

Big Data Analytics →

Big Data (huge Data)

↳ quantity is very high.

- * It is collection of data that is huge in volume
- growing exponentially with time

Big Data Analytics →

process of examining large data sets containing

a variety of data types to uncover hidden
patterns, customer preferences.

Challenging →

- * challenges → (Amazon story)
 - ① Insufficient understanding & acceptance of big data
 - ② Confusion while selecting Big Data tools.
 - ③ paying lots of money.
 - ④ Data integration. (Data ko alag alag jagah se (different) collect karna)
 - ⑤ Data Security (Data ko jis tarah store karna)

- * Problems with Traditional large ~~scale~~ Scale Systems →
 - Majority of data comes in the form of semi structured or unstructured data.
↳ e.g., image/video
 - Big data bahut hi teji ke sath grow hoتا है jisse traditional system nahi handle kar sakte. Vo Keval steady data ko hi handle kar sakte hai.
 - Data is so expensive to store in traditional systems. Data is filtered & aggregated and large volume of data is thrown out but data ko minimize करने se accuracy reduce hoti hai.

* Sources of Big Data →

- Social Networking sites (Instagram)
- E-commerce sites (Amazon)
- Weather station (Satellite)
- Telecom company (Jio/Airtel)
- Share Market (daily transaction)

* 3V's of Big Data

Velocity

Variety (photos, videos, gif, pdfs etc)

Volume

* Types of Big Data

Structured

store, access,
process
in fixed format

e.g., Tabular data

unstructured

(Photos
audios and
videos)

Semi-
structured
data
(xml file)
(log file)

Simple Questions

Working with Big Data

* Google file system (GFS)

- GFS is a scalable distributed file system.
- GFS provide fault tolerance.
- Runs on inexpensive hardware.
- Delivers high performance to the large number of clients.
- files are organised hierarchically; in directories and identified by path name.
- supports (CRUD) (Creating, reading, updating, deleting).

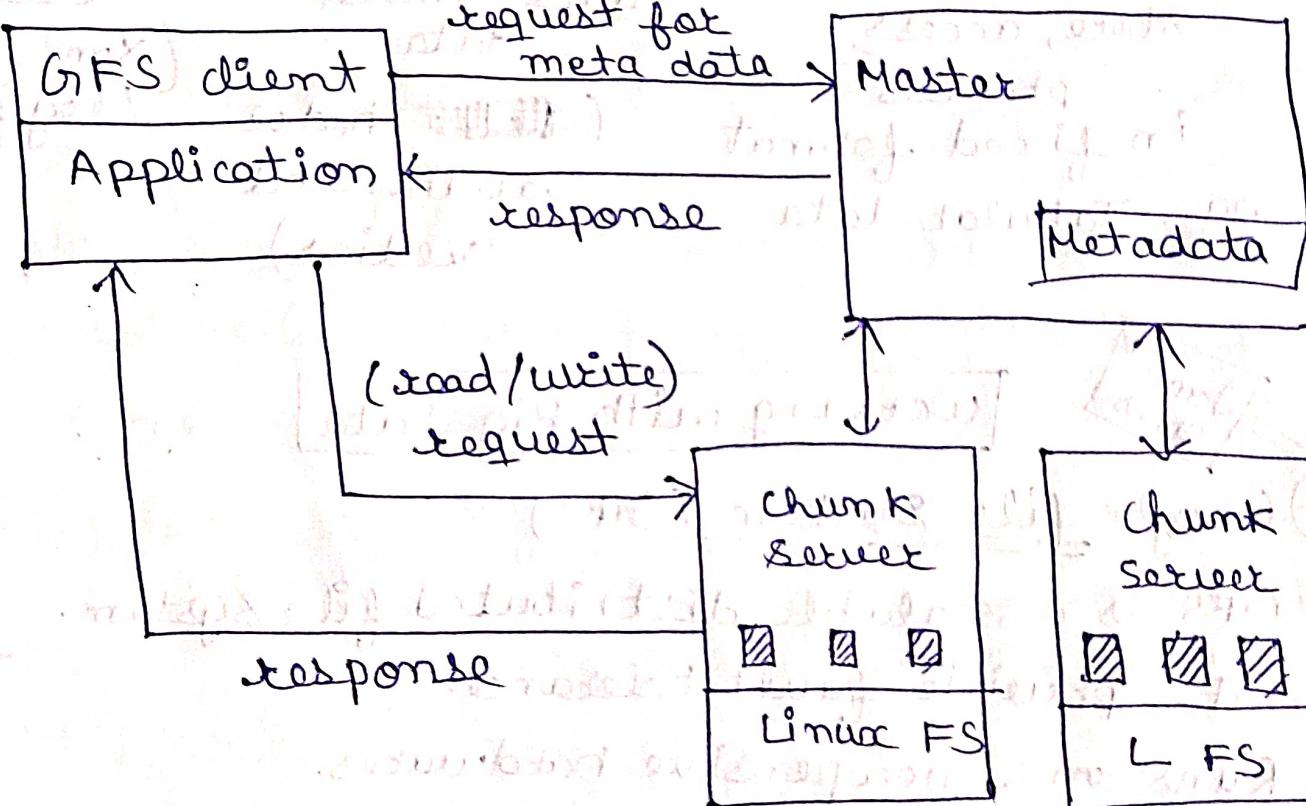
* Snapshot operation

file / folder ki copy banta hai low cost mai.

* Record operation

ek time par ek se jayada person ek file mai data append kar sakte hain

GFS Architecture →



* Hadoop Distributed file system →

→ based on GFS.

→ designed to run on large clusters

(thousand of computers)

* Master / Slave architecture



Name Node Data Node
(contains (actual)
meta data) data

→ files split into several blocks and stored
in a set of data nodes.

