

DHARMSINH DESAI UNIVERSITY, NADIAD  
FACULTY OF TECHNOLOGY  
ONLINE SESSIONAL EXAMINATION

B.Tech (CCE) sem → 7

subject : Big Data Analytics

Roll no. → 142

signature → Pankaj

Date : 9/10/2020

Time : 10:00 to 11:15

Total pages. : 12

Q. 1

(a) Reasons of Apache spark becoming more popular

- Far more faster than other tools for large scale data processing
- can be used with Scala, Java, R & Python.
- provides powerful caching & disk persistence with simple programming API.
- deployment possible with Mesos, Hadoop via Yarn or spark's own clusters
- Machine learning
- Advance analytics
- lazy evaluation
- supports various data formats

(b) Replication strategies of Cassandra.

2 replication strategy

## ① simple strategy:

- move around trying to meet application factors
- stop when replication factor is met

## ② Network Topology Strategy:

- somewhat complicated
- different replication factor allowed for different data centers
- Inside a data center, replicas on different nodes are stored.

### (c) Spark Dynamic Memory Utilization

- We have to set it true to use it.
- With this, ~~it~~ once task is finished, it kills idle executors automatically.
- ~~If we set it false, i.e. don't use it,~~ these idle executors are kept as it is.
- Executors can be killed automatically after completion of jobs but driver has to be stopped. even after job completion, driver will still be there & eat resources

### (d) ways deleted to improved performance via storage techniques of Hive

- By partitioning
- use of indexing
- vectorisation
- bucketing

(3)

(4)

Positioning → without hive sends all delta in directory & applies query on it which is slow & expensive.

Indexing → separate index table is created which acts as reference for original table.

(e) SerDe, and its working.

- it is short form of Serializer / Deserializer.
- it is used in hive for I/O
- It allows Hive to read delta from a table & write it back to HDFS in my custom format.
- It handles both serialization & deserialization and also interprets results of serialization as individual fields of processing.

(f) advanced delta types available with Hive & Pig.

along with primitive delta types, hive has

(4)

arrays : ARRAY < data-type >

maps : MAP < primitive-type, data-type >

Structs : STRUCT < column: data-type [comment  
col-comment], ... >

in string it has VARCHAR & CHAR & STRING  
in Date/Time it has TIMESTAMP & DATE

in Misc types it has BOOLEAN & BINARY

→ in pig along with primitive types,  
it has

charArray

byteArray

tuple

bag

map

(5)

## No 2

(a) Meekawa's quorum based approach

- For ensuring mutual exclusion in distributed systems.
- Here a site doesn't request permission from every other site but from a subset of sites which is called quorum

3 message types

① request

② reply

③ release

algo used for following

① entering critical section

② execute " "

③ release " "

⑥

⑦

(c) **I** File formats supported by Hive with their purposes.

→ TEXTFILE → famous input / output format used in Hadoop  
→ if table written in this format, JSON & CSV data can be loaded  
→ In this format, each line is considered as a record.

packages ~ org.apache.hadoop.mapred.  
TextInputFormat  
org.apache.hadoop.mapred.TextOutputFormat

→ SEQUENCEFILE

→ solves issue of working with file of size smaller than typical block size.

(7)

- ~~THE~~ SEQUENCE files are flat files consisting of binary key-value pairs
- they are in binary format which can be split or used for clubbing 2 or more small files to one file.

package → org.apache.hadoop.mapred.SequenceFileInputFormat  
org.apache.hadoop.mapred.SequenceFileOutputFormat

→ RCFILE → stores few Record columns file  
→ binary file format.  
→ for high compression setup on the top of the rows  
→ used when we want to operate multiple rows at a time

package → org.apache.hadoop.io.RCFileInputFormat  
org.apache.hadoop.io.RCFileOutputFormat

→ ORCFILE → optimized few columns  
→ stores data in optimized way than other formats  
→ reduces size of data upto 75%.  
→ increases speed of processing.  
→ stores data in groups (Stripes) along with file footer

package → org.apache.hadoop.hive.orc.OrcInputFormat

③

## (II) Hive crud operations

\* Create table.

e.g. creating one table.

CREATE TABLE STUDENT

STD\_ID INT,  
STD\_NAME STRING,  
ADDRESS STRING )  
CLUSTERED BY (ADDRESS) into 3 Buckets  
ROW FORMAT DELIMITED  
FIELDS TERMINATED BY ","  
STORED AS ORCtbl properties ('orc.  
Hive'='true');

\* insert

INSERT INTO TABLE STUDENT VALUES (101,  
"ABC1.", "NADIA") ;

\* Select

SELECT \* FROM STUDENT;

\* update

update STUDENT  
SET std\_id = 110  
WHERE std\_id = 101 ;

\* delete from table

DELETE FROM STUDENT WHERE std\_id = 110;

9

Q1

\* delete all rows

DELETE FROM STUDENT;

Q3

(c)(i) Mapreduce vs Spark in terms of processing.

- Spark is able to execute batch processing jobs b/w 10 to 100 times faster than the MapReduce. Although both the tools are used for processing Big Data.
- Mapreduce is used when we want linear processing of large dataset & no intermediate solutions are required.
- Spark is used when fast & interactive data processing is needed.
- Spark provides Graph processing, Real-time processing, machine learning & iterative jobs.
- mapreduce & spark can be used

(10)

together in same Hadoop cluster or  
spark alone can be used as a processing  
framework

→ spark was built on top of Hadoop  
MapReduce & it extends MapReduce to  
efficiently use more types of  
computations.

(II) Spark can work in following ways  
with Hadoop.

\* spark + HDFS

spark can run on top of Hadoop's distributed  
file system (HDFS) to leverage  
the distributed replicated storage.

\* spark + mapreduce

spark can be used along with MapReduce  
in same Hadoop cluster or we can  
use it alone as a processing framework

\* spark + Hadoop YARN

Spark applications can also be run on YARN  
(Hadoop NextGen)

\* spark is not intended to replace Hadoop  
but it can be considered as its extension.

→ MapReduce & Spark are used together where MapReduce is used for batch processing & Spark for real-time processing.

(b) different types of collections in Cassandra

-day

Types of collections :-

- ① set collection (unordered collection of unique values)
- ② list collection (ordered collection of values)
- ③ map collection (used to map one thing to others like key-value pairs)

examples

\* set collection → phone nos of employee.  
→ when query is returning values in sorted order.

ALTER TABLE employee\_by\_id ADD phone  
SET <text>; // add column

UPDATE employee\_by\_id SET phone = phone  
+ '12345' WHERE id = 1; ~~REORDER BY~~  
// add phone no. to column

\* list collection

→ retrieve in order in which they were entered.

(12) Alter Table employee\_by\_id ADD language  
list<text>; // add column

Update employee\_by\_id Set language = language  
& ['Hindi', 'Gujarati', English']  
where id=1;

#### \* Map collection

e.g. To store timestamp deleted info.

→ each element stored as timestamp column  
& operated individually

Alter Table employee\_by\_id ADD todo  
map<Timestamp, Text>;

Update employee\_by\_id Set todo='2020-10-10'

: 'Assignment', '2020-10-10'  
: 'work' } where id=1;

→ end of answer sheet