



Big Data Analysis - BDA class notes

Big Data Analytics (University of Mumbai)



Scan to open on Studocu

* Difference between Small and Big Data.

Small Data

① Structured (in rows & col)

② stores data in

GB / MB / TB

③ Increases gradually

④ SQL & Oracle

⑤ Single Node
(Centralized system)

Big Data

① Unstructured. (Any kind of data).

② $PB | EB | ZB$

$PB = 1000 \text{ of } TB$

$EB = 1000 \text{ of } PB$

③ Increases exponentially.

④ Hadoop, Spark

⑤ Multinode cluster.

* Imp

[Q] Big data characteristics.

- 3V's (Volume, Velocity, Variety) (Variety, Value)

1) Volume :- large amount of data.

2) Velocity :- speed of using data.

3) Variety :- ignore the missing.

4) Variety :- different types of data.

5) Value :- finding the current meaning (value) of data.

[Q] Types of Big data

(1) Structured

(Table format e.g. SQL, RDBMS etc.)

(2) Semi structured

(JSON, XML, YAML, NOSQL)

(3) Unstructured

3. Semi structured

XML → eg. text based markup language

`<programmer Details>`

`<FirstName> John </FirstName>`

`<LastName> Doe </LastName>`

`<coding platforms>`

`<coding platform = "FAV"> coded Era </coding platform>`

`</programmer Details>`

YML → eg. ~~yaml~~ Yet another markup language

Emp 1 : Yam 1

name : Yam 1

name : Mary Jacob

Job : Developer

skills : Java

* [Q] Why Big Data?

Traditional Data

① It is generated in

Enterprise level

② Its volume ranges from

GB to TB

③ Traditional DB system deals with structured data.

④ Generated per hour / per day / or more.

⑤ Size of data is very small

Big Data

① Generated outside the Enterprise level

② Its volume ranges from PB to ZB

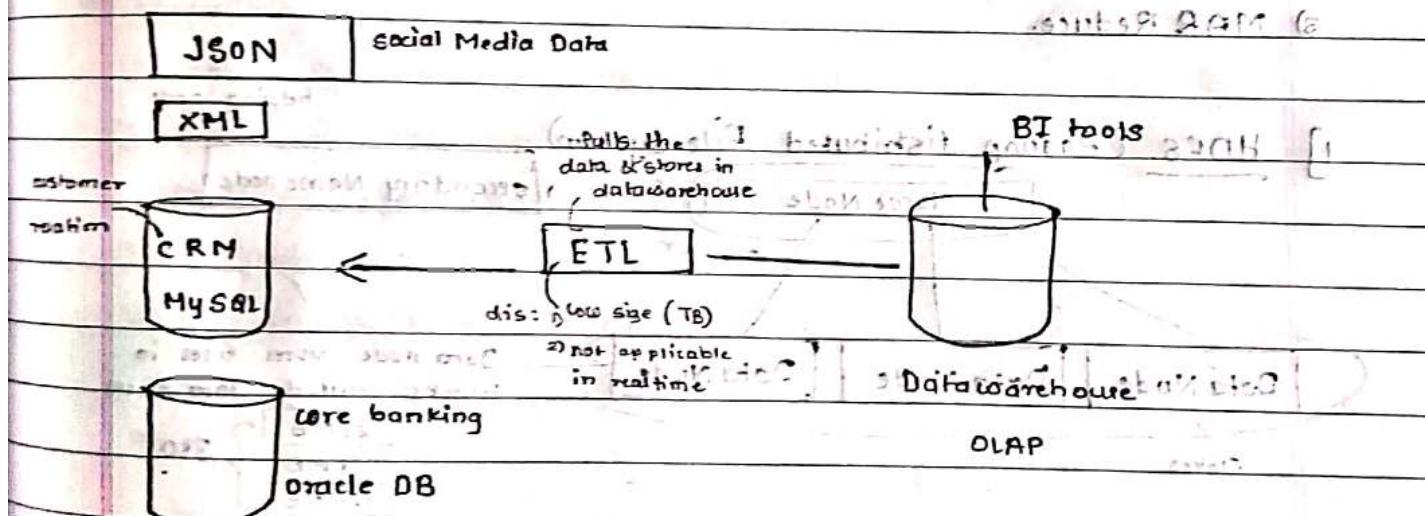
③ Big data deals with semi-structured, unstructured, etc.

④ Generated more frequently mainly per second

⑤ Size is more than traditional data

- | | |
|--|---|
| ⑥ Tradition data is stable and inter | ⑥ Big data is not stable & unknown relationship. |
| ⑦ Easy to manage & manipulate the data. | ⑦ Difficult to manage & manipulate data. |
| ⑧ Data sources include ERP Transaction data, CRM Transaction data, financial data, Web transaction | ⑧ Data sources social media, sensor data, video, audio, images. |
| ⑨ Traditional data source is centralized and managed in centralized form | ⑨ Big data source is distributed and managed in distributed form. |
| ⑩ Data integration is easy | ⑩ Data integration is difficult. |

⑩ Extract, Transform, Load



Transaction

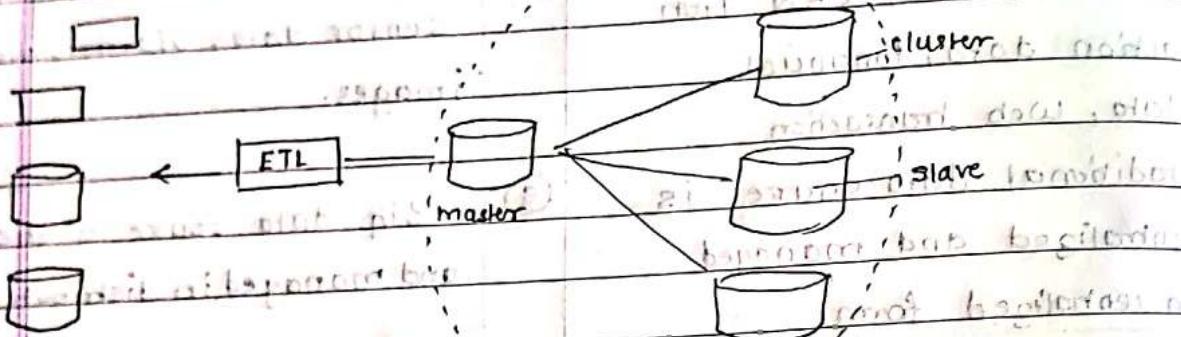
Traditional Enterprise architecture

Tools in ETL

1) informatica

2) oracle

- * MPP (Massive Parallel Processing)
- hadoop - cloud era
- Apache (using this download hadoop V2.6)
- hadoop uses Master/Slave Architecture



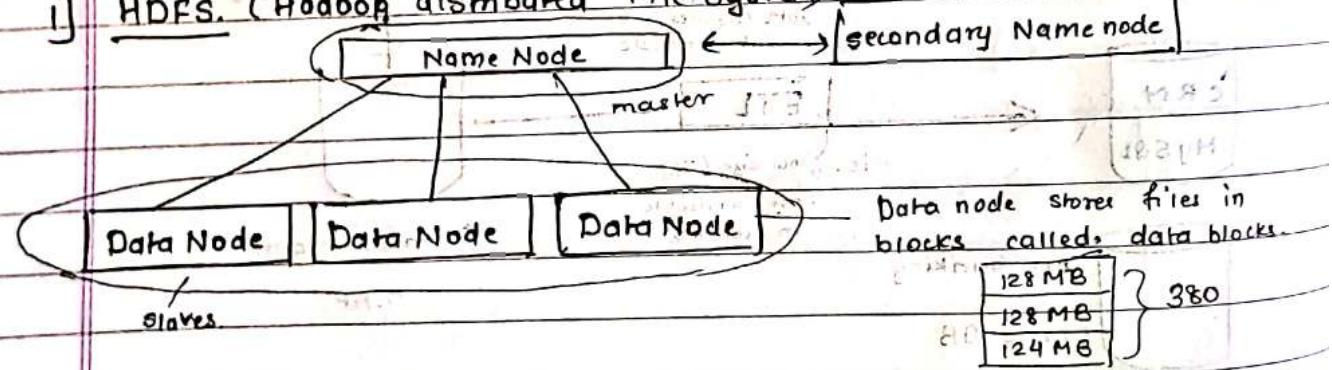
- * Hadoop is divided in components.

D) HDFS

E) YARN

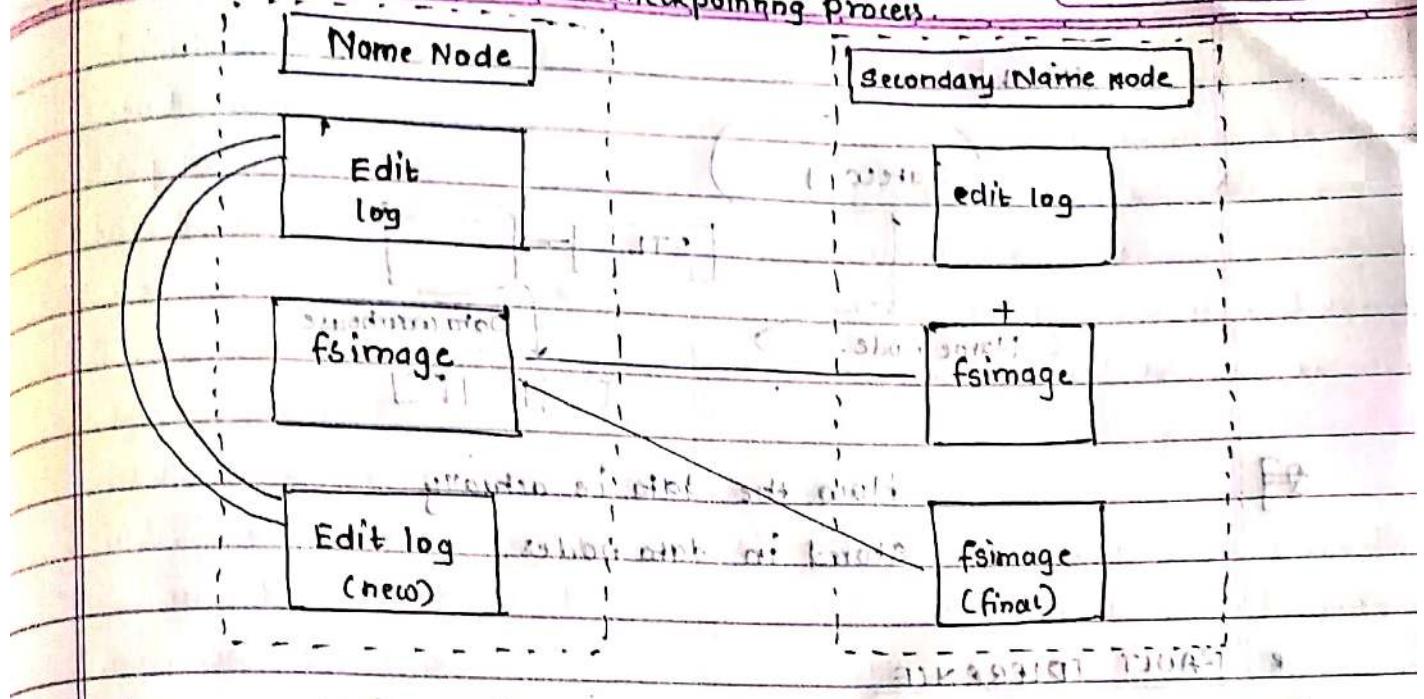
F) Map Reduce

- I) HDFS (Hadoop distributed File system)



- master manages the Data nodes and checks for the file related information.
- It checks for the data nodes be working or are either dead.
- Secondary Node is also called as helping node.