* Building Blocks of Hadoop

① Name Node (Master)

→ Stores metadata
→ kyuki hamare pass ek hi Name Node
hota hai to agar ye fail hua to pura
hadoop system fail ho jata hai.

② Data Nodes →

→ Client can communicate directly with these
nodes.

③ Secondary Name Node

→ assistant that monitors the state of the
cluster
→ har cluster ke pass ek secondary Name
Node hota hai.
→ ye Name Node se communicate kar
sakta hai HDFS metadata ka Snapshot lene
ke liye.

④ Job Tracker

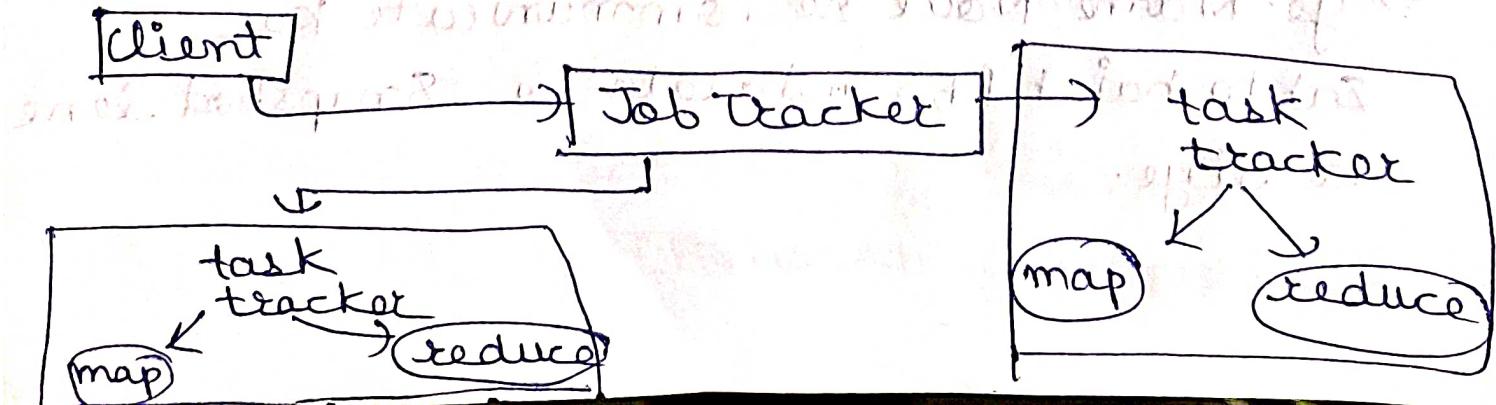
- Link between your application and hadoop
- hadoop cluster ke pass ek Job Tracker hota hai
- ye basically dekhta hai ki aur task ko

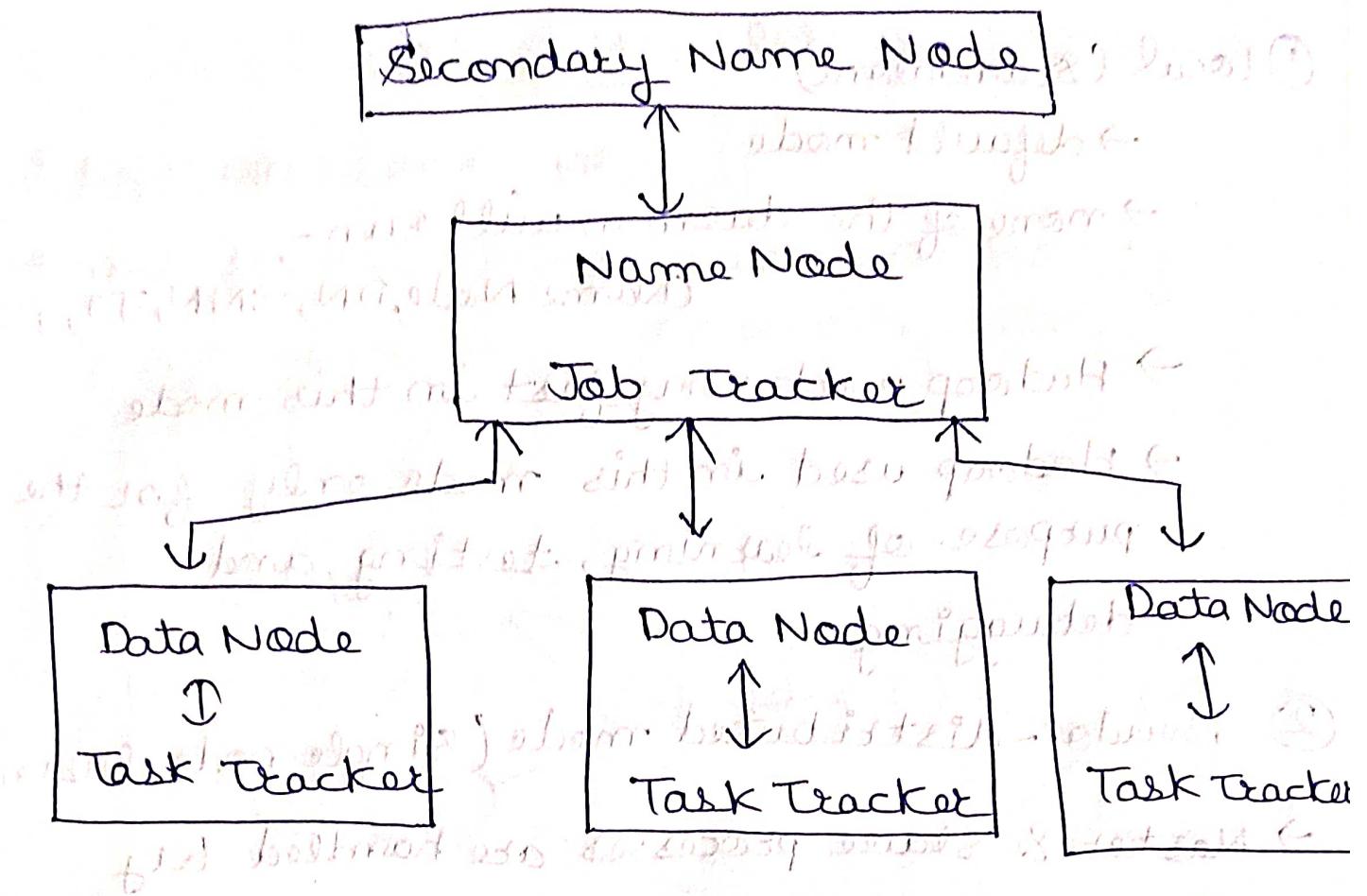
- ① Kaunsi file pehle process hogi
- ② Kaunsi file ko kis data node ko de raha hai
- ③ aur task ko monitor karata hai aur agar koi task fail ho jata hai to usse re-launch karata hai.

