



Date: _____

Pg-1

Syllabus: BR16 Class: BE Seat No: 7278650

Sem: VII Reg No - 2017C02 Sub: BDA

(Q1) MCQ

- 1) Option: B : High level of communication exists between the various nodes.
- 2) Option: A : Ordered
- 3) Option: D : MongoDB
- 4) Option: C : Data Node
- 5) Option: B : 8
- 6) Option: B : 2
- 7) Option: A : Cold start
- 8) Option: D : Biased reservoir Sampling
- 9) Option: C : Microsoft instant Messenger
- 10) Option: C : Periodic



Date: _____

Pg - 2

Syllabus: R16 Class: BE Seat No - 7278650
Sem: VII Sub: BDA

Q2 A

2) Hadoop Core Components

- ① Hadoop has a master-slave topology
- ② In this topology, we have one master node and multiple slave nodes.
- ③ Master node's function is to assign task to various nodes and manage resources. The slave nodes do the actual computing.
- ④ Slave node stores the real data whereas on master we have meta data.

Hadoop Distributed file System (HDFS)

- ① HDFS is based on Google file system (GFS)
- ② HDFS runs on clusters on commodity hardware.
- ③ The file system has several similarities with the existing distributed file systems.

HDFS follows the master-slave architecture and it has following core elements:

① Name Node/Namenode:

- i) It is a daemon which runs on master node of hadoop cluster.
- ii) There is only one namenode in the cluster.
- iii) It ~~can~~ contains metadata of all the files stored on HDFS which is known as namespace of HDFS.
- iv) It maintains two files i.e. Edit Log & FsImage.
- v) Edit log is used to record every change that occurs to file system metadata.
- vi) FsImage stores entire namespace, mapping of blocks to files and file system properties.



R16 / BE / Sem: VII / Sub: BDA / Seat : 7278650

Page No. _____

(2) Datanode :

- i) It is a daemon which runs on Slave machines of Hadoop
- ii) There are number of data nodes in a Cluster.
- iii) It is responsible for serving read/write request from the clients. It also performs block creation, deletion, and replication upon instructions from the name node.
- iv) It also sends a Heartbeat message to the namenode periodically about the blocks it holds.

(3) BLOCK :

- i) # generally the user data is stored in the files of HDFS
- ii) The file in a file system will be divided into one or more segments and stored in individual data nodes.
- iii) These file segments are called as Blocks.

(4) Mapreduce :

- i) Mapreduce is a Software framework
- ii) Mapreduce is a data processing layer of HADOOP
- iii) It is a software framework that allows you to write applications for processing a large amount of data.
- iv) In mapreduce a application is broken down into number of small parts are also called as fragments.
- v) These blocks can run on any node in the cluster.
- vi) Data processing is done by mapreduce.
- vii) Mapreduce core functions are
 - a) Read input
 - b) Function mapping
 - c) Partition, Compare & sort
 - d) Function Reducing
 - e) Write Output



R16 / BE / Sem-VII | Sub: BDA | Seat: 7278650

HADOOP Ecosystem Components

- (1) Core components of Hadoop Ecosystem is nothing but the different components that are built on the hadoop platform directly.
- (2) Hadoop ecosystems are :
- a) HDFS
 - b) MapReduce
 - c) HIVE
 - d) PIG
 - e) HBase
 - f) Zookeeper
 - g) Sqoop