Checkpointing process.



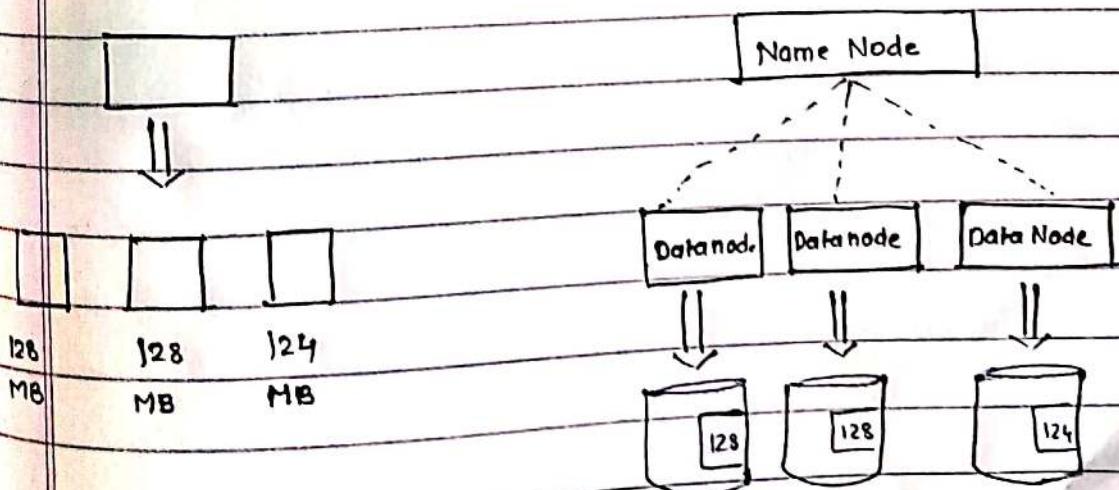
fsimage → file system image (at time of snapshot or part of node)

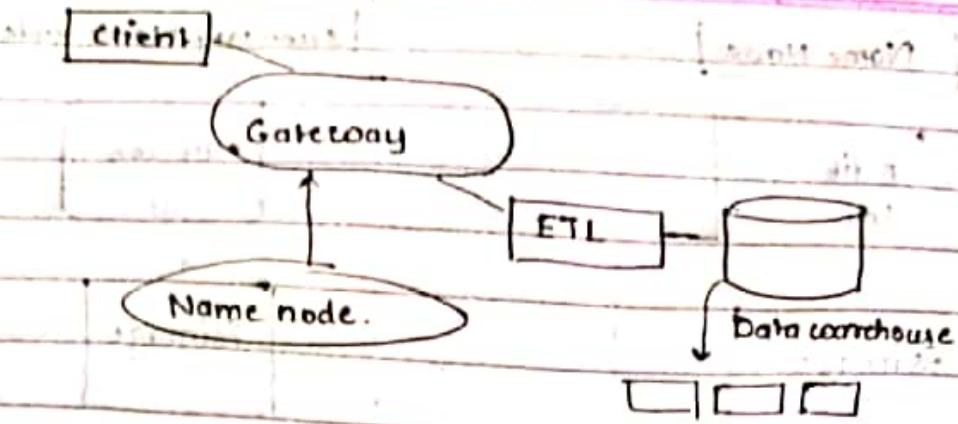
checkpointing → process of combining fsimage & Edit logs

- fs image cannot be continuously updated.

* Secondary Namenode:- (SNN)

- (1) process of combining fsimage & Edit logs (checkpoints)
- (2) SNN takes over the responsibility of checkpointing therefore making Name node more available.
- (3) checkpointing happens periodically (by default 1hr)
- (4) If Name node fails, from secondary name node we will recover the data.





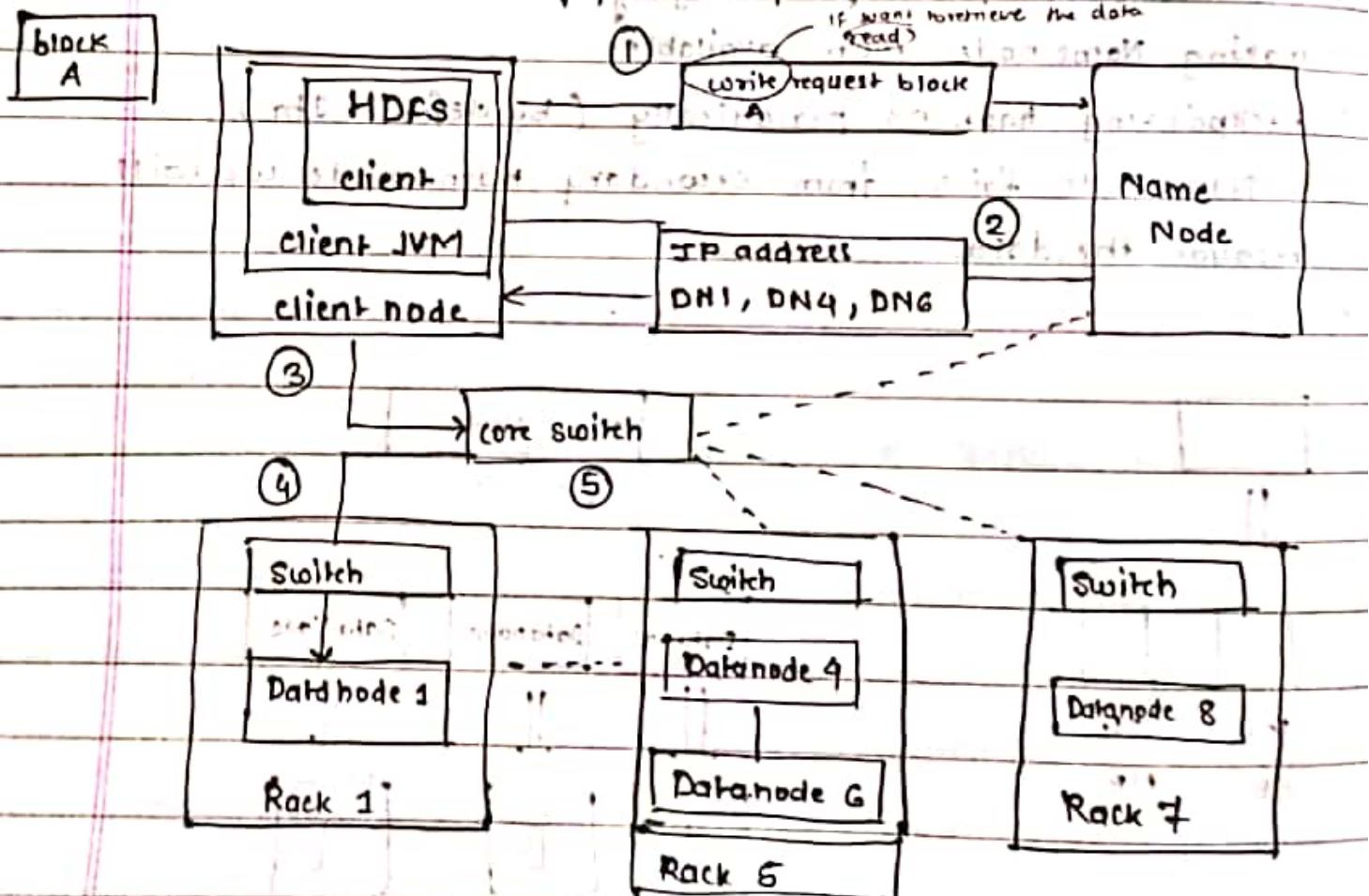
Q] How the data is actually stored in data nodes.

* FAULT TOLERANCE

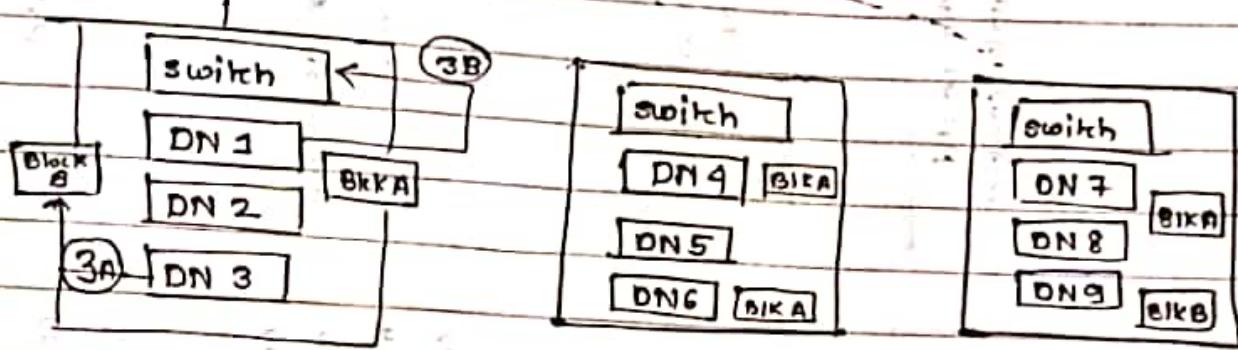
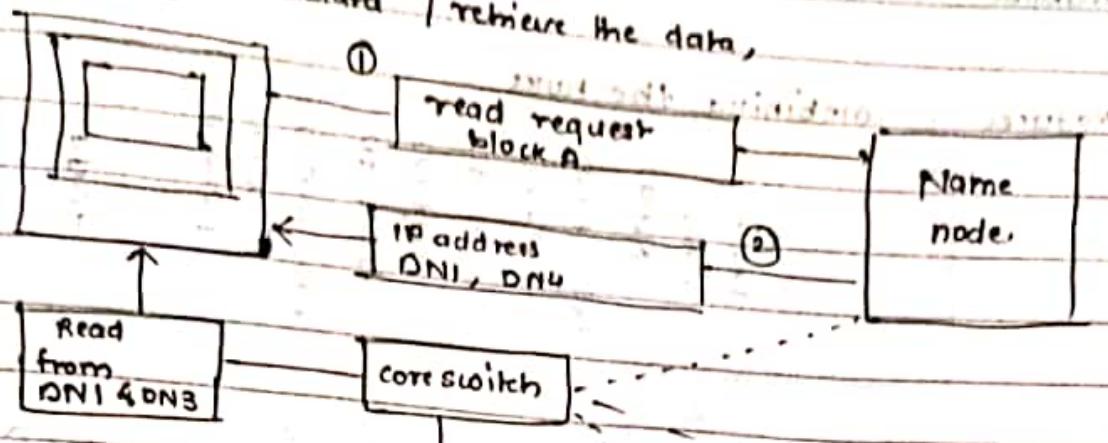
- * How Hadoop cope up with data node failure?
- Each data blocks are replicated (3 times by default) and distributed across different data nodes.

* Hadoop

writes mechanism pipeline setup.



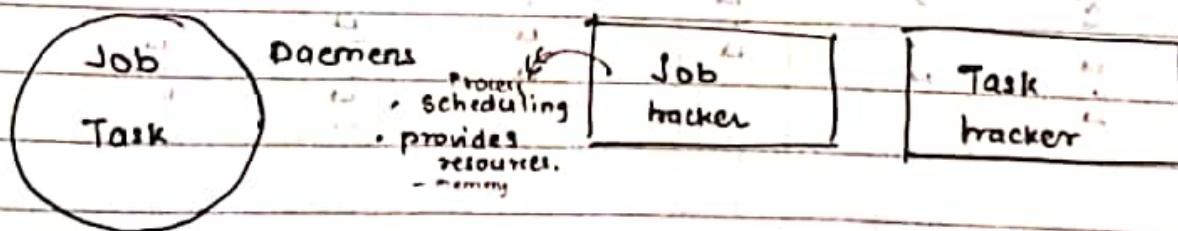
Data nodes are placed randomly for reading the data / retrieve the data.



