

## Session - 2

- Basics of storage.
- Code in ~~hd~~<sup>hard disk</sup> is known as program.
- But when it is in running state then, it is called process.
- Ctrl + c ≈ console ≈ Terminal
- Ctrl + c vs Ctrl + z
  - ↳ Terminates ↳ pauses the prog in bg.
- # fg 2 → forward ground (resume the pro)
  - ↳ job number
- # firefox &
  - ↳ run the cmd in bg.
- # touch file.txt
  - ↳ creates an empty file
- # mkdir file
  - ↳ creates a directory
- # ls ~~ll~~ → long list
  - ↳ field switch
- # whoami
  - ↳ displays the current user
- # useradd gaurav
  - ↳ adds a user
- # passwd gaurav
  - ↳ adds the password for the user
- # Ctrl + Alt + F1 / F2 / F3
  - ↳ changes the desktop/screen.

# vi my.txt

↳ esc → cmd / control mode

↳  $i \rightarrow$  insert mode

→ esc ;

↳ :  $w \rightarrow$  written or same

↳ :q → exit the file

↳ :wq → saves and exits the file.

# df -h → converts bits to Mbs.

↳ maximum limit of data we can store.  
→ Shows h.d.d

## Big Data → Major problem

## Volumes & sub-problem

Velocity  $\rightarrow$  cost high

$B_{\text{IT}} \rightarrow B_{\text{Yttr}} \rightarrow kB \rightarrow MB \rightarrow GB \rightarrow PB \rightarrow EB \rightarrow ZB \rightarrow VB$

← Brontobryas ↑

- If we create hard disk of size 4B issue will come of input / output (new problems will emerge)

# free - m

↳ shows memory (ram)

# watch date

→ runs the cmd every 2 sec for optim  
watch -n 1 file -m

# vi

↳ esc

↳ yy → copies whole line

↳ p → paste data

↳ np → pastes n times.

- concept

- ↳ distributed storage

- ↳ split data

Servers  
↓  
OS  
↓  
[ Ram / CPU / HD ]

{ More powerful  
than set of PCs.

• Cost ↑.

• Reliable ↑.

Name pc hardware → commodity

- Master : Receiving the storage.

- Slave : & contributing the storage.

Master - Slave relationship.

- All are connected through network. NameNode (NN)

Topology  
(center)



• similarly we can connect RAM, CPUs, etc  
↳ distributed computing

**HADOOP** → software

↳ solve the problem of  
big data through the  
concept of distributed  
storage.

\* self-work

doc / Article / Post / Blog  
regarding big data.

↳ Q9. solved by hadoop.

↳ what is hadoop  
what does hadoop do  
what is hadoop used for  
what is hadoop + yester  
what is hadoop + future

↳ what is hadoop  
what is hadoop + future

↳ what is hadoop  
what is hadoop + future

## SESSION - 3

- Prompt is known as shell prompt.

# vim my.txt

# cat my.txt → used to read a file.

# source my.txt

↳ runs all the cmd in my.txt.

- Bash shell Scripting

# echo hello → prints the text.

# espeak -g hello → speaks the text.

TASK : Find a cmd which will send the data to whatsapp, sms, email and print the data.

# echo `date` or # echo \$(date)

special symbol that will differentiate bt cmd and string  
(Shift + console / back quote)

# rpm -qf /usr/lib/date

↳ query file

↳ displays the name of the software.

# man cmd

↳ opens the manual for the cmd.

TASK: espeak -ng ; Find the cmd to change the language of the cmd.

# rpm -a -f 'which firefox'

# rpm -e firefox  
↳ erase

# init 0  
↳ shut downs the OS.

# rpm -q hadoop

- similar like Hadoop  
↳ S3, Ceph, Ganesha
- HDFS protocol is used in Hadoop  
Master  $\xrightarrow{\text{HDFS}}$  Slave  
↳ Hadoop Distributed File System

# java -version

- winscp - program used to transfer doc from windows to linux

• Software

Drag & drop.

in ~~Recycle~~ treat.

(i) hadoop

(ii) JDK.

# rpm -i ~~jk~~ jdk<TAB>

# java -v (Hotspot) ORACLE

# rpm -i hadoop <TAB> ← fails

# rpm -i hadoop --force.

# rpm -q hadoop

# hadoop \*version || check

## SESSION - 4

- Basics of Binary
  - ASCII - basic knowledge
  - In networking we use bits not bytes
  - ASTERISK - software used to create our own conference servers.
  - SPoF - single point of failure.
  - Horizontal scaling & Vertical scaling
    - (SPoF)
    - (SPoF)
      - ↳ replica
  - Valid ip address.

# jps → does the system is master or slave

# cd /etc/hadoop/  
↳ configuration file of hadoop.

# Vim hdf5-site.xml

Master <property>  
<name>dfs.name.dir</name>  
<value>/nn</value>  
</property>

# mkdir mn ~~(not a dir)~~ || in 1 dir.

# rpm -i -v jdk-1.8.0\_121.rpm  
↳ verbose (shows all the bg process).

Slave =>

# cd /etc/hadoop/

# vim hdfs-site.xml

// Same as master

' <name> ~~dfs~~ dfs.data.dir </name>  
~~dfs~~

b

2022-08-20 20:20:45

07-08-2022 20:20:45

2022-08-20

2022-08-20

2022-08-20 20:20:45

2022-08-20 20:20:45

2022-08-20 20:20:45

2022-08-20 20:20:45

2022-08-20 20:20:45

2022-08-20 20:20:45

2022-08-20 20:20:45

2022-08-20 20:20:45

2022-08-20 20:20:45

2022-08-20 20:20:45

2022-08-20 20:20:45

2022-08-20 20:20:45

2022-08-20 20:20:45

## SESSION - 5

DATE: 18/Sept  
PAGE NO.:

- IP address consists of 4 octets, each of 8 bits. 32 bits/ 4 bytes in total.
- ifconfig in Linux  
ipconfig in Win.
- MAC address - 6 bytes.

