Ans1)

Architecture of Hadoop :-

Architecture of HDFS :-



Data Nodes

Name - node :-

- Maintains block IDs
- Maintains block health reports
- Runs on system memory
- Single point of failure

Secondary Name Node :-

- NOT a back-up
- Performs memory intensive tasks.
- Runs on system memory
- Mainly a helper for Name Node

Data Node :-

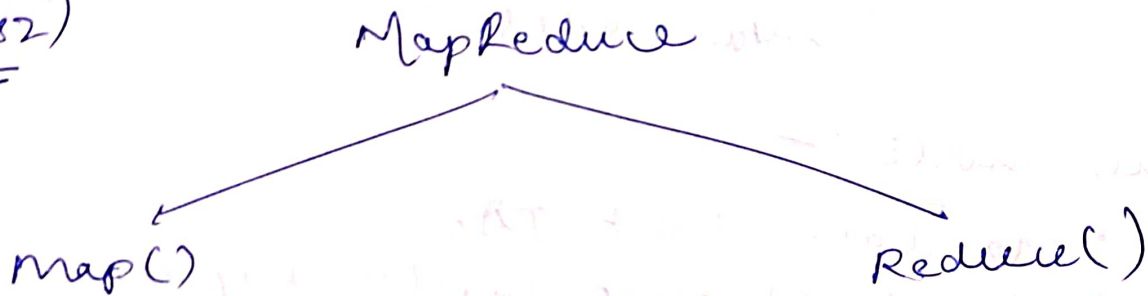- Stores broken down data.
- Slave in nature

Read operation :—

- Namenode keeps track of block IDs & health reports.
- At the time of reading, pointer is shifted to the required block Id & data is read. Once a block is exhausted, pointer shifts to the next block.

Write operation :—

- Data is broken down into chunks & stored on blocks.
- NameNode maintains a <key, value> pair of the data & block IDs it is associated to.

Ans2)

Map Reduce

map()                    Reduce()

- Map → splitting
  Receives the data & breaks it into <key, value> pairs or tuples.

- Reduce → shuffling, reorganization
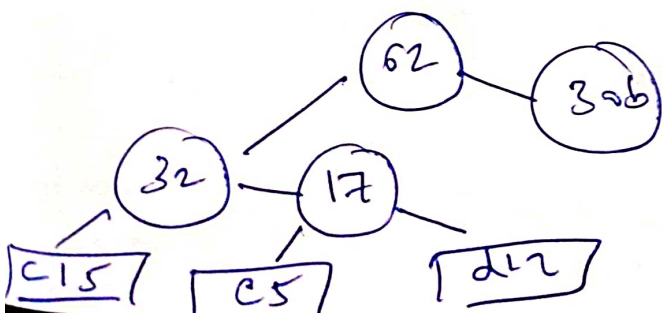  Tuples are further broken down to reduce the size of data while
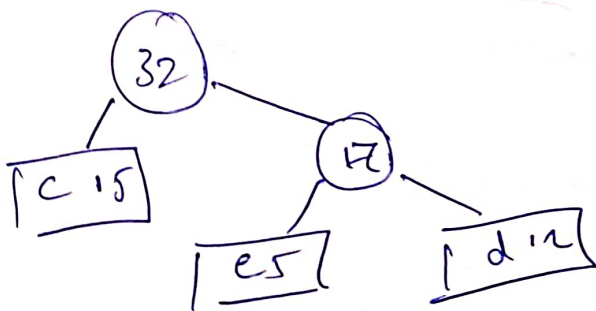
maintaining integrity. ③

Ex — Huffman coding or compression
algorithms.

frequency of characters in a file is given:-

| a | b | c | d | e |
|---|---|---|---|---|
| 40 | 30 | 15 | 12 | 5 |

On reducing :—

| e | d | c | b | a |
|---|---|---|---|---|
| 5 | 12 | 15 | 30 | 40 |



| | | b | a |
|---|---|---|---|
| c | | 30 | 40 |
| 15 | | | |



| | b 30 | a |
|---|---|---|
| | | 40 |

```
           102  1
        0        
                  ┌──────┐
       62         │ 40 a │
     0     1      └──────┘
   32      ┌──────┐
  0   1    │ 30 b │
┌────┐     └──────┘
│ C15│  17
└────┘ 0    1
      ┌────┐ ┌──────┐
      │ e5 │ │1 d12 │
      └────┘ └──────┘
```

Now the frequency of characters is
Same but the memory needed to
to store them has been
reduced.

Ans3) ~~NoSQL Replacement~~ p

   NOSQL schemas :-

① key-value

② Document based

③ Column based

④ Graph based

## Key-value :-

| ID | Name | dept |
|----|------|------|
| SQL 1 | Om | IT |

key value :  1 : { Name : Om
                    dept : IT }

- Key-value dictionaries are used to store data.
- Best for shopping carts.

## Document based :-

| SQL : | ID | Name | dept |
|-------|----|------|------|
|       | 1  | Om   | IT   |

document : JSON {
                 ID : 1
                 Name : Om
                 dept : IT
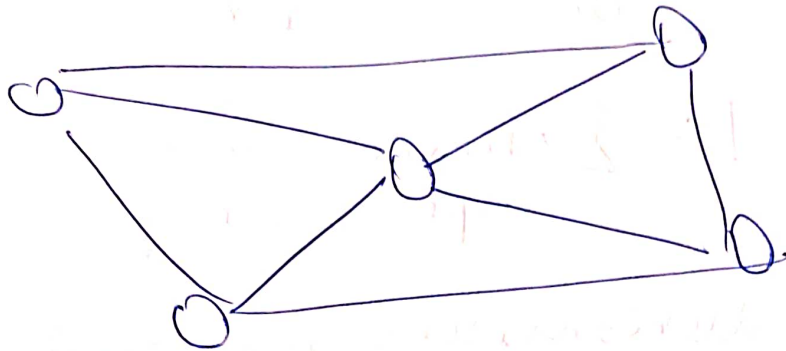                }

- Best for content management service

## column based :-

- Best for data warehousing
- Reduces the size of data using Vectorization.

## Graph based :—

- Best for social media setg



• stores entity - relationships

**Ans4)** Replication management
policy :—

↳ When a rack goes offline, ~~the~~ all
the nodes on that rack will be
unavailable.

↳ Not more than 1 replica of block
is stored on the same node, not
more than 2 replicas of a node
are stored on the same rack.

↳ default replication is 3-fold.

Rack separation

↑

Node separation

↑

Block separation

**Ans 3)** 4Vs :-

① **Velocity :-**

(a) On-going transactions
(b) DB side balance management
(c) Hashing the user data to maintain data integrity.

② **Volume :-**

(a) Loan book
(b) Customer details
(c) Transaction details

③ ~~Visibility~~ **Visibility :-**

(a) Bank's Loan book should have limited visibility
(b) Account details should be visible to account owner only.
(c) Internal-employee data must not be visible to anyone who is not in a managerial position.

④ **Variety :-**

(a) Customer details are stored as strings.
(b) Verification records are stored as PDF.

(c) Transaction details are stored as data dictionary.

(8)

Om Shree
2006077
IT - 02