⑤ Task Tracker

- ye task karata hai aur constantly Job Tracker se communicate karata hai in the form of heart beat. aur agar ek time interval ke baad Job Tracker ka heart beat nahi milta to vo uss task ko fail mark kar ke launch Kar deta hai.

Architecture





* Introducing and configuring Hadoop

→ we need to configure several xml files)

↳ configuration directory of Hadoop-Hom
hadoop-default.xml

→ Hadoop cluster can be configured by one of the three modes:

- ① Local (stand alone) (default mode)
- ② Pseudo-distributed
- ③ fully-distributed

① Local (Standalone)

→ default mode

→ none of the daemon will run -

(Name Node, DN, SN, JT, TT)

→ Hadoop work very fast in this mode

→ Hadoop used in this mode only for the purpose of learning, testing and debugging.

② Pseudo-distributed mode (single node cluster)

→ Master & slave processes are handled by single system.

→ All the processes inside cluster will run independently to each other.

→ All daemon will run separately on various JVM (Java Virtual Machine)

→ used for development & debugging purpose

③ fully distributed mode (Multi Node cluster)

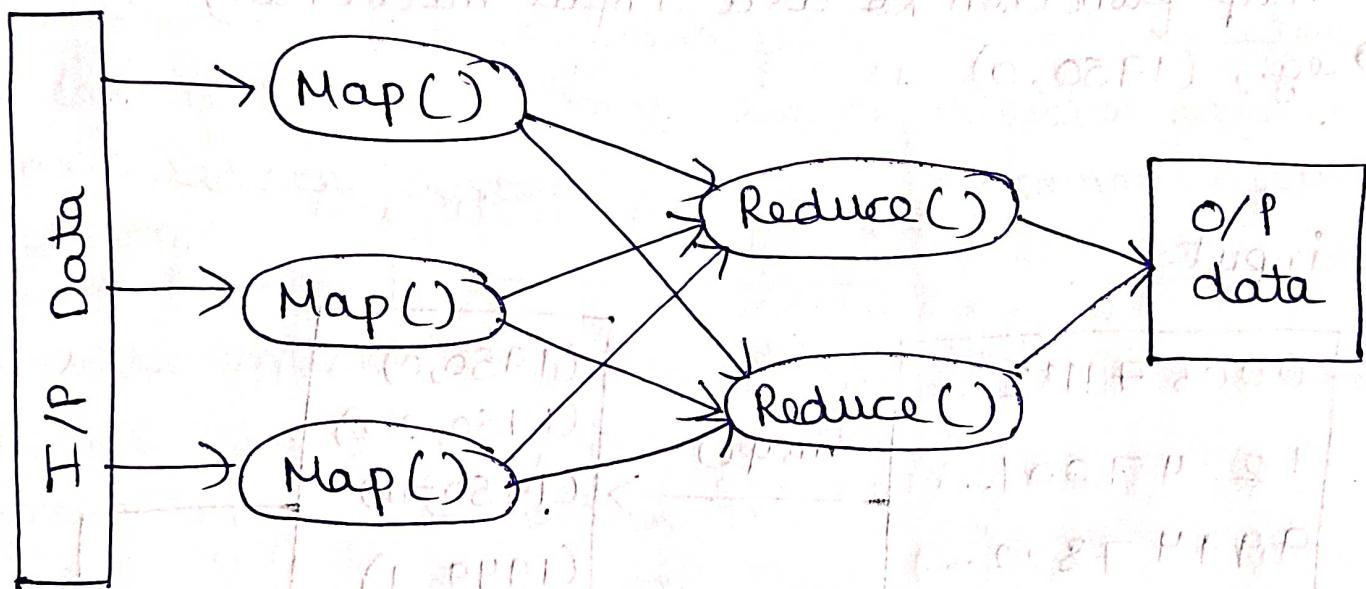
→isme bahut saare nodes hoti hei kuch Master Daemon ko chalati hai aur kuch Slave Daemon ko chalati hai

→ ye production mode mai use hota hai

writing Map-Reduce program →

① Software framework

→ used for writing app which process big amount of data



→ Single Master (Job Tracker) and one Slave (Task Tracker) per cluster.

