

AJAY KUMAR GARG ENGINEERING College, GHAZIABAD.

Dept of MCA

Solution ST-2

Course: MCA

Session: 2017-18

Subject: Big Data

Max Marks: 50

Semester: 5

Section: MCA182

Sub Code: NMCA  
E-44

Time: 2 hour.

## Section A

Q 1) What is the purpose of Hadoop Pipe.

Ans 1. Hadoop Pipes: Apache Hadoop provides an adapter layer called pipes which allows C++ application code to be used in Map Reduce programs.

# It allows C++ code to use Hadoop DFS and Map/Reduce.

Q 2) What is the advantage of using materialized view.

Ans 2. A Materialized View is a database object that contains the result of a query.

# For Example, it may be a <sup>local</sup> copy of data located remotely, or may be subset of rows and/or columns of a table or join result, or may be a summary using an aggregation function.

# It is a database query that contains the



result of pre computed query.

Q 3 Differentiate between Scale up and Scale out.

Ans 3 In vertical scaling: the data resides on a single node and scaling is done through Multicore i.e. spreading the load between the CPU and RAM resources of that machine. It is often limited to the capacity of a single machine, scaling beyond a capacity becomes too difficult.

# Horizontal scaling: is often based on Partition of the data i.e. each node contains only a part. It allows you to combine the power of multiple M/C into a virtual single M/C with combined power of all of them together.

Therefore we conclude that Horizontal scaling is achieved through partitioning and vertical scaling is achieved through multi core support.

Q4) what is Codec? Name the interface for its Implementation.

Ans 4. Codec is the implementation of a Compression - Decompression Algo. In Hadoop a codec is represented by an implementation of the Compression Codec interface.

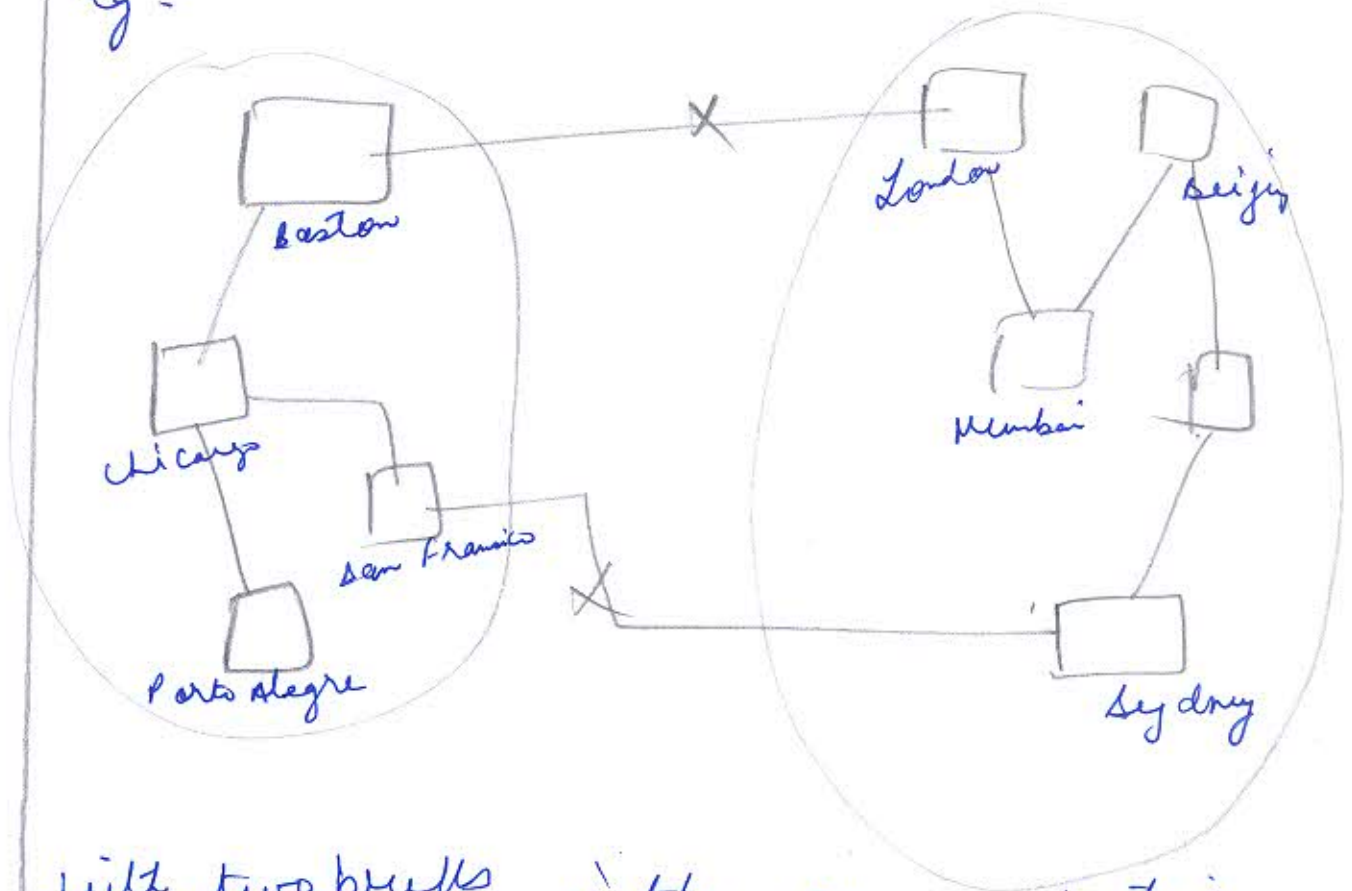
For eg: GZIPCodec encapsulates the compression & decompression algorithm for gzip.

