

(1)

DHARMSINH DESAI UNIVERSITY, NADIAD.  
FACULTY OF TECHNOLOGY  
ONLINE SESSIONAL EXAMINATION

B.Tech (CE) sem → 7

subject : Big Data Analytics

Roll no. → 142

Signature → Pankuri

Date : 28/8/2020

Time : 10:00 to 11:15

Total pages : 10

Q. 1

(a) responsibilities of namenode.

- manages & maintains datanodes
- receives block report & heartbeat from all datanodes.
- records metadata about blocks

(b) checkpointing.

- process involving merging of fsimage along with latest editlog for creating new fsimage for namenode to have latest metadata
- This task is performed by secondary name node or standby name node.

(c) components of resource manager

① components - scheduler

- allocates resources to applications & does not

(2)

offers guarantees about restarting failed tasks.

→ it has pluggable policy like capacity schedulers & fair schedulers

(3) application manager

→ accepts job submissions

→ negotiates first confines for executing application specific Application Masters

→ provides service for restarting Application Masters on failure.

(d) IoT deals with million of devices that collect & operate very large

~~structured~~ unstructured data. This is mainly unstructured data.

→ Big Data helps IoT to process this large data ~~on~~ on real-time basis and storing it ~~on~~ using different storage technologies.

(e) schema on read

→ This is a strategy for analysing the data in ~~on~~ tools like hadoop in which data is applied to schema when it is pulled out from a storage location, rather than when it is stored.

→ Data is loaded as it is and can be read in different manners.

(3)

example : hive , impulse

→ It is used when raw / atomic delta is being stored or we need flexibility on consumption of delta.

(f) sharding

- process of storing delta in multiple machines.
- helps mongoDB meeting the demands of delta growth
- As delta size increases, single machine may not handle it
- sharding solves this issue.

policies :

① vertical scaling

→ refers to adding more resources to server as on demand.

② horizontal scaling.

→ mongoDB supports horizontal scaling through sharding.

→ here dataset is divided & loaded over multiple servers.

Q.2

(b) Task:

- execution of Mapper or Reducer on a slice of data
- also called Task-In-Progress (TIP)
- This means processing the data is in progress on mapper or reducer.

Task attempt:

- particular instance of attempt to execute task on a node.
- it is possible that a machine can go down.
- if it happens, task is rescheduled on some other node. But this cannot be done infinite time.
- default value of task attempt is 4.
- if a task fails 4 times, job is considered as failed job.
- this value can be changed for huge or high priority job.

Failed task:

- A task that is failed more than no. of task attempt times.

~~decided~~

(S)

→ it is a task attempt that is completed but with an unexpected status value. i.e. task that generate ~~err~~ errors.

killed task →

→ a task attempt which is a duplicate of a task attempt which was started as a part of speculative execution, & is killed by Hadoop because of below reasons.

- it didn't report progress during timeout
- scheduler needed slot for other pool or queue.
- its result is not needed because it has completed on other place.

speculative execution :-

- execution of same task is started in multiple boxes.
- The first task to finish is winning & other tasks are killed duplicate.

