

Q. Interview Question:

Let me know the configuration of Hadoop cluster:

Configuration of one m/c :

40TB (HDD)

512GB (RAM)

256 cores (CPU)

What will be the total capacity / configuration of 1200 m/c.

It will be

40TB (HDD) }
512GB (RAM) } x 1200.
256 Cores (CPU) }

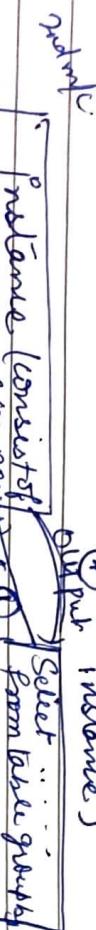
2nd problem

Problem: → Complete data is travelling from Storage Database mfc to instance mfc over network

Database Server → consists of 2 physical mfc (storage database)

lets look with example

(1) output instance)



we are trying a query from this mfc.
we are trying a query from this mfc.

Storage

1TB
100ms

1024x1024 → 10485.76 sec.

100 → 17476.2 MB

= 99 hours

= 4 days

How long it will take to transfer 1TB data with 100ms

(2) look for particular table
Storage dbf. → any table is maintained by this storage database.

(taking care of storage)

High level of overview of Database

instance works on your behalf which ever table you are looking for. This instance is going to look for the particular table in storage database. Once the table is scanned ~~the~~ instance is going to RDB of instance where process is going to happen. final off goes to the mfc asking that query.

Q) what could be the possible problems:

$$1024 \times 1024 \times 1024 = 100$$

$$100 \text{ days}$$

$$124 \text{ days}$$

[it is going to take 124 days just to transfer of data]

DBA coming to you asking for please remove archive table, temp data.

~~problem~~ Data getting transferred over the network.

Solution:

That takes much time which is not acceptable

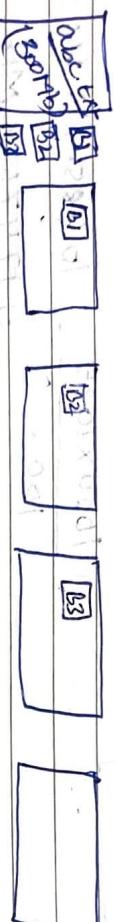
so the solution provided by Hadoop is :-

How hadoop is solving a problem :- which uses

a horizontal scaling

Before talking about that we have to see many things

[Data Blocked]
Size: 128 MB



File is going to be divided

into small - 2 blocks
Size of every block = 128 MB.
So ABC total will be divided into 3 blocks

B1 → 128 MB.

B2 → 128 MB.

B3 → 44 MB

Only data comes to hadoop it remains there

flower because Hadoop is horizontal scaling system. You don't remove the data whenever you are shortage of data we just add more

System:-

CDR → Data is processed for different purpose
→ increasing revenue of [margin / profit] data

when you process the data we write SQL query.

Result is higher line of query than any one can

written :- (SQL Query) (higher line of query) → 150 lines of

(huge query) → It is stored in line of query into any file

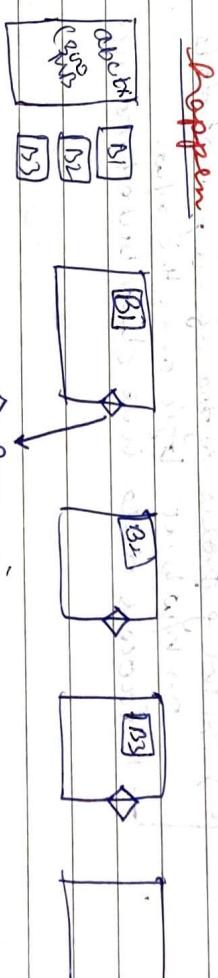
file

What would be file size? of any query → it would be. → It would be few KB of data

Important Point: Processing ◇

Hadoop says In hadoop, here it writes the data only once whenever point of done in case

of processing. → Instead of transferring the data it transfers the processing logic wherever is the very small size. where block processing is going to happen.



↓ Processing
This is how it solve the problem

There is important fact that hadoop achieves:

B₁ is processed at 1 machine
B₂ is " " 2 machines

B₃ " " " 3 machines

at the same time

So we are achieving at the same time distributed processing parallel processing of data.

so it not only achieve the solution of problem

At 2nd problem But it also achieves very important thing that is distributed processing

At the end we say (Definition)

Hadoop → software free of cost (not only academic community driven but developer point of view)

→ solving the problem of existing system (vertical scaling by horizontal")

→ infinitely Scalable

→ solving problem of transferring not data on ip but logic over the ips.

process data into distributed manner

Q. If you go to the office with new solution (Hadoop)

what do you think the your company will accept the new sof. or they will ask some more question

What will be their concern? (costing)
(horizontal Scaling)

Are these machines (in horizontal Scaling) free of cost means their concerned regarding hardware cost In case of Hadoop commodity hardware is sufficient for wider hardware installation.

It is not necessary to have highend machines costing of commodity is almost nominal for company point of view.

So Definition :
it is not only cheap but also

→ parallel processing
→ infinitely Scalable.

