



BDA unit-1 - Hand written Notes

big data analytics (Jawaharlal Nehru Technological University, Anantapur)

UNIT - I

15/3/98

Introduction to Big Data

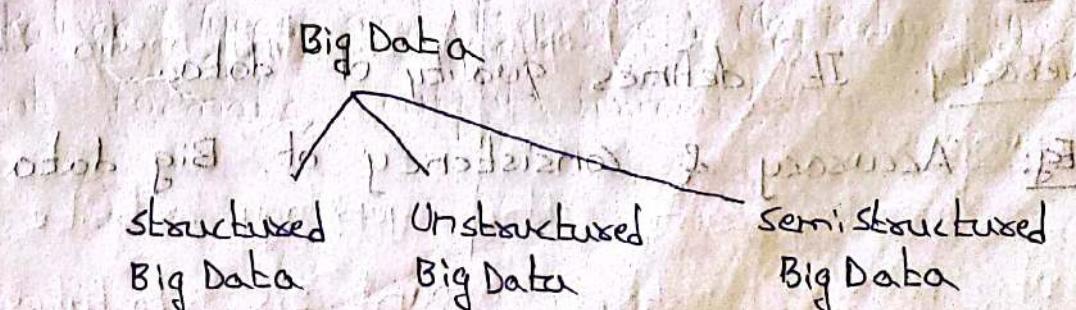
What is big data:

Data: Qualities, characters (or) symbols on which operations are performed by a computer which may be stored and transmitted in form of electrical signals & record on magnetic optical recording media like CD (or) pendrive is called data.

- Big data is also data with huge size, growing exponentially with time.
- It is collection of large datasets that can't be processed using traditional & current tools efficiently.

Examples of Big Data

- * Social Media - Facebook generates 500 TB + of data per day.
 - * NYSE (New York Stock Exchange) generates one TB of data per day.
 - * Google Maps
 - * Medical Industry
 - * Transportation
- There are 3 types of big data as follows:



→ Structured Big Data : A data that can be stored & processed in fixed format.

Eg: Employee table in database

→ Unstructured Big Data: A data that can be stored without any fixed format.

Eg: Google search (80% is U.D)

→ Semistructured Big data: A data with combination of both structured & unstructured Big data.

Eg: XML (Extensible Markup Language)

Characteristics of Big Data:

These are 4 characteristics called as 4 V's

- i) Velocity
- ii) Volume
- iii) Variety
- iv) Variety

Volume: Defines size (or) quantity of big data

Eg: bits < bytes < KB < MB < GigaBytes (GB) < TerraBytes (TB)
< PetaBytes (PB) < ExaBytes (EB) < ZettaBytes (ZB)
< YottaBytes (YB)

Velocity: It defines speed of big data.

Eg: Kbps, Mbps

Variety: It defines type of Big data.

Eg: Structured BD, Semistructured BD, Unstructured BD

Versatility: It defines quality of data.

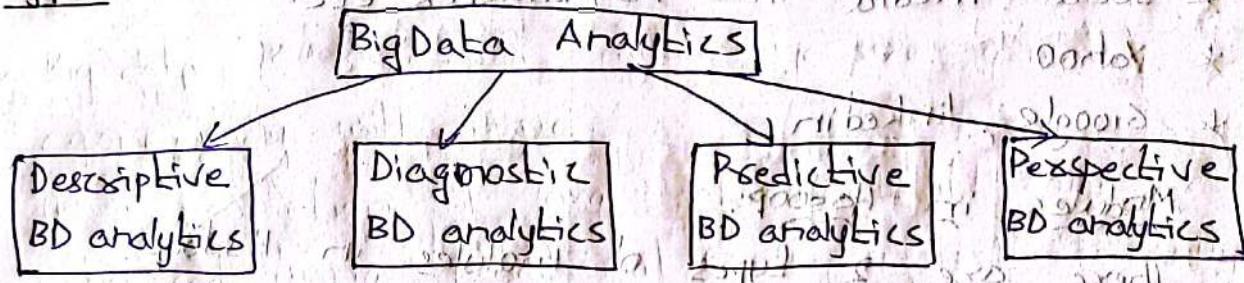
Eg: Accuracy & consistency of Big data

Designations for Big data specialisation:

1. Big data developer
2. Big data engineer
3. Hadoop developer
4. Big Data Analyst

Big Data Analytics: The process of analysis of large volumes of diverse data sets using advanced analytical techniques refers as Big Data analytics.

Types:



Descriptive analytics: It simplifies data & summarise past data into readable form.

Eg: Dow chemical company utilises its past data to increase facility utilization across offices & labs.

Diagnostic analytics: It gives a detailed & indepth of data & finds route cause of problem.

Predictive analytics: It makes use of historical (past) & present data to predict future data.

Prescriptive analytics: It allows business to determine best possible solution to problem.