# ifconfig enp0s3 4.3.2.1

↳ changes the ip address of a system

# ping www.google.com

↳ checks the connectivity.

# nslookup www.google.com

↳ gives the ip of the web site.

• DNS - domain name system

↳ gives the name-to-ip addresses.

# bc

↳ binary calculator

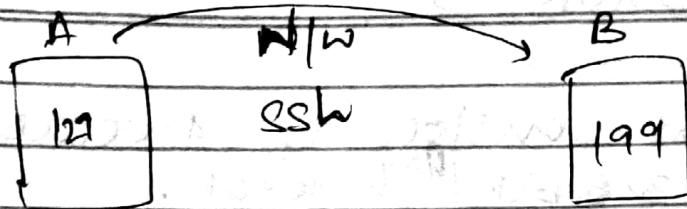
base = 10 ↳ input base

obase = 2 ↳ output base

• protocol == process

win ↳ RDP (Remote desktop protocol)

linux ↳ ssh (secure shell)



# \$ ssh 192.168.0.199

# init 6

↳ reboots the system

- For android - linux connection

↳ Juice Ssh - (app)

# putty

↳ connects windows to linux using ssh.

# ping 192.168.0.129

↳ from slave to master to check the connectivity.

- we never connect with the ip address. Through we can check the connectivity. To connect we use some external software.

~~Master~~ windows slave

Master

# jps

↳ checks the process of master node.

# vim core-site.xml

```

<property> → fs.default.name
<name> fs.name.default </name>
<value> local </value>
</property> ← local hdfs://
                  → hdfs://192.168.0.129:9001
  
```

# netstat -tnlp

↳ checks how many ports are running  
in the linux system

\* 909/cupsd is used by printer

# hadoop namenode -format

- process == service == daemon

# hadoop-daemon.sh start nodename

# jps // check the conn.

# hadoop dfsadmin -report

↳ shows how many slaves are connected.

• the first number probably is

slave.

# jps

# vim core-site.xml

```

<properties>
  <name>fs.default.name</name>
  <value>hdfs://192.168.0.129:9001</value>
</properties>
  ↑MASTER (ip)
  ↗ SAME IN ALL MASTER & SLAVE.
  
```

# jps

# hadoop-daemon.sh start datanode

The service will not start.

master

# systemctl stop firewalld.

slave

# hadoop-daemon.sh start datanode

# Jps . WORKS!!!

master

# hadoop dfsadmin -report.

AWS cloud  $\rightleftarrows$  Hadoop

- Launched an instance - Redhat.
- Use putty to interact with the OS.
- change .pem to .ppk using puttygen
- Run putty
- Upload JDK & Hadoop via <sup>winscp</sup>  
 → (winscp) SCP protocol used to move files  
 → ssh - remote login (putty)
- Connecting private key in winscp.  
 Advanced... → Advanced... → SSH(Auth)
- OOM - Out of Memory  
 Hadoop is heavy if not provided with sufficient memory.

# cat /proc/sys/vm/drop\_caches  
 $\xrightarrow{0}$  each no. has its own meaning.  
# echo 3 > /proc/sys/vm/drop\_caches  
# echo 3 > /proc/sys/vm/drop\_caches  
 $\xrightarrow{\text{cleans the cache mem.}}$

- Configure hdfs-site.xml

some change in core-site.xml

slave

<property>

<name> fs.default.name </name>

~~localhost~~

<value> hdfs://11.15.207.72:9001</value>

<property>

↓ public ip of master

only  
Master

value> hdfs://~~11.15.207.72~~:9001

~~localhost~~ open  
ip

EC2 → pub ip  
→ pr ip

Master  
only #

hadoop namenode -format

↳ format the disk

# jps

Master  
only #

hadoop-daemon.sh start namenode

↳ start the master node

# hadoop ~~dfsadmin~~ -report

b) dfsadmin

## SESSION - 9

- when we stop & start the instance, the ip address changes.
- logical storage != physical storage
- logical storage ≈ virtual storage
- Hadoop 1.2.1 ; by default size of strip is ~~is~~ 64 Mb
- # ip of master: 50070 ← url  
↳ shows the report in web.
- configured client.

- In client we do not have to start the service through home to be configured in order to ~~be~~ connect to the master node.

Master

# vi core-site.xml ↳ to ch hdfs-site.xml

```

<property>
<name>fs.default.name </name>
<value> hdfs://192.0.31.143:9001 </value>
</property>           ↳ master pub ip

```

client# hadoop fs -ls /

# hadoop fs -put lw.txt  
copy file name directory

# hadoop fs -rm /zzz.txt

# hadoop fs -cat /zzz.txt  
no space

# hadoop fs -D dfs.block.size = ~~33554432~~ 33554432  
- put zzz.txt  
changes the block size

• It does not support mb, only Bytes is supported. we can change block size per size

# ~~yum~~ install bc

# rpm -q bc

# bc

$32 \times 1024$

$32 \times 1024 \times 1024 \approx 32M$  in bytes.

## SESSION -10

#! hadoop-daemon.sh start datanode/nodemaster  
 ↳ starts the service temporarily.

- Port no: 50070 → web ui master
- Port no: 9001 → service port of hadoop 1.

# cd /etc/rc.d/  
 # cat rc.local

13.233.51.