Q. What will happen if commodity box fails?

Because we are storing the data on the commodity box if one data box fails then what happens if commodity box fails goes down.

example

abc	def	ghi
ijk	lmn	opq

abc	def	ghi
ijk	lmn	opq

abc	def	ghi
ijk	lmn	opq

Hadoop

Hadoop keeps multiple copies of block on different machine. (By default it replicates at 3 times) This process is called **Data Replication**.

We can increase/decrease this replication factor as per my requirement.

So if machine (two) fails it will be lost of B1, B2, B3 block. Only our still 2 replications are already present.

Incase Google - Google never says i have lost your Email. Reason is Google is maintaining a replication factor 5. Even in the worst case scenario . Data doesn't get lost.

Although data gets replicated 3 times but because of commodity we are using the focus on reliability. It is economical also as we are using commodity h/w.

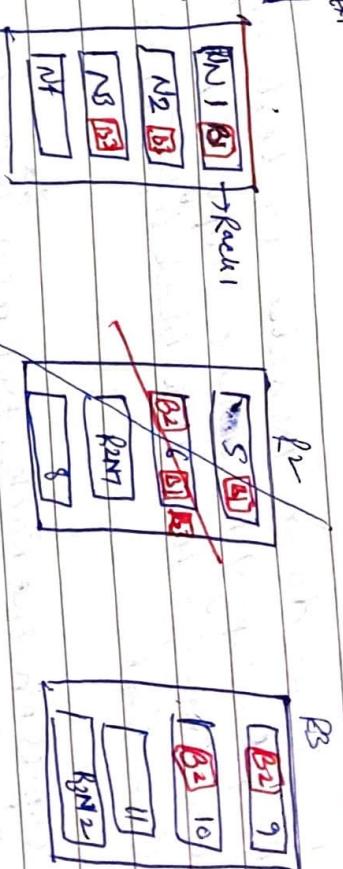
RACK AWARENESS AND REPLICATION



Rack Awareness

Q what happens when you store any file?

This is 3 Rack, 12 Node hadoop cluster. We are making first copy of B1 available at R1.



Rack Awareness say 2nd will go to other rack. 3rd will go to same but on different rfc.

3 copies on the same Rack - 1 another Rack
3 machines will be consumed.

lets talk about [B2]. First copy is available on R2 N6 and 2nd one R3 N9, 3rd copy on R2 N10

Note is B3 block, first copy will place at R2 R2 N6, 2nd will go to R1 N2 3rd copy is going to R1 N3.

losing one rfc [R2 N6] is not a problem because we can get the info/data from other rack or rfc.

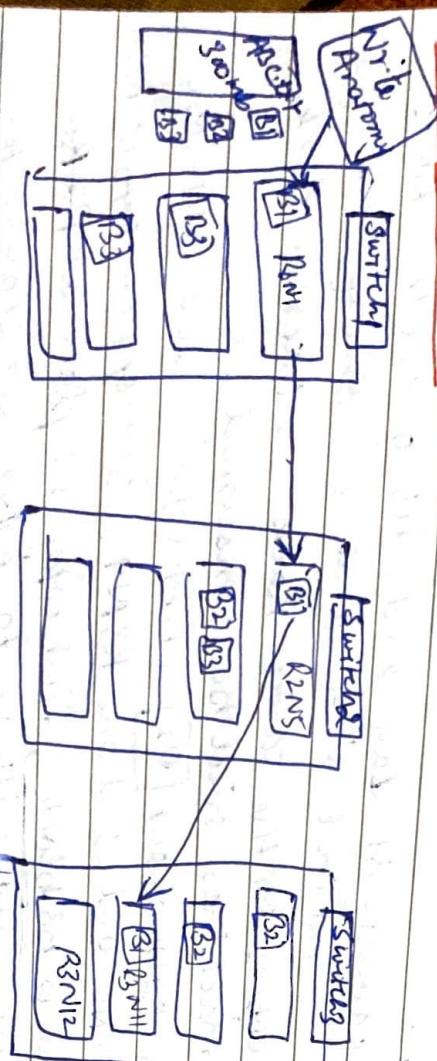
Q losing Rack is also no problem but chances are minimal

If Rack 2 goes down then also we lost 2 copies of B1, One copy of B2, and one copy of B3. Since we can get the data block info [B1, B2, B3] from other Racks.

Two copies on one Rack. Not more than one Rack.

Q. Why it is not happening and 2nd on Rack 2 & 3rd copy on Rack 3? Is this not a reliable solution?

Solution / Answer:



Whenever we want to write data there is a process called write anatomy.

Write Anatomy will write into 1st m/c then it will write into 2nd m/c then it will write into 3rd m/c.

If two m/c on the same Rack talk to each other this is called Intra Rack m/c

If two m/c on the diff. Rack wants to talk then m/c of first Rack will talk to the switch of 2nd Rack. Then switch of 1st will talk to switch of 2nd Rack then that switch of 2nd Rack connects to intra particular m/c communication:

RINI → RIN2 (IntraRack)
RINI - R2NS (InterRack)

Scenario

Now for write Anatomy process If B1 at RINI wants to communicate to R2NS (After replication) then R2NS wants to replicate at R3N11. There are Inter Rack communication whereas Scenario - 2

If two Copies is stored at one Rack this time there is only one Inter Rack communication whereas in Scenario - 1. 2 inter Rack Communications were there which was costly.