* MAP Reduce

- It stores data in distributed manner. the result is gained parallelly manner.

Class → Mapper class & Reducer class.



To check network, there will be delay.

$$\therefore PD = \frac{\text{Distance}}{\text{Time}}$$

Mapper & Reducer

This document is available free of charge on

Mapper - Data nodes are divided, into tasks.

Reducer - combining the tasks

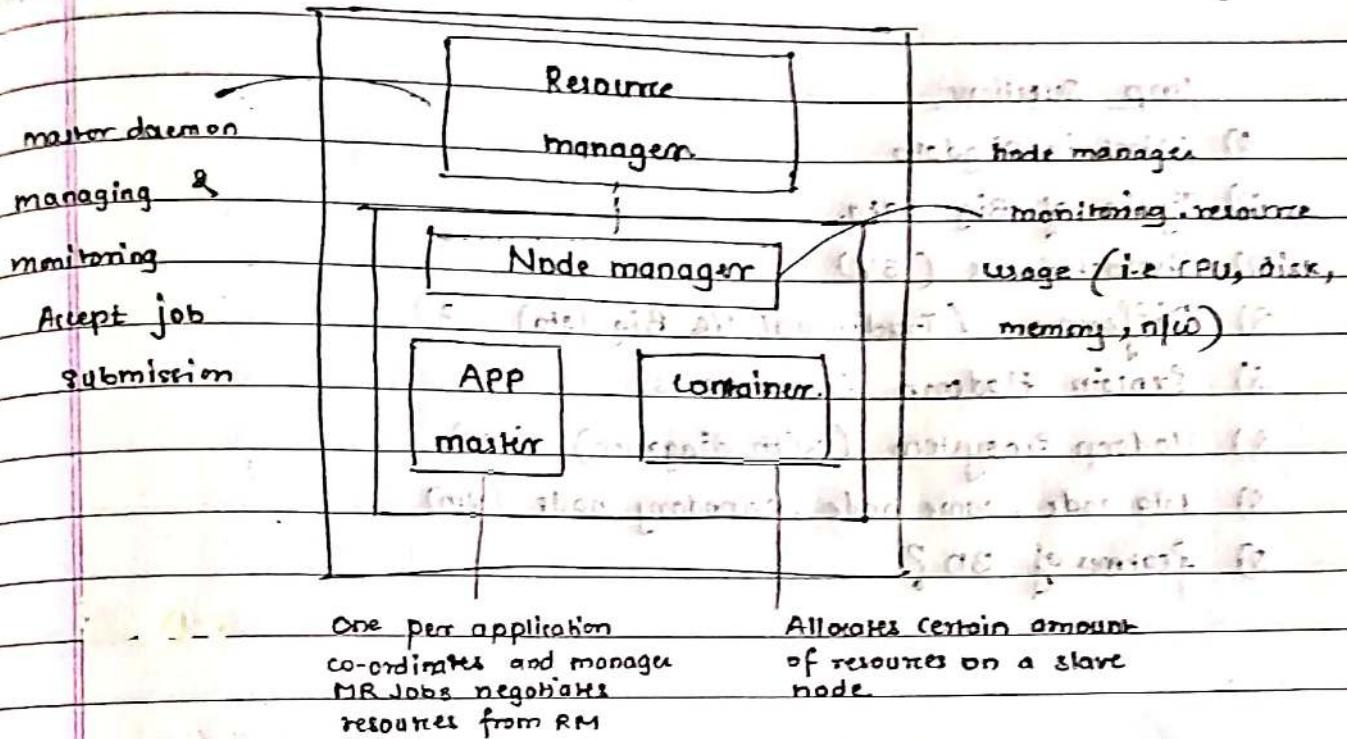
Reduced 8
mappe.



3) YARN (Yet another resource manager) (hadoop 2.0)

* Hadoop Ecosystem

Apache Hadoop Yarn management (Yet another resource negotiator)



* ECOSYSTEM

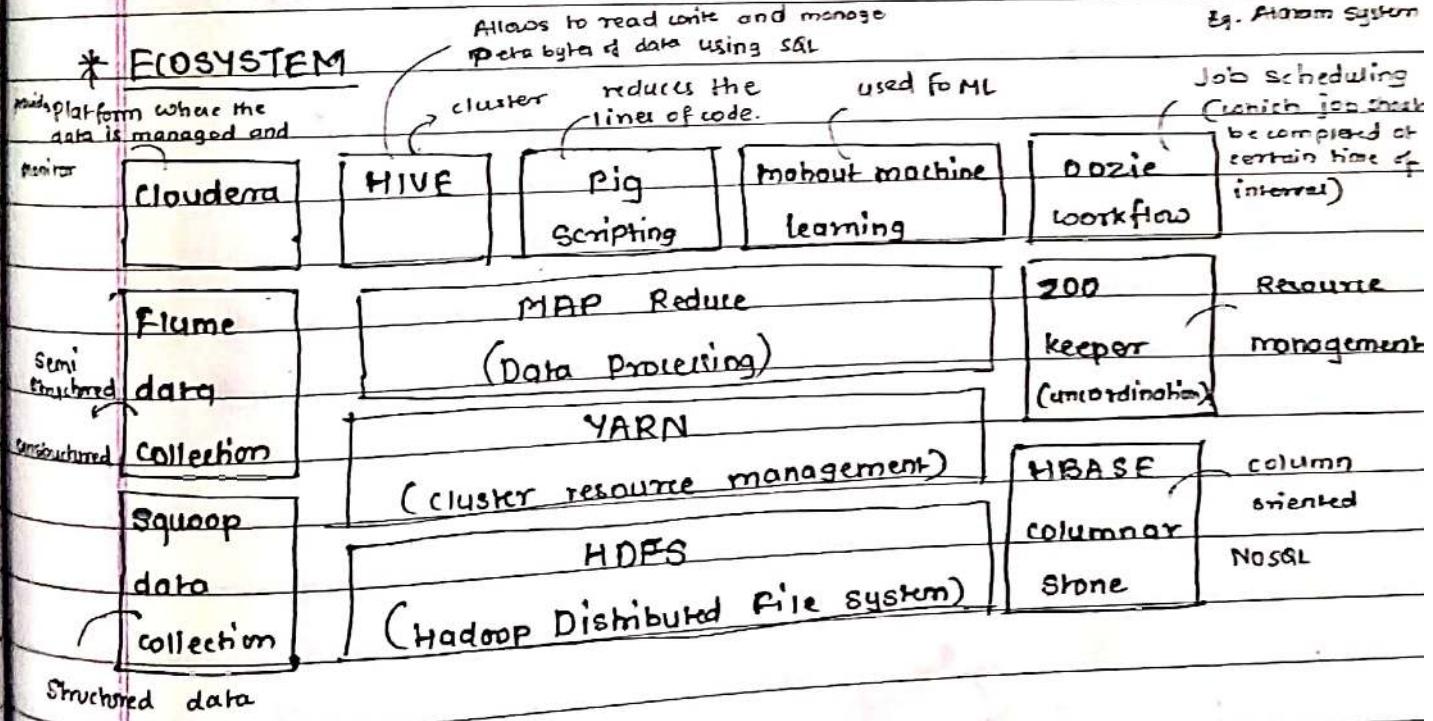
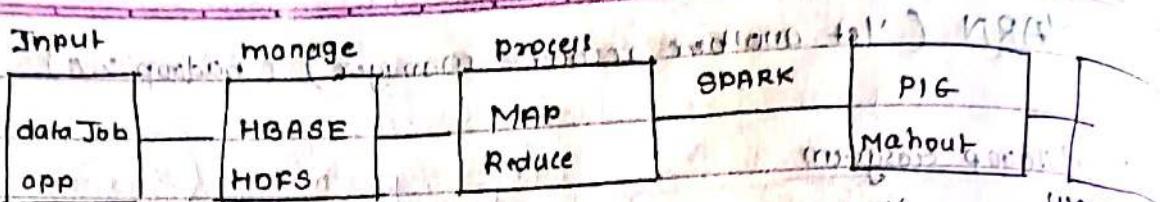


fig. Ecosystem (Hadoop)



User can
Search data
Using HDFS

Imp Questions:

- (Q) What is Big Data
- (Q) Types of Big data
- (Q) Characteristics (5V)
- (Q) Difference (Traditional Vs Big data) (2)
- (Q) Explain Hadoop (Components)
- (Q) Hadoop Ecosystem (with diagram) (5m)
- (Q) data node, name node, secondary node (2m)
- (Q) Features of BD?

18/7/23
Tuesday.

LAB (BDA)

- vi & nano editors.

for vi editor, editor's working

! vi filename (enter) : edit mode

write the content and

press esc, :, wq! (enter) [you will come out of the editor]

create file touch filename.txt

instead of hadoop fs -ls we can use

hdfs dfs -ls

(option) command

when creating a directory,
`drwxr` is seen where (d) is prefixed.

* Basic of MAP Reducer

Mapper

$((1, 1, 1), (1, 1))$ \rightarrow $i = 1, j = 1$
 $((1, 1, 1), (0, 0))$ \rightarrow $i = 1, j = 2$
 $((1, 1, 1), (1, 0))$ \rightarrow $i = 1, j = 3$

$((2, 1, 1), (0, 0))$ \rightarrow $i = 2, j = 1$
 $((2, 1, 1), (1, 0))$ \rightarrow $i = 2, j = 2$

* matrix Vector multiplication by Map reducer

$$\begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \underset{2 \times 2}{\times} \begin{bmatrix} 2 & 0 & 1 \\ 5 & 1 \end{bmatrix} \underset{3 \times 2}{\rightarrow} \begin{bmatrix} a & b \\ c & d \end{bmatrix} \underset{2 \times 2}{=} \begin{bmatrix} ((B1, A2 \times 2), 0, 1) \\ (A, B1, A2 \times 2, 0, 1) \end{bmatrix}$$

Whenever we are multiplying two matrices, there is a rule, the total no. of columns in the first matrix should be equal to the total no. of rows from the second matrix. And if they are equal

$$\begin{bmatrix} M_{00} & (M_{01}, 0, 1) (0, 0) \\ M_{10} & (M_{11}, 0, 1) (1, 0) \end{bmatrix}$$

Steps for two matrix mul. using MAP reducer.

- * ① for matrix 1 $\rightarrow (A, 0, 1) (1, 0)$
 $((1, 0, 1) (0, 0))$
 $((i, k), (A, j, Aij))$
value

Key

for matrix $2 \rightarrow B$

$((i, k), (B, B_j, B_{jk}))$

key

value.

for matrix A $(i, j) (A, A_{ij}, A_{ij})$

$A_{00} = 1 ((0, 0), (A, 0, 1))$

$A_{ij} ((0, 1), (A, 0, 1))$

$A_{01} = 2 ((0, 0), (A, 1, 2))$

$A_{ij} ((0, 1), (A, 1, 2))$

$A_{10} = 3 ((1, 0), (A, 0, 3))$

$A_{ij} ((1, 0), (A, 0, 3))$

$A_{11} = 4 ((1, 0), (A, 1, 4))$

$A_{ij} ((1, 1), (A, 1, 4))$

for matrix A, key, value $(i, j) (A_{ij}, A_{ij})$

Operations on map, addition, subtraction with arithmetic operators, multiplication with multiplication operator, division with division operator.