Q5. Explain Partition Tolerance through an Example

Ans 5. Partition Tolerance means that the cluster can survive communication breakages in the cluster that separate the cluster into multiple partitions unable to communicate with each other.  
# This situation is also called split brain.



eg.



With two breaks in the communication lines, the network partitions into two groups.

\* clusters have to be tolerant of network partitions.

## Section B

Q6

Compare and Contrast various Compression Algorithms in Hadoop.

Ans 6

Compression bring two major benefit:

- It reduces the space needed to store files
- It speeds up data transfer across the network or to or from disk.

Various Compression formats & respective features is as following:

| Compression Format | Tool  | Algo    | Filename Extension | Splitting |
|--------------------|-------|---------|--------------------|-----------|
| Deflate            | N/A   | Deflate | .deflate           | No        |
| gzip               | gzip  | Deflate | .gz                | No        |
| bzip2              | bzip2 | bzip2   | .bz2               | Yes       |
| lzo                | lzo   | LZO     | .lzo               | No        |
| Snappy             | N/A   | Snappy  | snappy             | No        |

#

All compression Algo exhibit a space/time trade off: Faster compression and decompression speed usually comes at the expense of smaller space savings.



H GZIP is general purpose compressor, and in the middle of speed/time trade off.

- BZIP2 compresses more effectively than GZIP, but is slower.

- BZIP2's decompression speed is faster than its compression speed, but it is slower than other formats.

- LZ0 & Snappy, both optimize for speed and are around an order of magnitude faster than gzip, but compress ~~less~~ less effectively.

- Snappy is slightly faster than LZ0 for decompression.

Q7) What is the use of Serilization? What are Methods of Serilization in Hadoop.

Ans. 7. Serilization is the process of turning structured objects into a byte stream for transmission over a network or for writing to persistent storage.

Serialization is used for:

- ① IPC (Interprocess Communication)
- ② Persistent storage.

Methods of Serialization in Hadoop are.

- Writable
- AVRO.

Writable is compact & fast but not so ~~easy~~ easy to extend or use from language other than Java.

Avro: supports many language ~~other than~~ including <sup>Java</sup> ~~Avro~~ for serialization.

Q8 Differentiate between Schemed and Schema based Data Base.

Ans- ~~Schemed database~~ NOS.

