

## Introduction to NOSQL

### NOSQL

- ① open source, non-relational, distributed database
- ② it can handle all 3 forms of data
- ③ does not support ACID properties
- ④ no proper fixed schema for the table

### Types of NOSQL databases

- ① key value (key value pair)

Ex: frame: Dhan Ladnu: Nayak

- ② schemas - Document form of data

i) column - data is stored in column

ii) graph - consists of node & links

btitle: —

author: —

publisher: —

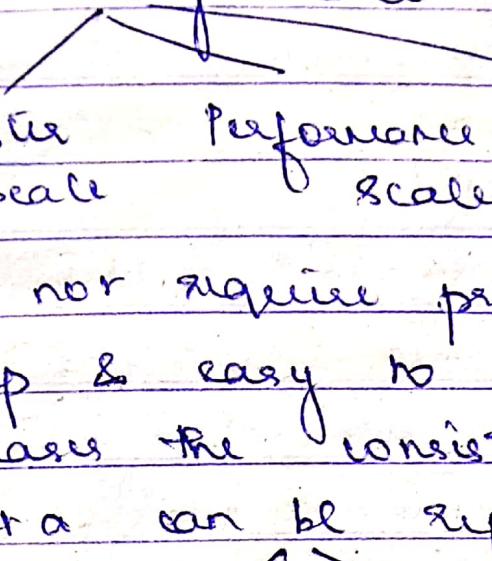
}

### Need of NOSQL

- ④ relational DB have monolithic but the NOSQL architecture have scale out architecture
- ⑤ it can store large amount of structured, unstructured & semi-structured data
- ⑥ Dynamic schema - (insert | update | delete are easy; changes to appl in real-time is easy)
- ⑦ Auto Sharding  
The availability data is distributed into different arbitrary servers so that there is no downtime if any server crashes
- ⑧ Replication  
because of the replication of data into all the arbitrary servers that are present downtime will not be present

## Advantages

- ① can easily scale up & down.

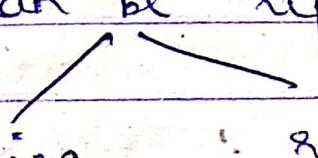
  
cluster      Performance      data  
scale            scale            scale

- ② do not require proper schema

- ③ cheap & easy to implement

- ④ releases the consistency of data

- ⑤ Data can be replicated & partitioned

  
Sharding      replication

## Disadvantages

- ① There are no joins

i.e. inner join, outer join etc & it's group by clause

- ② NO ACID properties.

- ③ integration of the application with other applications which support SQL is not possible

# use of NoSQL in industry / Application



## NoSQL

- accepts ACID, OLAP features

<u>SQL</u>	<u>NoSQL</u>	<u>NoSQL</u>	<u>New SQL</u>
1) ACID properties	Yes	No	Yes
2) OLAP/OLTP properties	Yes	No	Yes
3) Schema rigidity	Yes	No	Maybe
4) Dataformat flexibility	No	Yes	Maybe
5) Scalability	Vertical	Horizontal	Scalable
6) Distributed computing	Yes	Yes	Yes
7) community support	Huge	Growing	Slowly growing

## Mongo db

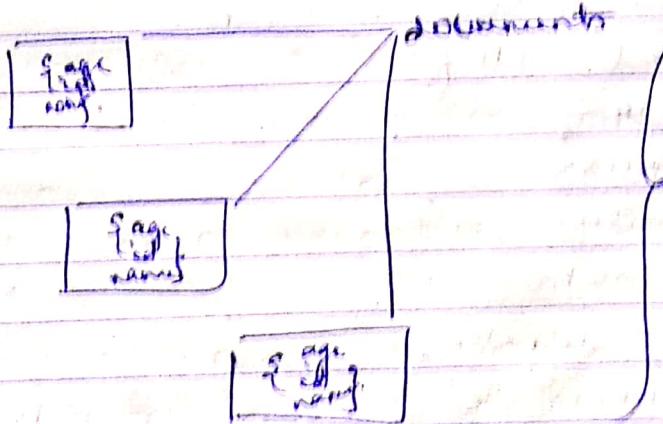
- cross platform
- open source
- distributed
- Multinational
- Document Oriented data store
- NO SQL

using a single script object mutation  
open standard  
complex

name : abc  
name :  
contact no :