~~# open vim~~

13.233.254.248

# yum whatprovides vim

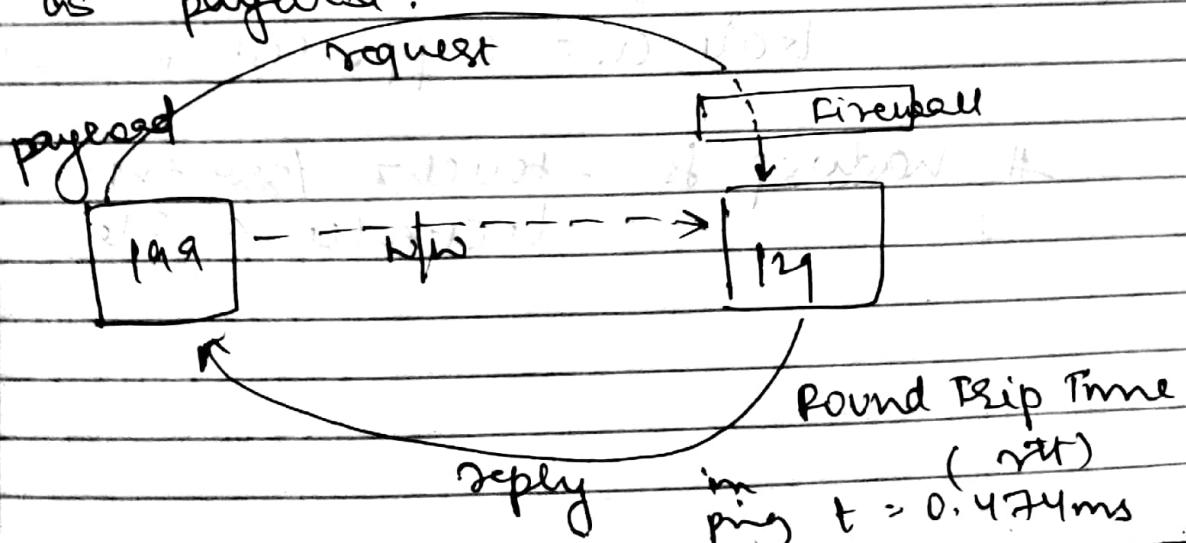
# vim rc.local

hadoop-daemon.sh start namenode

# chmod +x /etc/rc.d/rc.local

- For client we only need to update the core-site.xml file.

- Information sent on network is known as payload.



64 bytes is the size of payload attached with the packet.

• ICMP : Internet Control messaging Protocol

\* • ICMP Tunnel

• PING : Packet Internet Grouping

# ping -c 2 192.168.0.129

↳ command (it will stop after 2 pings)

# alias vimal="date"

↳ gives a nickname to cmd.

# useradd -s /usr/bin/cal tom

↳ tom is restricted only to cal cmd.

# hadoop fs

↳ available and supported in hadoop

# hadoop fs - touchz 1889.txt

↳ creates a file

## SESSION - 12

- System (CS) crashes after year 2038.
  - ↳ proof ; date cmd.
  - ↳ memory overflow
- When more space is used than required than given memory then the extra data will automatically discarded by the computer.

Dec 13, 1901 15:46:00  
— 0 sec

- Time is stored in seconds.
- Jan, 18. 2028 22:14:07
  - All 32 bits are full
  - ↳ last date
- Dec, 13, 1901 15:46:00
  - ↳ resets
- Difference bet 32bit & 64bit OS.

- Master (Hadoop) ; SG → Rule in inbound traffic  
Type - All
  - ↳ client can do all
  - ↳ ssh
  - ↳ ping
  - ↳ seq, all.

# vim core-site.xml

<value> hdfs://ip(master):9001 </value>

- We only have one public ip but they can have ~~three~~ dif. private ips.

# rm -rf /dn

↳ deletes the file.

# mkdir /dn

- Delete & recreate the folder (Quick ~~fix~~ troubleshooting)

- If the master mode is in safe mode we can't make any changes in it.

Master # hadoop dfsadmin -safemode get  
↳ displays the mode.

" # hadoop dfsadmin -safemode leave.  
↳ changes the state (Safe mode is off)

20212821

Hadoop has a replication factor of 3.

- Block size : Slices the file (64MB def.).
- Replication : Makes copy of file (3 by def.).

Task : Create a cluster

6 slaves

1 client

When you upload a file

Installation of Py environment.

(Install Python 3.8.5)

(Install pip 20.0.2)

(Install numpy 1.19.2)

(Install pandas 1.1.3)

(Install matplotlib 3.3.2)

(Install scikit-learn 0.24.2)

(Install tensorflow 2.4.1)

(Install pyarrow 1.0.1)

(Install opencv 4.5.1)

(Install keras 2.4.3)

(Install tensorflow-gpu 2.4.1)

(Install tensorflow 2.4.1)

## SESSION - 14

- Basics of networking

# tcpdump -i enp0s3 -n

↳ lists all the ICMP packets

↳ monitor/ packet capturing

# pgrep firefox

↳ shows PID (process ID) of the program.

- clear text : Text anyone can read  
↳ has no encryption.

# tcpdump -i enp0s3 icmp -n

↳ filters the icmp packages.

- But ssh will be failed.

Set. cmd

# tcpdump -i enp0s3 tcp port 22 -n

↳ inside step; port : 22

(SSH).

- Basics of Python (Introduction)

↳ crossplatform language.

SOO10 → data transfer port.

## SESSION - 15

- Only & Only Client decides the block size & replication of a file.

as well as client.

- Hardware address == Mac Address

Master #

vi hdfs-site.xml

<property>

<name>dfs.replication </name>

<value> 2 </value>

</property>

↳ If it fails. Replication factor is still 3.

client #

vi hdfs-site.xml <name> dfs.replication

<value> 4 </value>

</property>

client #

vi hdfs-site.xml

<name> dfs.block.size </name>

<value> 1024 </value>

↳ changes the block size

- Replication & Size & Memory

- REPL does not save the file

# edit lw.py

# python3 lw.py

• IDE - Integrated Development Environment

- Jupyter Notebook

- Mostly in DS & ML we use Jupyter notebook

- alt + enter

shift + enter : Runs the prog.

~~ctrl + enter~~

esc + dd : deletes the selected cell

- Basics of data type in Py

- int

string

bool

list

tuple

- : - slice operator

- list can only perform row-wise operation.

- Array can perform column-wise operation.

- programming file in which the code is written is known as module.

```
#import numpy
numpy.array(name)
```

- columns are accessed by (, )

#  
#put setaf 3      1 → bl 7 → white  
# changes the color of terminal

- If we want to print the content exactly we use triple quotes (" ").