Example: $\{((1, 0), 1), ((1, 1), 2)\} + \{((1, 0), 3), ((1, 1), 4)\}$

$\begin{bmatrix} 2 & 1 \\ 5 & 1 \end{bmatrix}$ for matrix B,

$B_{00} = 2 ((0, 0) (B, 0, 2))$

$B_{ij} ((0, 1) (B, 0, 2))$

$B_{01} = 1 ((0, 1) (B, 0, 1))$

$B_{ij} ((0, 0) (B, 0, 1))$

$((i, i), (j, j))$

$$B_{10} = 5 \quad \begin{pmatrix} (0,0) & (B, 1, 5) \\ (1,0) & (B, 1, 5) \end{pmatrix}$$

$$B_{11} = 1 \quad \begin{pmatrix} ((1,0), (B, 1, 4)) \\ ((1,1), (B, 1, 1)) \end{pmatrix}$$

*2 Combiner Task

used to combine key value pairs to reduce network congestion.

$(0,0)$, $(A,0,1)$, $(A,1,2)$, $(B,0,2)$, $(B,1,5)$

$(0,1)$, $(A,0,1)$, $(A,1,2)$, $(B,0,2)$, $(B,1,1)$

$(1,0)$, $(A,0,3)$, $(A,1,4)$, $(B,0,2)$, $(B,1,5)$

$(1,1)$, $(A,0,3)$, $(A,1,4)$, $(B,0,1)$, $(B,1,1)$

*3 Reducer Task

a) for $(0,0)$

$(A,0,1)$

$(B,0,2)$

1×2

$(A,1,2)$

$(B,1,5)$

2×5

$$2 + 10 = 12$$

Element selected at $F_0 = 12$

b) for $($

every job does some combiner task which is not yet given

partitioned partitioned minor and major

minibatch artifacts

$$\delta = \Delta \cdot \frac{\epsilon}{\delta}$$

minor

major

$(s,1), (s,2), (s,3)$

$(s,4), (s,5), (s,6)$

$(s,7), (s,8)$

* Relational algebra operations using map reduce.

- selection - difference
- projection - natural join
- union - Grouping & Aggregation
- intersection

(1) Selection Algorithm

$\text{map}(\text{key}, \text{value}) :$

$\text{emit}(\text{key}, \text{tuple})$

MAP worker 1

A	B	A	B
1	2	1	2
3	1	3	5

Map worker 2

A	B	A	B
2	3	1	1
3	5	2	1

Big data divides entire data in various parts each part will be process separately & parallelly.

* Selection conditions

e.g. $B \leq 2$

MW1

key	value
(1, 2)	(1, 2), (1, 2)
(3, 1)	(3, 1)
(1, 2)	(1, 2)

MW2

key	value
(2, 1)	(2, 1)

Map worker checks for duplicacy of the data.

- * A hash function will be applied.

MW1

key	value
(1, 2)	(1, 2), (1, 2)

MW

key	value
(3, 1)	(3, 1)

MW2

key	value
(1, 1)	(1, 1)

key	value
(2, 1)	(2, 1)

- * To reduce task swapping is done.

- Swapping will be done to discard redundancy that can be caused by duplicate.
- It will swap the adjacent tables with the adjacent workers.

RW1

key	value
(1, 2)	(1, 2), (1, 2)

RW2 reduced.

k	v	k	v	k	v
(3, 1)	(3, 1)	(3, 1)	(2, 1)	(2, 1)	

- * To join two tables

RW1

key	value
(1, 2)	(1, 2) (1, 2)
(1, 1)	(1, 1)

RW2

key	value
(3, 1)	(3, 1)
(2, 1)	(2, 1)

*

RW1

A	B	A	B
1	2	3	1
1	1	2	1

Map worker checks

for duplicacy of the data.

- * A hash function will be applied.

MW1 visible after 3rd MW2

key	value	key	value
(1,2)	(1,2), (1,2)	(3,1)	(3,1)

(1,2) (1,2) sum

MW2

key	value	key	value
(1,1)	(1,1)	(2,1)	(2,1)

- * To reduce task Swapping is done.

- Swapping will be done to discard redundancy that can be caused by duplicate.
- It will swap the adjacent tables with the adjacent workers.

RW1

key	value	k	v	key	v	k	v
(1,2)	(1,2), (1,2)	(1,1)	(1,1)	(3,1)	(3,1)	(2,1)	(2,1)

RW2 reduced.

- * To join two tables

(1,1) RW1 (1,1)

(key) value (1)

(1,2) (1,2) (1,2)

(1,1) (1,1)

(key) value (2)

(3,1) (2,1)

(2,1) (2,1)

*

RW1

RW2

A	B (1,2)
1	2
1	1

A	B
1	2
2	1

* Projection Algorithm

$\text{Map}(\text{Key}, \text{Value})$

for tuple in value:

 if ts = tuple with only attributes in s:

 emit(ts, ts)

$\text{Reduce}(\text{key}, \text{value})$

 emit(key, key)

project B : C

Eg

Name
age
Salary

Map worker 1					
A	B	C	A	B	C
1	2	3	2	2	3
3	1	4	4	3	4

Map worker 2					
A	B	C	A	B	C
5	7	2	3	2	5
7	1	1	4	1	0

Step 1: (1,2) MW1 (1,3) (1,1) (1,4) MW2 (1,1)

key	value	key	value
(2,3)	(2,3) (2,3)	(7,2)	(7,2)
(1,4)	(1,4)	(1,1)	(1,1)
(3,4)	(3,4)	(2,5)	(2,5)
		(1,0)	(1,0)
		(1,1)	(1,1)

Step 2: Hash function will be applied, to divide entire table in two parts.

for Map worker 1		MW1	key	value
(2,3)		(2,3) (2,3)	(3,4)	(3,4)
(1,4)		(1,4)		

for map worker 2.

MW2			
key	value	key	value
(7, 2)	(7, 2)	(2, 5)	(2, 5)
(1, 1)	(1, 1)	(1, 0)	(1, 0)
(T1)		T2	
(swap pair) time			

Step 3: Swapping is done to remove the redundancy, i.e. present in both Map workers.

(Reducer worker 1)

RW1 :

key	value	key	value
(2, 3)	(2, 3)	(7, 2)	(7, 2)
(1, 4)	(1, 4)	(1, 1)	(1, 1)

(Reducer worker 2)

RW2 :

key	value	key	value
(3, 4)	(3, 4)	(2, 5)	(2, 5)
		(1, 0)	(1, 0)

Step 4: Merge after combining.

(E, 2) (E, 5)

(E, 1) (E, 3) (E, 4)

(E, RW1) (E, RW2)

key	value
(2, 3)	(2, 3)
(1, 4)	(1, 4)
(7, 2)	(7, 2)
(1, 1)	(1, 1)

(E, E) (RW2)

(1, 1) key value

(3, 4) (3, 4)

(2, 5) (2, 5)

(1, 0) (1, 0)

Steps

(1, 1) (E, RW1) (E, RW2)

RW2 (1, 1)

B	C
2	3

B	C
3	4

1 SW 4

2 SW 5

V S 7 2 V S

V S 1 0 V

(1, 1)	(1, 1)	(1, 1)	(1, 1)
(1, 1)	(1, 1)	(1, 1)	(1, 1)

* Union Algorithm

Map (key, value):

for tuple in value:

emit (tuple, tuple)

reduce (key, value)

emit (key, key)

Q Map worker 1

A	B	A	B
1	2	1	2
3	1	2	1

Map worker 2

A	B	A	B
2	3	1	1
4	5	5	2
1	1	1	1

Step 1:

MW1

MW2

key	value
(1, 2)	(1, 2) (1, 2)
(3, 1)	(3, 1)
(2, 1)	(2, 1)
(2, 1)	(2, 1)

key	value
(2, 3)	(2, 3)
(4, 5)	(4, 5)
(1, 1)	(1, 1)
(2, 1)	(2, 1) (2, 1) (2, 1)

Step 2:-

MW1 (1, 2)

MW2 (2, 1)

K	V	K	V
(1, 2)	(1, 2) (1, 2)	(2, 1)	(2, 1)
(3, 1)	(3, 1)		

K	V	K	V
(2, 3)	(2, 3)	(1, 1)	(1, 1)
(4, 5)	(4, 5)	(2, 1)	(2, 1)

Step 3:-

RW1

RW2

K	V	K	V
(1, 2)	(1, 2) (1, 2)	(2, 3)	(2, 3)
(3, 1)	(3, 1)	(4, 5)	(4, 5)

K	V	K	V
(1, 1)	(1, 1)	(2, 1)	(2, 1)
(2, 1)	(2, 1)		

Step 4: Merge & combine.

	RW1	RW2	
key	value	key	value
(1,2)	(1,2)(1,2)	(2,1)	(2,1) (2,1)
(3,1)	(3,1)	(2,1)	(2,1)
(2,2)	(2,2)	(1,1)	(1,1)
(4,5)	(4,5)	(2,1)	(2,1)

Step 5 (Union)

	RW1	RW2	
A	B	A	B
2		2	
3	1	1	1
2	3	4	5
(1,2)(4,5)	5 (1,2)	(2,1)(2,1)	(2,1)
(1,2) (1,2)		(1,2) (1,2)	

* Intersection Algorithm

Map (key, value):

```
def map(self, key, value):
    for tuple in value:
        self.addtuple(tuple)
```

```
    emit (tuple, tuple)
```

```
reduce (key, values)
```

```
if values == [key, key]:
```

```
    c
```

Step 1: MW1

K	V
(1,2)	(1,2)(1,2)
(3,1)	(3,1)
(2,1)	(2,1)

MW2

K	V
(2,3)	(2,3)
(4,5)	(4,5)
(1,1)	(1,1)
(2,1)	(2,1)

Step3:-

MW1

K	V
(1,2)	(1,2)(1,2)
(3,1)	(3,1)

MW1

K	V
(2,1)	(2,1)
(2,3)	(2,3)

K	V
(1,3)	(1,3)
(4,5)	(4,5)

MW1

K	V
(1,1)	(1,1)
(2,1)	(2,1)

K	V
(1,1)	(1,1)
(2,1)	(2,1)

Step4:-RW1

K	V
(1,2)	(1,2)(1,2)
(3,1)	(3,1)

RW2

K	V
(2,1)	(2,1)(2,1)
(1,1)	(1,1)

Step4:-RW1

K	V
(1,2)	(1,2)(1,2)
(3,1)	(3,1)
(2,3)	(2,3)
(4,5)	(4,5)

RW2

K	V
(2,1)	(2,1)(2,1)
(1,1)	(1,1)

Condition:- If length of values is greater than 1 then emit the key part of it.

A	B
1	2

Output

Input

(1,2)

(1,1)(1,1)

(3,1)

(1,1)(1,1)

(2,3)

(1,1)(1,1)

(4,5)

(1,1)(1,1)

* Difference Algorithm

map (key, value)

if (key == R)

for tuple in value:

emit (tuple, R)

else

for tuple in value:

emit (tuple, s) *with space after emit*

reduce (key, value):

if values == [R]:

emit (key, key)

Step 1.

MW1

A	B	AB
1	2	1 2
3	1	2 1

MW2

AB	A	B
2	3	1
4	5	1 1

Step 2 *using MW1 info to get difference in MW2*

Key	Value ranges starting at key	Value ranges
(1, 2)	[T ₁ , T ₂]	(T ₁) <i>subset</i>
(3, 1)	(T ₁)	(T ₂) <i>subset</i> (1)
(2, 1)	(T ₂) <i>beginning</i> (1, 1)	(T ₂) <i>overlap</i> - (2, 1) <i>overlap</i> (1)

Step 3: Apply hash function.