* Weather Dataset →

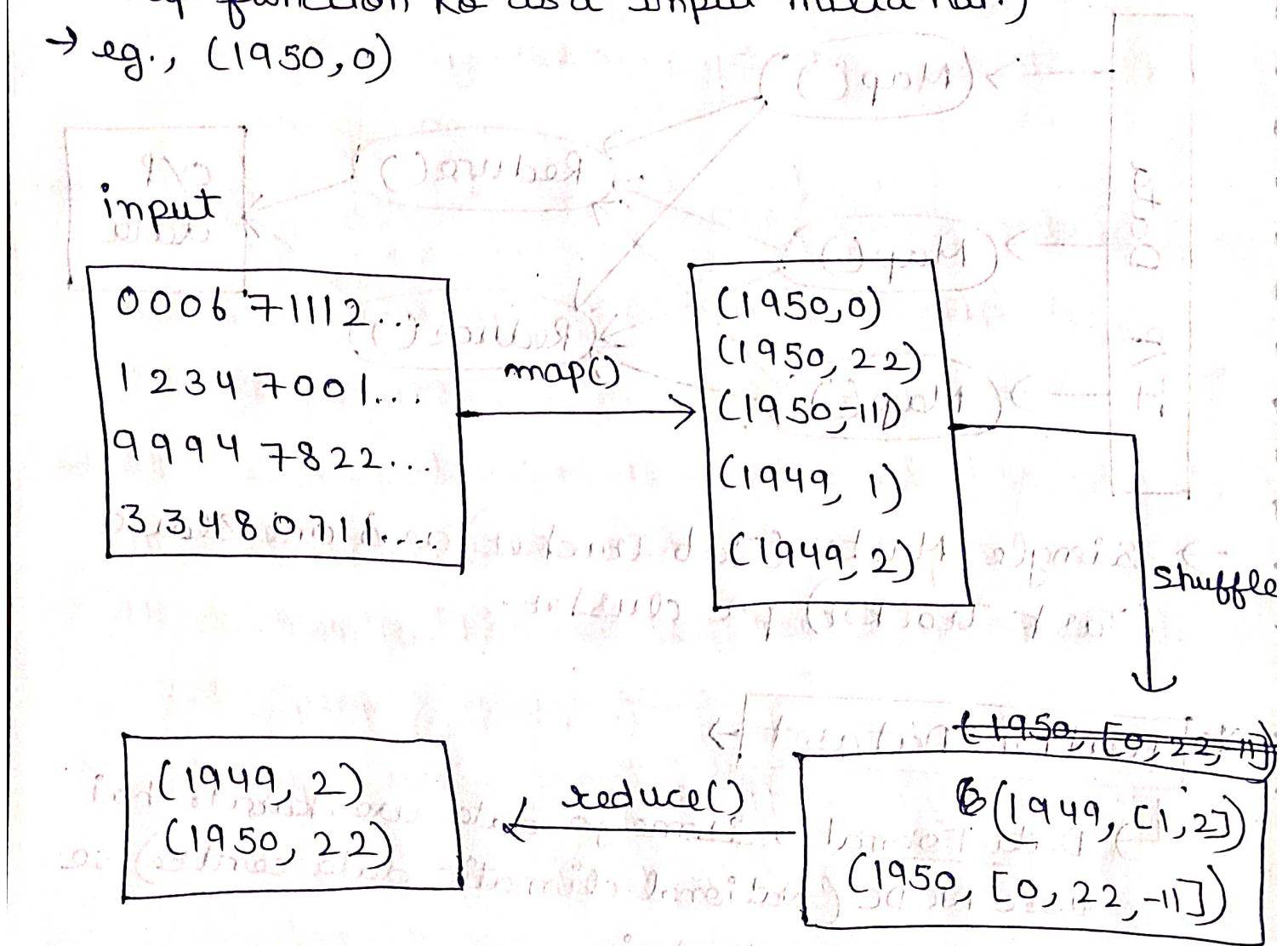
↳ Data Format →isme jo data use karate hai usse NCDC (National climatic data centre) se collect kiya jata hai

→ line oriented ASCII format

→ Each line is a record.

→ Data files are organized by Date & weather station

- Analyzing the Data with Hadoop →
- Weather Data set → Map-Reduce Program → Hadoop
- I/P → Raw data from NCDC.
- Pull out year and the air temperature in map function. (matlab ka year ka air temperature map function ko as a input milta hai.)
- eg., (1950, 0)



* Old and New Java Map Reduce API →

Difference	② New API	Old API
① Mapper & Reducer	It uses Mapper & Reducer as class	as interface

② Package	New API is in org.apache.hadoop.mapreduce package	old API is in org.apache.hadoop.mapred package
③ user code to communicate with map-reduce System	use context object to communicate with map-reduce System	uses Job Conf, output collector, reporter object to comm. with map-reduce system
④ extra control Mapper & Reducer execution	It allows both mappers & reducers to control the execution flow by overriding the run() method	controlling mappers by writing a Map Runnable but no equivalent exists for reducers
⑤ Job control	Job control is done through the Job class	Job control was done through Job client
⑥ O/P file name	Map → part-m-nnnnn reduce → part-r-nnnnn	for both part-nnnnn
⑦ reduce() method pass value	Java.lang.Iterable	Java.lang.iterator

Basic Program of Hadoop Map-Reduce

→ Driver code → (initialize the Job) and (control mapping, reduce). Ki kaha par O/P file store hogni)

//code

Configuration config = new Configuration();

Job job = new Job(config, "my word count prog.");

Job.setJarByClass(WordCount.class);

Job.setMapperByClass(Map.class);

Job.setReducerByClass(Reduce.class);

Job.setOutputKeyClass(Text.class);

Job.setOutputValueClass(IntWritable.class);

Job.setInputFormatClass(TextInput.class);

Job.setOutputFormatClass(TextOutputFormat.class);

Path output-path = new Path(args[1]);

FileInputFormat.addInputPath(job, new Path(args[0]));

FileOutputFormat.setOutputPath(job, new Path(args[1]));

for setting up the path (matlab O/P file kaha

store hogni)