- If - else loop

```
if int(exp) == 1:
    print("no job")
elif else int(exp) == 1:
    print("job")
else:
```

- Python does not support switch-case.

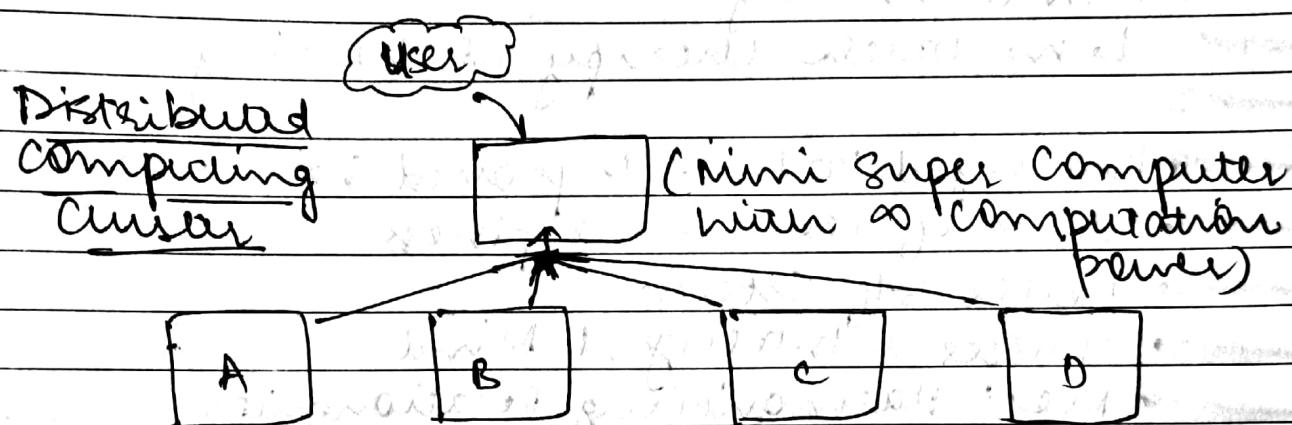
## SESSION: 143 Hadoop

## ★ Problem

↳ Large RAM (Same PB)

↳ No SPOF

This all can be solved by Big Data Concept (Super Computer)



- Sharing of Computing (CPU, RAM)

- Storage → HDFS (Hadoop)

- S3 (Cloud)

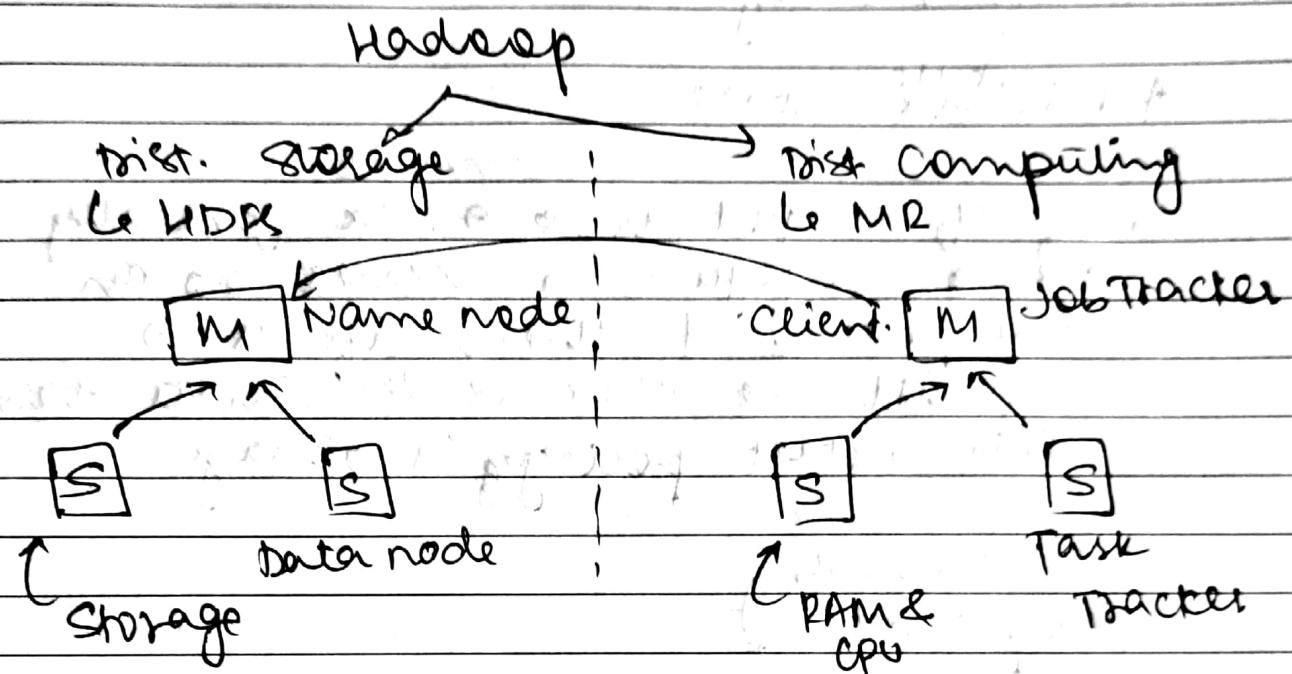
- Computing → MR Cluster (Map Reduced cluster)

- In this setup, we create some chunks/blocks of data and send it to the memory of all the slaves/nodes.

## ★ Product for distributed computing:

Hadoop (Map Reduced Cluster)

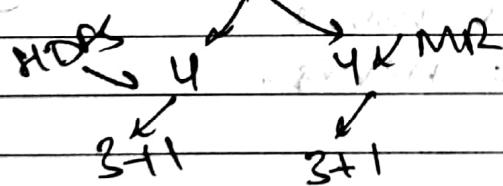
↳ Only for distributed storage but we can take its internal concepts & with some changes we can enhance it -



• we need HDFS cluster for MR because,  
for computing we need data.  
The large data is stored in the distributed storage.

