

## \* What is Big data -

Data which is so large, fast or complex that it's difficult to process using traditional methods is Big data.  
Generally - TB's to 10's of PB's.

Data which is large, fast or complex that is difficult to process using traditional methods is Big data.

Data which is large, fast or complex that is difficult to process using traditional methods is Big data.

Data which is large, fast and complex that is difficult to process traditionally.. using traditional method is Big data

Data which is large, fast and complex which is difficult to process using traditional methods is a big data.

## \* Veracity :-

Data and algorithm should be same.  
Then only we will get the accurate result.

Now it will work in Amazon.

① User click on particular photo.

② Then amazon will take that data in that is photo and run the algorithm with that only and with this help it will recommend the user similar type of the product with the help of algorithm and data.

So Amazon makes more profit by suggesting this product.

\* There are 310 million active users on Amazon.

• 100% of shopping cart items are bought by 310 million users.

• 31 million → 31 billion users.

• Average price range of ~~1000~~ is 1000 to 10000.

\* Value: -

• Any user can search about the product.

• Is the raw data.

i) Amazon takes that raw data and suggest the product to the user and user buy that product and amazon makes money with this.

• Amazon has 310 million users.

• Shopping of: Shopping cart items are bought by 310 million users.

• Web page is a bottom bottom page.

• Amazon is a bottom bottom page.

## Hadoop.

Hadoop is open source framework of tools designed for storage and processing of large scale data (big data). *free to use to all.*

Hadoop is open source framework of tool  
designed for storage and processing of large-  
scale data (big data)

Hadoop is open source framework of tools designed for storage and processing of large scale data (big data).

Hadoop is open source framework of tools which is designed to storage and processing the big amount of data (that is big data).

### I.M.P. points

i) It is created by Doug cutting & Mike Cafarella in the year ~~2006~~ April 4<sup>2005</sup>, 2005.

ii) In 2006 only yahoo! funded this Hadoop.

b) Why its symbol is elephant?

i) because the main author Doug cutting his favorite toy was the elephant. Hence he given the symbol of elephant.

ii) and in 2006 Apache software foundation taken the Hadoop from yahoo!

\* Who are the major users of Hadoop

i) Netflix

ii) Amazon

iii) Microsoft

iv) IBM

\* Two major components of Hadoop

Hadoop

HDFS

↓

Map Reduce

HDFS :- Hadoop Distributed File System

\* HDFS :-  
 i) Designed to store and manage huge data in efficient manner.

↳ Distributed storage.

Designed to store and manage huge data in efficient manner.

i) HDFS is the Hadoop distributed file system.  
 Designed to store and manage the large data in efficient manner.

ii) For the large data only we uses the HDFS system & it's not suitable for small data.

centralized storage → storage on the one place  
distributed storage → storage on multiple nodes  
→ this can original, duplicate, or replicated

PAGE 11  
DATE 11/11/19

## \* MapReduce

i) MapReduce is the massive parallel processing technique used to process data.

MapReduce is the massive parallel processing technique which is used for processing large amount of data parallelly, so we get the output quickly.

Input :- It takes input in the form of list  
Output :- It gives the output in the form of list.

ii) Most of Java language is used.