Key	Value over K	Value over K	Value over K	Value over K
(1, 2)	[T ₁ , T ₂] (2, 1) [T ₂]	(T ₁)	(T ₁) <i>subset</i> (1, 1)	(T ₂)
(3, 1)	[T ₁]		(T ₁) <i>subset</i> (2, 1)	(T ₂)

RW2

Step 3:

RW1

K	V	K	V
(1,2)	[T ₁ , T ₂]	(2,1)	[T ₂]
(3,1)	[T ₁]	(1,1)	[T ₂]
(2,3)	[T ₁]	(2,1)	[T ₂]
(4,5)	[T ₁]	(3,4)	

Step 4: eliminate the keys that are present both in [T₁] & [T₂] as well as only from T₂: (1,1), (2,1), (3,1)

A	B	$\{a\} = \{1, 2, 3\}$
3	1	$\{1, 2, 3\} \cap \{1\} = \{1\}$
2	3	
4,5	5	

* Hadoop Limitation

(1) Security problem.

- there is no inbuilt security in hadoop. instead we can use SPARK technology to provide security. (AES) Advance Encryption Standard.

(2) Issue with the small files.

- Hadoop is typically java based. Each block size is 128 bits MB.

(3) Slow processing speed.

- Hadoop SPARK overcomes this limitation.

(4) Support for batch processing

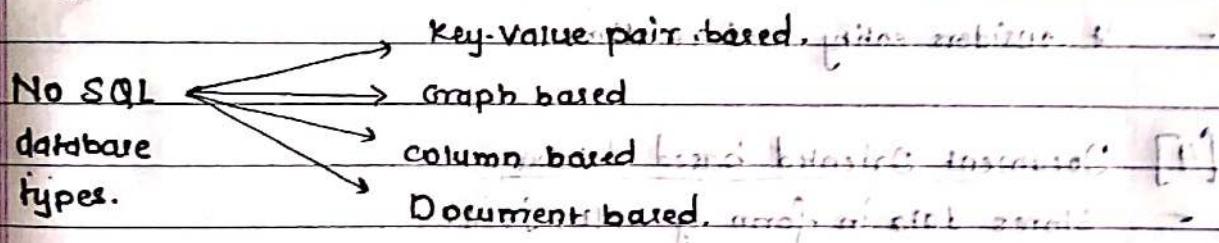
- The batch wise processing. by using hadoop SPARK, the continuous processing is possible. (flink) it's also useful for live streaming of data.

(5) Not easy to use.

- basic hadoop is java oriented and requires hardcode. instead HIVE Spark can be used for implementation.

3. No SQL

- No SQL stands for 'Not only SQL'.
 - why use No SQL?
 - DBMS: RDBMS can be used for tabular data storage, but, the data now-a-days is in semi structured, unstructured and structured form. Therefore, NoSQL can store data 'not only' tabular form but also can store semi-structured & unstructured data.
 - who invented:- NoSQL is invented in 1998 by Karl Storaasli.
- * NoSQL handles various types of data:
- Semi-structure
 - Unstructure
 - Rapidly changing data.
 - Big Data.



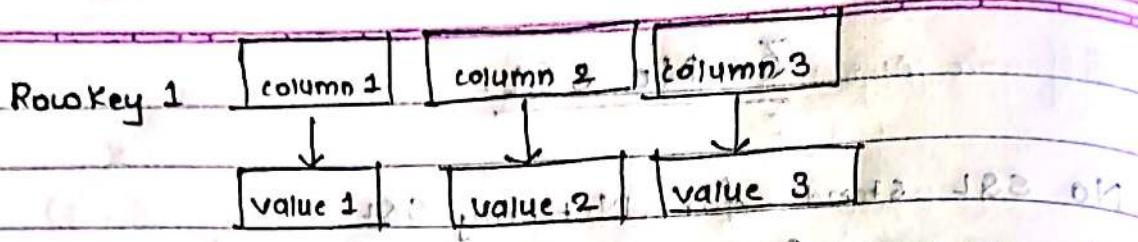
[1] Key-Value pair based:

Eg	key	value
	Name	sakshi
	DOB	30-10-2002

Eg - Redis, Riak are the databases that follows key-value pair.

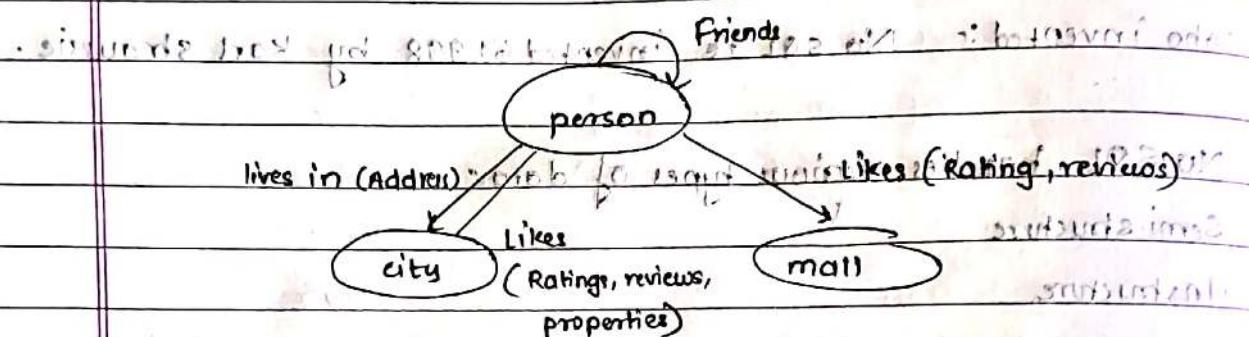
[2] Column based oriented based:

- Eg. HBASE, Cassandra



- Hbase and Cassandra follows column oriented based database.

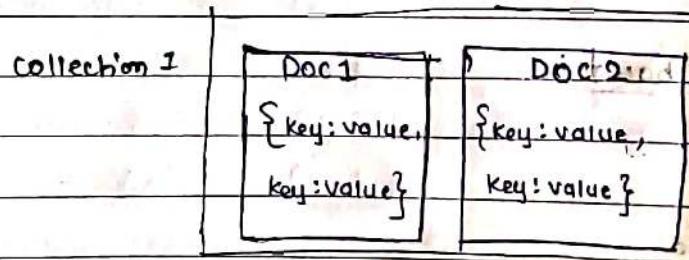
[3] Graph based database: neo4j, OrientDB, spatial Data



- It can be used in multirelational database.
- It considers entity and relationships.

[4] Document Oriented based database

- Stores data in form of XML or JSON.



e.g. MongoDB, CouchDB are databases which follows document based database.

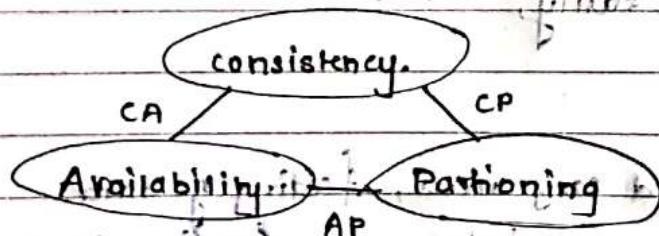
[Q] Difference between SQL & NoSQL.

SQL

NoSQL

- | | |
|--|--|
| (1) Relational Database Management system. | (1) Non-Relational/Distributed database system. |
| (2) Fixed or static (or predefined) schema. | (2) It has dynamic schema. |
| (3) Not suited for hierarchical data storage. | (3) Best suited for hierarchical storage (using trees). |
| (4) Vertically scalable. | (4) Horizontally scalable. |
| (5) It follows acid ACID properties. (Atomicity, consistency, Isolation and Durability). | (5) It follows CAP (Consistency, Availability, Partition tolerance). |

* CAP Theorem



Eg

A : 1000	B : 2000	C : 3000	D : 4000	E : 5000	F : 6000	G : 7000	H : 8000
----------	----------	----------	----------	----------	----------	----------	----------

the nodes connected together have the replicated copies stored

therefore if some node fails, then other nodes will have

(Network, where nodes fail in the replicated copies, hence have some data).

no work is stopped.

for partitioning, (horizontal arrangement, so to distribute the work)
 for availability - The system works even if network fails.

- In NoSQL, only two condition is satisfied at a time.
 either, CA (MariaDB, SQL server)
 or AP (CouchDB, Cassandra)
 or CP (MongoDB, Redis)

* Business Rivalry Drivers

parameters:

1) Volume

2) Velocity

3) Variability

4) Agility

[1] Volume:

- need to scale up

- Data is extending continuously.

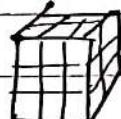
[2] Velocity:

- slow response time

- RDBMS

[3] Variability:

- OLAP server where
be stored easily,



3D, where variety of data can

[4] Agility:

- Quick and easy access, feeding of data
- The time of retrieving and feeding data must be instant.
- RDBMS find difficult if data contains or divides into many subgroups which decreases response time.

* NoSQL features

- Most of the NoSQL database are open source.

- It is schema free / schema less.

- They doesn't require object oriented mapping & data normalization.

- Most of the NoSQL database executes in distributed fashion and they offer auto-scaling.
- NoSQL not follows relational model and table with flat fixed column records.
- It offers heterogeneous structure of data in the same domain.
- It is generally used to store big data and real time data.
- It follows CAP theorem (Any two at once). conditions

* Features of MongoDB.

(1) Adhoc queries:

- These are types of queries which are not known while structuring the database. Whenever we are designing a database we don't know in future which types of queries may come.

(2) Aggregation:

- We can batch process the data and get single result after doing some sequence of task on group of data.
steps in, input → match → group → sort → output.

(3) Schemaless / schema-free:

- It doesn't care about pre-defined schema or structure.

(4) Grid File System (GFS):

if file contains multiple members in different parts
chunk chunk chunk chunk.

(5) Document oriented.

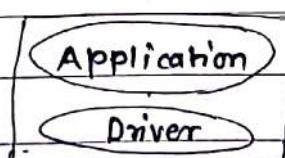
structured for documents

(6) Sharding:

means data partitioning

- High availability of data.

→ scalability

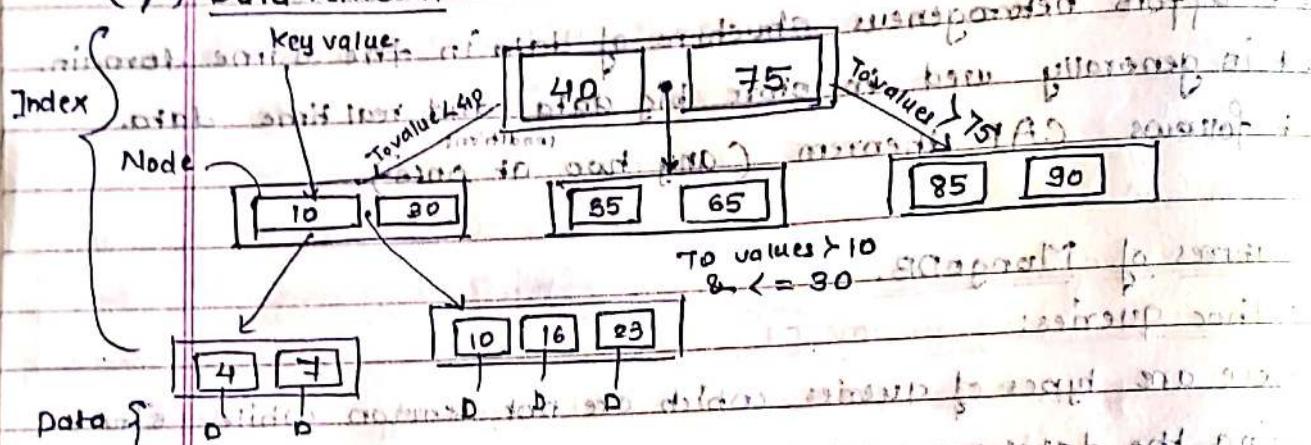


Shard 1
Primary
Secondary
Secondary

Shard 2 Shard 3 Shard n

Sharding separates very large db's into smaller, faster, more easily managed called data sharding. It is often used here when documents are too large for a single database.