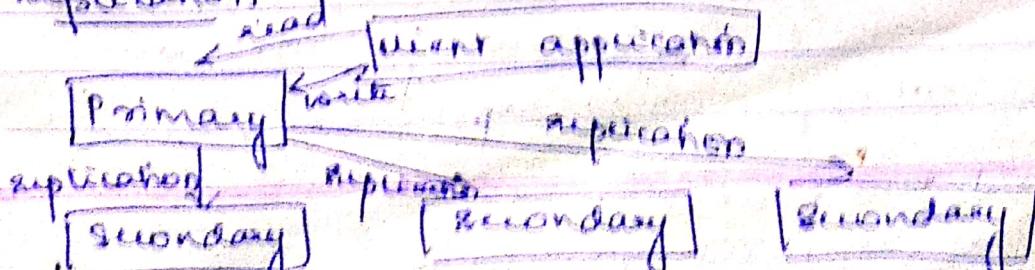
object 2  
keys to other data

- ② writing & querying structure
- here there are collection - which have data (structured)
  - one mongo db server can have multiple data bases.
  - Accessing of a collection is done through fields
  - Also multiple documents can be available



- ③ support for dynamic queries
- ④ support to store binary data using GridFS
- The meta data stored in form of logs
  - The data is divided into chunks & stored in file

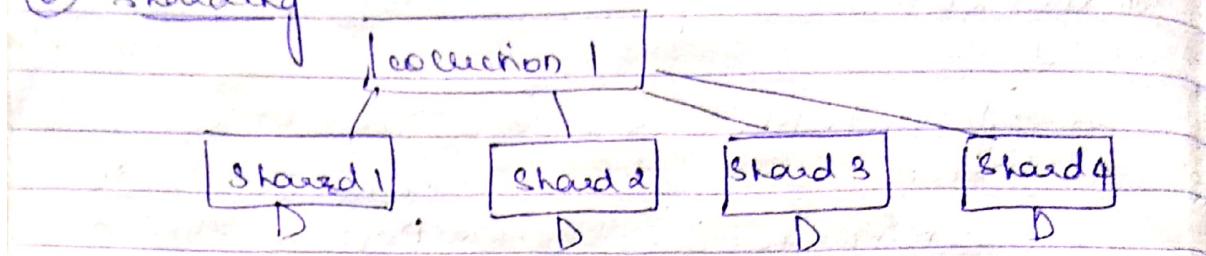
### ⑤ replication



When there are multiple copies of data, the data

MUST receive the updated data. The data should be updated in all copies.

#### ④ Sharding



Large data set is divided into multiple shards or servers.

Each shard will have some data to be stored and processed. Therefore easy to process & maintain.

#### ⑤ updating data in place

use the storage less the performance.

### Apache Cassandra

feature - highly scalable.

- distributed interface

#### ① peer to peer network

→ it has multiple peers.

→ meritable datastructure

→ easier to access data

#### ② gossip & failure detection

or will come to know the failures in network

#### ③ Partitioner

#### ④ replication factor

no. of copies of rows that is available

simple strategy

network topology strategy

also add the additional note  
to existing network

#### ⑤ anti-entropy & read repair

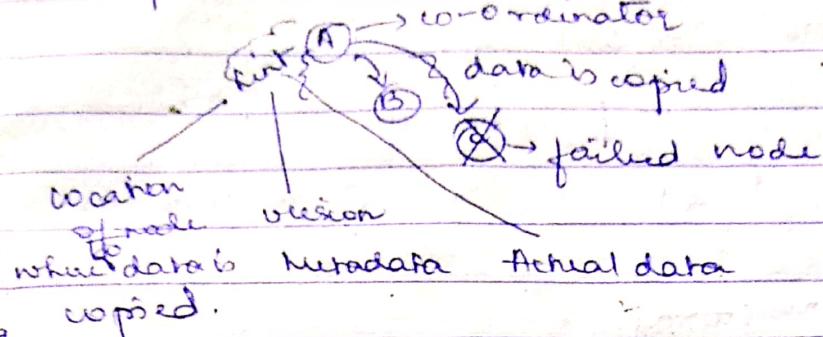
updated data is compared  
with all replicated data