## \* Features of Hadoop

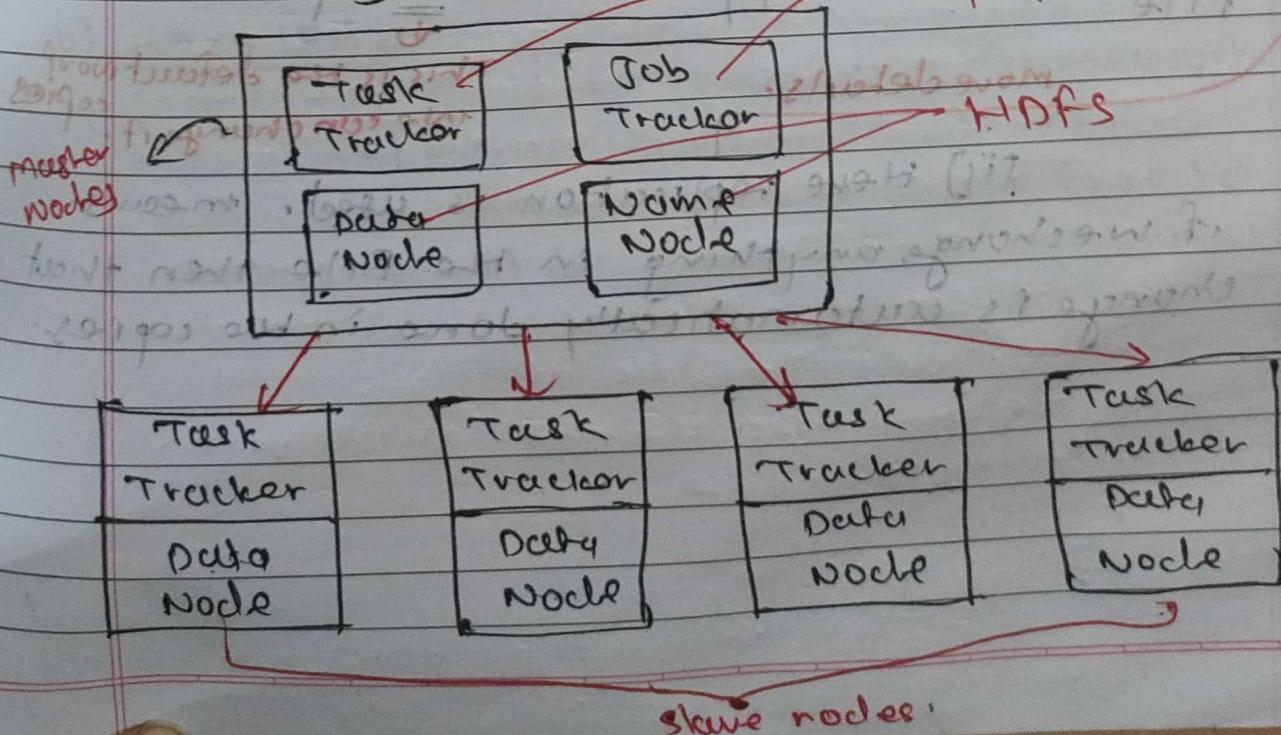
i) Fault Tolerance.

ii) Highly scalable.

iii) Easy programming.

iv) Huge and flexible storage.

v) MapReduce



## Replication

If we change anything in file then it will be updated automatically on each node where it is present.

## Duplication

If we change any file anything in file then either it will not get change or update on each node where the same file is present on node.

### Fault Tolerance

i) All the nodes in Hadoop are the hardware if any node get damage because of any reason.

ii) If suppose one node get damage then the information inside that data will get loss hence Hadoop what do, hence Hadoop checks 3 copies of each file.

more details  
27/11

This is the default No. of copies.  
we can change it.

iii) Here replication is used, means if we change anything in the file then that change is automatically done in two copies

2/20T	2/20T	2/20T	2/20T
rabbit	rabbit	rabbit	rabbit
elephant	elephant	elephant	elephant
about	about	about	about

## Advantages of Hadoop

### i) Highly scalable:

270H i] In the Hadoop all the slave nodes are highly scalable means we can scale that slave nodes in the thousands of number.

### ii) Easy programming:

i] Here as a programmer we have to know the Java programming language, and programmer should know the logic means able to create the logic implementation will become easy.

### iii) Huge and flexible storage:

i] the data we are storing that can be structured or un-structured

ii] Before storing the data there is no need need to preprocess the data.

### iv) Low cost:

i] It is an open source framework

ii] Hence it is free to use all the

iii] because because of the free to use everyone it is low cost.