(*) Data retrieval



(*) High performance

- MongoDB shows high availability

- It has better query response for indexing and replication

(*) Replication

- This feature creates copy of document into different machines. One can have primary node and its replicated node where primary node cannot work for some reason then replicated or secondary node takes care of it.

* Components of MongoDB

(1) Collection.

(2) Database.

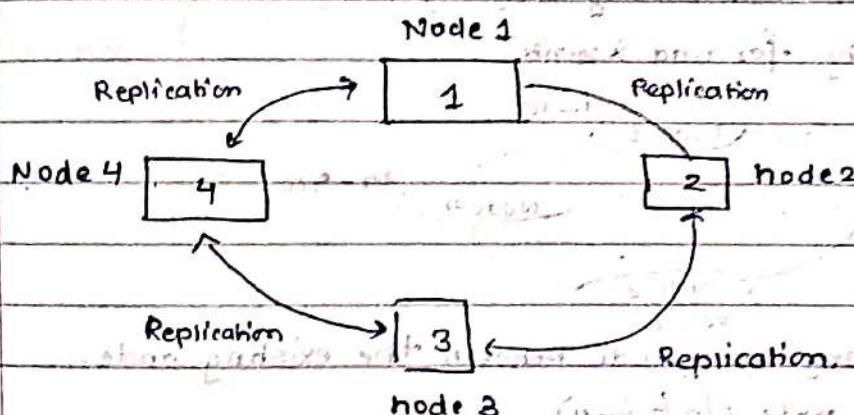
(3) document

* Cassandra

- It is a key-value as well as column oriented database.
- It is an open source database.
- You can go for 3rd party support.

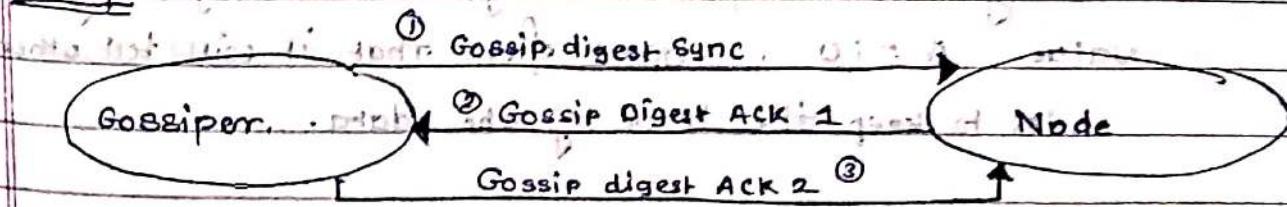
- It is also schemaless / Schema free.
- It is scalable & distributed.
- data is distributed across multiple nodes called as a cluster.
- It is horizontally Scalable
- fault tolerant.
- If any node is failure in cluster it will not affect client.
(Read & write still occurs even if it fails)
- high performance
- It handles multiple petabytes of data.
- It uses best optimized algorithm for read & write.
- Eg. Netflix uses cassandra for recommendation system.
- It supports ACID properties.

* Architecture of Cassandra



- Note:- It uses Gossip protocol which allows node to communicate with each other or to find faulty node in structure.
- Data replication in Cassandra.
 - It will check for failure detection and used for communication.

① Gossip



To detect failure **Phi accrual** is used.

- ② * Snitches (latency)
- It comes under the Cassandra architecture.
 - Note: which node provides minimum latency for read & write operation.

① Simple Snitching:

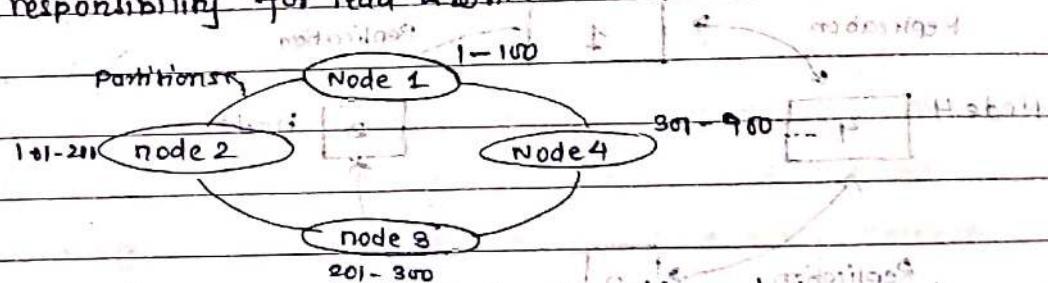
- It identifies the node in a single datacenter which of the rack provide minimum latency for read & write operation.

② Dynamic Snitching:

- which of the node will provide least latency for read & write operation, for this gossip protocol will be used in a multiple cluster.

③ * Rings and tokens

- co-ordinator node takes care of assigning the / communication. takes responsibility for read & write.



- co-ordinator is any one node between the existing node. (here, it can be node 1/2/3/or 4)

- e.g. Read & write operations

- Write:
- In this, all the request go to the co-ordinate, multiple node.
 - It will take hash of the key, for eg., hash value = 50.
 - Co-ordinator node sends the request to the node that is dealing with the partition range (token range) by to save value A = 10, apart from that it will tell other nodes to keep the track of the data.

* Q7 write short note on Cassandra Architecture & Cassandra MongoDB.

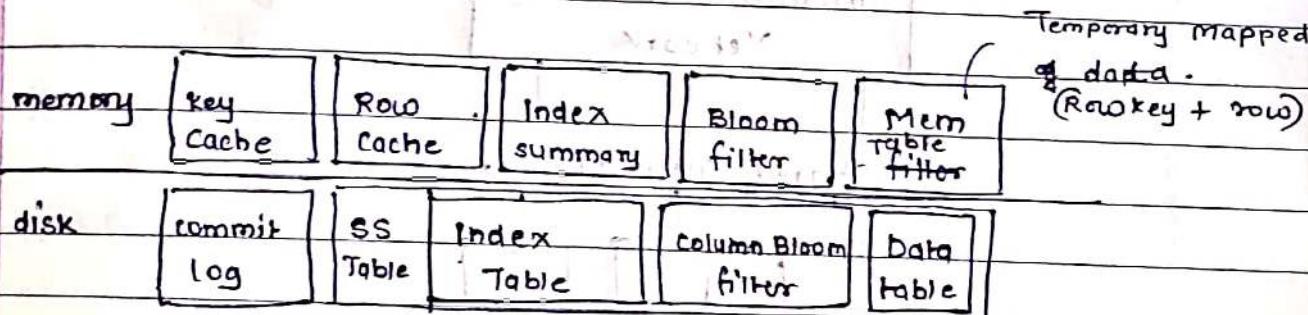
Ques	1	1
Page		

Read :

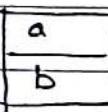
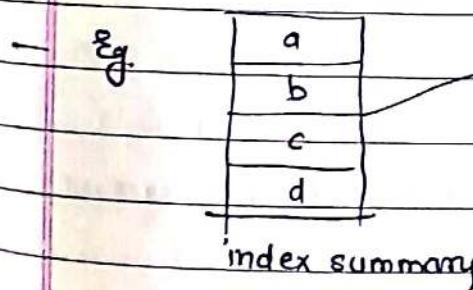
- If none of the nodes responded values of $a=10$ & one node is responded with value $a=5$, in this case co-ordinator will revert to client with value $a=10$.
- The co-ordinator node will tell that node to update the value i.e., $a=5$ to 10 , such type of mechanism is called read repair.

* Components

- Node - data collection of related nodes
- Data center - one or more data centers
- Cluster - sequence of actions
- commit log - memory
- mem table
- SS Table
- Bloom filter - (probabilistic structure)
 - value is (0 to 1) range
 - tells if data is present or not present.

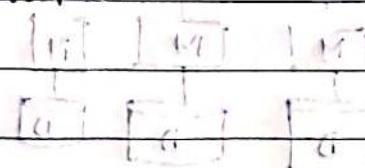


- SS \rightarrow sorted string tables \rightarrow (Index table, column Bloom, Data table).



index table

- Column bloom filter: It will tell whether the column is present inside key-value pair.



key cache: stores primary key of row vs offset (location)
in data file

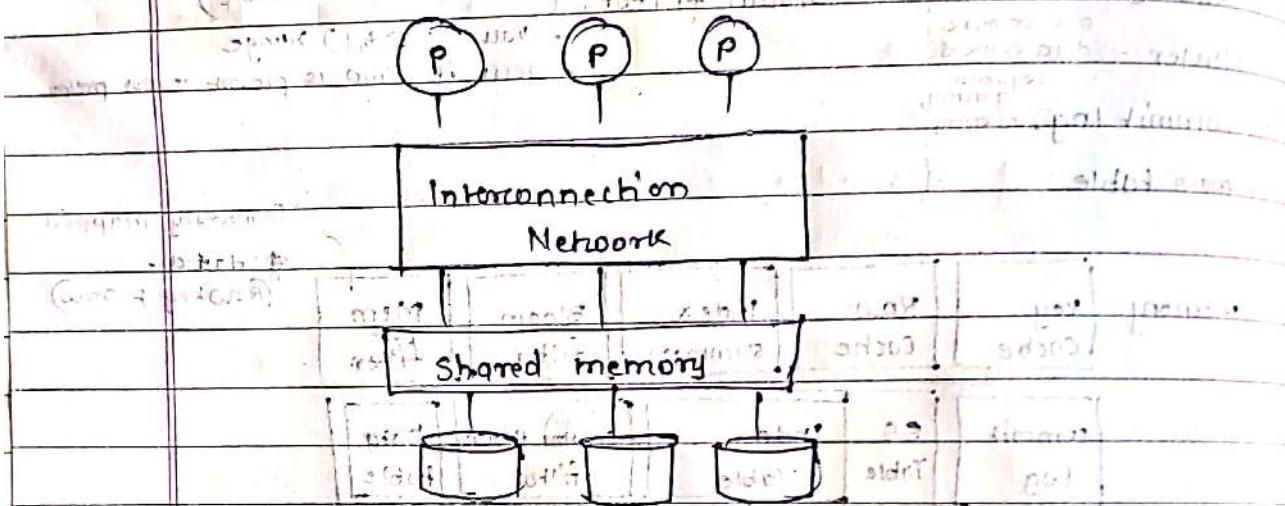
row cache: holds the list of rows.
- a complete row can be stored in a single node.

* Analyzing Big Data with shared nothing Arch.