Q.

Sharmi is writing a file ABC.txt into HDFS now Atuljeet wants to read this file today which was written by Sharmi yesterday. How Atuljeet will be able to read?

Is this a data which
Project wants to read

Q:
This info

abc.txt

data: B1, B2, B3

Information is data

Arrow data. It is

Called **METADATA**

abc.txt
B1, B2, B3
B1: 1, 2, 3
B2: 2, 3, 4
B3: 1, 2, 4

(This info has to be available
Somewhere.)

Thruonly you will be able to

read B1, B2, B3

Imp Point: This metadata is a critical info
And where it will be kept. If it is not
there you will not be able to read the
data, reconstruct data etc so this metadata
is highly critical.

If you deposit money

: Accno, Bankname, Branch address] If you loose
this info then is equivalent to ~~you lost~~

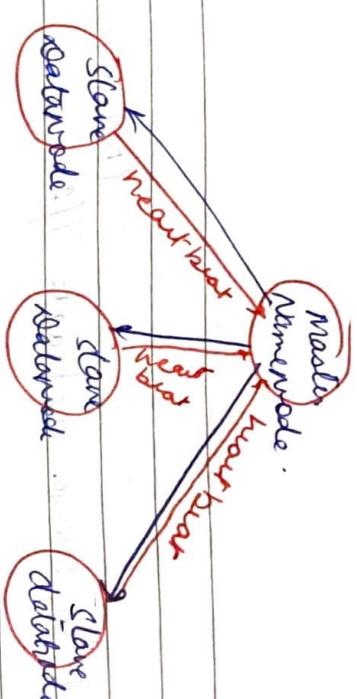
Money. This info is highly critical.

So in hadoop what happens is this
Highly critical if metadata is stored in
RAM of that m/c That m/c is referred
as NameNode

Hadoop is something which works as
Master Slave architecture.

all the m/c coming under this master. assume

slave m/c coming under master.



All the Slave m/c are separate and they
continuously sending signals [acknowledgment]
to the master m/c after 3 sec.

⇒ Hey master, I am alive and you can give me command

This signal is called **heartbeat**

This heartbeat is sent after 3 sec. After every
10 heartbeat means interval of 30 sec. What

is going to happen is :

These slave m/c send block report to the master
This particular m/c will send it
a block report that block
is having B1, B2, B3

Block Report

from slave to master.

Whenever you are keeping any data or metadata is generated. That metadata is kept in the RAM of master m/c (NameNode). CAN I afford to lose metadata iff.

METADATA is highly critical iff.

Q Can I afford to lose this metadata iff. On

Commodity m/c
It should not happen - It should be high end m/c.

Commodity m/c is not applicable for master m/c. Master m/c keep critical iff.
NameNode - BMW, Lamborghini
DataNode - TATA NANO!

Master m/c has to be high end m/c and slave m/c should be commodity m/c.

Q: Does my master m/c never breakdown
Probability of failure is less. But still

You cannot append, update the file.
In hadoop you can do these 2 activity

- 1) Add a new file
- 2) delete an existing file.

Can you afford to lose this metadata iff.
There is a big requirement for the backup of NameNode :->

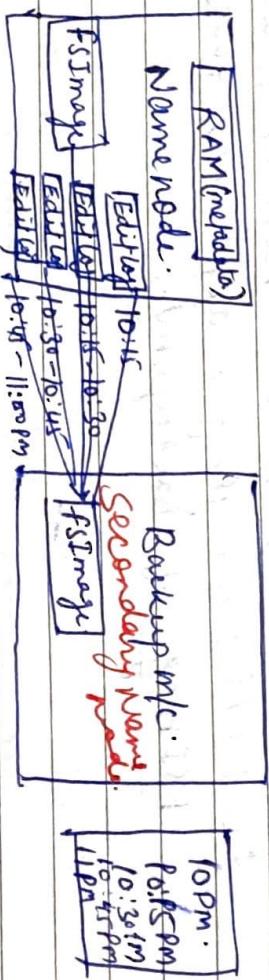
m/c.

Now at 10:15 PM, all the changes which

Backup of NameNode in Hadoop Generation 1

first version of hadoop that was released called HDFS 0.1

- 1. Add a new file
- 2. delete an existing file



NameNode having huge amount of RAM for keeping metadata iff

happen to cluster are captured on Editlog. means either you added or deleted a file those changes were captured and made available to Editlog.

finally at 11:00 pm. all changes were "made available to one more Editlog". At 11 all the Editlog are going to applied on these FSTimage or Backup mfc.

Q Q Q all the Editlog changes are applied on FSTimage. what will be op?

If we apply all Editlog on top of FSTimage then this FSTimage will be in proper sync with the metadata present in the RAM of NameNode.

At 11 both metadata & FSTimage both will be same. This is how backups maintained.