review test whole out approach

- labor working at

review test whole out approach

- labor working at

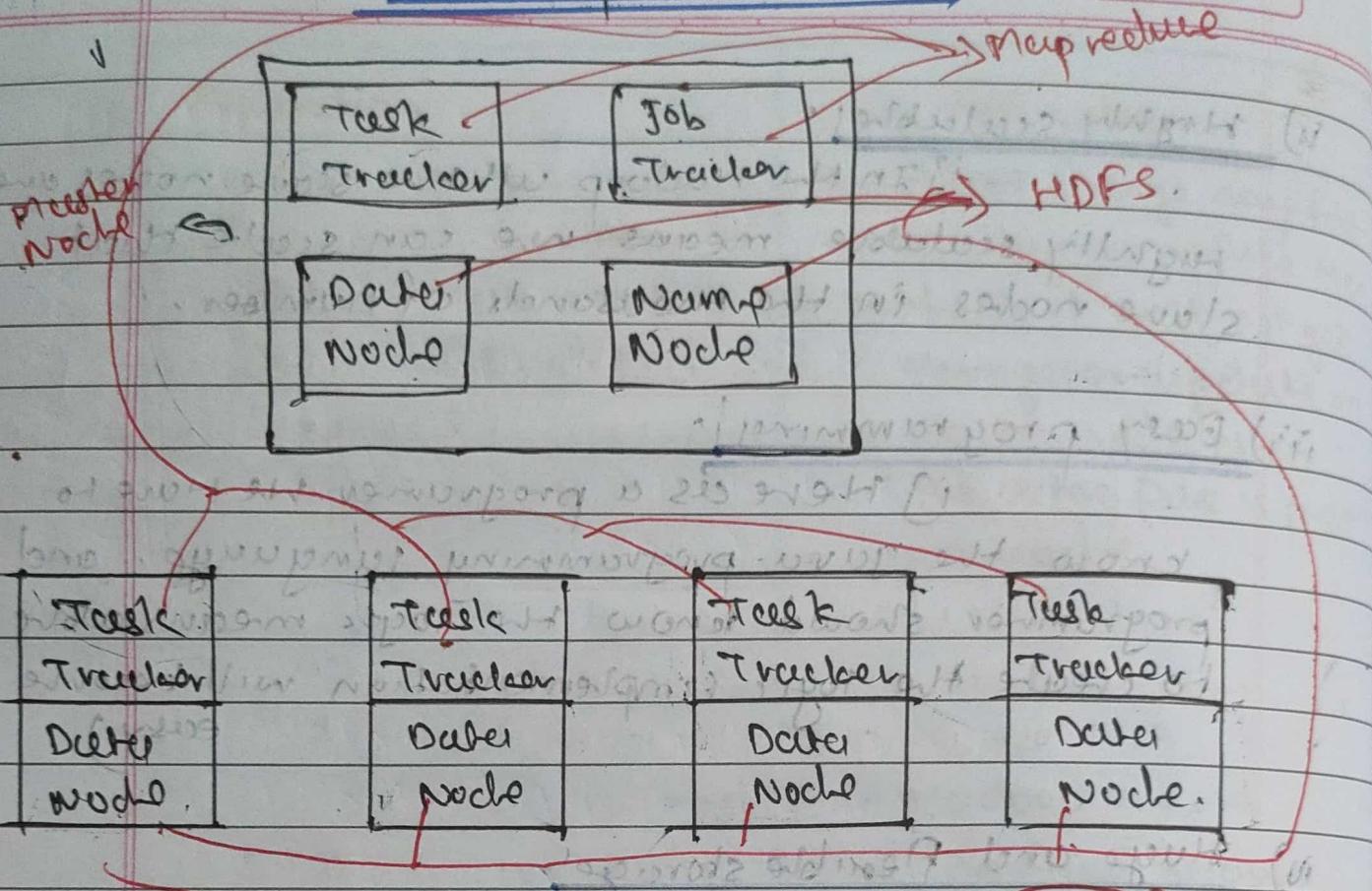
Disadvantages

270H following are some of the

- this are reasons of disadvantages

## Hadoop Architecture

PAGE: / /  
DATE: / /



### ① Task Tracker: working of basic term

It processes the small ~~process~~ piece of data given to that particular node.

It processes the small piece of data given to that particular node.

### ② Data nodes:

It manages the data that received on the particular node.

It manages the data that received on the particular node.

### ③ Map Reduce:

It is massive parallel processing technique to process the data.

### \* Map Reduce:-

It is massive parallel processing technique to process the data.

\* suppose Task is given find the one youtube channel who have more number of video.

i) This takes biggest task first taken by Job tracker. It divides into smaller task and passes to the all slave node.

ii) then that slave nodes process process that data parallel and give the output.

### \* HDFS:- (Hadoop distributed file system).

Designed to store and manage the large amount of data in efficient way.

### \* Job Tracker:

i) Breaks the bigger task into smaller task and forward it to the Task Tracker.

### \* Name Node:

It keeps information about the which part of the data is having with which node.

## Hadoop ecosystem &

Coordination (ZooKeeper)	Scripting (pig)	Machine Learning (mahout)	Query (Hive)	MapReduce (map reduce)	Distributed processing (map reduce)	Hadoop Distributed file system (HDFS)	MapReduce (Scoop)
ZooKeeper	JVM	Machine Learning	Hive	MapReduce	Distributed Processing	HDFS	Scoop

① HDFS ② MapReduce ③ Flume ④ Hive ⑤ HBase ⑥ Mahout ⑦ Pig ⑧ Scoop

② MapReduce ③ Mahout

③ Flume ④ Hive ⑤ HBase ⑥ Mahout

④ Hive ⑤ HBase ⑥ Mahout ⑦ Pig ⑧ Scoop

### ① HDFS :- (Hadoop distributed file system).

It have majorly two components

i) Data Node, ii) Name Node.