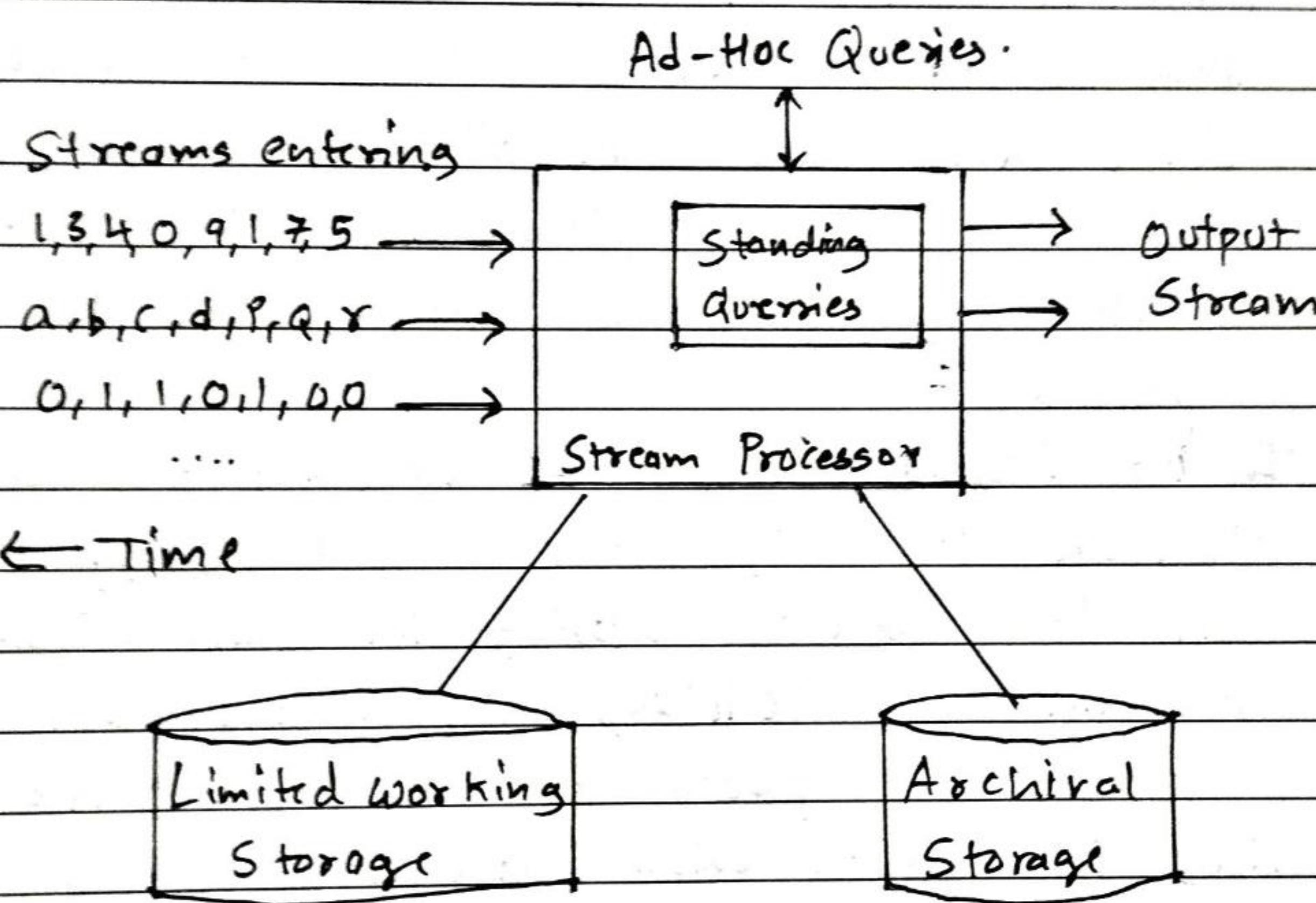
R16/BE/Sem-VII/2017C02/Sub:BDA/Seat:7278650

Page No. 65

Q2B

⇒ Data Stream Management System.

- ① For handling and controlling the data streams, we require a concrete, standardized model or framework so that data stream will be processed in a properly defined manner.
- ② Block diagram of data stream management system.



- ③ The data-stream management System architecture is very similar to that of conventional relational database management system architecture.
- ④ The basic difference is that processor block or more specifically a query processor is replaced with specialized block known as Stream processor.
- ⑤ The first block in the system architecture shows the input.
- ⑥ The number of data stream generated from different sources will enter into the System.



Date: _____

Pg - 6

R16/BE/Sem-VII/2017CO2/Sub: BDA/Seat: 2278650

~~B6~~

- ⑦ Every data stream has its own characteristics such as:
- Every data stream can schedule and rearrange its own data items
 - Every data stream involved in heterogeneity heterogeneity i.e. in each data stream, we can find different kinds of data such as a numerical data, alphabets, alphanumeric data, graphics data, textual data, binary data or any converted transformed or translated data.
 - Every data Stream has different input data rate.
 - No uniformity is maintained by the elements of different data streams while entering into the Stream processor.
- ⑧ In the Second block of architecture, it is abstracted that, there are two different Subsystems exist One of which will take care of storing the data stream and other responsible for fetching the data stream from Secondary storage and processing it by loading into main memory.
- ⑨ Hence the rate at which ~~Stream~~ enters into the system is not the burden which is involved in the stream processing.
- ⑩ the third block represents the active storage or working storage for processing the different data streams.
- ⑪ The working Storage area may also contain Sub-streams which are integral part of main core stream to generate result for a given query.
- ⑫ Working Storage basically a main memory but situation demands then data items within stream can be fetched from Secondary Storage. Working Storage has limited size.
- ⑬ The fourth block of system architecture is known as archival ~~Storage~~ Storage



