marks-[An/C,2])

SS		) sve	eve	
Data Access	ad Cell - level	schema Level	Better Write Schema Level	Row Level
e R/W	d Better Re	Query	Better Wr	Scripts
Architecture R/W	Master base	Masterless Query	Metastore	Multi-Query Scripts
Data model	Columns are grouped Master based Better Read Cell - level	CASSANDRA Nested Relational	Data Warehouse Tool Metastore	Data is partitioned
Hadoop	HBASE	CASSAND	HIVE	PIG

Complete Marked out of 6

```
The following code is used to create table in HBASE.

HBaseAdmin admin= new HBaseAdmin(config);

HColumnDescriptor []column;

column= new HColumnDescriptor[2];

column[0]=new HColumnDescriptor("columnFamily1:");

HTableDescriptor desc= new HTableDescriptor(Bytes.toBytes("MyTable"));

desc.addFamily(column[0]);

admin.createTable(desc);
```

Refine the HBASE code with suitable logic to add two more columns
Refine the HBASE code using GET() method to return the corresponding record with the help of Key
Refine the HBASE code using PUT() method to return the corresponding record
(CO:5; PO:4.1.1) (6 Marks-[An/C,3])

HBASE code using GET() method to return the corresponding record with the help of Key

Get get = new Get(Bytes.toBytes("row1"));

Result r = htable.get(get);

byte[] b = r.getValue(Bytes.toBytes("cf"), Bytes.toBytes("attr")); // returns current version of value

Refine the HBASE code using PUT() method to return the corresponding record

Put put = new Put (Bytes.toBytes(row));

put.add(Bytes.toBytes("cf"), Bytes.toBytes("attr1"), Bytes.toBytes(data));

htable.put(put);

11/15/21, 8:39 PM

Marked out of 2 Complete

"Create Keyspace Command in cassandra is used to create keyspace in Cassandra.Replication factor is the number of replicas of data placed on different nodes"

Construct a keyspace and determine the replication factor to ensure no single point of failure. Sometimes, the server can be down, or network problem can occur, then other replicas should be able to provide service with no failure.

(CO:5; PO:4.1.1) (2 Marks-[AP/C,2])

Create keyspace KEY\_NAME with replication={'class':'SimpleStrategy','replication\_factor':3};

Information

"Hive is a data warehousing system to store structured data on Hadoop file system. Hive Buckets in hive is used in segregating of hive table-data into multiple files or directories. It is used for efficient querying. Partitions is a way to organizes tables into partitions by dividing tables into different parts based on partition keys"

Based on this answer the following questions

Marked out of 2 Complete

Query for Create Dept\_bucket with column names such as first\_name, job\_id, department name, department id, salary and country is given below

create table Dept\_bucket first\_name string , job\_id int , dept\_name string , dept\_id int ,salary string ,country string );

Construct four buckets for the "country" attribute.

(CO:5; PO:4.1.1)(2 Marks-[AP/C,2])

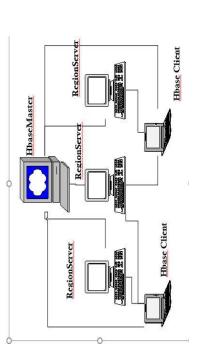
CREATE TABLE Dept\_bucket PARTITIONED BY (string ,int, string ,int ,string ) CLUSTERED BY (first\_name ,job\_id ,dept\_name,dept\_id,salary, country | SORTED BY (country[ASC]) INTO 4 BUCKETS;

Question **14** Complete

11/15/21, 8:39 PM

Marked out of 4

Abstract the roles of the HBase master, Region Server and HBase Client given in the architecture below.



(CO:5; PO:4.1.1)(4 Marks-[U/C,1])

## HBase master

- HMaster is the implementation of a Master server in HBase architecture.
- It acts as a monitoring agent to monitor all Region Server instances present in the cluster and acts as an interface for all the metadata changes.
  - In a distributed cluster environment, Master runs on NameNode. Master runs several background threads.

## Region Server

- When Region Server receives writes and read requests from the client, it assigns the request to a specific region, where the actual column family resides.
  - However, the client can directly contact with HRegion servers, there is no need of HMaster mandatory permission to the client regarding communication with HRegion servers.
    - The client requires HMaster help when operations related to metadata and schema changes are required.
      - HRegionServer is the Region Server implementation

## **HBase Client**

Client wants to write data and in turn first communicates with Regions server and then regions

11/15/21, 8:39 PM

Marked out of 2 Complete

Determine the PIG Grunt Queries for the following operations.

To specify the number of tuples from the relations.

To extract the data from file system to relation. To combine the data in a single relation,

To arrange the relations in sorted form

• LIMIT - To specify the number of tuples from the relations

(CO:5; PO:4.1.1)(2 Marks-[U/C,1])

• LOAD - To extract the data from file system to relation0.

• JOIN - To combine the data in a single relation.

**ORDER BY** -To arrange the relations in sorted form

¥

18CS005 BigDataAnalytics 24-09-21 PREVIOUS ACTIVITY

11/15/21, 8:39 PM