"HDFS is designed to store the large amount of data in efficient way"

#### i) Data Node,

It manage all the data that received on particular node.  
iii)

(i) Name Node: (each file have 3 copies it track that also).

i) It keeps all the information about the which part of the data is given to the which slave node.  
ii) It maintains some record or information.

② Map Reduce: (Input - List, Output - List)

MapReduce is a massive parallel processing technique used to process data parallel.  
MapReduce is a massive parallel processing technique used to process data parallel.  
MapReduce is a massive parallel processing technique used to process data parallel.

It have major two component

i) Job Tracker

ii) Task Tracker

i) Job Tracker : ii) It also set the priority of time to nodes  
i) Job tracker breaks the bigger task

into smaller task and forward it to the

task tracker got no fluid window

i) Task Tracker:

i) Job Task Tracker process the the small piece of task to generate the output

### ③ Flume

i) It is the data injection tool in HDFS.

means we can inject the data with the help of flume, that data can be.

Unstructured, structured, or semi-structured data.

ii) It is also more reliable and fault

~~Tolerant~~

### ④ Apache Hive

i) It is an open-source data warehouse system for querying data.

ii) Analyzing data.

iii) summarizing data.

This is large dataset stored in H-files.

iv) It uses HQL = Hive + SQL

v) It is SQL like query language

vi) Feature 1 - Highly scalable

### ⑤ Hbase

i) It is a distributed column oriented database which is built on top of Hadoop file system.

ii) It is designed to provide the quick random access to huge amount of data.

### ③ Flume

i) It is the data injection tool in HDFS.

means we can inject the data with the help of flume , that data can be.

Unstructured, structured, or semi-structured data.

ii] It is also more reliable and fault

~~(fault-tolerant)~~

### ④ Apache Hive

i) It is an open-source data warehouse system for i) querying data

ii] Analyzing data.

iii] summarizing data.

This is large dataset stored in H-files.

iv) It uses HQL = Hive + SQL

v) It is SQL like query language

vi) Feature 1 - Highly scalable.

### ⑤ HBase

i) It is distributed column oriented database which is built on top of Hadoop file system.

ii) It is design to provide the quick random access to huge amount of data.

## ⑥ Mahout :-

i) With the help of Mahout it becomes easy to implement the machine learning algorithms in the Hadoop ecosystem.

ii) It has three techniques:-

i) Recommendation

once classify as spam it go in spam folder  
ii) Classification

iii) Clustering. → One type of things come together and form their group is known as clustering.

• previous

• next

## ⑦ Pig :-

i) This is also data processing tool.

ii) It has its own language that is Pig Latin. It is scripting language.

iii) It also has its own component that is Apache Pig Engine. → It accepts the Pig Latin scripts.

N) 1 line in pig Latin = 100 line in Java

## ⑧ Sqoop :-

i) It imports structured data from RDBMS to HDFS in Hadoop

ii) It exports structured data from HDFS to RDBMS in Hadoop

## ⑨ Zookeeper :-

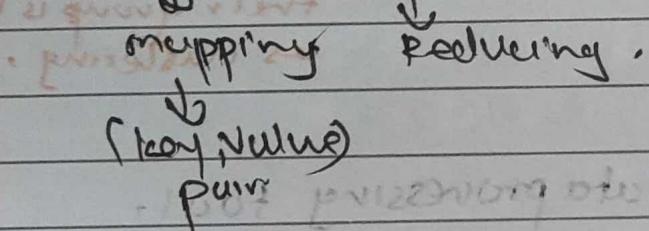
i) With the help of Zookeeper each component in the Hadoop ecosystem they can communicate with each other. ex. between the problem realme

i) And perform or complete the task in less time, together with others.

ii) It also takes care that no one component can get interruption, in between the task processing.

iii) It is vulnerable to failures.