NOSAL is another name for Schemed database while SQL is Schema based Data Base. Following are the differences between the two based on following parameters.



|   | <u>Entity</u>     | <u>SQL</u>   | <u>NoSQL</u>   |
|---|-------------------|--|--|
| ① | Scaling           | Vertical scaling   | Horizontal scaling   |
| ② | Schema            | Fixed  | Dynamic  |
| ③ | Development Model | Monolithic   | Open source  |
| ④ | Consistency       | Follow ACID<br>(Atomicity -<br>Consistency<br>Isolation<br>Durability) | Follow BASE<br>(Basically Available,<br>Soft State<br>Eventually consistent) |
| ⑤ | Model             | Relational   | Non-Relational   |
| 6 | Development       | Since 1970   | Since 2000   |

Q9  
Ansa Discuss various Distribution Model with their advantages and disadvantages. Various Distribution Models can be:

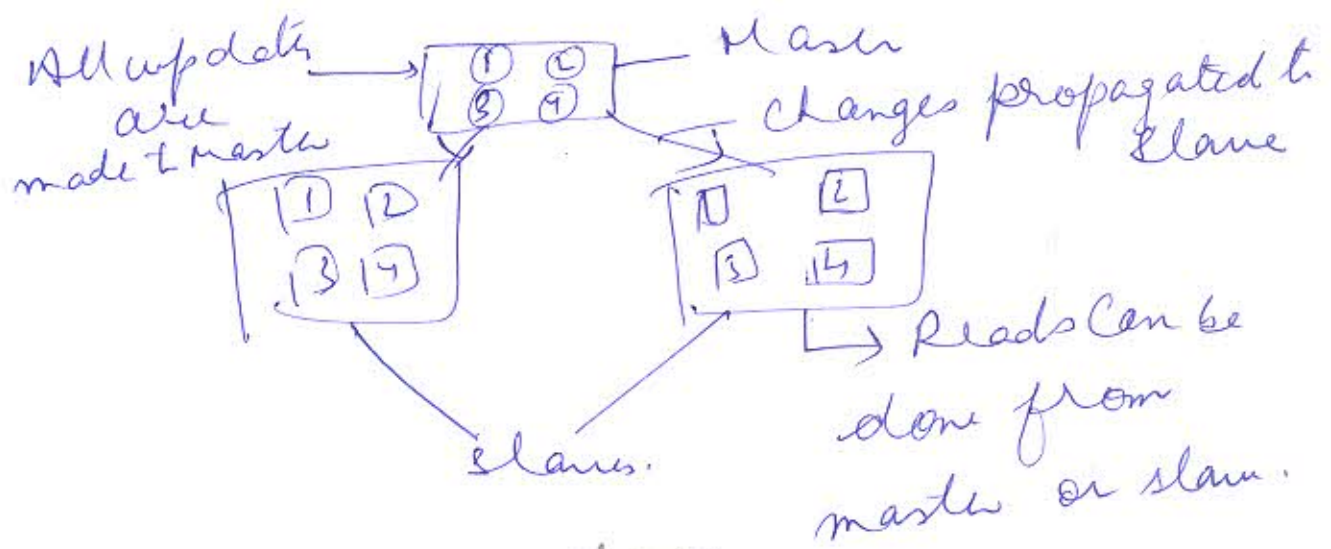
Replication & Sharding:

Replication puts same data and copies it over multiple nodes while Sharding puts different data on different nodes.

They can broadly be.

- ① Master slave
- ② Peer-to-peer replication
- ③ Sharding
- ④ Combining sharding & replication

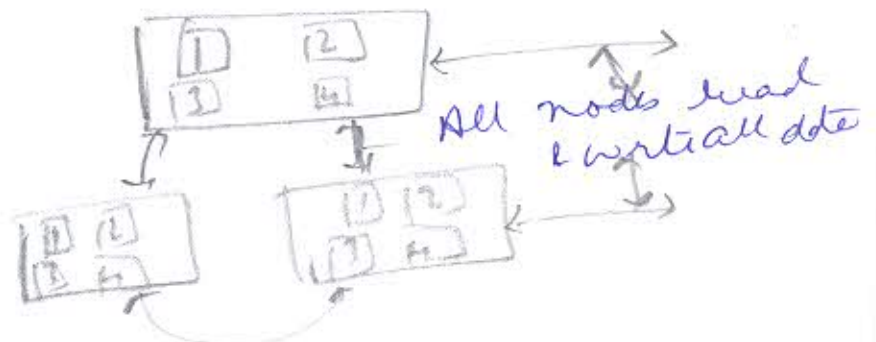
In master slave: one node is designated as master other nodes are slave



Advantages Master slave.