Driver class → add main-class

Submit the job

Mapper Code

```
public static class Map extends Mapper  
<LongWritable, Text, Text, IntWritable>  
{  
    public void map(LongWritable key, Text value,  
                    context) throws IOException, InterruptedException  
    {  
        String line = value.toString();  
        StringTokenizer token = new StringTokenizer(line);  
        while (token.hasMoreTokens())  
        {  
            value.set(token.nextToken());  
            context.write(value, new IntWritable(1));  
        }  
    } // void map close  
}  
} // class close
```

* Reducer Code →

```
public static class Reduce extends Reducer  
<Text, IntWritable, Text, IntWritable>  
{  
    public void reduce(Text key, IntWritable value,  
                      context) throws IOException, InterruptedException  
    {  
    }  
}
```

```

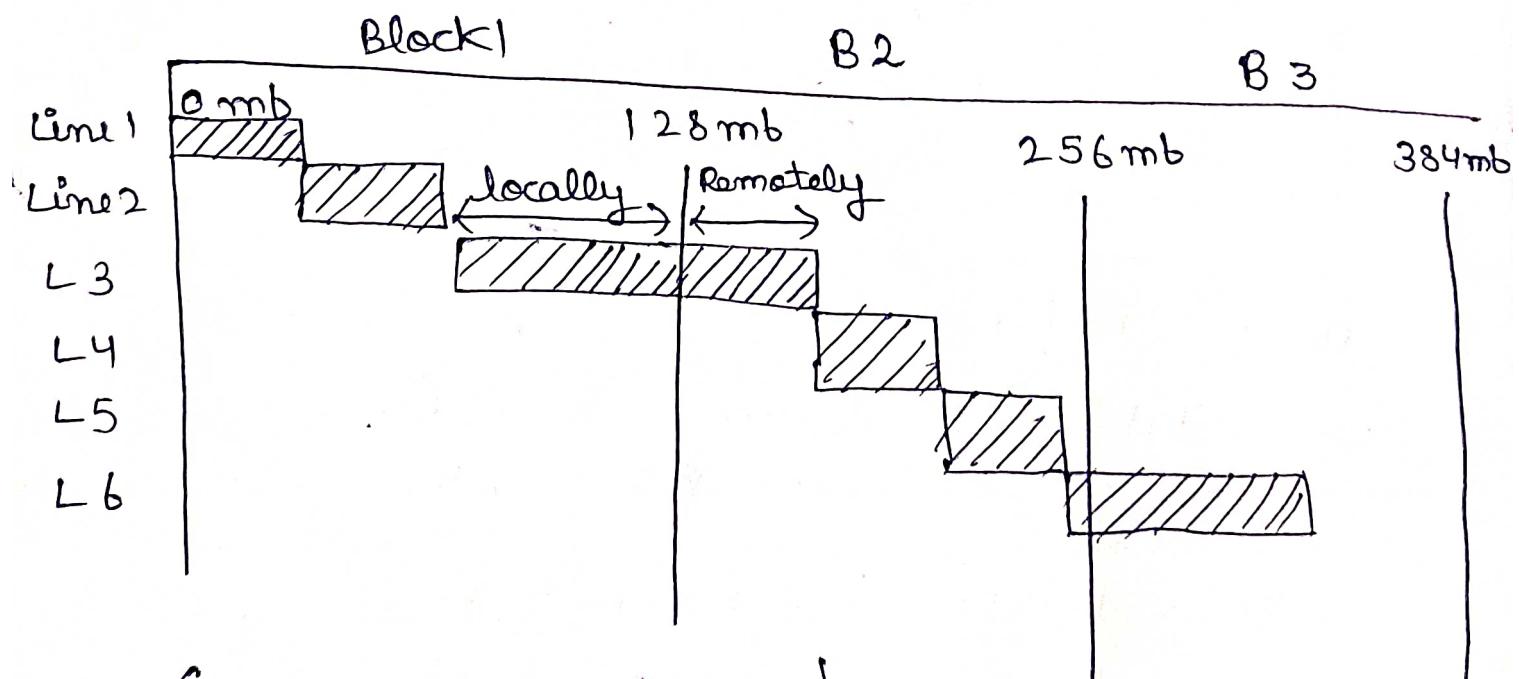
int sum = 0
for ( IntWritable x : values)
{
    sum += x.get();
}
context.write(key, new IntWritable(sum));
}
}

```

④ Record Reader → (This comes before mapper)

↓
Data/Input files

- Loads data & convert into key value pair
- Its Instance is defined by the InputFormat
- Text InputFormat - Default
- Key associated with each line is its byte offset.



(Data set = 300 mb
6 lines each of 50 mb)

* Combiner → It is used to summarize the map output records with the same key.
(Semi Reducer)

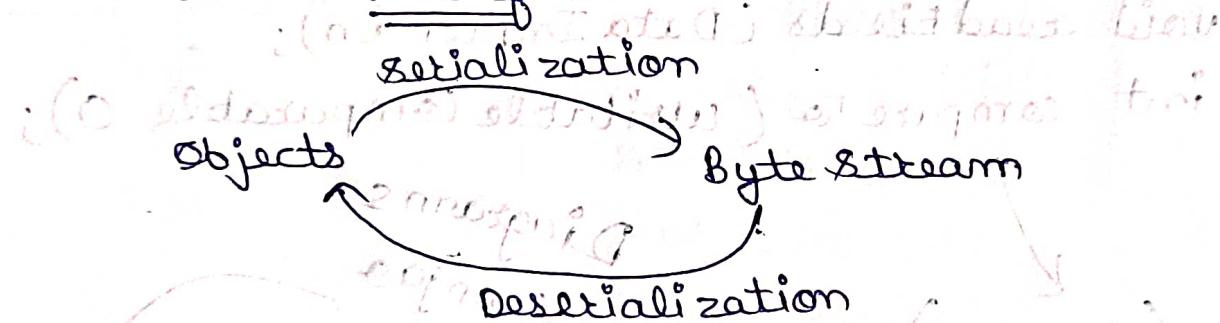
→ See (Map - Reduce Basic diag)

* Partitioner → Controls the partitioning of the keys of the intermediate map O/P.