## \* MAP-Reduce



Shuffling → similar record group

Reducing → Add the final result

Map Reduce = matching of entities

2018 most web browser request 15%

google at 17% of

2704 most web browser requests 15%

google at 17% of

most requests to 194.169.254.11

most part of 2018 google was the second

2nd most website along with alexa

## \* Relational database.

student ID	Name	Location	Gender	college
18	Tallu	Nerul	M	Don-Bosco
27	Gullu	LST	M	K.J.Somayaji
36	PUPPU	Thane	M	R.A.I.T
61	Pipa ke Puri	Seawoods	M	Don-Bosco

## \* NO-SQL Database.

Types of database in NO-SQL - i] key value store

- ii] wide column store iii] Document DB
- iv] Graph DB .

```
{ "student_id" : "gg",
  "name" : "Bebhu",
  "Hobby" : "Drinking",
  "Branch" : "Computer Engineering"}
```

y

```
{ student_id : "6g",
  "Name" : "PUPPU",
  "Location" : "Dadar",
  "Hobby" : "Killing kids For Fun"}
```

y

# RDBMS VS NO-SQL

PAGE: 11  
DATE: 11/11/2023

RDBMS	NO-SQL
i) Have fixed or static predefined schema.	ii) Have dynamic schema.
ii) Not suited for hierarchical data storage.	ii) Best suited for hierarchical data storage.
iii) Vertically scalable.	iii) Horizontally scalable.
iv) Follow Acid property.	iv) Follow CAP (consistency, availability, partition tolerance) (cap theorem).
v) RDBMS supports transactions. (also complex transaction with joins).	v) NoSQL databases don't support transactions (support only simple transactions).
vi) RDBMS manage only structured data.	vi) NoSQL databases can manage structured, unstructured, semi-structured data.
vii) Relational databases have a single point of failure with failover.	vii) NoSQL database have no single point failure.
viii) Support powerful query language.	viii) Support a very simple query language.
ix) Support complex realme transaction.	ix) Support simple transaction.

102-04

2M809

Example - MySQL, Oracle,  
MS-SQL etc.

example MongoDB, BigTable,  
Redis, RavenDB etc.

IP file.

Matrix	j	k	l	m
A	0	0	0	0
A	0	1	2	1
A	0	0	3	2
A	1	1	4	3
B	0	0	0	0
B	0	1	0	0
B	0	0	0	0
B	1	0	0	0

Now we have to make the map function,  
so we will now convert above table.

## Matrix multiplication with 2 MapReduce step

DATE: / /

1 step

2 step

i) Map

ii) Reduce.

\* There are two types of matrix multiplication.

i) One step  $\rightarrow$  1 time map, 1 time reduceii) Two step  $\rightarrow$  2 time map, 2 time reduce

## i) One step Matrix multiplication

Assume there is A matrix of size  $2 \times 2$ Also assume there is B matrix of size  $2 \times 2$ .

1st condition:

We can do matrix multiplication when only,

the columns of first matrix is equal to the rows

of second matrix.

\* Assume the below matrices are suitable for us

o  $j \rightarrow$ k  $\rightarrow$ 

$$\text{matrix } A \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix} \quad 2 \times 2$$

$$\text{matrix } B \begin{bmatrix} 5 & 6 \\ 7 & 8 \end{bmatrix} \quad 2 \times 2$$

(B)  $\rightarrow$  No. of columns of matrix B  
 This should have to be same  
 Then only matrix multiplication  
 is possible.

i  $\rightarrow$  No. of rows in matrix A.j  $\rightarrow$  No. of columns of matrix Bj  $\rightarrow$  No. of rows of matrix Bk  $\rightarrow$  No. of columns of matrix B.

$$\begin{array}{c} \text{No. of rows} \\ \downarrow \\ 2 \times 2 \end{array} \quad \begin{array}{c} \text{No. of columns} \\ \downarrow \\ 2 \times 2 \end{array}$$

$$2 \times 3 \quad 3 \times 2 \quad 2 \times 2$$

102-001 2V 201809

PAGE: / /  
DATE: / /Most Imp question.  
Matrix Multiplication with 4 Map Reduce step  
DATE: / /Example - MySQL, Oracle,  
MS-SQL etc.example MongoDB, BigTable,  
Redis, RavenDB etc.

Now? To make the input file. The will make the four columns i.e. ip file.

Matrix A	i	j	k	Value	MapReduce File
A[0][0]	0	0	0	10	0
A[0][1]	0	1	0	10	1
A[0][2]	0	2	0	10	2
A[1][0]	1	0	0	10	3
A[1][1]	1	1	0	10	4
B	2	0	0	10	5
B[0][0]	0	0	0	10	6
B[0][1]	0	1	0	10	7
B[0][2]	0	2	0	10	8

Now we have to make the map function,

so we will now convert above table.

i) Map

ii) Reduce.

There are two type of matrix multiplication.