∴ Min 9 system we require

1. 1 Name node  
2. 4 Data nodes  
3. 1 Job tracker  
4. 1 Task tracker

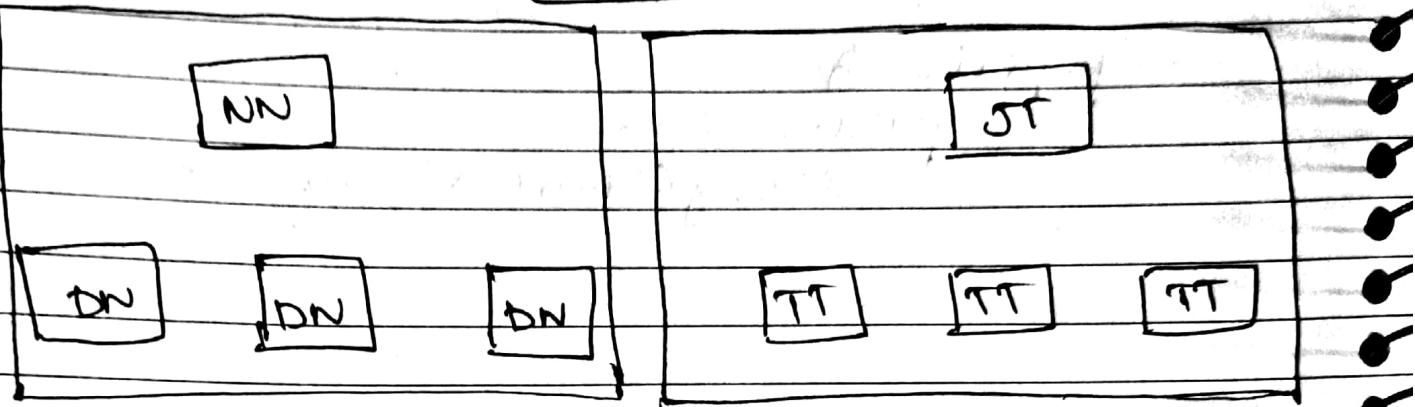


SESSION: 145

Hadoop

Distributed Computing Cluster

Client



Requirements: JDK & Hadoop software  
in every node.

(ec2 instance)

- Create one hadoop-template, and winscp JDK & Hadoop.

# cd /home/ec2-user/

# ls

↳ JDK &amp; hadoop

# rpm -ivh <jdk & hadoop>,  
<hadoop> --force

# init 0

- Now we will clone the hadoop-template

# Actions&gt;Image &amp; Template&gt;

Create a Image

↳ Image name: hadoop

&gt;&gt; Create Image

- Create 9 clones of hadoop-template AMI.
  - ↳ Allow all traffic
  - ↳ Public IP enabled
  - ↳ Security Group: allow-all-port-for-hadoop
    - # ec2 create-security-group --group-name allow-all-port-for-hadoop
    - ↳ Security group rule
    - ↳ Name: allow-all-port-for-hadoop
    - ↳ VPC: default
    - ↳ Add rule
      - ↳ All traffic to anywhere
      - ↳ All traffic anywhere

Create SG.

- ↳ SG: allow all traffic
  - Launch Instances.

~~NN~~

```

# sudo su -
# cd /etc/hadoop/
# vim hdfs-site.xml
<property>
  <name>dfs.name.dir</name>
  <value>/nn</value>
</property>
# mkdir nn # mkdir //nn
# hadoop namenode -format
# hadoop-daemon.sh start namenode
  
```

# vim core-site.xml

```
<property>
<name>fs.default.name</name>
<value>hdfs://0.0.0.0:9001</value>
</property>
```

# start service

```
# netstat -tulp | grep java
↳ 50070
↳ 9001
```

~~DN~~ # sudo su -

```
# cd /etc/hadoop/
# vim hdfs-site.xml
<property>
<name>dfs.data.dir</name>
<value>/dn1</value>
</property>
```

\* vim core-site.xml

```
<property>
<name>fs.default.name</name>
<value>hdfs://<master_ip>:9001</value>
</property>
```

• we can also put port ip  
bc all are in same ip

# mkdir /dfs

# hadoop-daemon.sh start datanode

~~# hadoop dfsadmin -report~~

## JOB TRACKER

~~# sudo su -~~

# cd /etc/hadoop/

# vim mapred-site.xml

<property>

<name> mapred.job.tracker </name>

<value> 0.0.0.0:9002 </value>

</property>

← No protocol used

# hadoop-daemon.sh start jobtracker

# netstat -trep | grep java

l : 9002

l : 50030 ← Map reduced cluster

url = ip : 50030/jobtracker.jsp

~~# sudo su -~~

# cd /etc/hadoop/

# vim mapred-site.xml

<property> job <!-- job will also be same -->

<name> mapred.\*\*.tracker </name>

<value> ip of ST:9002 </value>

</property>

# hadoop daemon.sh start tasktracker

Starting tasktracker on port 50030

1377415826

tasktracker starting up on port 50030

- Data is stored in Data Lake.  
  ↑  
  hdfs / s3
- Hadoop is the produce that come under the umbrella of data engineer.

~~nn - hadoop~~

- Security → Inbound rule.
  - ↳ Delete all inbound traffic rules  
~~we don't want~~
  - ↳ Add: all traffic sg-0d7f  
  ↑ (Launch wizard: 47)  
  security gp
- Now all the instances with this sg will be able to send the traffic to namenode.
- Now we will attach this security group (Launch wizard 47) to all the nodes and change their inbound traffic to sg47.  
(same as nn.)
- \* we can also change all these sg from by going to the security group.