### ⑥ writes in cassandra

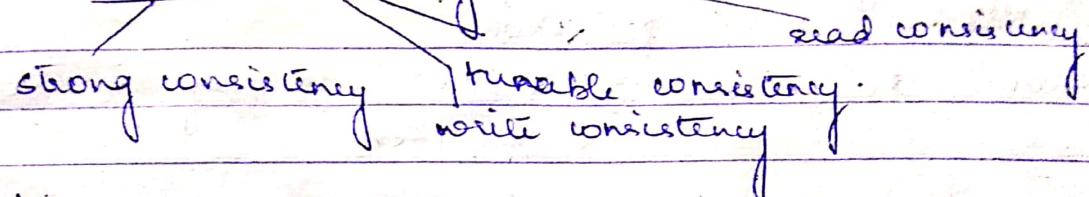
One write is updated to log. It's stored in mutable & written to immutable. Total data is scattered & replicated. [SStable = sorted string table]

### ⑦ hinted handoffs

- multiple nodes in network



### ⑧ tunable consistency



weak

## SUMMARY

### Hadoop

#### ① Every day

Facebook - 2.7 billion

Google - 2.4 petabyte

#### ② Every minute

FB - 2.5 million

Twitter - 3,00,000

Instagram - 2.20,000

Youtube - 12 hours

Email - 20,00,000 message

Amazon - \$800,000

### Why Hadoop

- handling large amount of data
- different categories of data.

## feature

### ① low cost

Since open source & use hardware to store large amt of data

### ② computing power

use distributed computing

↑ no. of nodes ↑ computing process

### ③ Scalability

more network, add extra nodes to improve performance.

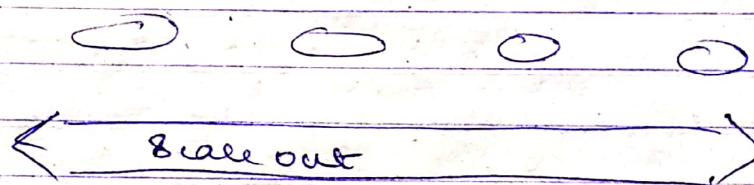
### ④ flexibility in storage

storage of any kind of data, is possible

### ⑤ Intact data protection

How data is protected. In case of failure other nodes take up the task

## Hadoop Framework



chunk 25-30

35-40

30-35

30-35

25-30

40-45

43-50

- multiple nodes together - chunk

- there can be repetition of chunks for backup

- resources are parallelly utilized.

- any no. of chunks can be added

## Why Hadoop & not RDBMS?

- not suitable for advanced analytics - DBMS such as unstructured data

- in DBMS scaling out of data & huge computation is very costly

<u>Paranutes</u>	<u>R B M S</u>	<u>Hadoop</u>
① system	Relational Database Management System	node based flat structure
② Data	Structured	unstructured - video, Audio, images - different format
③ Processing	OLTP (online transaction processing)	Big data processing analytical files
④ choice	consistency & relationship b/w data	Big data & no consistency
⑤ Processor	high End processor & costly hardware	cluster & multiple nodes
⑥ cost	£10,000-14,000/terabyte bytes of data	hardware & NC card \$4000/terabyte

### Distributed computing challenges

① hardware failure → in case of hardware fail → how node access data.