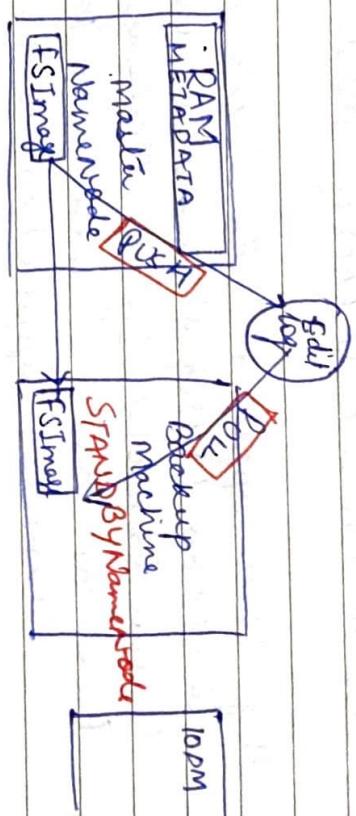
This Backup mfc is called Secondary NameNode.

Q What happens if what happens if at 10PM

Q What happens if name node fails then what will happen.

As all the Editlog are on NameNode. If it will fail then it will lose some info (lose of some data) is there in Hadoop 1. This was problem in Hadoop generation 1.

The hadoop 2 solution to this problem is given:
Solution: By Hadoop Generation 2
Backup in Hadoop Generation 2



I started doing backup at 10 PM & complete metadata if hadoop is available as FSTimage in hadoop of master mfc.

This FSTimage is made available to Backup mfc.

After that any of the change which is going to

happen to cluster those changes are being captured in terms of Editlog.

These Editlog are placed at shared location

classmate

Editing made available at Shared location ^{and changes}

are going to made available in Edit log or

Shared location into the real time basis and

again on the Real time basis those changes are

Captured and directly being applied to FSImage

to Backup machine. It is called Pull mechanism.

It works as PUSH & PULL mechanism.

PUSH

At any time FSImage is sync with

Metadata present at RAM of Master machine

Same way we have Super Manager available

Any change happening to the cluster that change

is captured on Editing at Shared Location On

the real time basis by Push Mechanism and

Once again with pull mechanism. Even real time changes are applied on the FSImage

Backup m/c

This procedure is called the Backup Mechanism

of Hadoop 2.0.

Backup m/c is called Standby Namenode

At any given point of time there are

2 nodes one is active Namenode &

Standby Namenode Both of them having

are having the real time Metadata

Information available at any point of time

There are 2 masters : (At any given point of time there are 2 masters)

Q. What will happen if there is conflict between 2 managers.

A. If there is conflict between 2, whom to follow, not to follow.

If this is the case how this problem is sorted out : There is

always a Super manager - Super. Manager will resolve the conflict.

In Hadoop 2.0

ZooKeeper

→ Zookeeper ensure at any

given point of time all the Slave

Machine are reporting to One master i.e. active Namenode.

If there is any problem comes to Master

keeper immediately will communicate to

All Slave m/c that to report to new

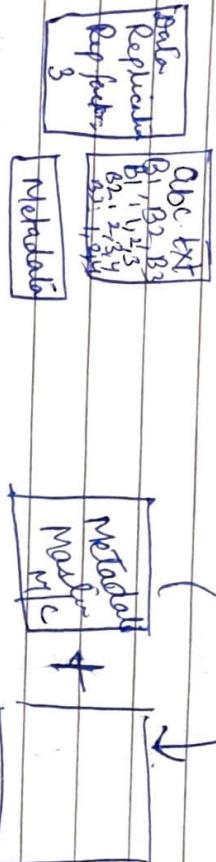
master (Slave m/c)

The whole process does not require any downtime. This is called High availability

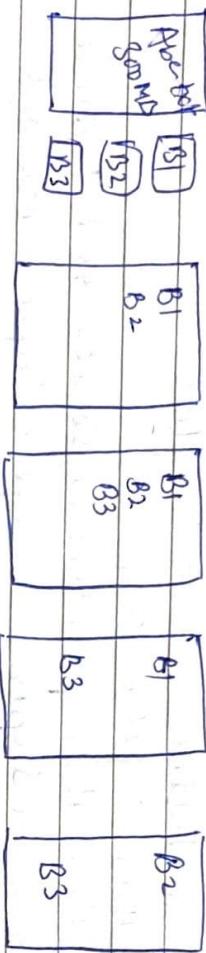
in Hadoop 2.0.

High Availability : So we achieved HA

FEDERATION



And add 2nd master machine. Once again there is no limit for master also machine also. This concept is called as Federation. This only comes with Hadoop 2.0.



4 Node Hadoop cluster

If we are running short of RAM, Harddisk then we add more Slave (commodity m/c)

because of Horizontal Scalability.

Similarly in case of Master m/c, if unresponsive metadata is stored on RAM master m/c (RAM)

Q. What if RAM of master m/c gets filled up.

→ Will you be able to store any file on hadoop? → NO.

If no space available in the RAM of master m/c we will not be able to store any file in hadoop. We will never be able to reconstruct the data.

If RAM of master m/c gets filled up we will ~~witdraw the RAM with the help~~

of horizontal scalability

classmate

classmate

- Q. who is taking the decision of Metadata? (Master)
 Q. No. of Blocks will be decided on which based on the basis of BlockSize and size of file
 Q. what would be the criteria for selection of 3 machines for B1: 5, 2, 4

federation: is to apply horizontal Scalability for hadoop - master machine.

Create One Folder
To create directory

Command:

hdfs dfs - mkdir

There is no command cd available in HDFS.