- ① It is Read Intensive
- ② Read Resilient.

Peer to Peer.



Peer to Peer

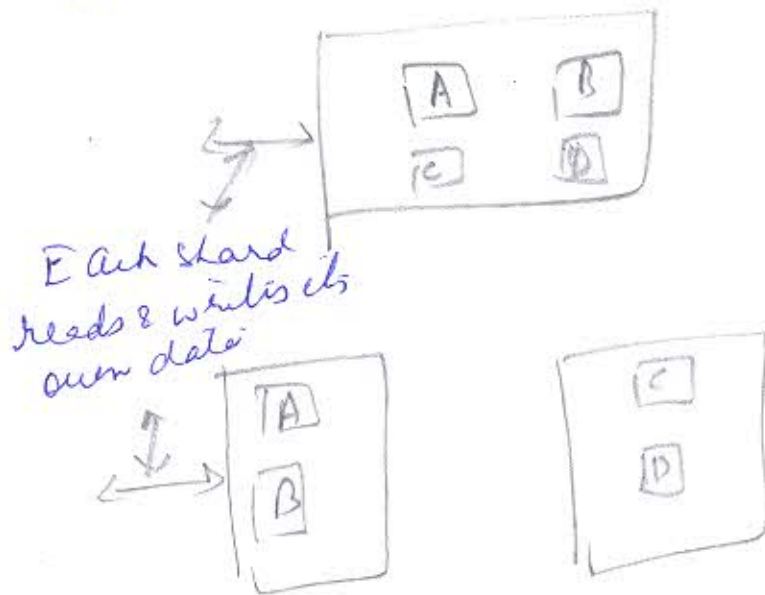


Peer to Peer Disadvantage is :

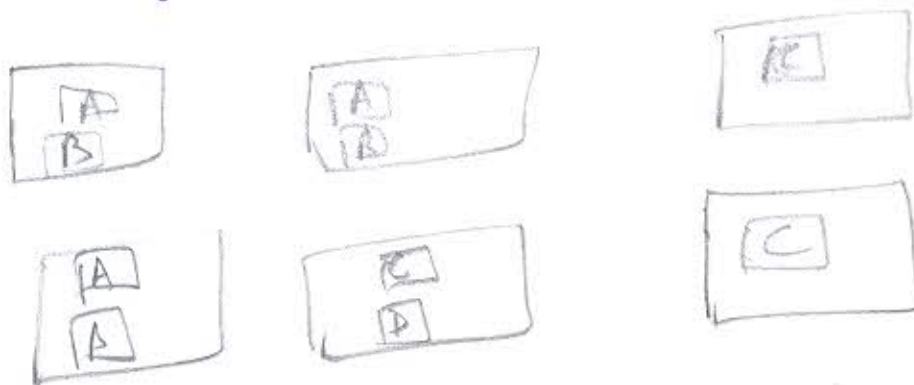
1) Write - write conflict.

# In sharding different people are accessing different part of datasets.

# Different parts of data are put onto different servers!



Combining sharding & replication



# Thus sharding & Replication are orthogonal strategies that can be.

Q10) What is the advantage of AVRO File Based data structures over sequence file format?