R16 | BE | Sem-VII | 2017/02 | Sub: BDA | Seat-7278650

~~Page~~

- (14) As name indicates, the block is responsible for maintaining the details of every transactions within the system.
- (15) It is also responsible to maintain edit logs.
- (16) Edit logs are nothing but updation of data.
- (17) The fifth block is responsible for displaying or delivering the output stream generated as a result of processing done by the stream processor. Usually by taking the support of working storage and occasionally by taking support of archival storage.

DBMS

- (1) DBMS refers to Database management System
- (2) DBMS deals with persistant data
- (3) In DBMS random data access takes place
- (4) It is based on query driven processing model ie. pull based model
- (5) The data update rate is low
- (6) It does not provide real time service
- (7) DBMS uses Unbounded disk store, means unlimited secondary Storage

DSMS

- (1) DSMS refers to Data Stream management System
- (2) DSMS deals with stream data.
- (3) In DSMS Sequential data access takes place
- (4) It is based on Data driven processing model ie. push based model
- (5) The data update rate is high
- (6) It provides real time service.
- (7) DSMS uses bounded main memory means limited main memory.



Date: _____

R16/BE/Sem-VII/2017/02/Sub: BDA/Seat: 7278650 Pg-98

(Q3A)

- 1) Hadoop is an open source Software programming framework for storing large amount of data and performing computation.

① Features of Hadoop

- a) Low cost
- b) High computing power
- c) Scalability
- d) Huge & Flexible Storage
- e) Fault Tolerance & Data Protection.

② Limitations of Hadoop

- a) issue with small files.
- b) slow processing Speed
- c) Latency
- d) No real-time data processing
- e) Support for batch processing only
- f) Not ~~not~~ Caching.
- g) Not ease of use
- h) No Delta iteration.

3) Key Value Store Database

- ① It is one of the most basic types of NoSQL databases.
- ② The kind of NoSQL database is used as a collection, dictionaries, associative arrays, etc.
- ③ Data is stored in ~~as~~ key/value pairs.
- ④ It is designed in such a way to handle lots of data and heavy load
- ⑤ Key Value pair storage database store data as a hash table where each key is unique, and the value can be a JSON BLOB, String etc.



Date: _____

Pg-90

R16/BE/Sem-VII/2017C02/Sub: BDA/Seat: 7278650 P2

⑥ Example : ~~Redis~~ → Cassandra, Hypertable

- a) Azure Table Storage (ATS)
- b) DynamoDB

Document Database.

- (1) Document oriented No-SQL database stores data as key-value pair but the value part is stored as a document.
- (2) The document is stored in JSON or XML formats.
- (3) Every document contains a unique key used to retrieve the document.
- (4) Key is used for storing, retrieving and managing document oriented information also known as semi structured data.

(5) Examples:

- a) MongoDB
- b) Couch DB .

Q3 B

2)

How dead ends are handled in Page rank

→ A dead end is a web page with no links out. The presence of dead ends will cause the Page Rank of some or all the pages to go to 0 in the iterative computation, including pages that are not dead ends.

(3) Dead end can be eliminated before undertaking a page rank calculation by recursively dropping nodes with no arcs out.

(3) Note that dropping one node can cause another which linked only to it to become a dead end. So the process must be recursive.

(4) There are two approaches to deal with dead ends.

a) We can drop the dead ends from the graph and also drop their incoming arcs. Doing so may create more dead ends which also have to be dropped recursively. However eventually, we wind up with a Strongly Connected Component (SCC) none of those nodes are dead ends.

Recursive deletion of ~~of~~ dead ends will remove part of the out-components, tendrils, and tubes but leave the SCC and in-component, as well as parts of any small isolated components.

b) We can modify the process by which random Surfers are assumed to move about the web. This method which we refer to as 'taxation' also solves the problem of spider traps.