HDFS: is Virtual file System. You can imagine it. Think of it you cannot live w/o it. You cannot go to virtual location.

WRITE ANATOMY: Transferring some data

from local file system to HDFS. You want to write Data to HDFS.

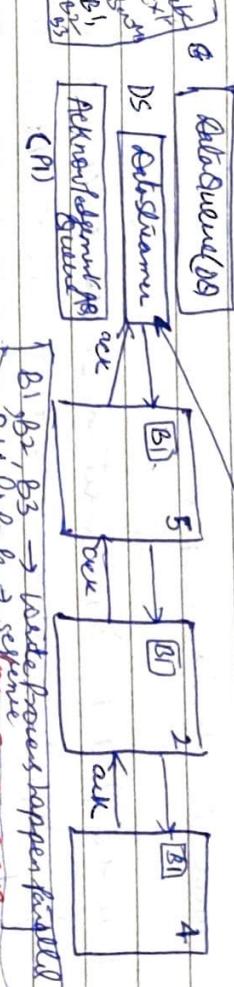
WRITE Anatomy
HDFS dft - Put localFilePath HDFSfilepath

When you fire this command you can login to any machine & from there you fire the command.

first point of contact is master node. This master node is going to prepare Metadata info for the file.

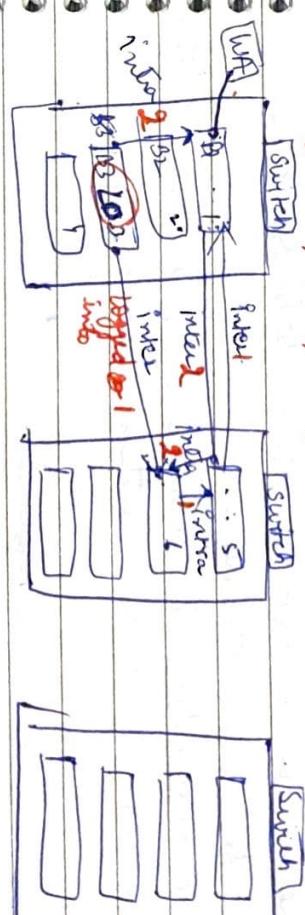
Ans: 3 concept of Rack Awareness will decide whether block information will go to classmate

[HDFS dfs - put localFilePath HDFSfilePath]



Q. what could be the BlockID for sequence

Why it is 5, 2, 4 only why not any other seq. 2, 4, 5 or 2, 5, 4 etc.



by says.

Intersect communication. Inter block is Path 1: 2 inter + 1 intra

Path 2: 1 inter + 2 intra

This sequence is decided on the basis of inter/intra rack communication

classmate

Name node connected with data streamer.

What DS does

Data streamer further divides B1 to

Small -2 Packets of 16 MB = Total of 3 Packets (P₁, P₂, P₃, P₄) -

All the packets become part of Data Queue. It will set first packet to Acknowledgement Queue

then Data Streamer connect to first m/c 5 and it is going to write a file on the local file system of m/c name it as B1 after that m/c 5

Connect to 2 same procedure

Once the packet B1 is written successfully on all m/c info & send ack to 2 from 2 to 5 it will

Send confirmation to Data Streamer Yes Packer(B)

is written successfully on all the m/c's.

Data Streamer remove B1 from Acknowledgement

Queue & put packet 2 from data Queue to

Acknowledgement Queue

Same procedure is done for all the packets (P₁-P₄)

In this way Block B1 written successfully

In Not-happy Path Scenario:

Q. How the detection of failed machine is done?

As the machine which is getting fail will

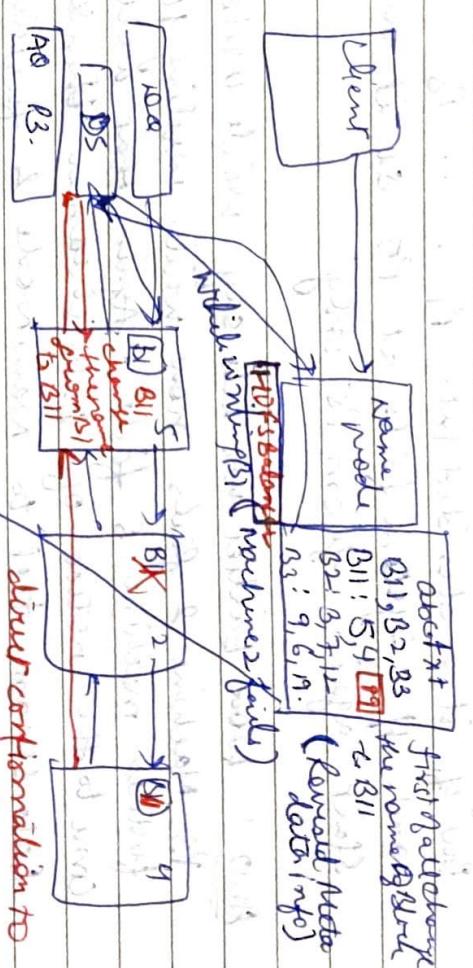
not send acknowledgement packet because

the receiving m/c bcos of not getting

classmate

Acknowledgment understand that something happen to m/c [failure occurs]

This information is passed to Data Streamer. Once Data Streamer get to know it will put packet P3 to Data Queue then it will connect to namenode.