Ans 10) Sequence file is a file format consisting of binary key/value pair

# It provides writer, Reader & Sorter classes

# There are three different formats for sequence file depending on the compression types specified:

- Uncompressed
- Record compressed
- Block compressed format

Sequence File Header

3 Bytes (3E0) + 1 Byte (version)

Text - key class name

Text - value class name

Boolean - Is compressed

Boolean - Is Block compressed

Compression code class name

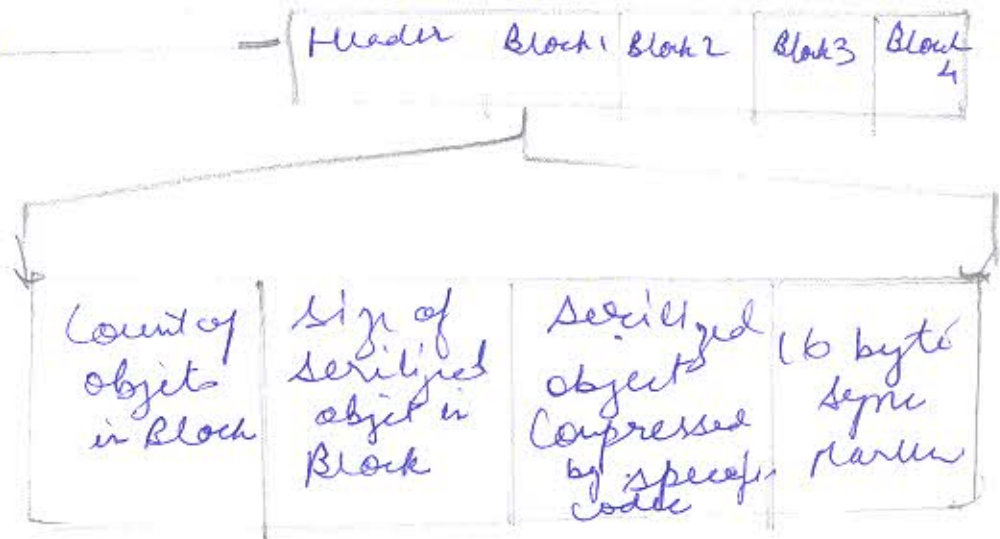
Metadata

Sequence Name



# Avro File

4 bytes  
ASCII  
0's followed by  
meta data.  
  
16 byte sync  
marker



~~Big using Avro~~

- # Avro schema are usually written in JSON
- # Supports splittability
- # Its Row Based, Compact binary format.