- 1) One step.  $\rightarrow$  1 time map, 1 time Reduce
- 2) Two step.  $\rightarrow$  2 time map, 2 time Reduce.

i) One step Matrix multiplication:

Assume there is  $A$  matrix of size  $2 \times 2$   
 & also assume there is  $B$  matrix of size  $2 \times 2$ .

1st condition:

We can do matrix multiplication when only  
 the columns of first matrix is equal to the rows  
 of second matrix; i.e. first and second

Assume the below matrices are equal to each other

$$\begin{matrix} A & \xrightarrow{i} \\ \begin{matrix} 1 & 2 \\ 3 & 4 \end{matrix} & \xrightarrow{j} \\ \begin{matrix} 5 & 6 \\ 7 & 8 \end{matrix} & \xrightarrow{k} \\ B & \xrightarrow{o} \end{matrix}$$

This should have to be same  
 then only matrix multiplication  
 is possible.

 $i \rightarrow$  No. of Rows in matrix A. $j \rightarrow$  No. of columns of matrix B $j \rightarrow$  No. of Rows of matrix B $k \rightarrow$  No. of columns of matrix B.

$$\begin{matrix} 2 \times 2 & & \\ \downarrow & & \\ 2 \times 2 & & \end{matrix} =$$

$$\begin{matrix} 2 \times 2 & & \\ \downarrow & & \\ 3 \times 2 & & \end{matrix} =$$

realme

Shot on GT Master

## Map Function

i, k	value	Note:- To fill the value format is $(Matrix, j, value)$
0	(A, 0, 1)	
0	(A, 1, 2)	
1	(A, 0, 3)	
1	(A, 1, 4)	
0	(B, 0, 5)	
0	(B, 1, 6)	
1	(B, 0, 7)	
1	(B, 1, 8)	

Now while making the map function we forgot that the  $j$  for matrix A is no. of columns and for matrix B it is no. of rows.

So here we have to transpose the map function for matrix B i.e. that is when  $i = j$  do not do anything else transpose the matrix where column will become row.

Hence above map function is wrong.

Correct one is 10 numbers to an  $\infty$ .

1st row point at 2nd row 1st col  $\leftarrow i$

2nd row point at 2nd row 2nd col  $\leftarrow j$

map function (one one, after  $\text{map}(\text{value}) + (\text{key})$ )

i/k	value	$(0,0)$
0	$(A, 0, 1)$	
0	$(A, 1, 2)$	$(1,0)$
1	$(A, 0, 3)$	
1	$(A, 1, 4)$	
0	$(B, 0, 5)$	$(0,1)$
0	$(B, 0, 7)$	
1	$(B, 1, 6)$	
1	$(B, 1, 8)$	$(1,1)$

✓ correct  
MAP function.

This above one is correct map function.

Now before Reduce function we have to see the shuffle [group]: "

\* shuffle [group]: ~~number of tasks~~  $(x \times 3) + (0, 82)$

(i, k)

$(0,0)$   $(A, 0, 1)$ ;  $(A, 1, 2)$   
 $(B, 0, 5)$ ;  $(B, 0, 7)$

$(0,1)$   $(A, 0, 1)$ ;  $(A, 1, 2)$   
 $(B, 1, 6)$ ;  $(B, 1, 8)$

$(1,0)$   $(A, 0, 3)$   $(A, 1, 4)$   
 $(B, 0, 5)$   $(B, 0, 7)$

$(1,1)$   $(A, 0, 3)$   $(A, 1, 4)$   
 $(B, 1, 6)$   $(B, 1, 8)$

Now Reduce 1.: without going forward

$$(0,0) \quad ((1 \times 5) + (7 \times 2)) = 19 \quad 11$$

$$(0,1) \quad ((1 \times 6) + (2 \times 8)) = 22 \quad 0$$

$$(1,0) \quad ((3 \times 5) + (4 \times 7)) = 43 \quad 0$$

$$(1,1) \quad ((3 \times 6) + (4 \times 8)) = 50 \quad 1$$

Note! Vertical "multiplication"  
Horizontal "Addition".

Hence Resultant matrix =  $\begin{bmatrix} 19 & 22 \\ 43 & 50 \end{bmatrix}$

In this way we do matrix multiplication using map Reduce.

$$(5, 4, 1, 4) : (1, 0, A) = (1, 0)$$

$$(8, 1, 8) : (2, 1, 2) =$$

$$(11, 1, A) : (2, 0, A) = (0, 1)$$

$$(F, 10, 8) : (2, 0, 8) =$$

$$(H, 1, A) : (2, 0, A) = (1, 1)$$

$$(8, 1, 8) : (2, 1, 8) =$$