While writing Block B1 node 2 fails so namenode change the Block B1 to B11 along the Block B1 to B11

Question 1 Why Block B1 is divided to P₁, P₂, P₃ - 8

Q. 2. what happen if further m/c fails.

Q. 3. If block B11 is only replicated to 2 m/c now the problem of under replication is solved

Q. 4. Why B1 name changes to B11

Q. 5. What will happens to failed m/c classmate

What will happen

Q.S. w/c no. 2 and about the corrupted Block

Ans1: Packets P1, P2, ..., P8

Lets we have not divided it & we have written 100 MB of data when we start

again we have to start with 0. But if it is divided into P1, ..., P8 and if system fails at P5 so it will start again from P5 only.

Ans2: Earlier lets say while writing P3

Machine 2 fails now while writing P6 further machine 5 fails. Again that will come to the notice of Data Streamer then same Data stream tell to NameNode & further NameNode change the name of the block.

If you are able to write successfully a block even on a single machine. That is called a successful write

It will be underreplicated but still it is called successful write.

Ans: 3 Under replicated

Write Anatomy complete even if the blocks
if is written even on one machine

After failure we are writing one 2 or 1 machine but it is called underreplicated. How to solve it

there is one more soft process running with the name HDFS Balancer running under NameNode

Q. What does HDFS Balancer do?

HDFS Balancer check metadata information on pair of time at regular interval then it will come to know that Block 11 is underreplicated then it will inform to master node. Then master node decide that it should be placed at P9. Machine. So Reid copy will be replicated to P9 w/c hence problem of under replication will be solved.

Ans: 4 why the name of the block B1 is changed to B11

After coming back HDFS Balancer will never come to know that this Block B1 is corrupted Block.

Ans: 5 what happened to machine 2 and corrupts

Block.

As soon we are losing commodity of any file
it can go out we don't have any constraint
it cannot come back its own. So the hadoop
administration is going to indirect the file
again to cluster.

Machine 2 after coming back will send

heartbeat after 3 sec. to hadoop master.

If 10 heartbeat file 2 will say that

it is containing Block B1.

Then HDFS Balancer will check the complete

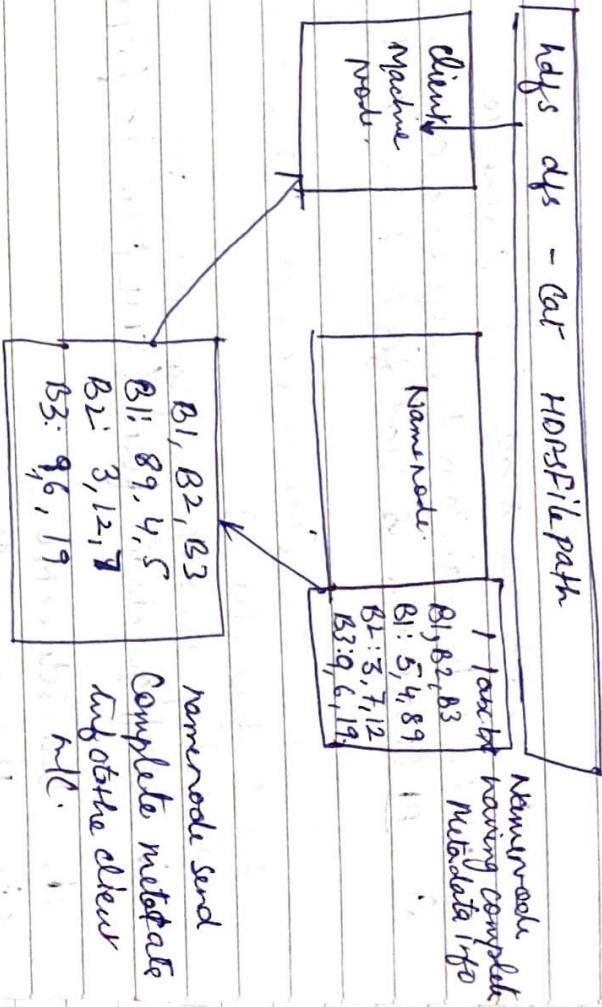
int metadata information & check the entry

for B1. After checking it will come to

know that it is half written block or

corrupted Block. HDFS Balancer delete the

Block B1 from machine 2



Q.4. What is reason for changing of sequence?

B1 is present on node 1, 2, 6.
and you are logged into node 3
from which file you will

read := 1

As each 6 internode is small as compare to cost of inter

race

so cost of internode is small as compare to cost of inter

Point to point communication is less expensive

thus internode communication.

nearest block is selected for reading

B1 is read from 89

B2 : is read from 3

B3 : is read from 9

READ ANATOMY: Read the data from HDFS

Q.1 Who you think the first point of contact of NameNode
Block are available on diff. hlf.

Q.2 If I want to read B1, I should read from all

the hlf on one is enough. → One is enough.