Why Big Data importance

- | | |
|--|-----------------------------|
| * cost reduction | * Time saving |
| * Faster & better decision making | * Improved customer service |
| * Increased productivity (Performance) | * Fraud detection |
| | * Innovation |

Meet Hadoop:

- * It is an open source framework from apache.
- * It is used to store, process & analyze the data which is in very large in volume.
- * It is written in Java.
- * It is used for offline processing/batch processing.

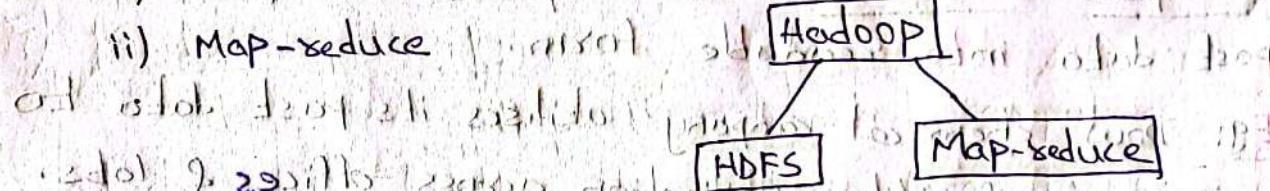
Applications:

- * Social media like fb, twitter etc.
- * Yahoo
- * Google, linkedin

Modules in Hadoop:

These are 2 types of modules in Hadoop:

Distribution File System (HDFS)



Advantages of Hadoop:

- * Resilient to failure
- * Fast to access data stored in distributed manner
- * Scalable
- * Cost effective

Data: Qualities, characters (or) symbols on which

operations performed by computer which may be

stored & transmitted in form of electric signals

record on magnetic optical recording media like CD.

Examples:

Ancestry.com, the genealogy website (study of families) stores around 12.5 petabytes per day.

→ Large Hadron Collider near Geneva, Switzerland
will produce about 15 petabytes of data per year.

- Data Storage & Analysis:
- * In 1990's, one drive would store 1370 MB data & had transfer speed of 4.4 MB/sec, within 5 min, all data can be sent.
 - * Over 20 yrs later, 1 TB drives are the norms with a speed of 100 MB/sec and takes 2½ hours to load all data.
 - * Since, it takes long time to read all data on single drive and written is even slower.
 - * To reduce time, we use multiple disks at once but it seems wasteful & cost effective.
 - * Hence we can store 100 data sets, each of size 1 TB and provide shared access to them with shared analysis types.

Problems in Data Analysis:

These are 2 problems:

- i) H/w failure,
sol: Replication
- ii) Combined data, solution is map reduced.

Hadoop provides - reliable data storage (HDFS), Efficient analysis (Mapreduce).

Comparison with other systems:

- i) RDBMS

Type of DBMS that stores data in row based table structure.

Eg: MySQL

This document is available free of charge on

Grid RDBMS

1) Data size in TB.

2) It is an interactive (online) & offline processing (batch).

3) It reads the data & writes many times.

4) static structure schema.

5) High integrity.

6) scaling is non linear.

MapReduce/Hadoop

Data size in Petta Bytes.

It is offline processing.

Data is write once, read many times.

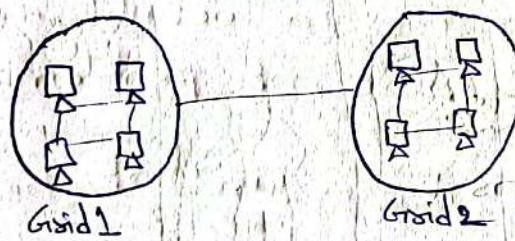
Dynamic structure schema.

less integrity.

scaling is linear.

[Grid Computing]

It is a group of networked computers which works together to perform large tasks.



Grid Computing

- 1) Data access is slow.
- 2) low level programmed.
- 3) Highly expensive.