input split :-

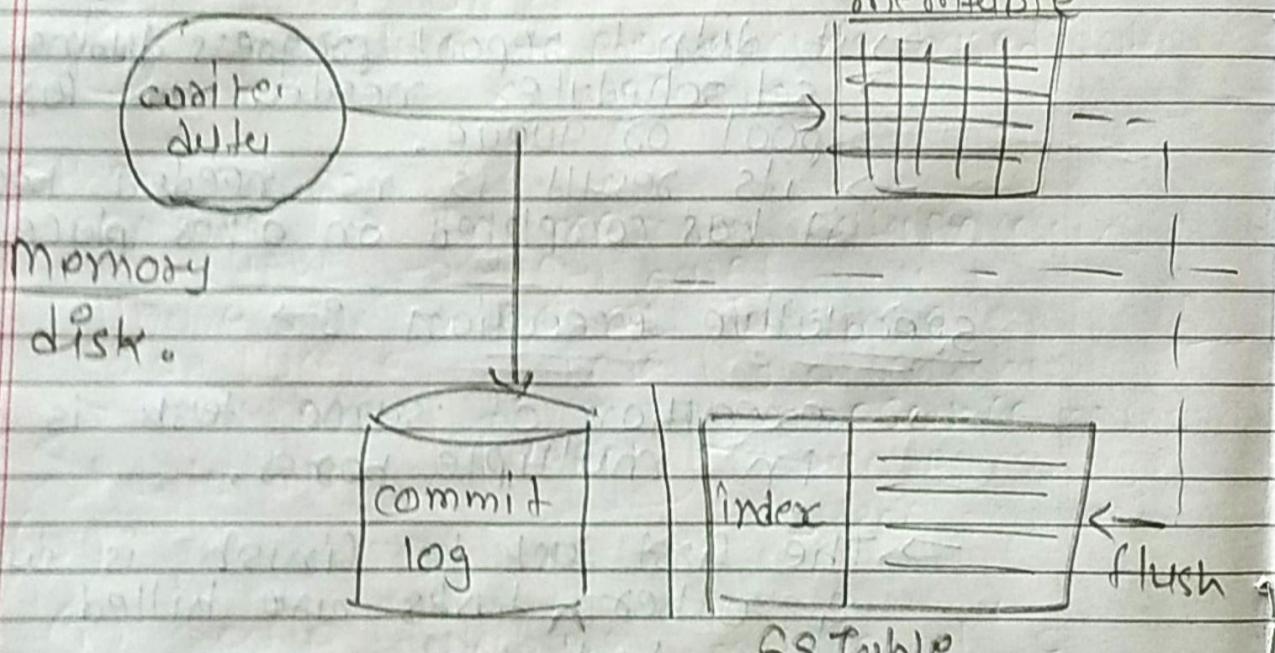
- it is created by input format
- logically represents data processed by

(6)

individual mappers.

- one map task created for each split.
- split divided in records & each record is processed by mappers.
- its length is measured in bytes.
- size is approx. equal to block size.

(c) write operation in cassandra.



- when a write operation is completed, it is firstly written in commit log.
- Then, it is pushed in memory resident data structure called memtable having a predefined threshold value.

(2)

- A node responds ~~with~~ successful message to coordinator only when write is completed in memtable & commit log.
- When number of objects reaches a threshold, at that time, the contents of memtable are flushed to disk in a file called SSTable (Sorted Staging Table)
- Flushing is a non-blocking operation.

(8)

Q-3

(a) Hadoop eco system & its deleted components.

- Hadoop ecosystem is a platform providing various services to solve big data problems.
- An open source framework.

→ There are 4 major components in Hadoop.

- ① HDFS (Hadoop Distributed File System)
- ② MapReduce, Yet another Resource Negotiator
- ③ YARN
- ④ Hadoop Common.

MapReduce → programming based data processing.

Other components.

Spark → In memory data processing

PIG, HIVE → query based "

HBase → NoSQL database

Solr, Lucene → indexing & searching,

HDFS

namenode  
datanode

YARN

resource manager  
node manager  
app manager

(b)

(i) create database 'bdaspecial' & collection 'word cloud'.

use bdaspecial

to create database.

show dbs

to check if database is created.

(9)

db.createCollection("wordcloud")  
(to create collection)

show collections

(to check if it is created)

(ii) import records from csv file

mongoimport -d bduspecial -c wordcloud --type csv --file myfile.csv --headerline  
collection → database name  
→ file address → to use 1st line  
as field names.

(iii)

find max salary per group of branches only for year 2020

(c) (2) steps

① map fn → take delta & converts into another set of individual elements

input → set of delta (comma separated words)

(10)

output  $\rightarrow$  another set of data  
(key, value.) e.g. (car, 1)  
(bus, 1)

reduce fn

takes map's o/p as i/p  
& combines tuples into smaller  
set of tuples

i/p = o/p of map fn

o/p = (key, value)

e.g. (car, 5)

(bus, 3)

{ (iii)

[map fn.]

public void map( Long Writable key, Text  
value, Context con)  
throws IOException.

{  
String line = value.toString();

String[] words = line.split(" "));

for (String word : words)

{  
Text okey = new Text(word.toUpperCase()  
case(' ').join());

IntWritable ovalue = new IntWritable();  
con.write(okey, ovalue);

end of answer sheet