② how to process gigantic data  
replication factor = 2

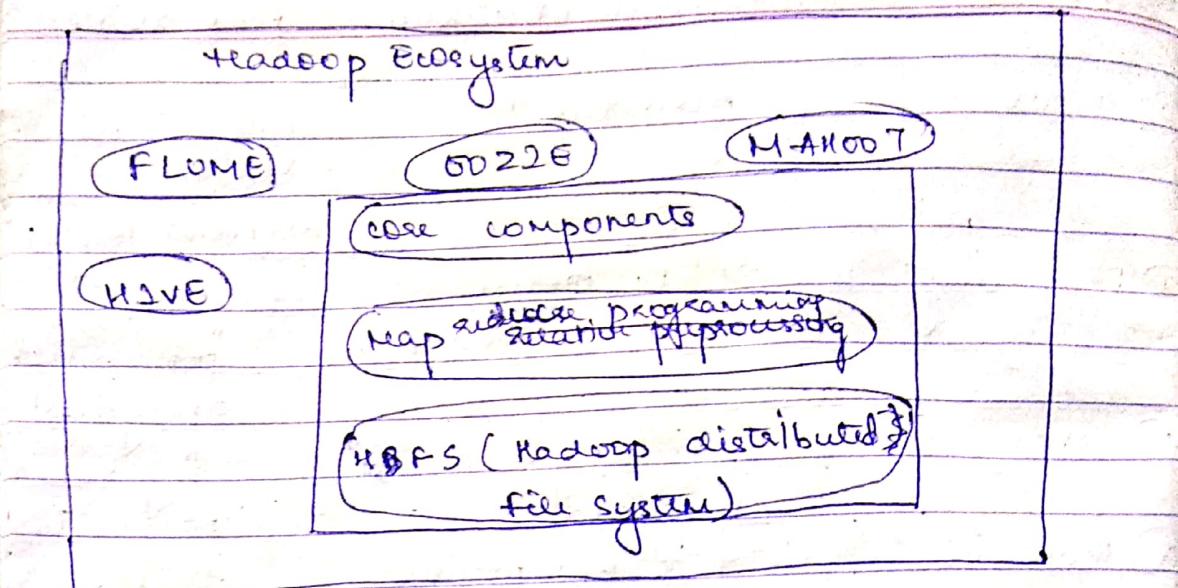
### Hadoop Overview

- ① faster data processing
- } 2 basic tasks,
- ② massive data processing

### Key Aspects of Hadoop

- open source - free to use & contribute
- framework - tools program & apps
- Distributed - among multiple nodes
- massive storage
- faster processing - Quick response

# Hadoop Basic components



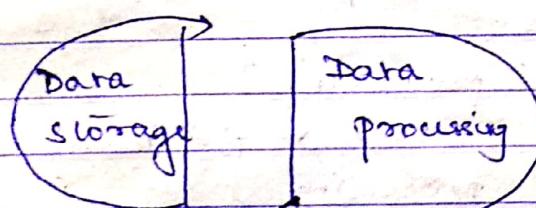
## ① HDFS:

- distributed system • data distributed across various nodes
- redundant
- storage component

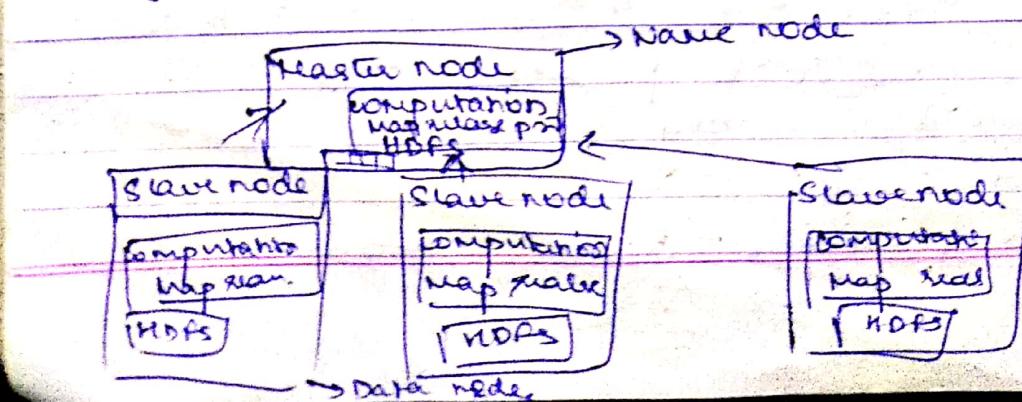
## ② Map Reduce programming

- computational framework
- task & distributed across nodes
- parallel processing of data

# Hadoop conceptual layer



- master slave architecture
- high level architecture of Hadoop



Master HDFS - keeps location of data in slave  
- distributes data among slave

Master MapReduce - computations & distributed among slaves equally

use case of hadoop

→ clickstream data - analysis on different data based on click

→ customer relationship management - knows behavior of particular customer

→ scalability - increase storage.