→ configure() and getPartition()

Chapter → 3

* Writable Interface →



→ Hadoop uses its own ~~serial~~ serialization format known as Writables.

→ Two Method → Data output (write Byte stream)
Data Input ~~ff read~~

Program →

Package org.apache.hadoop.io;

import java.io.*;

public interface Writable

void write(DataOutput out);

void readFields(DataInput in);

{
String s;
}

Note → If writable is not present in Hadoop, then it uses the serialization of Java which increases the data overhead in network.

* suitable comparable and compatible →

* Sub interface of writables

public interface WritableComparable extends
Writable, Comparable

```
void write (DataOutput out);  
void readFields (DataInput in);  
int compareTo (Writable compar
```

Diagram 5

④ returns $(-1 \ 0 \ 1)$

$(K_1, V_1) \xrightarrow{\text{Map}} (K_2, \text{List}[V_2])$

$(K_2, \text{List}[v_2]) \xrightarrow{\quad} \text{Reduce} \rightarrow (K_3, v_3)$

writable compatible
is used here

Program

* Writable class →

→ Hadoop provides classes that wrap the Java primitive types and implement the Writable Comparable & Writable interface.

→ provided in org.apache.hadoop.io package

↳ set() → store wrapped value

↳ get() → retrieve " "

Java primitive	writable Implementation	bytes
byte	Byte Writable	1
short	Short Writable	2
int	Int Writable	4
long	Long Writable	8
float	Float Writable	4
double	Double Writable	8
boolean	Boolean Writable	1

~~Text~~

Some Writable class

Text → Writable for UTF-8 sequence

feature → indexing charAt()

Iteration

Mutability (change ho sakti hai string)

Byte Writable →

↳ Byte Writable ki value ko change kar sakte hai by calling its set() function.

Syntax

```
ByteWritable b = new ByteWritable  
(new byte[] {3, 5});
```

* Null Writable →

↳ used as placeholder & store empty value.

* Object and Generic Writable →

Object

- Java Primitive

- String

- enum

- null

* useful when a field can be of more than one type

Generic

better than object writable because in object writable we append the class declaration as a string into the output file in key value pair.

* writable collections →

Array writable classes

→ Array writable (1 D)

→ 2 D Array writable

Map writable class

→ Abstract Map writable

→ Map writable

→ Sorted M.W

Chapter 4

Pig

→ tool / platform used to analyse big data

→ Represent big data as data flow

analyse

* why do we need Apache pig if hadoop exist?

→ It is suitable for those coders who are not good in Java.

→ Pig Latin is used here and it is SQL like language (easy to learn)

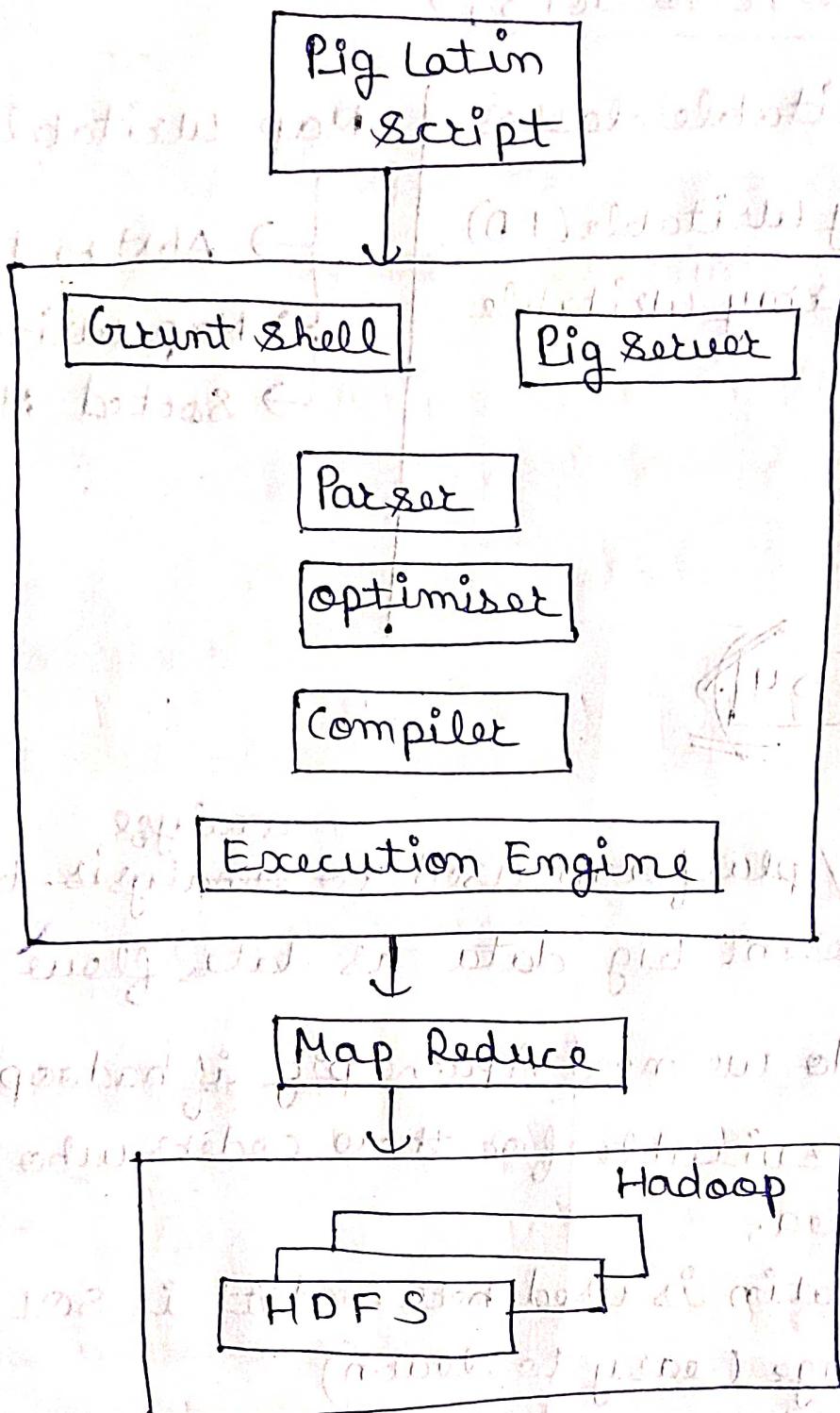
→ provides many built in operators (e.g., joins, filters)