Q.3. Why NameNode send complete metadata info

B1. if machine 89 goes down, to avoid
roundtrip so that it can read from classmate

Write is parallel & Read is in sequence

Q. B1, B2, B3 are read ⁱⁿ Parallel or in sequence

B1, B2, B3 are read into sequence

We have to maintain the sequence as well

So B1 block info is read first then B2 & then B3

CONFIGURATION: Where you will go change

the configuration

In hadoop mostly configuration are maintained

in configuration files (XML file) in property

- (i) core-site.xml
- (ii) hadoop-env.sh
- (iii) hdfs-site.xml
- (iv) mapred-site.xml
- (v) yarn-site.xml

Output file path

Any configurable properties are maintained in these files

JAR explanation: Processing logic (It contains classes A1, A2, A3, A4)
If Input file path is put in local file system it will be single file. But for distributed processing it should be present in HDFS.

Output file path: Meaningful info should be available on highly reliable hardware. HDFS is reliable storage.

PROCESSING IN HADOOP GEN-L

CLASS 4

commands

Q. What is jar file.

Jar file: (Java archive file.) compiled version of your program written in Java programming language. With the help of Java compiler it get compiled and converted into Byte code (jar file). This Byte code is given to JVM. Then JVM convert into machine level code which is understood by the operating system.

Java is platform independent language.

We have three file over here : where these files

should be placed?

local filesystem / HDFS.

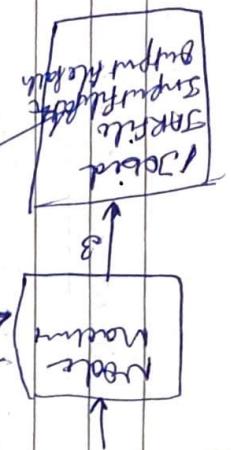
(no class)

✓

Inputfile Path

✓

B1, B2, B3 : are processed in parallel with the help of Mapper logic



Reducer logic : is applied to intermediate O/P. The O/P of Mapper logic is stored in local file system (emp directory) & local file system.

[Mapper + Sort and Shuffle + Reducer]

Map-Reduce Framework

Failure Scenario

Consolidated reasons

- 1) **Mapper failure** (due to networking issue, socket issue, cooling issue, data issue (Name instead of id), whatever may be the issue if mapper fails)
- Mapper failure is reported to job tracker through task tracker.

Job tracker when come to know - JT true

very simple approach to run it again

If Mapper fails : Then Task tracker runs again & it start processing work again from scratch.

Sort & shuffle cannot start till the time

#3rd scenario Task tracker failure:

off of all the mappers will be available.
So sort & shuffle will wait for the off from failure mapper machine.

→ **Mapper failure result into delay of job**

2) Reducer failure : Once again it is a logic. Reducer can also fail because of same reasons. (Programming issue, network issue, data issue)

Then Task Tracker report to JT that

Reducer has failed.

JT again will take same simple approach to ask (JT) to run reducer program

Again

→ Reducer failure result into delay of jobs

Q What will happen if mapper & Reducer fail again?

Same procedure happen

JT tell JT start → JT → Delay of Job

about again (Mapper failure happen) Reducer

→ Task tracker will also result into the delay of jobs.

4th Scenario Job Tracker failure:

All the jobs will be terminated. All the system resources will be released.

Computer Cluster become stand still.

Master failure result into comprehensive failure.

Job tracker is a single point of contact. In case

of JTF failure.

Or biggest problem in hadoop

Central processing : Single point of failure

Redundancy

(JTF will try 5 times) I think this is a good classmate

It will come to know to JT because JT knows that Task tracker is not working so B1 processing is not going on. So it will contact to NameNode where metadata info

B1 is present so if one node fails it knows that where on the other machine B1 is present so earlier B1 was running on 89

machine (no node failing of 89). One more node is going to be selected & that job process starts from scratch on the new selected machine.

Job tracker is the biggest problem (single point of failure) in Hadoop 1.0 Gen.

Q why the job tracker is failing? How that

problem can be solved out

The job tracker is heavily engaged doing almost all things

1) It's the job tracker → first (single) point of contact

2) It's job tracker who is maintaining the queue

3) Job tracker is looking to the ~~Speculative~~ folder of client

4) It is the job tracker who contact to name node for looking to the complete metadata information

Job tracker who select the m/c on the basis of Data Locality

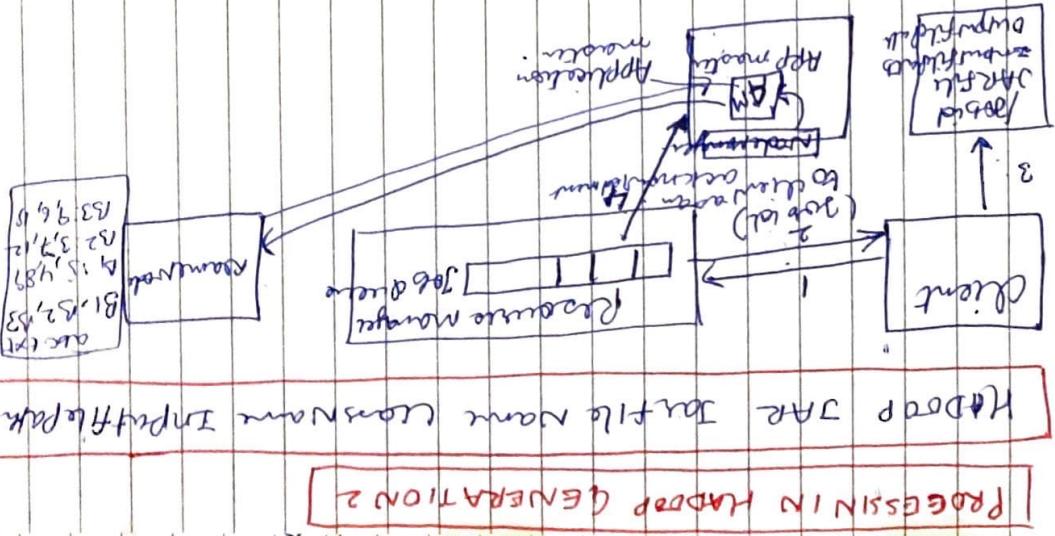
5) At the ~~task tracker to task~~ Job tracker connecting to TT asking for giving resources for data processing