Configure 1nn, 2dn, 1gt, 2tt

# hadoop job -list-active-trackers  
↳ lists all the active task trackers.

# vim /etc/hadoop/core-site.xml  
↳ we can make this site permanent by putting the start cmd here.

\* and make it executable chmod +x

~~client~~ # cd /etc/hadoop  
# vim core-site.xml

↳ same as dn.  
<name> fs.default.name </name>  
<value> hdfs://<ip>:9001 </value>  
</property>

# cd

# vim twitter.txt

# hadoop fs -put twitter.txt /

# hadoop fs -ls /

# hadoop fs -cat /twitter.txt

# vim /etc/hadoop/mapred-site.xml  
same as jobtracker.

153

# hadoop job -list

↳ NO jobs running

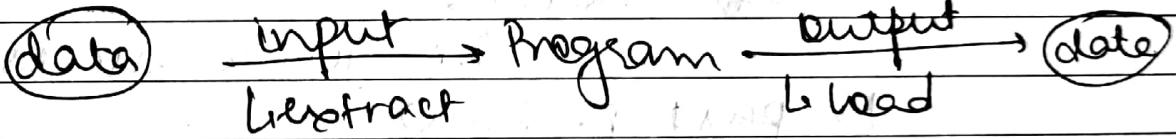
# cd /usr/share/hadoop/examples

↳ hadoop-examples-1.2.1.jar

# cd ..

# hadoop jar /usr/share/hadoop/hadoop-examples.jar wordcount

for job | /twitter.txt | output.



• In hadoop map reduce they always read the data from hdfs cluster storage.

• We also need to tell job tracker where the name node ~~is~~ is.

jt # cd etc/hadoop/

# vim core-site.xml

~~host~~

<same as dns>

~~check~~ ~~hadoop job~~

# hadoop-daemon.sh start jobtracker

At \_\_\_\_\_

~~client~~

# hadoop fs -ls /  
↳ new /tmp dir is created

# hadoop jar /hadoop-example  
~~wordcount~~ /twitter.txt /output

# hadoop fs -ls /  
↳ /output

# hadoop fs -ls /output  
↳ ~~tmp~~ /SUCCESS

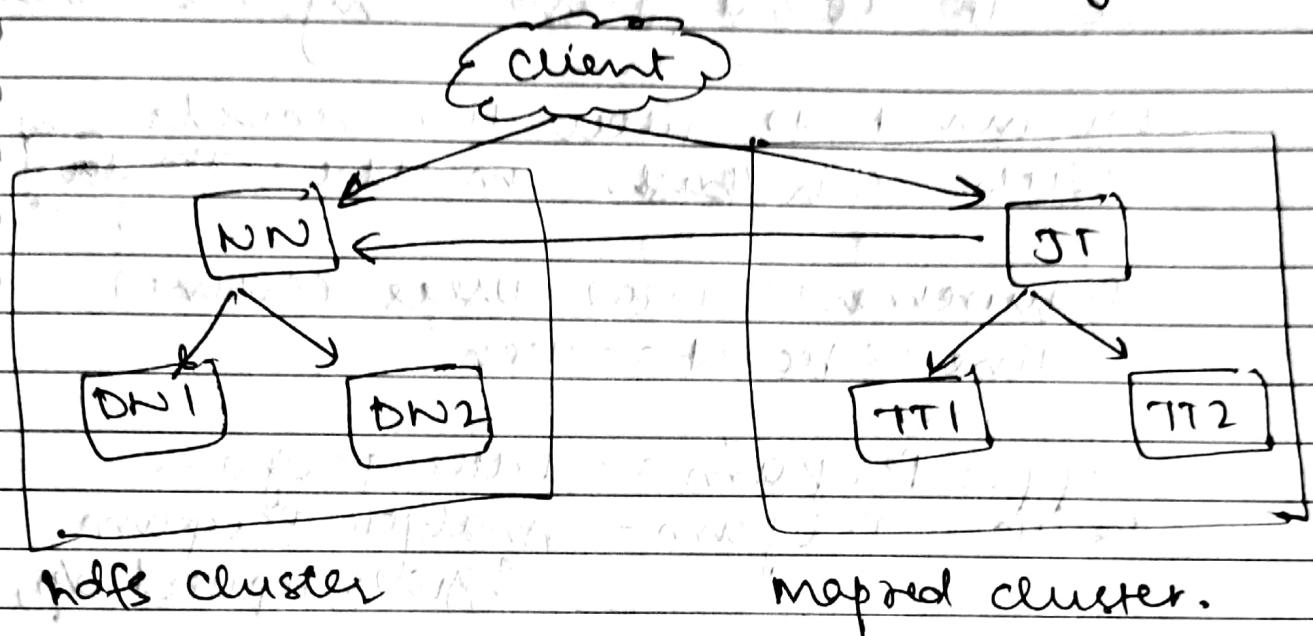
↳ i.e job is success.

↳ part-r-00000  
↳ O/p is stored here.

# hadoop fs -cat /output/part-r-00000

# Session 148

## Hadoop Map Reduce Programming



- Cluster is created and made permanent.

~~client~~

```
# hadoop fs -ls /
# hadoop job -list-active-trackers
↳ 2 Task Trackers.
```

```
# cat /etc/passwd/
↳ lists all the users
↳ jack:x:1001:1001::/home/jack:/bin/bash
    ↳ : → field separator/delimiter
    ↳ 7 columns
    ↳ User → User name
    ↳ : → User id
    ↳ : → Group id
    ↳ : → Comment field
    ↳ : → User's home directory
    ↳ : → User's login shell
    ↳ This user can
    ↳ log in in bash shell
```

- If  $uid \geq 1000$  ; user is general user  
 $\therefore$  limited power.  
& root has  $uid = 0$

- /etc/passwd/ is raw data  
bc. it doesn't have any metadata
- we want to filter the records  
with /bin/bash in the ~~for~~  
field

Requirement: Total users (count)

line == record == user

- Mapper Program == Filter program
- Reducer Program == Analytic Program.  
(gives required o/p)

```
# cat /etc/passwd | grep bash  
↳ filter bash (Mapper)
```

```
# cat /etc/passwd | grep bash | wc -l  
↳ 6 (count) (Reducer)
```

```
# hadoop fs -put /etc/passwd /  
# - - - - ls /  
↳ /passwd
```

\* In Hadoop, there is a streaming API  
where we can use languages other  
than Java, ex: Python, Scala, Shell,  
etc.