→ provides nested data types (tuples, map, bags)

* Admitting Pig Architecture →

Language
(Pig Latin)

Compiler
(converts Pig Latin to executable code)



→ Going with Pig Latin Application flow →

A = LOAD 'data-file.txt' { stream definition }

B = GROUP;

C = FILTER ...; } transformation functions

~~DUMP~~ DUMP B;

STORE C INTO 'Result'

Working through the ABC's of Pig Latin →

→ keep it simple

→ Make it smart (coders ko Jayada di maag na lagana Pade development mai)

→ Don't limit Development (mai, kya kya)

kyunki iska main point hi yahi tha ki jisse Java mali aata vo isse use kare

* Looking at Pig data types & Syntax

Two data Types

Scalar

(contains single value)

Complex

(contains other types)

for eg., Atom - for eg., "Lavesh", 1, 1.2

tuple → sequence of field

Bag → collection of non unique tuple

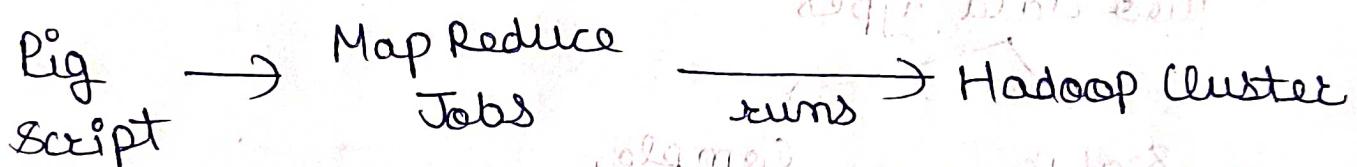
map → collection of key value pairs

* Evaluating Local & distributed Modes of Running Pig Script

Local Mode →

- useful when we have small data set
- useful for developing & testing Pig Latin
- Hadoop is not required
- Pig Program runs in context of local JVM

Distributed / MapReduce / Hadoop Mode →



* Checking out the pig Script Interface →

Script

Pig Script ko likhne ke tareeke
pig — Eg., Flight Data.pig

Grunt commands are interpreted by PL compiler
& executed in the order determined by
Pig optimizer.

Grunt

↳ It is a shell in which we can write small pig codes.