Here we modify the calculation of page rank by allowing each ~~random~~ random Surfer a small probability of teleporting to a random page rather than



Date: _____

R161 BE1 Sem-VII | Sub: BDA | Seat No - 7278650 Pg - 11
~~Page No.~~

following an out-link from their current page.
The iterative step where we compute a new vector estimate of Page Rank v' from the current page Rank estimate v and the transition matrix ' M ' is

$$v' = \beta M v + \frac{(1-\beta)}{n} e$$

Where, β is the Chosen Constant usually in the range of 0.8 to 0.9.

e is a vector of all 1's with the appropriate number of nodes in web graph.

n is the number of nodes in web graph.

If we do not dead ends in the graph then
Probability of introduction of New = Probability of not to user to a given web page Choose out-link for Current page by same user.

Another possibility is user will not be able to move to any page as $(1-\beta)e/n$ term is independent of $\sum v$ i.e. when we don't encounter with dead ends then,

$$\sum v < 1 \text{ but } \sum v = 0.$$

Q4 B

- 2) Stream Clustering Algorithm
- ① Stream Clustering Algorithm is also known as BDMO Algo.
- ② BDMO Algo follows the concept of 'counting ones' method, which means that there is a window of length N on binary stream and it counts the number of 1s that comes in the last k bits where $k \leq N$.
- ③ The BDMO algorithm uses the bucket with allowable bucket sizes that forms a sequence where each size is twice of the previous size.
- ④ In the algorithm the number of points represents the size of bucket.
- ⑤ It ~~does~~ does not consider that the sequence of allowable bucket size starts with 1 but consider only forming a sequence such as 2, 4, 6, 8, ... Where each size is twice the previous size.
- ⑥ For maintaining the buckets, the algorithm considers the size of bucket with the power of two.
- ⑦ In addition, the number of buckets of each size is either one or two that form a sequence of non-decreasing size.
- ⑧ The buckets that are used in the algorithm, contains the size and timestamp of the most recent points of the stream.
- ⑨ Along with this bucket also contains a collection of records that represents the clusters into which the points of the buckets have been partitioned.
- ⑩ This record contains the number of points in the cluster, the centroid or clustroid of the cluster, and other parameters that are required to merge and maintain the clusters.



Date: _____

Pg - 13

R16/BE/Sem-VII/Sub: BDA | Seat No - 7278650

Priya

(1) The major steps of the BDMO algorithm are as follows:

- Initialising buckets.
- Merging buckets.
- Answering queries.

(2) Initialising Buckets.

- The algorithm uses the smallest bucket size that is P with power of two.
- It creates new bucket with the most recent P points for P stream elements.
- The time stamp of the most recent point in the bucket is the time stamp of the new bucket.
- After this, we may choose to leave every point in a cluster by itself or performing clustering using an appropriate clustering method.
- For the initialisation of the bucket using selected clustering methods, it calculates the centroid or clusteroid for the clusters and counts the points in each cluster.
- All this information is stored and becomes a record for each cluster.
- The algorithm also calculates the other required parameters for the merging process.

(3) Merging Buckets.

- After the initialisation of the bucket, the algorithm needs to review the sequence of a bucket.
- If there happens to be a bucket in time stamp or more than N time units prior to the current time then nothing of that bucket is in window.

R16/BE/Bem-VII/Sub:BDA/Seat No. 207278650

Pg. - 14
B

- 3) In such case algorithm drops from the list.
- 4) In case we have created 3 bucket size of P, then we must merge the oldest two of the three buckets.
- 5) In case the merger can create two buckets of size $2P$, this may require us to merge buckets of increasing size recursively.
- 6) For merging two consecutive buckets, the algorithm needs to perform following steps.
 - a) For merging the size of bucket should be twice the sizes of the two buckets to be merged.
 - b) The timestamp of the merged bucket is the timestamp of the more recent of the two consecutive buckets.
 - c) In addition it is necessary to calculate the parameters of the merged clusters.

(14) Answering Queries:

- 1) A query in the Stream-computing model is a length of a suffix of the sliding window.
- 2) Any algo takes all the clusters in all the buckets that are at least partially within the suffix and then merges them using some method.
- 3) The answer of query is the resulting clusters.
- 4) For the clustering of the streams the stream-computing model finds out the answer to the query.
- 5) During the initialisation, the k-means method is used and for merging the buckets the timestamp is used hence the algo is unable to find the set of ~~best~~ buckets that covers the last in point.
- 6) After this the algo generates the answer & in response to query as the centroids or clusteroids of all the points in selected buckets.