Map Reduce

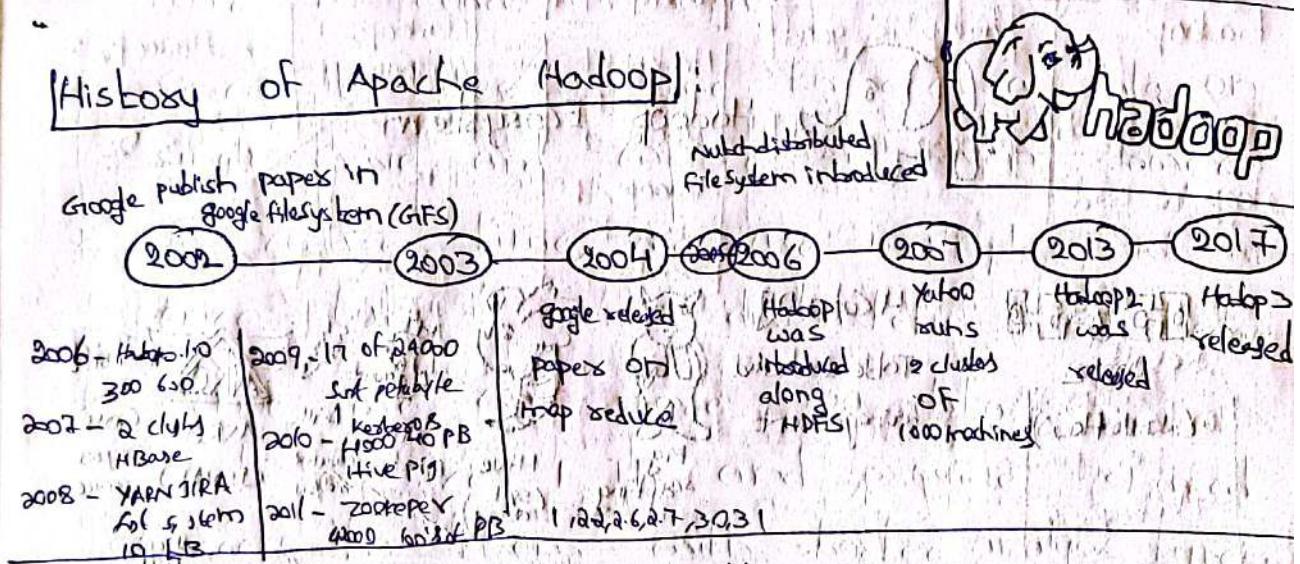
- 1) Data access is fast.
- 2) High level programmed.
- 3) less expensive.

iii) Volunteers (Computing),

It is a type of distributed computing where people donate their unused computers to perform large tasks.

Volunteer computing	Map Reduce
1) No data locality	with data locality
2) Multiple data centers	Single data center
3) Low bandwidth	High bandwidth

History of Apache Hadoop



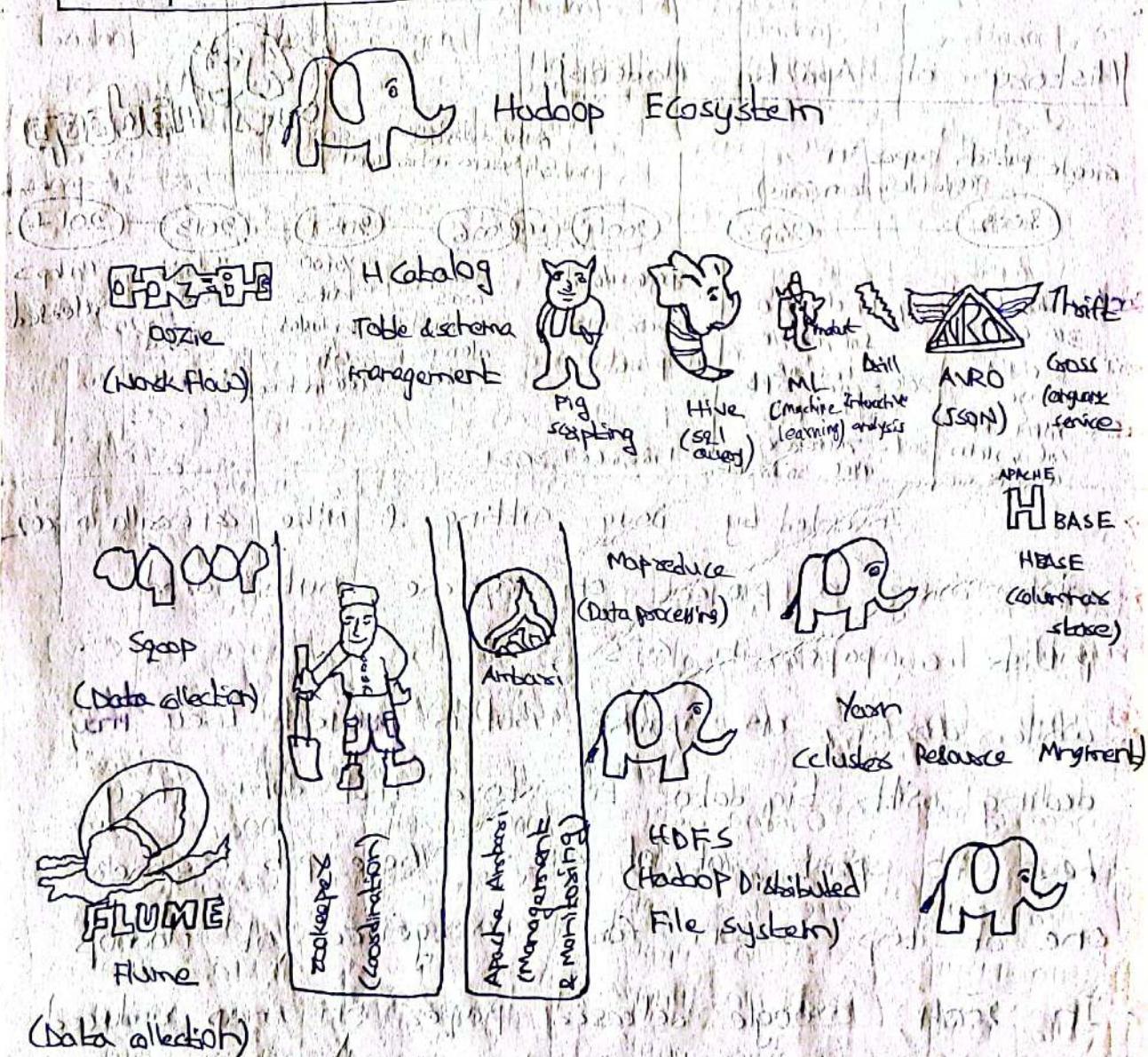
- Invented by Doug cutting & Mike Cafarella in 2002
 - while doing a project on apache hdfs and publish a paper in GFS
 - while working on apache hdfs project they may dealing with big data to store that data (they have to spent lot of cost). These prob becomes one of imp reason for emergence of Hadoop
 - In 2004, Google released paper on map, reduce
 - In 2005, NDFS was introduced
 - In 2006, Doug cutting quit Google and joined Yahoo and on basis of NDFS, he introduced HDFS & released Hadoop 1.0.