- shared memory system. (e.g. UNIX Filesystem)

IPC

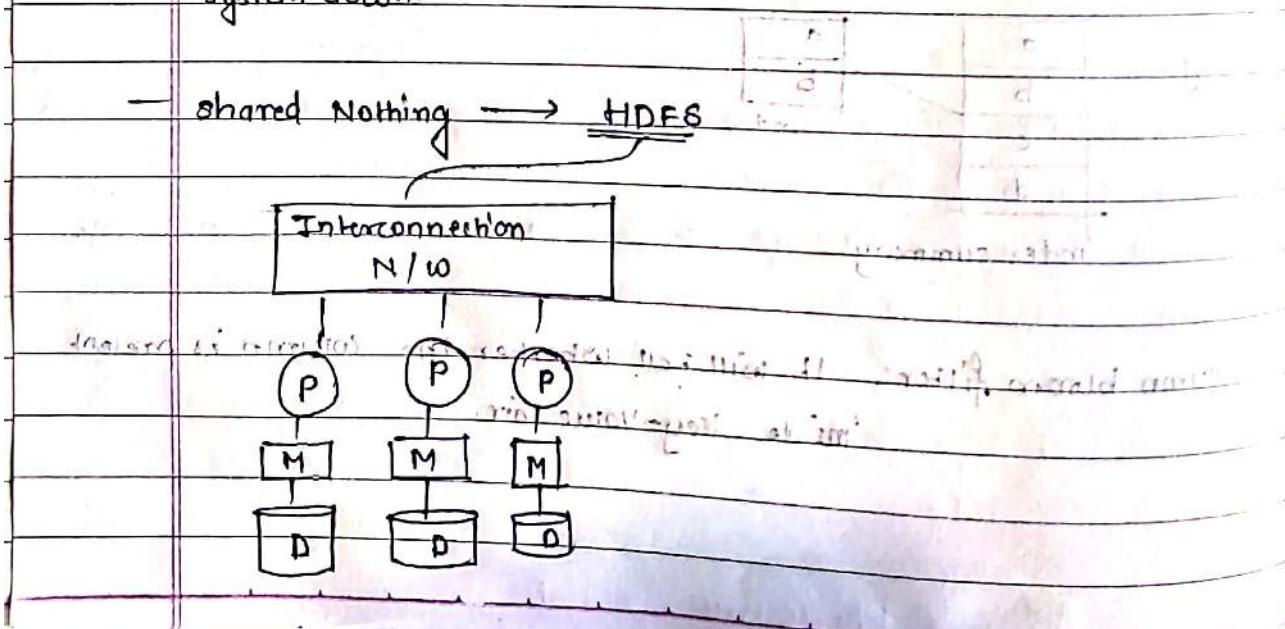
- Interprocess memory communication system.



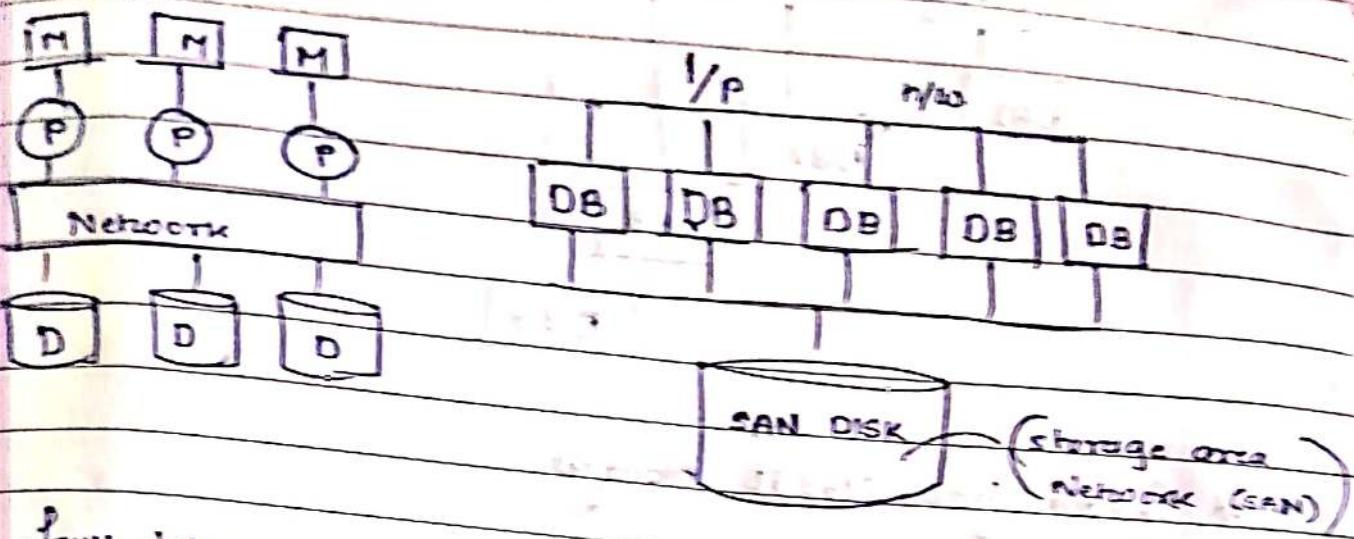
Disadvantage:

- like bus topology if one line fails, entire system fails
- System down.

- shared Nothing \rightarrow HDFS



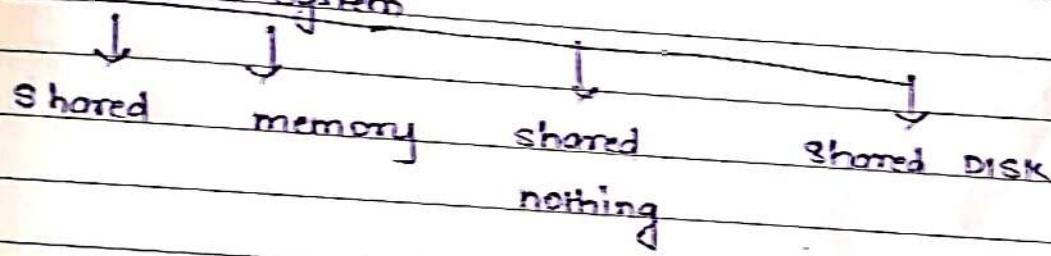
Shared Disk → ORACLE



- fault tolerance:

If a processor or its memory fails, the other processor can take its task since the database is present on disk and are accessible to all processor.

* Hierarchical System



→ Distributed Model.

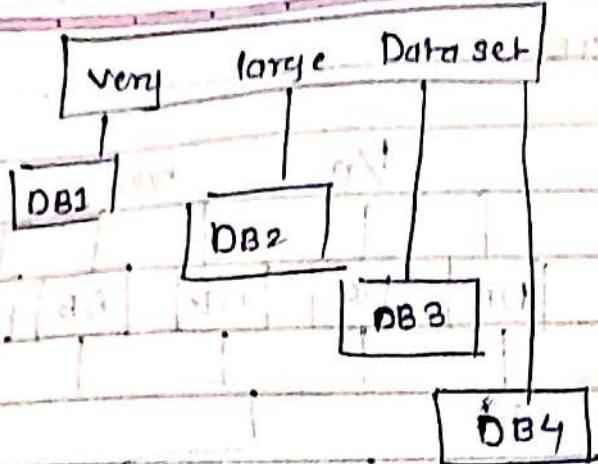
Master-slave Peer-to-Peer.

MongoDB

Cassandra

→ Sharding

- different database present in different server hence implementing sharding



→ Explain how NoSQL manages Big Data?

1 — No SQL limits

Types of



CQL → Astra (database)
cassandra query language.



Dynamo db uses key-value



Datamodels → mongoDB → BSON (binary) (Json like)
cassandra (CQL)

Neo4J (Cypher QL)

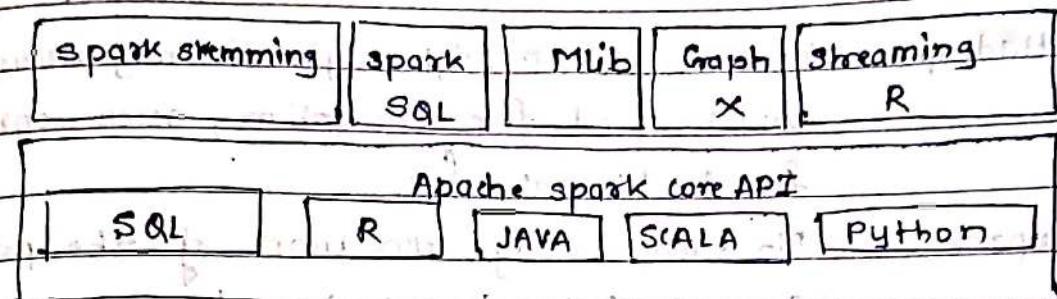
DynamoDB

consistency: Ø → cassandra → Neo4j → DynamoDB

MDB

3. Apache Spark

* Basic Spark Architecture.



- 100 times faster in memory (RAM)

- Requires more RAM, hence cost is more.

* Spark Vs Hadoop.

1) Data processing :

- Hadoop : - Batch processing, historical data
 - data which is collected over a period of time and process at later stage.

- spark : - Base on string processing
 - Real time information delivery.

2. Performance :

- Hadoop : slow speed of process (processing speed)
- spark : - Much faster processing speed
 - It stores & processes everything in memory.
 - suitable for Machine learning, IoT sensors, Security analysis.

3. Ease of Use :

- Hadoop : - It provides add-on to support user but doesn't have interactive mode.

- spark : user friendly API's (SQL, R, Java, Scala, Python) which provides interactive modes for developers.

4. Security :

- Hadoop : Kerberos, LDAP (lightweight directory access protocol) will be used for encryption, access control, etc for traditional file.
- spark : - No much security, uses more of the passwords & much likely to get hacked.

5. Price :

- Hadoop : less costlier.
- spark : more costlier.

6. Job schedulers :

- Hadoop : oozie, workflow, piggyback, spark
- spark : No external job scheduler is required.

7. Scalability :

Hadoop : Yahoo has 40 K nodes.

Spark : It has around 7000 nodes.

8. Latency :

- Hadoop : High latency computing web framework
- spark : low latency computing framework.

9. Fault tolerance :

- Hadoop : Job tracker, task tracker. (task tracker that provides heart beat. If heart beat is missed, all the pending operation will be rescheduled to other job tracker)
- Spark : RDD (Resilient distributed dataset) for building blocks for fault tolerance.
- they operate in a parallel.

10. compatibility

Hadoop : compatible.

spark : (compatible with each other.
(spark & map reducer)

MODULE 4: Mining Data Stream

* Characteristics of data stream

- (1) Huge volume of infinite continuous data.
 - (2) Fast changing require real time & fast response.
 - (3) single scan → entire data is streaming continuously
 - (4) System stores only summary (knowledge from the data received till now).
- e.g. IOT sensors, financial market, amazon website.