→ Business analysis - Apache, Hive.

→ HDFS (Hadoop Distributed File System)

◦ Storage components

◦ Distributed file system

◦ Optimized for file throughput

◦ Modelled after google file system

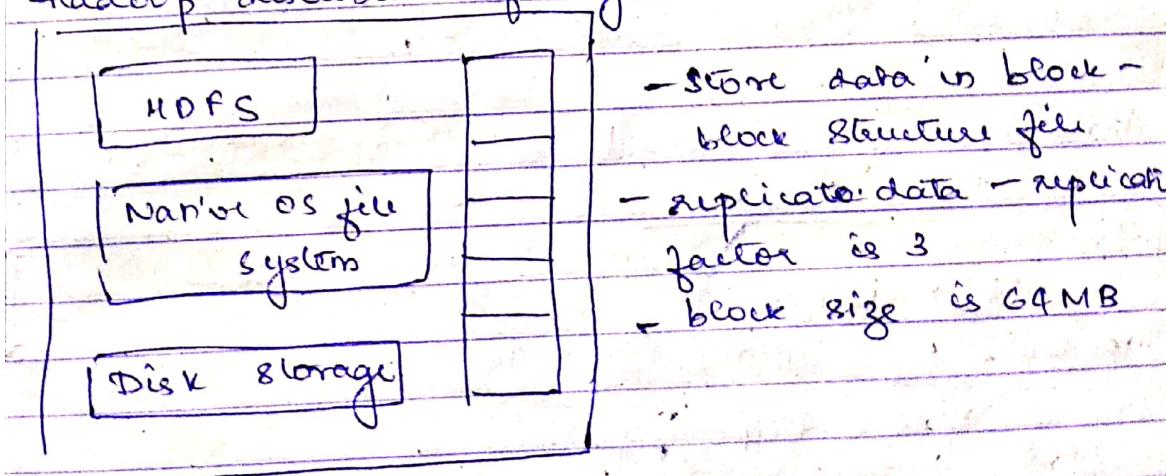
◦ replicate the file any no. of times

◦ replicate data on failure

◦ read / write

HDFS Daemons

Hadoop distributed file system

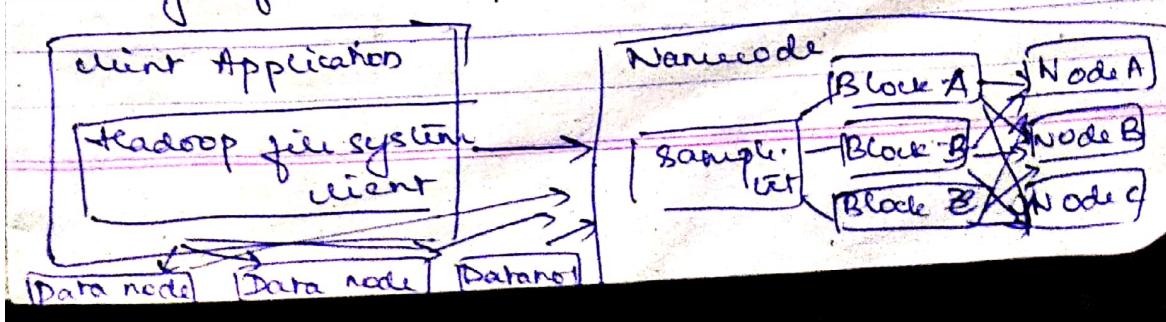


- store data in block -  
block structure file

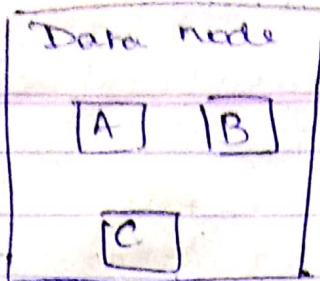
- replicate data - replication factor is 3

- block size is 64 MB

A large file is splitted into blocks! They are.