# cd /usr/share/hadoop/contrib/streaming  
↳ hadoop-streaming-1.2.1.jar  
(streaming API program)

# cd

# hadoop jar /usr/share/hadoop/contrib/  
streaming/hadoop-streaming-1.2.1.jar  
~~grepped~~ grep bash  
-input /passed -mapper "grep bash"  
-reducer "wc -l" -output  
"myuser01"

# hadoop fs -cat myuser01/part-\*

↳ G

## Map Reducing Programming in Py

python.

s = "vimal:tw:123"

s.split(":") [1]

↳ tw

s.split(":")

['vimal', 'tw', '123']

• This will split ":" (separator)

s = "vimal:x:1002:1002::/home/vimal/bin/"

s.split(":") [2]

↳ 1002 (string)

bash

```

if int(user[2]) >= 1000 and user[6] == "/bin/bash":
    print("User is a general")

```

```

fh = open("/etc/passwd", "r")
for i in fh:
    print(i)
    # removes the empty lines.
fh.close()

```

```

fh = open("/etc/passwd", "r")
for i in fh:
    print(i.strip())
fh.close()

```

```

# mapper.py
import sys
sys.stdin
fh = sys.stdin

```

```

# fh = open("user.txt", "r")
for i in fh:
    # removes the empty lines.
    i = i.strip()
    user = i.split(":")
    # user[2] = uid
    # user[6] = shell
    uid = user[2]
    shell = user[6]
    if int(uid) >= 1000 and
        shell == "/bin/bash":
        print(i)

```

fh.close()

- In filtering algo we always need to traverse the whole program.
- ∴ Best Case = Worst Case = Avg Case =  $O(n)$

# cat /etc/passwd | python mapper.py

# vim reducer.py

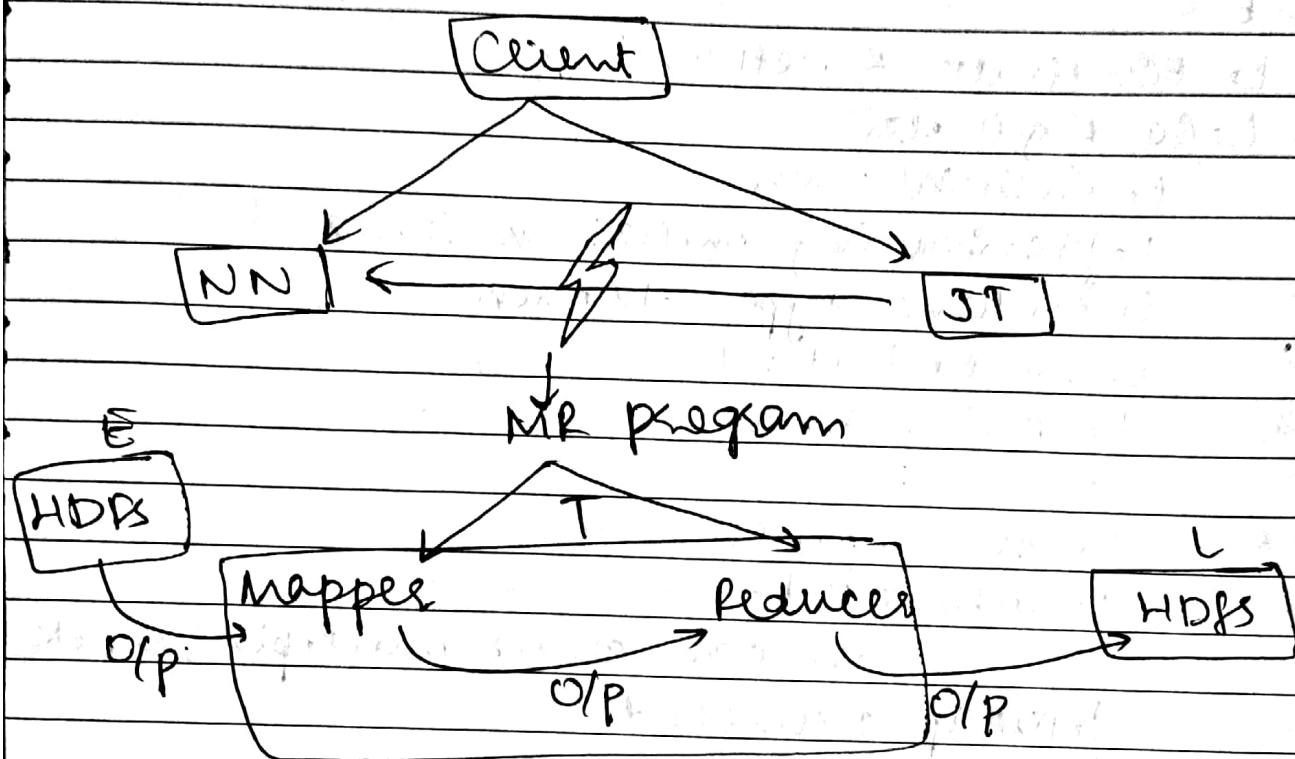
```
import sys
fh = sys.stdin
j = 0
for i in fh:
    j += 1
    print(j)
EOF
```

# cat /etc/passwd | python mapper.py | python reducer.py

```
# hadoop jar /usr/hadoop/lib/hadoop-mapreduce.jar
 -input "/passwd" -file "mapper.py" python
 mapper.py" -file "reducer.py" python
 reducer.py" -output "/op2"
```

SESSION: 152

Hadoop



- ETL process

- Mapper program always run in parallel ( $3TT \rightarrow 3$  mapper program run in parallel).

```
# cat db.txt | tr -s '\n' ''\n'
```

↳ all the words come in one line

ex: this is w

→ this

is

w

```
# cat db.txt | tr -s ' ' ''\n'|
```

sort | uniq -c

↳ 1 from

3 linux

3 word

→ For word count

↳ mapper.py

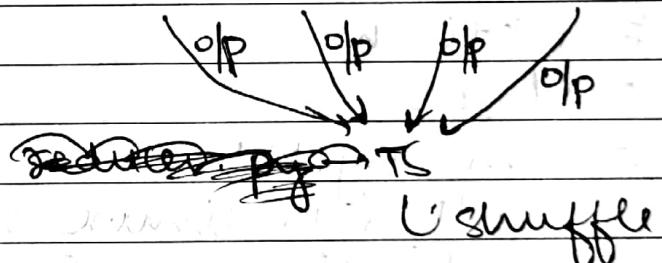
- word is the key

- ∴ split it to each word to diff. line.