* DBMS Vs. DSMS

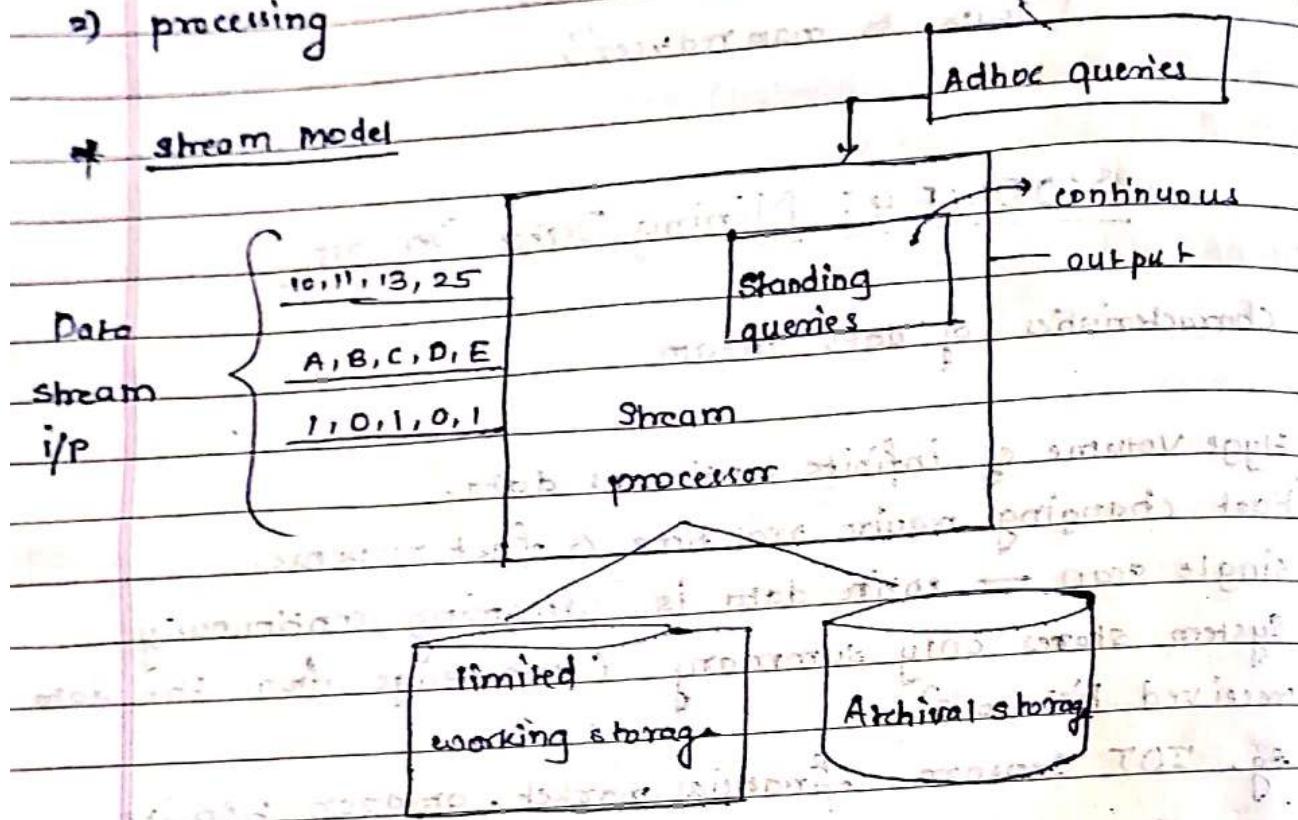
- | | |
|--|----------------------------|
| (1) Dealing with persistent & static data. | Dealing with stream data. |
| (2) Random Access of Data. | sequential access of data. |
| (3) Query provides exact answers. | Approximate answers. |
| (4) One time queries | continuous. |
| (5) No real time requirement | real time. |

* challenges

1) Storage

2) processing

* stream model



→ standing Q : Queries are asked to the stream (continuous)
e.g. Alert me whenever current value exceeds 50 mA

→ Adhoc Queries : Queries can be asked only one time in stream / It is not continuous.

e.g. what is the average of current values captured so far.

Ques *

* Bloom Filter. (filtering data stream)

False Positive : True negative.

- 1) It is used to check whether an element belongs to set (whether element is present or not).
- (2) - two categories : false
 - probably says that element belongs to set (False Positive) it won't give assurity that element is 100% present
 - Accurately says that the element ~~is not~~ does not belong to set (only true Negative).

* Two steps in bloom filter:

(1) Insertion.

(2) Presence of element.

Ex. Insert element 10 and 7

Given: $h_1(x) = x \bmod 5$

$$h_2(x) = (2x + e) \bmod 5$$

(array size) $m = 5$

Comment on presence of element 14 & 15

Element to insert	$h_1(x)$	$h_2(x)$	Bloom filter															
10	0	01	<table border="1" style="display: inline-table; vertical-align: middle;"> <tr><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>1</td><td>1</td><td>0</td><td>0</td><td>0</td></tr> <tr><td colspan="5" style="text-align: center;">0 1 2 3 4</td></tr> </table> <small>initial</small>	0	0	0	0	0	1	1	0	0	0	0 1 2 3 4				
0	0	0	0	0														
1	1	0	0	0														
0 1 2 3 4																		
7	02	0	<table border="1" style="display: inline-table; vertical-align: middle;"> <tr><td>1</td><td>1</td><td>1</td><td>0</td><td>0</td></tr> <tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td></tr> </table> <small>new</small>	1	1	1	0	0	0	1	2	3	4					
1	1	1	0	0														
0	1	2	3	4														
14	4	4	<table border="1" style="display: inline-table; vertical-align: middle;"> <tr><td>1</td><td>1</td><td>1</td><td>0</td><td>0</td></tr> <tr><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td></tr> </table> <small>new</small>	1	1	1	0	0	0	1	2	3	4					
1	1	1	0	0														
0	1	2	3	4														

To find for 14 & 15

Accurate says doesn't exist index 4, the value is exists (T-False) which means element is not present in

(4) 15

Condition:

If any of indices is zero, it is present [for 0 & 1 position which depicts that element is present in the set but in reality it is not].

It only provides the probability of presence (False positive).

$$[Q] h_1(x) = x \bmod 5$$

$$m=5$$

$$h_2(x) = (2x+3) \bmod 5$$

Insert 9 & 11

Element to insert	$h_1(x)$	$h_2(x)$	Bloom filter										
1) 9	4	$(2 \cdot 9 + 3) \bmod 5 = 1$	<table border="1"> <tr> <td>0</td><td>1</td><td>0</td><td>0</td><td>1</td> </tr> <tr> <td>0</td><td>1</td><td>2</td><td>3</td><td>4</td> </tr> </table>	0	1	0	0	1	0	1	2	3	4
0	1	0	0	1									
0	1	2	3	4									
2) 11	3	0	<table border="1"> <tr> <td>1</td><td>1</td><td>0</td><td>0</td><td>1</td> </tr> <tr> <td>0</td><td>1</td><td>2</td><td>3</td><td>4</td> </tr> </table>	1	1	0	0	1	0	1	2	3	4
1	1	0	0	1									
0	1	2	3	4									
Check for 15	0	3	<table border="1"> <tr> <td>1</td><td>1</td><td>0</td><td>X</td><td>1</td> </tr> <tr> <td>0</td><td>1</td><td>2</td><td>3</td><td>4</td> </tr> </table>	1	1	0	X	1	0	1	2	3	4
1	1	0	X	1									
0	1	2	3	4									
(3) 16	1	0	<table border="1"> <tr> <td>1</td><td>1</td><td>0</td><td>0</td><td>1</td> </tr> <tr> <td>0</td><td>1</td><td>2</td><td>3</td><td>4</td> </tr> </table>	1	1	0	0	1	0	1	2	3	4
1	1	0	0	1									
0	1	2	3	4									
			FP										

- * Flajolet Martin Algorithm: used for counting distinct values from the data in input string.
 - It is one of the approximate algorithm.

[Q] I/P stream Given: 1, 3, 2, 1; 2, 3; 4, 3; 1, 2, 3, 1;
 hash function $h(x) = (6x + 1) \bmod 5$

- calculate the function C.bash_fn) for following string.

$h(1) = 2$ $h(3) = 4$ $h(1) = 2$

$h(3) = 4$ $h(1) = 2$

$h(2) = 3$ $h(1) = 1$, if $h(2) = 3$ $h(1) = 1$ if $h(2) = 3$

$h(1) = 2$ $h(3) = 4$ $h(1) = 2$ if $h(3) = 4$

$h(2) = 3$ $h(1) = 2$ if $h(2) = 3$ $h(1) = 2$ if $h(2) = 3$

$h(3) = 4$ $h(1) = 4$ if $h(3) = 4$ $h(1) = 4$ if $h(3) = 4$

$h(4) = 0$

Step: Binary representation of the given decimal number.

2 010 → 1 0110 000 → b (cse)

$$4 \div 100 \Rightarrow 2 \text{----} 0.04 \leq 100 \Rightarrow 2 \text{ (ANS)}$$

$$3 = 011 \rightarrow 0 \quad 0100^2 = 010 \rightarrow 1(E)$$

$$2 \equiv 010 \rightarrow 1 \quad \text{and} \quad 3 \equiv 011 \rightarrow 0$$

$$g = 0.11 \rightarrow 0 \quad 100^4 \approx 100 \rightarrow 2 \cdot (0)$$

$$4 \stackrel{?}{=} 100 \rightarrow 2 \dots \quad 000\ 2\bar{1}0\ 010 \rightarrow 1 \text{ (c)}$$

0.0110 \rightarrow 88 + (2)

con't The count of trailing zeros. (right side zero)

Step: Count the count of trailing zeros. (right side zeros from
(right side zero only after 1's). i.e. binary).

Step: To find out distinct elements.

Trailing zero's = (n) - 2 (maximum trailing zero is 2)

$$\text{distinct value} = (R) = 2^r$$

$$= 2^2 = 4$$

1,2,3,4

- Approximate no. of unique objects in a string or data in one pass.
- If the string contains n elements with m of them unique, this algo runs $O(n)$ time & need $O(\log(m))$ memory. This algo for approximating the no. of distinct element in a stream.

Example: Stream: 4, 2, 5, 9, 1, 6, 8, 7

$$h(x) : (3x + 7) \bmod 82$$

8
32, 16, 4, 2, 1
32, 16, 8, 4, 2, 1

$$h(4) = 19$$

$$h(1) = 10$$

$$h(2) = 13$$

$$h(6) = 25$$

$$h(5) = 22$$

$$h(3) = 16$$

$$h(9) = 2$$

$$h(7) = 28$$

Step: Binary.

trailing zero

$$h(4) = 19 = 0,1,0,1,1 \rightarrow 01011$$

$$h(2) = 13 \rightarrow 0,0,1,1,0,1 \rightarrow 001101$$

$$h(5) = 22 \rightarrow 0,1,0,1,1,0 \rightarrow 010110$$

$$h(9) = 2 \rightarrow 0,0,0,0,1,0 \rightarrow 000010$$

$$h(1) = 10 \rightarrow 0,0,1,0,1,0 \rightarrow 01010$$

$$h(6) = 25 \rightarrow 1,0,1,1,0,0,1 \rightarrow 1011001$$

$$h(3) = 16 \rightarrow 0,1,0,0,0,0 \rightarrow 010000$$

$$h(7) = 28 \rightarrow 0,1,1,1,0,0 \rightarrow 011100$$

Step. maximum = 4. (it's given max skip size 4)

$$2^r = 2^4 = 16.$$

25 Explain Junction with its types.

BOA

- 1 Give all characteristic of big data.
- 2 Explain SV's of big data. (Value, Velocity...)
- 3 Types of big data (Short note)
- 4 Difference between Traditional & big data.
- 5 Explain Various tools used in big data
- 6 Explain component of Hadoop ecosystem.

module 2

- 7 What is hadoop ; component of hadoop ; features of hadoop.
- 8 Difference between hadoop 1.x & 2.x.
- 9 List down hadoop Limitations.
- 10 What is the use of Job & Task Tracker [map reduce & HDFS] — IMP'S
- 11 Explain concept of map reduce using an example (using word count example).
- 12 what is map reduce? Explain in detail different Phases with word count e.g. map, reduce, shuffle, sort, partition.
- 13 Short note on HDFS
- 14 List down Relational algebra Operation in Map Reduce & Explain
- * 15 Explain matrix vector multiplication by map Reduce.

module 3

- 16 What is no sql? List down Advantages.
- 17 Difference between SQL & NOSQL.
- 18 Explain NOSQL Database types.
- 19 Explain Following based Architecture :- Share nothing, Share bus.
- 20 Explain following Distribution Model :- Master Slave, Peer-to-Peer
- 21 Short note :- Cassandra, MongoDB