## Section

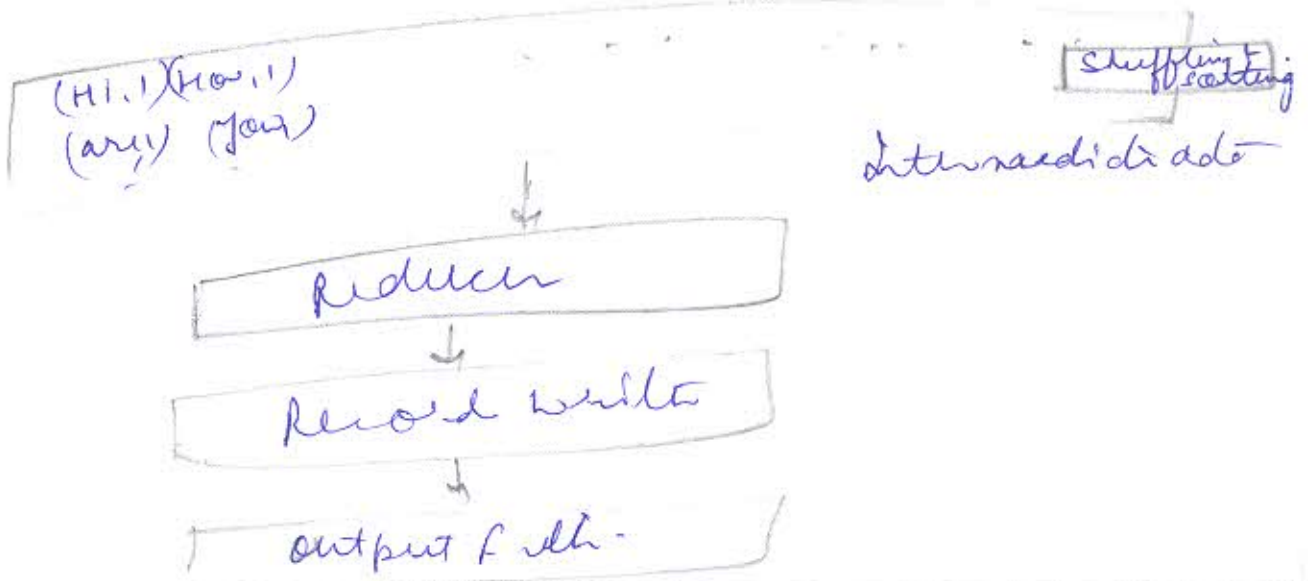
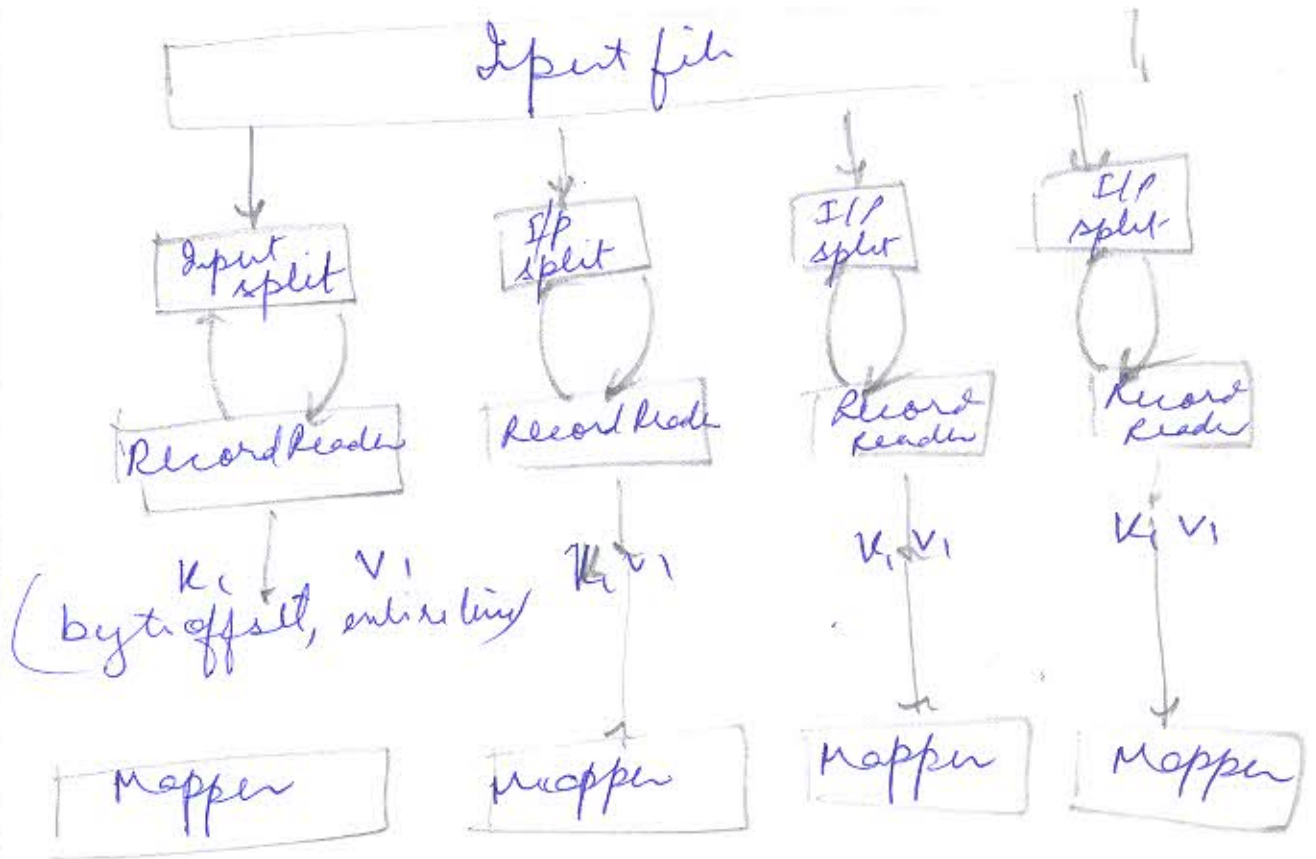
Q11 ~~map reduce~~ what are the stages involved in Map - Reduce? Explain giving an example

Ans 11 Map Reduce is a programming paradigm that uses key-value pair.