6) Job tracker Any failure scenario (mapper, Reducer, task tracker) job tracker is contacted

J.T is not only handling one job but handling many jobs parallelly

It is the heavy engagement of Job tracker which become the reason for its failure.

who is going to be the first point of contact
→ Reserve Manager



Map-Reduce Framework

Mapper :

Takes input : key, value
Key : line number
Value : complete line
Output : key, value

Sort and Shuffle :

Input : key, value (output of mapper)
Output : key, list of values.

Reducer :

Input : key, list of values
Output : key, value

Sort & Shuffle

① Sort & shuffle occurs on the basis of key → ② shuffle the values &

Output :
Name : Supriya Name : Kausik
Age : 19 Age : 21
NAME : {Supriya, Kausik}

O/P of shuffle : list of values

How the requirement is solved using user simple map-reduce:

Example: →

Cid, temp, date

1, 23, 1

2, 31, 1

3, 44, 1

4, 19, 1

1, 25, 2

2, 30, 2

3, 43, 2

4, 18, 2

1, 24, 3

2, 34, 3

3, 44, 3

4, 19, 3

Likewise we have all the cities across the globe.
we have data to last.

or Sales etc.

Process this particular file

Max. Temperature so far

we want to explore more temp.

city wise so far.

Input to mapper		Output of mapper		Reducer flow
(split entire)		Key value		Op
2, 31, 1	K . Value	1, 23, 1	City value	S
3, 44, 1	K . Value	(2, 31, 1)	City value	S
4, 19, 1	K . Value	(3, 44, 1)	City value	S
1, 25, 2	K . Value	(4, 19, 1)	City value	S
2, 30, 2	K . Value	(5, 43, 2)	City value	S
3, 43, 2	K . Value	(6, 44, 2)	City value	S
4, 18, 2	K . Value	(7, 44, 2)	City value	S
1, 24, 3	K . Value	1, 24, 3	City value	P
2, 34, 3	K . Value	2, 34, 3	City value	P
3, 44, 3	K . Value	3, 44, 3	City value	P
4, 19, 3	K . Value	4, 19, 3	City value	P

Op of Reducer

K, Value
1, 25
2, 34
3, 44
4, 19

Loop through values
and find max temp.

Q. Can i get the max temp for city from the mapper?
for aggregate values (sum, max, avg), we can not get
from mapper.

Example 2: Word Count

② Now If there are 93 blocks how many
mapper required (3 mappers)

③ If former 1 Mapper

key-value

classmate

classmate

Description about word count program.

Program starts with main method.

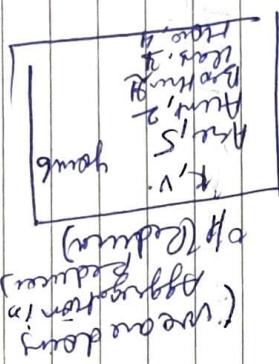
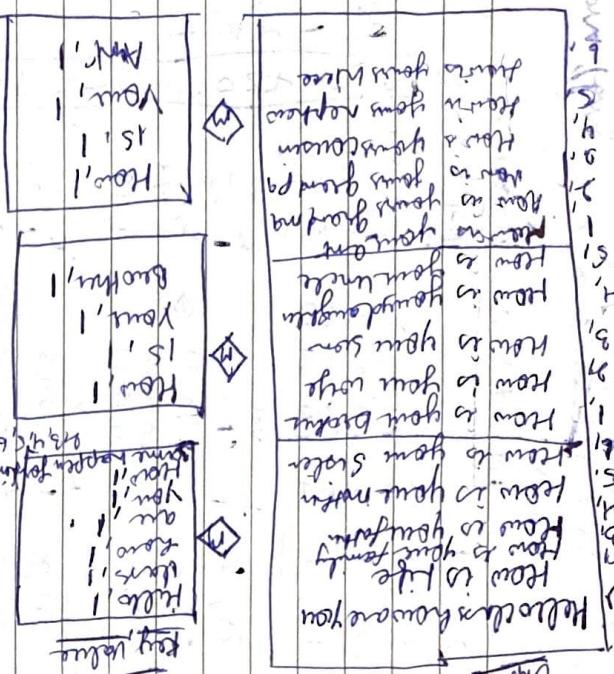
Q why the object of configuration is required?

We discussed about the configuration files, their

properties & values

Either we give the value of the properties

we don't want to give each & every properties
of the configuration files..



classmate