→ He gave named his project Hadoop after his son's toy elephant.

In 2008, Hadoop becomes fastest system to sort 1 TB of data on 900 hadoop clusters within 20 seconds!

→ In 2018, 3.1 released.

Hadoop ecosystem / Hadoop stack / Hadoop structure:



Yet Another Resource Negotiator

HDFS (Hadoop Distribution File System), which is used to store large data sets.

Map Reduce: It is a programming based data processor.

YARN: Yet Another Resource Negotiator used to manage the resources.

SPARK: Used for in memory data processing.

PIG, HIVE: Used for query based processing of data services.

Eg: MySQL

HBASE: Used for NoSQL database.

Mahout, Spark MLlib: Used for formal algorithm libraries.

ZooKeeper: Used to manage the clusters.

Oozie: Used for job scheduling.

Solr, UCB: Used for searching an index.

Analyzing the data with Hadoop:

Map reducing is a programming model to analyse and process data with Hadoop.

Hadoop can run MapReduce written in various language like Python, Java, C++.

Eg:

Map Reduce program works in 2 phases.

1. Map phase

2. Reduce phase

In Map phase it has 2 sub-phases: 1. I/P splitting

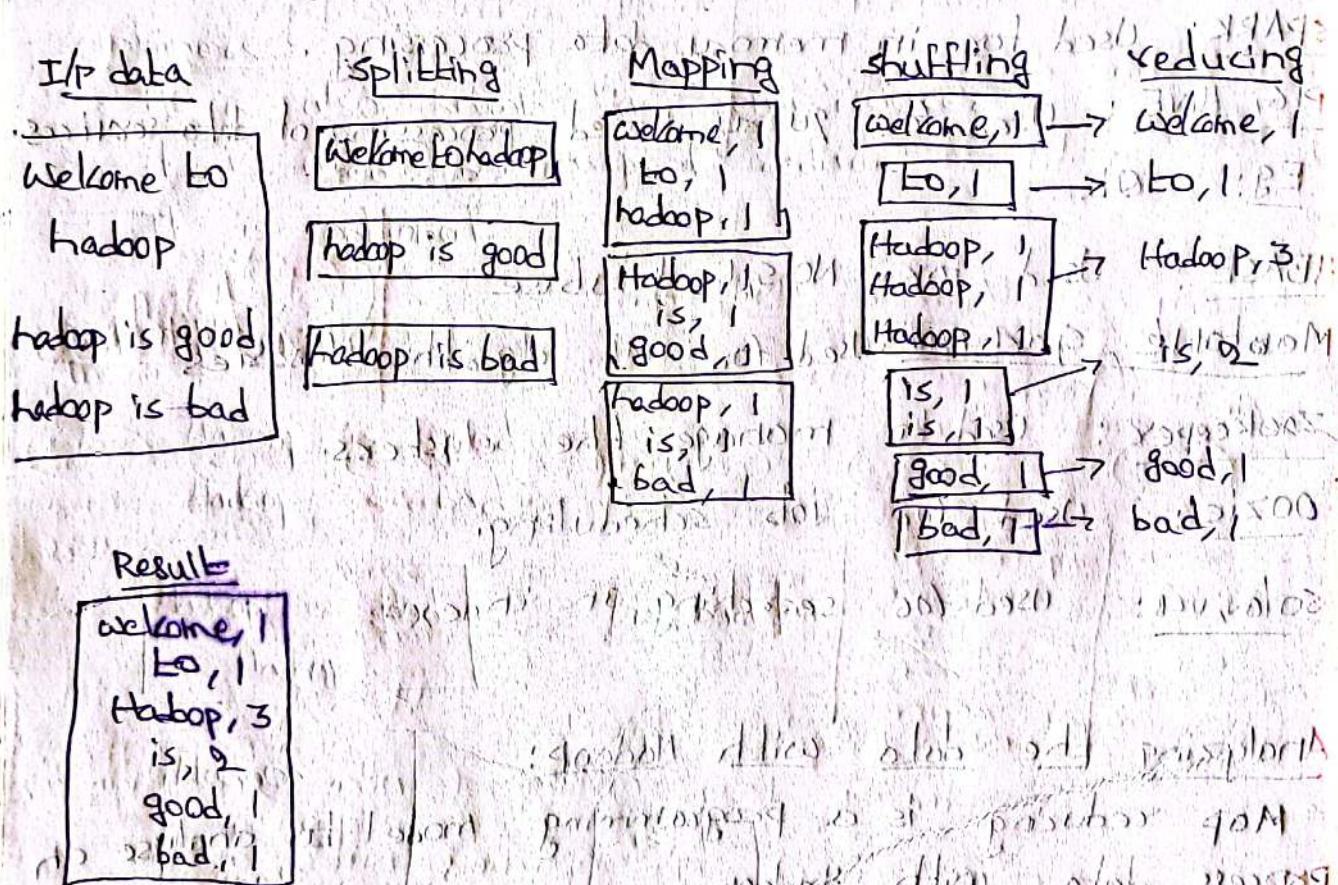
2. Mapping

→ I/P to map reduce job is divided into fixed size pieces called I/P splits.

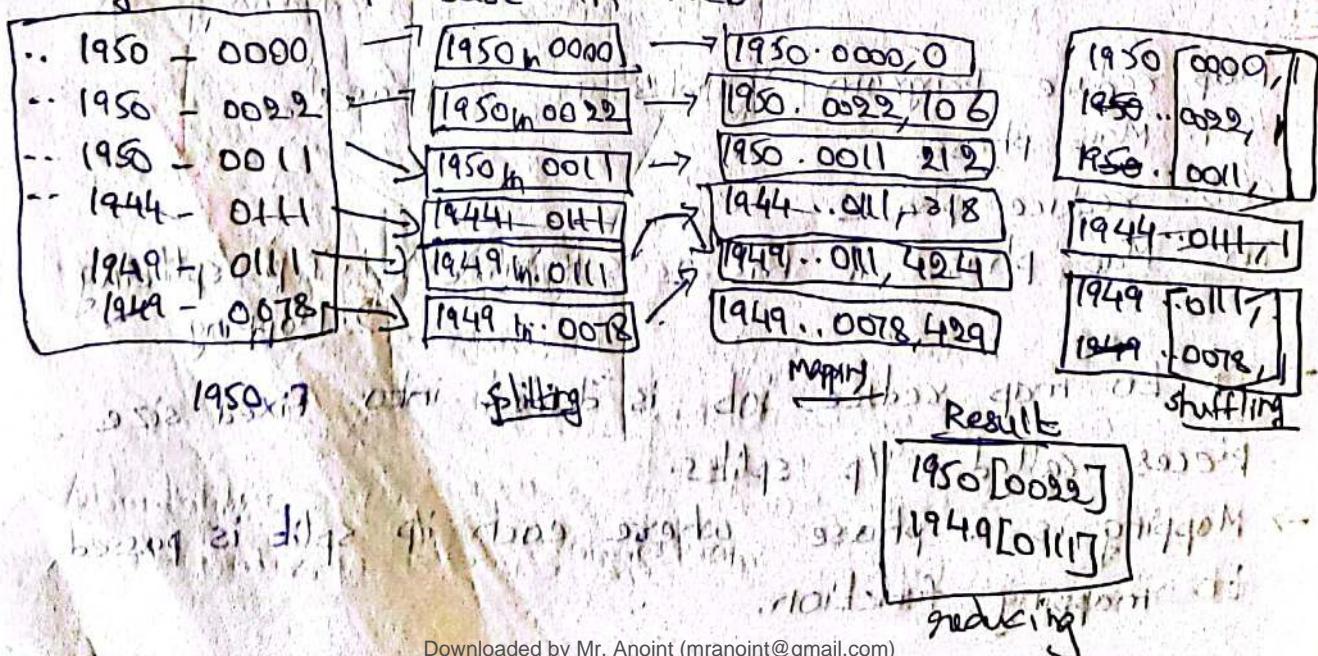
→ Mapping is 1st phase where each I/P split is passed to mapping function.

Eg: (word, frequency)

- Reduce phase is again subdivided into 2 phases.
- I - shuffling II - Reducing
- shuffling phase is used in to console data relevant parts!
- Reducing phase is used to aggregate data.