Each data node has following



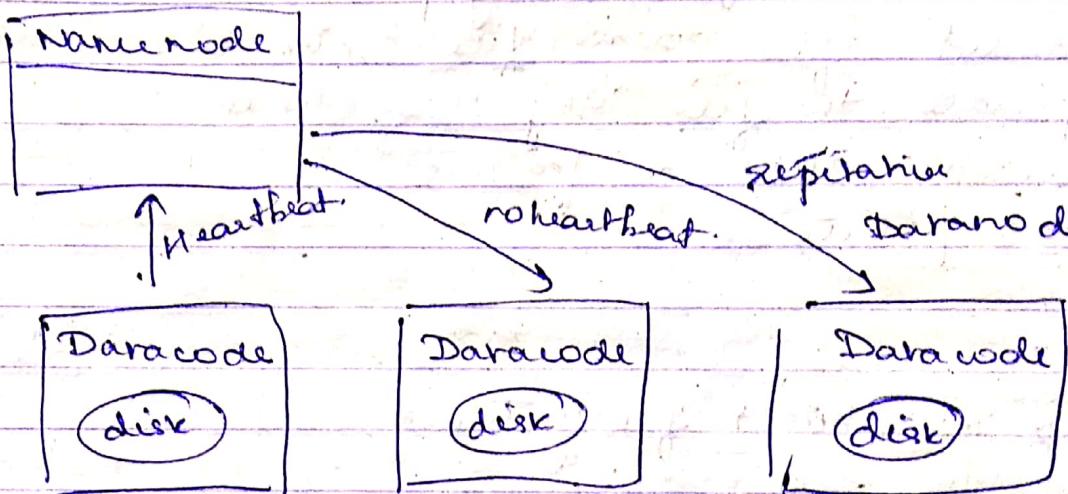
Name node manages various file operations that are done & manages file systems

delete  
read write update

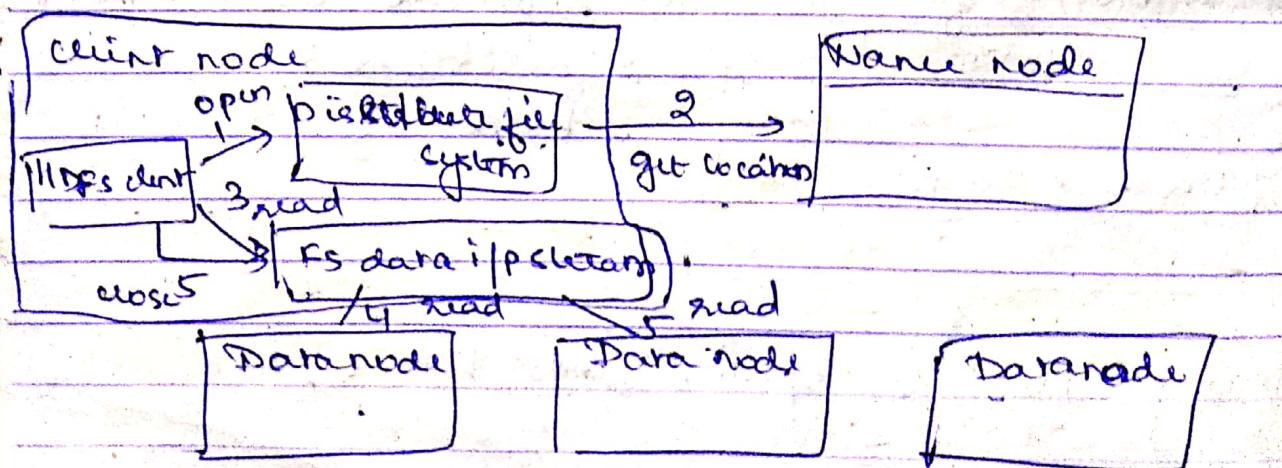
FS image - Interface to map between input file & blocks.