Embedded

→ we can use other languages also to execute Pig Latin statements

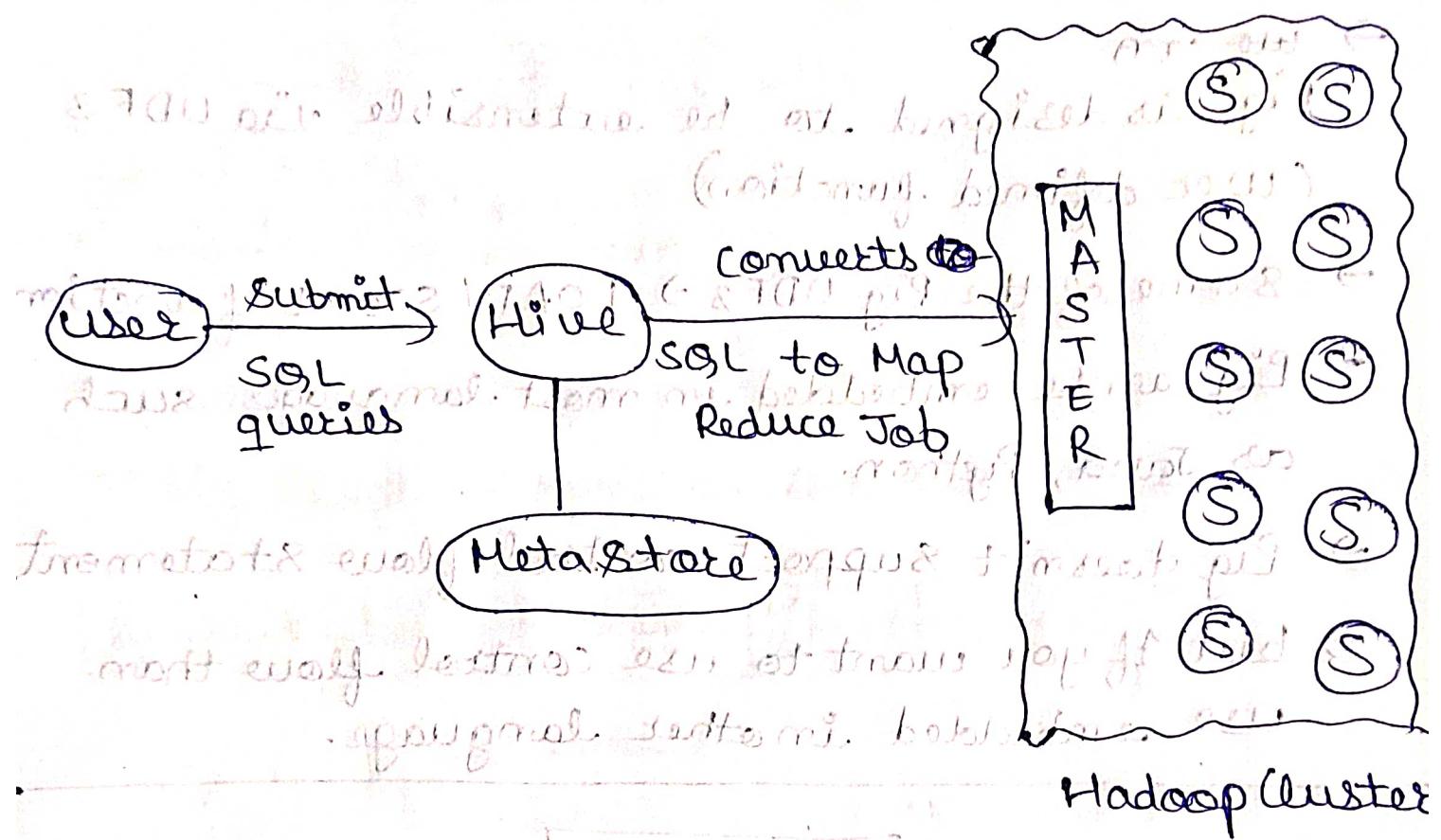
* Scripting with pig latin

- we can
- Pig is designed to be extensible via UDFs (user defined function)
- Some of the Pig UDFs → LOAD/STORE function
- Pig can be embedded in most languages such as Java, Python.
- Pig doesn't support control flow statement
- but if you want to use control flow then use embedded in other language.

Chapter → 5

Hive

- It is a open source data warehouse system for querying, analyzing large data set stored in Hadoop file
- { Analyze , querying, data summarize }
- It uses HiveQL
- $\text{SQL} + \text{Hive} = \text{HiveQL}$
- Highly Scalable



Master bigger task ko sub task mai divide
Karte slaves ko de data hai

Topic 2
Hive Data types → simple → complex

* Simple

Data types

tinyint

smallint

int

largeint

float

double

string

size

1 byte

2 byte

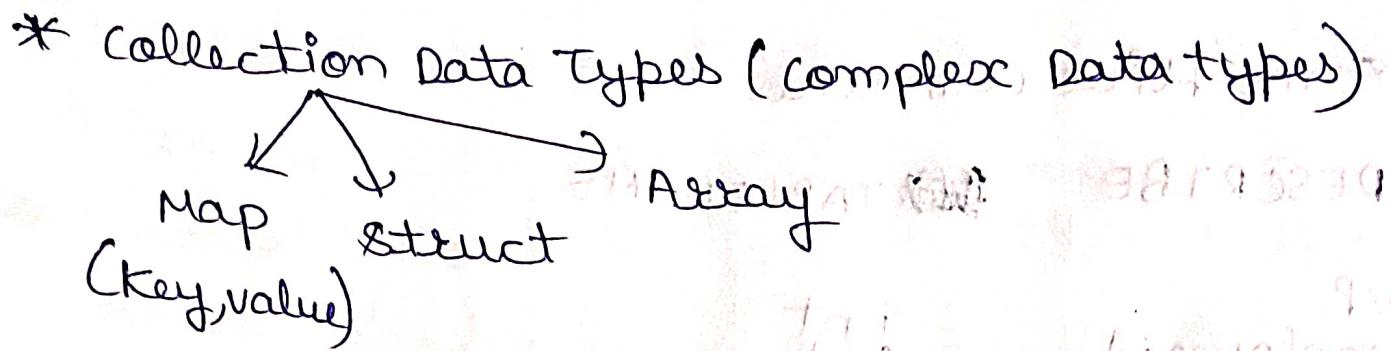
4 byte

8 byte

4 byte

8 byte

"John"



Topic → 3

* Different types of Hive Table

① Managed Table

CREATE TABLE TABLE NAME (col1 datatype,
col2 datatype, col3 datatype) row format
delimited fields terminated by " " STORED as
textfile;

② External Table

CREATE EXTERNAL TABLE -
— (Managed Table, Owned by HDFS
— — (External Table, not owned by HDFS)

* Difference

Managed Table

- ① Tab hum isse drop karenge tab data + meta data will be deleted

External Table

- ① only metadata is deleted.

Some Hive Commands

① DESCRIBE TABLENAME

O/P

employeeid int

int
long

employee name string

string

② Show tables → show all the tables

→

CREATE TABLE TABLENAME AS SELECT * FROM TABLENAME;

③ Display → SELECT * FROM TABLENAME;

④ Insert ↓

INSERT INTO TABLE TABLENAME VALUES
(108, 'Lavesh'), (109, 'Shubham');

⑤ drop table (delete the table)

drop table tablename;

⑥ WHERE clause

SELECT * FROM TABLENAME WHERE

condition

gives base condition

Database did not find condition

⑦ ALTER STATEMENT

→ Rename Table, add column, remove col

* Rename Table

ALTER TABLE TABLENAME RENAME TO
NEWTABLENAME

* add columns

ALTER TABLE TABLENAME ADD COLUMNS
(newcol1 datatype, newcol2 datatype);

* ~~RENAME~~ Replace / Remove columns

ALTER TABLE TABLENAME REPLACE
COLUMNS (col1 datatype, col2 datatype);

for eg.:

employeeid int

employeename ~~int~~
 int



→ ALTER TABLE TABLENAME REPLACE
COLUMNS (employeename string)

→ after that table will be

employeename string

* Truncate vs delete	
Deletes	Deletes
Rows not	Rows as well
Schema	Schema

* Delete

DELETE from TABLENAME WHERE cond

* Truncate

↓
delete rows not schema

TRUNCATE TABLE TABLENAME

* ~~DROP~~ (Deletes entire table)

DROP TABLE TABLENAME