Ex:-2 :- weather climate set before, data format, weather data set from NCDC (National Climate Data Set). Find the highest temperature in NCD.



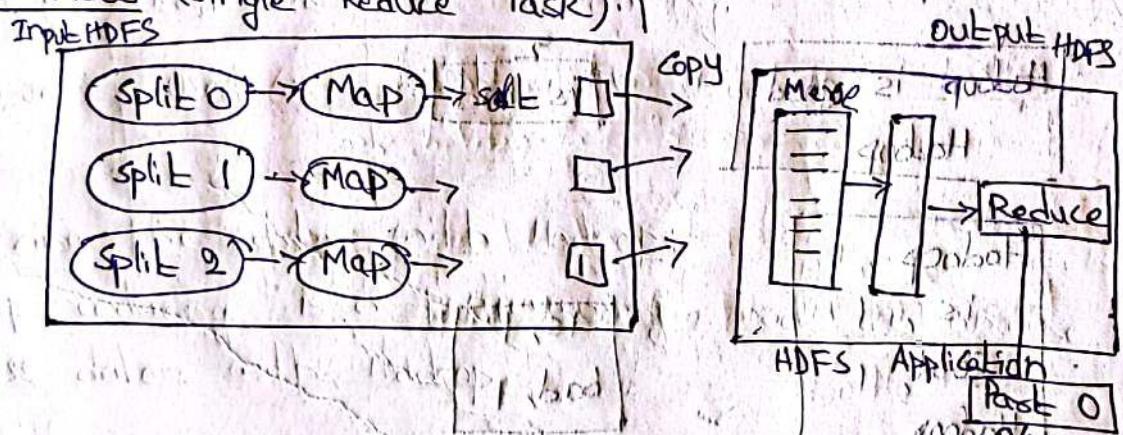
Scaling Out:

- To scaling out the Hadoop, we need to store the data in distributed File System. Typically HDFS.

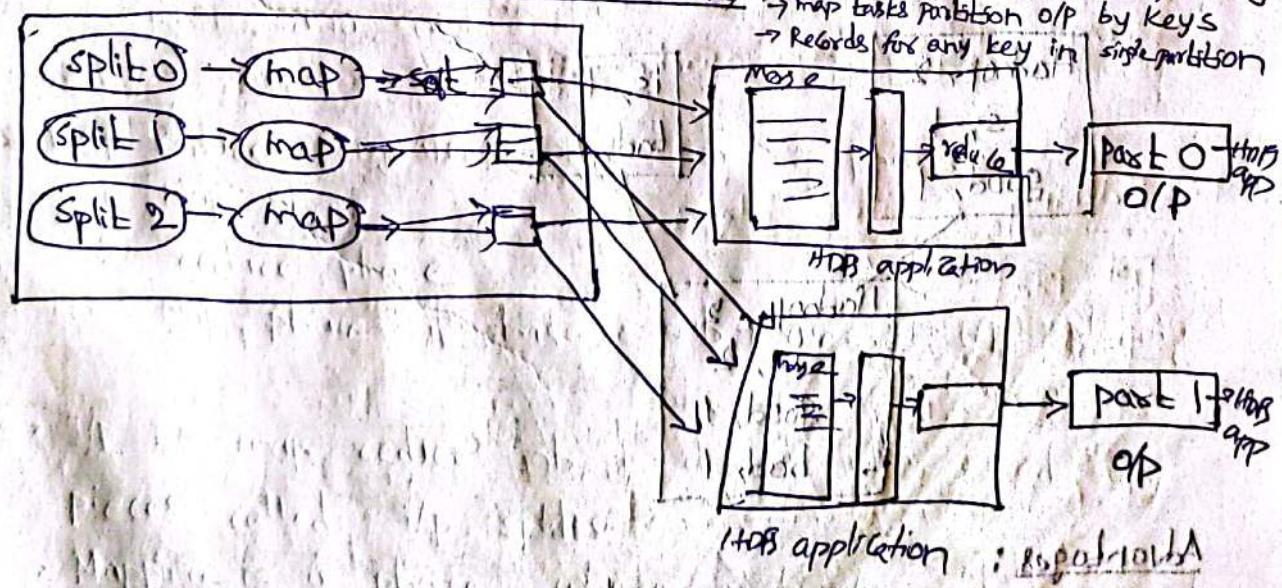
Data Flow in HDFS:

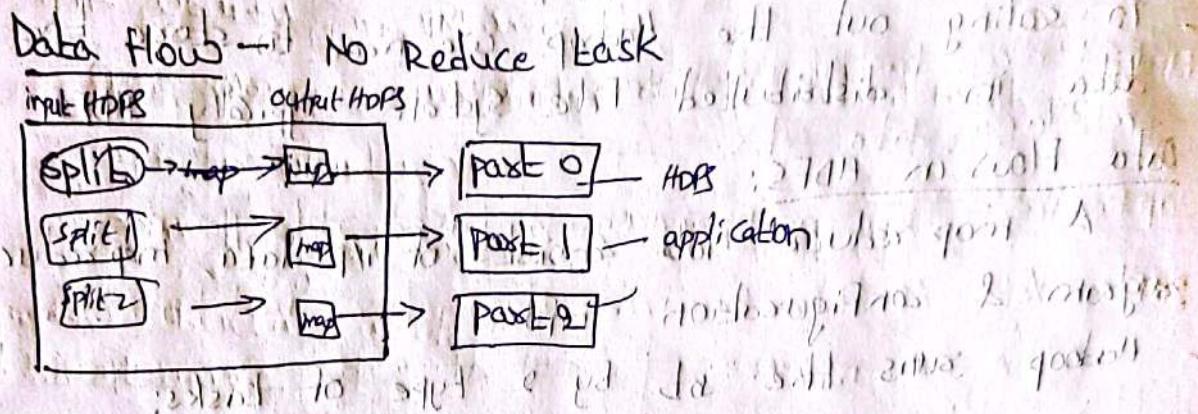
- A map reduce job consists of i/p data, map reduce program & configuration info.
- Hadoop runs the job by 2 type of tasks:
 - map task
 - reduce task
 - Hadoop contains 2 types of nodes
 - job tracker (single)
 - multiple task tracker (several)
 - Hadoop divides data into fixed pieces with size of 64 MB each.

Data Flow (single Reduce Task):



Data Flow (Multiple Reduce Task):

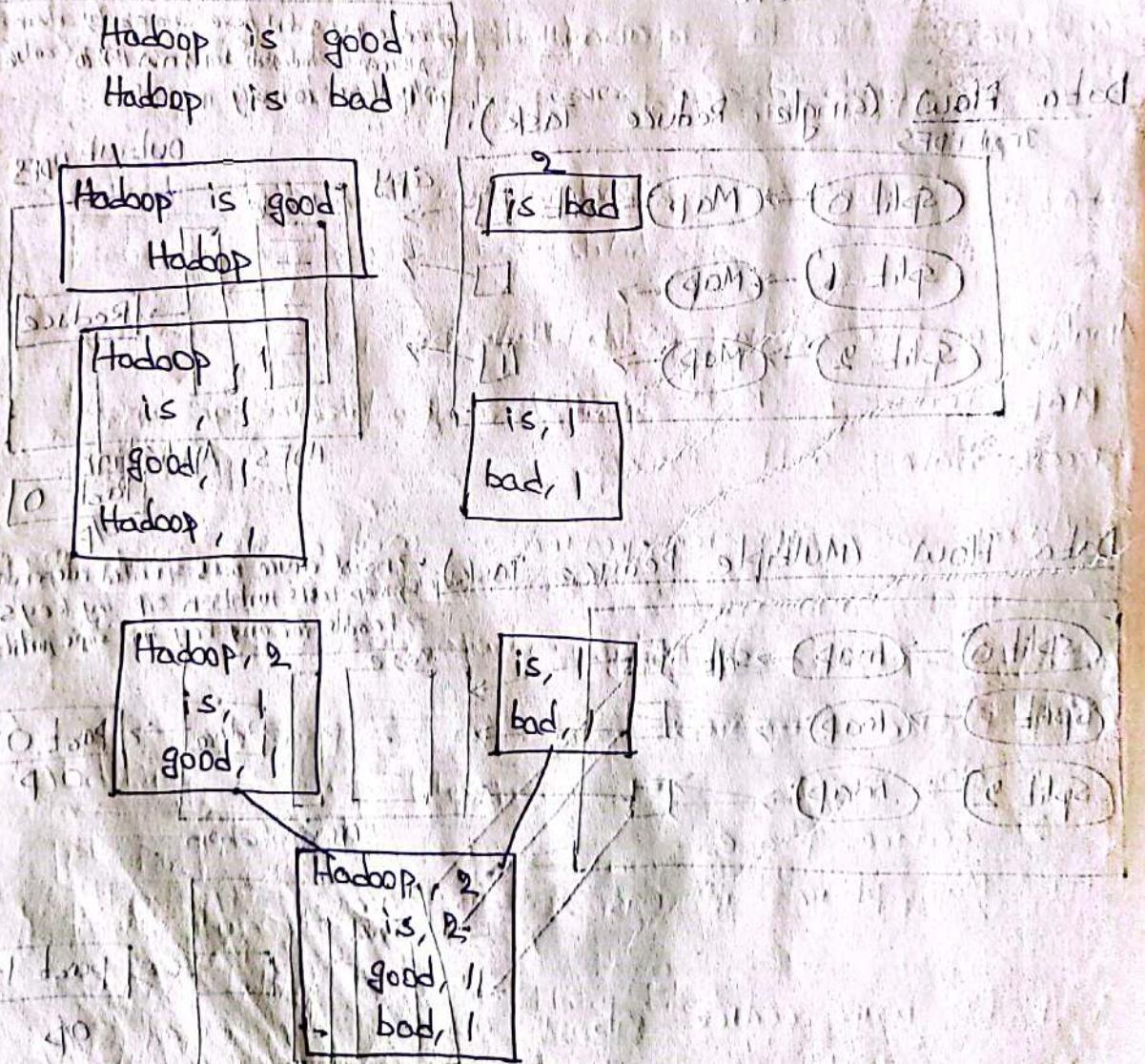




Combiner:

- It is also known as semi-reduces
- It is used w/o map & reduce
- It is optional

Data Flow:



Advantages:

Minimize N/w Congestion

Experiment 1: Install Apache Hadoop in VMWare

VMware installation for Hadoop

Aim: To install single node Hadoop cluster backed by the Hadoop Distributed File system on Ubuntu using VMWare

Different types of Linux OS:

raspbian, centOS, kali, ubuntu

Installing VMWare:

- i) Double click to launch VMWare-workstation-full-15 application.
- ii) Security warning panel & click on Run to continue.
- iii) Initial screen will appear, wait for processes to complete.
- iv) VMWare Workstation setup wizard open, click next.
- v) Select I accept terms within license agreement, click on next.
- vi) select directory during which you'd wish to install the appliance also select enhanced keyboard (drivers) checkbox.
- vii) leave it to defaults setting & click next.
- viii) Select both options desktop & start Menu programs folder & click next.
- ix) click install to start installation process.
- x) Installation in progress, wait for this to complete.

Install Ubuntu in VMWare:

- i) Open VMWare workstation & click on Create New VM.
- ii) Choose Installer disc image file (.iso): to make ws to detect that iso file is appropriate or not.
- iii) Fill info about full name, user name & PW.
- iv) Then click next & give your VM relevant name.
- v) Allocate size of Hard disk.
- vi) Run VM
- vii) Install Ubuntu 20.04 LTS Desktop.
- viii) To begin installation, click install Ubuntu.
- ix) choose keyboard layout
- x) This will take while to complete

3. Run Ubuntu in VMWare

Open VMWare and run Ubuntu 64 bit virtual machine

→ 1. login to Ubuntu OS

→ open terminal or use "ctrl+alt+t"

4 - check Ubuntu is updated

\$ sudo apt update

5. Installing Java:

\$ sudo apt install openjdk-8-jdk-4

→ checking Java version

\$ java -version

\$ javac -version

6. Installing SSH (Security Shell)

\$ sudo apt install openssh-server openssh-client

7. Create & setup SSH certificates:

\$ ssh-keygen -E rsa -P "" -f ~/.ssh/id_rsa

\$ cat ~/.ssh/id_rsa.pub > ~/.ssh/authorized_keys

\$ chmod 0600 ~/.ssh/authorized_keys

\$ ssh localhost

8. Downloading Hadoop:

\$ wget https://dlcdn.apache.org/hadoop/common/hadoop-3.3.1/hadoop-3.3.1.tar.gz

9. Editing 6/ imp files:

1st file

\$ sudo nano .bashrc

Add below lines in this file

Hadoop related imports

```
export HADOOP_HOME = /home/bigdata/hadoop-3.3.1  
export HADOOP_INSTALL = $HADOOP_HOME  
export HADOOP_MAPRED_HOME = $HADOOP_HOME  
export HADOOP_COMMON_NAME = $HADOOP_HOME  
export HADOOP_HDFS_HOME = $HADOOP_HOME  
export YARN_HOME = $HADOOP_HOME  
export HADOOP_COMMON_LIB_NATIVE_DIR = $HADOOP_HOME/lib/native  
export PATH = $PATH:$HADOOP_HOME/sbin:$HADOOP_HOME/bin  
export HADOOP_OPTS=-Djava.library.path=$HADOOP_HOME/lib/hadoop/
```

Now type:

```
$ source ~/.bashrc
```

2nd file:

```
$ sudo nano $HADOOP_HOME/etc/bigdata/hadoop-env.sh
```

Add below line in this file in end

```
$ sudo nano $HADOOP_HOME/etc/bigdata/hadoop-env.sh
```

```
$ export JAVA_HOME=/usr/lib/jvm/java-8-openjdk-amd64
```

(3rd file):

```
$ sudo nano
```

10 Marks

1. What is big data? Explain its details.
2. Why big data is important? Discuss in detail.
3. Explain (i) data (ii) Data storage & analysis.
4. Distinguish hadoop with other systems.
5. Elaborate history of apache hadoop.
6. Explain hadoop ecosystem.
7. Discuss about VMware.
8. List steps in hadoop installation.
9. How can you analyze data with hadoop?
10. Discuss scaling out in hadoop.