editlog - different transactions are dotted down

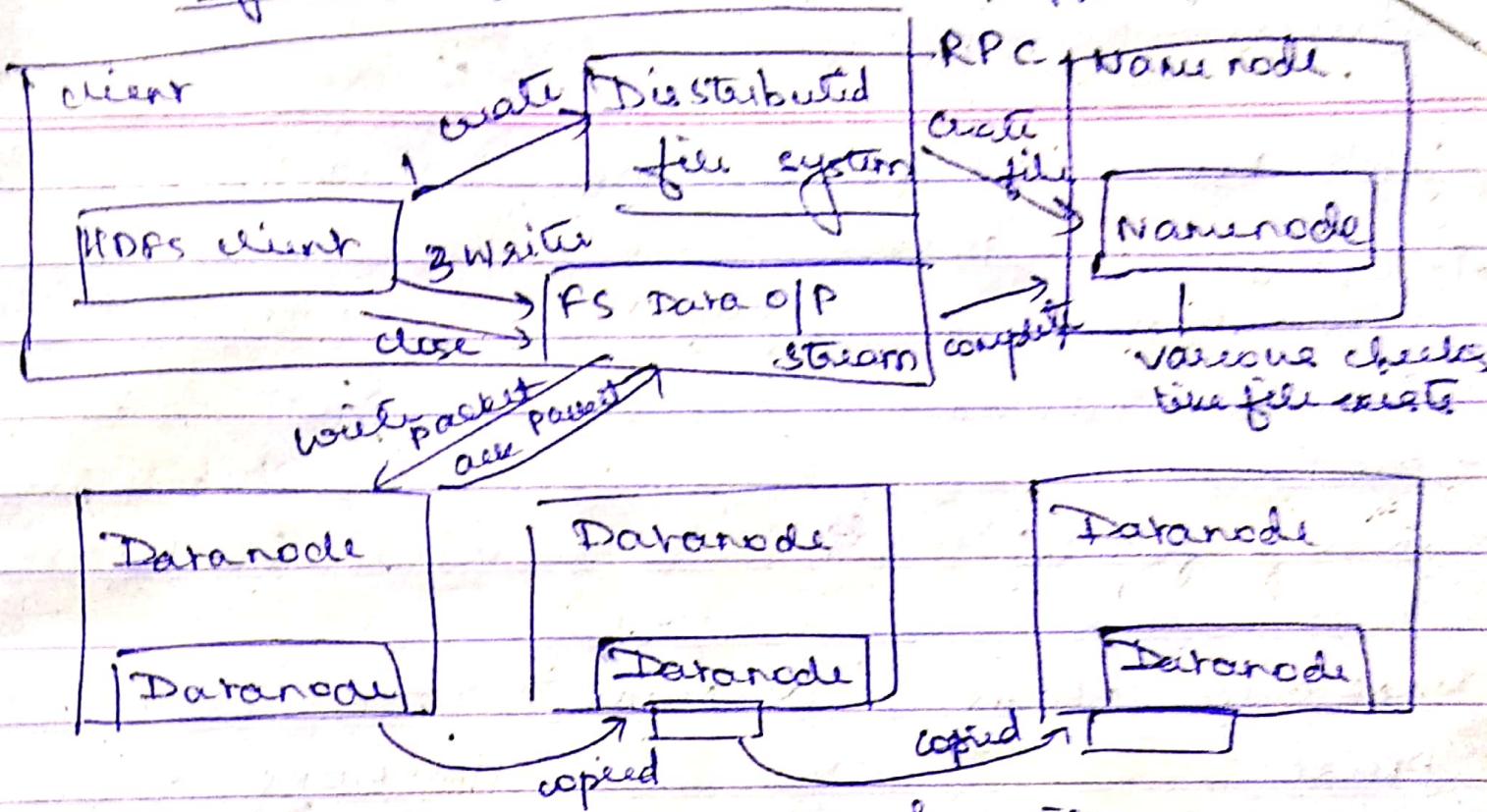
Communication between Name node & data node



Anatomy of File read

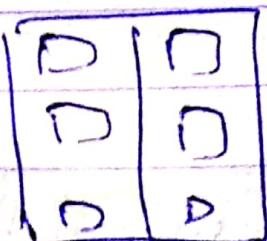


## Anatomy of file write



data is divided into packets & written into  
internal queue

## Replica placement strategy



rack1 rack2

## Special features of HDFS

- i) Data replication: Client will check for nearest node.
- ii) Data pipeline: Data is put to first node & then moved, & replica is in all nodes of pipeline.