File input (20 MB)

|                  |       |
|------------------|-------|
| Hi How are you   | 64 MB |
| How is your lab  | 64 MB |
| How is ur family | 64 MB |
| How is ur class  | 64 MB |
| What is time now | 64 MB |
| How are you      | 64 MB |
| What is the      | 64 MB |
| strength of      | 64 MB |





Q 12) Explain key value and graph data model in detail giving examples?  
Key value database:

A 12

# Basic data structure is map, where we can store value as string, a JSON structure along with key to reference that value

# eg are Riak, Redis, Memcached,

# eg.

| Car      | Attributes  |
|----------|---|
| key<br>1 | Make:<br>model:<br>color:<br>year:                            |
| 2        | Make:<br>model:<br>color:<br>color:<br>year:<br>Transmission: |

various key value pair Car:

(Business) key → value

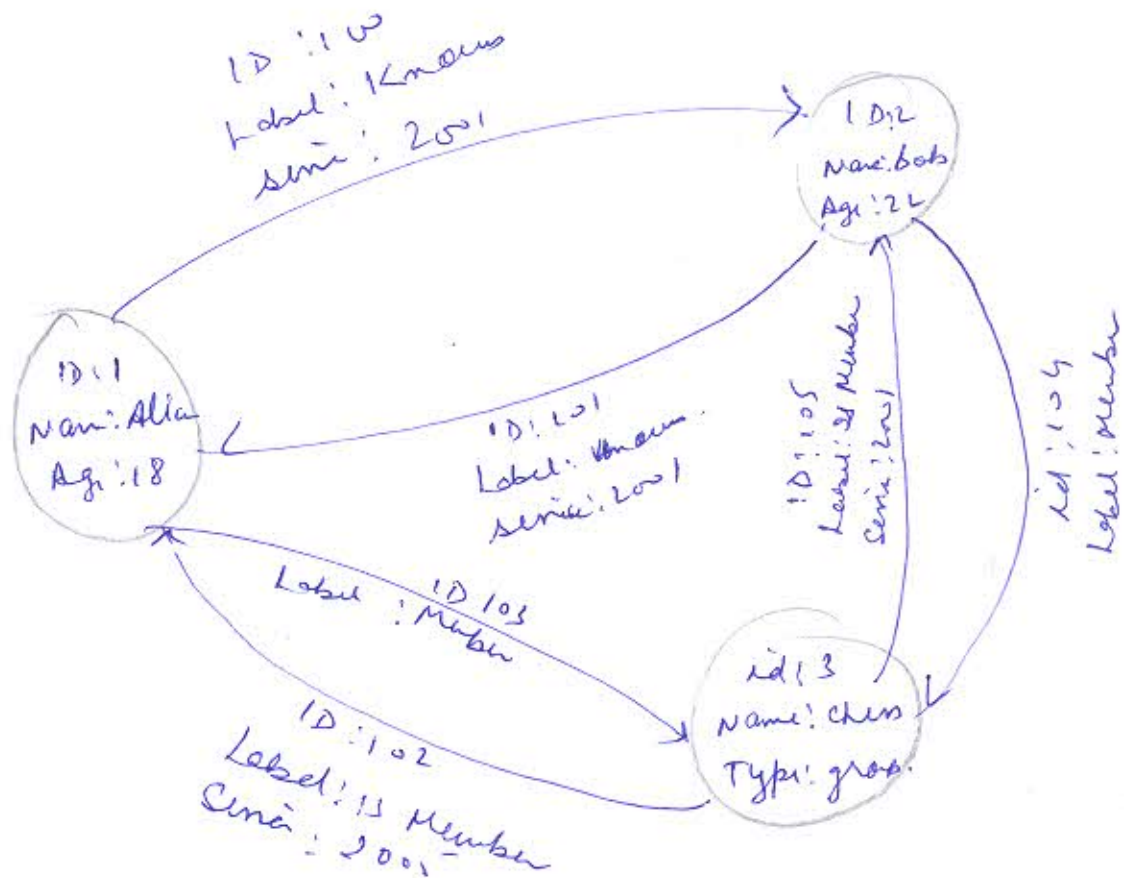
(twitter id) key → information about tweet & value

# graph model : uses graph data structure to store data.

# It is used when traversing relationships are core to the application like social

solid network connection

eg Neo4j & Graph Graph are its commercial products.



The node user: Bob is a vertex with a property Bob.

we also see Relations which are edges.

- Base of
- friend of
- known to
- married to.