Ex: linux 1  
word 1 } frequency  
hi 2 }  
this 3 }

- It is imp to sort the data first before running the cmd : `uniq -c` bec it will make the program more optimized.

Mapper.py → T1 T2 T3 T4



Reducer.py

↳ shuffle

- we already have a sorting algo in hadoop.

Hadoop : sorting algo (Mapper : quick sort)

(Reducer : merge sort)

ex: is 1

linux 1

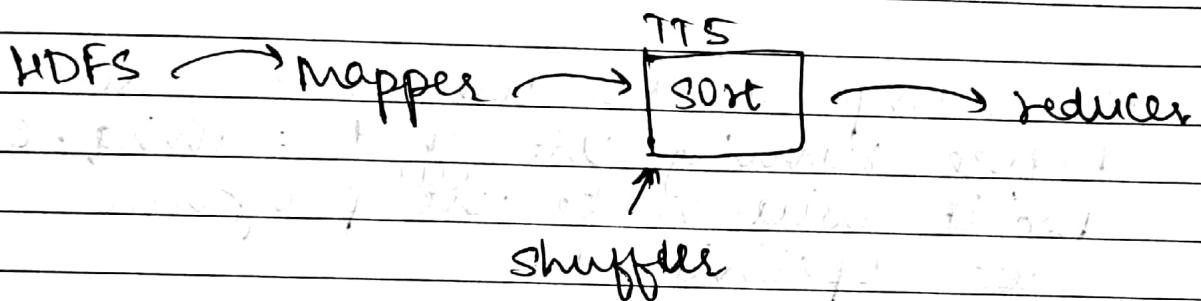
linux 1

linux 1

world 1

world 1

world 1



client

# mkdir ws

# cd ws

# vim my.txt

↳ this is linux world

hi linux world

welcome from linux world

hi

# hadoop fs -put my.txt /

# hadoop fs -ls /

# vim m.py

#!/usr/bin/python

import sys

~~@app.route~~

```
for line in sys.stdin:  
    lineList = line.split()  
    for word in lineList:  
        print word, 1 # without ()
```

# chmod +x m.py

# vim r.py

#!/usr/bin/python

```
import sys
```

```
for i in sys.stdin:  
    print i.strip()
```

# chmod +x r.py

```
# hadoop jar /usr/hadoop/contrib/streaming.jar  
hadoop jar -input /my.txt -file m.py  
m.py -mapper file r.py -reducer r.py -output /sl
```

# hadoop ps -cat

↳ automatically our data is sorted

(1st partition has min value, 2nd partition)

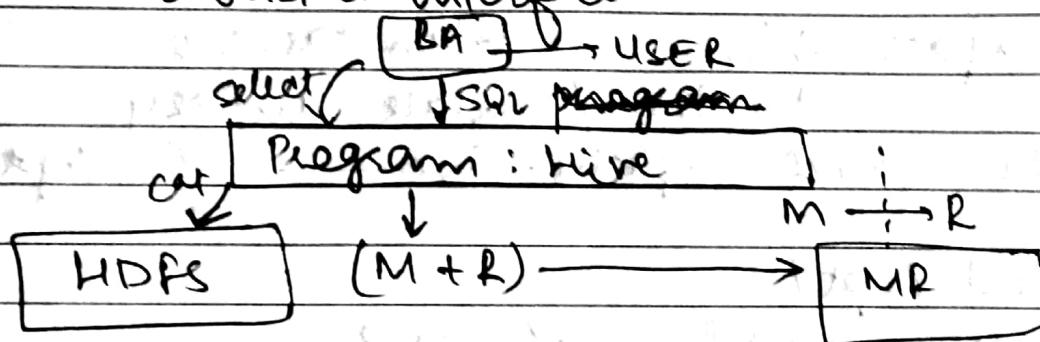
(2nd partition has max value, 1st partition)

SESSION: 154

Hadoop

Hive

↳ JUST a interface



↳ Generally we use  
MR program in  
Analytics.

where we can divide  
a problem in two  
parts:  $\oplus$  Mapper &  
Reducer.

- Nobody use hadoop directly. They always use intermediate program which automatically converts the program into Mapper & Reducer program.
- This intermediate program is Hive or a layer bet. client and hadoop.
- Hive ~~is~~ should be installed in the client.

~~client~~ # hadoop fs -ls /

# apache-hive-0.13.1-bin.tar.gz

↳ drive/uninstall directly

# tar -xvzf apache-hive-0.13.1-bin.tar.gz -C /opt/

# cd /opt/

# mv apache-hive /

# cd /hive/

↳ bin

All the cmd are here.

# export

# HIVE\_HOME=/opt/hive/

# export

# PATH=(\$HIVE\_HOME/bin):\$PATH

# vim /root/.bashrc

↳ export HIVE\_HOME=

export PATH=

- Generally we store raw data in hadoop but for analysis we need metadata/schema.

- Hive is a meta data service. It only creates/stores the schema of the data. It picks data from hdfs.

# mkdir /hivews

# cd /hivews

# ~~hive~~ hive

↳ initializes hive

hive> show databases;

> create database hw;

↳ creates a user in hdfs.

hdfs : /user/hive/warehouse/hw.db

> use ~~hw~~;

- Creating a table in ~~hive~~ hive means creating metadata + raw data.

> create table myuser(username string,  
password string, uid int, ejt int,  
comment string, homeDir string,  
shell string) row format  
delimited fields terminated by ":"  
lines terminated by "\n"  
stored as textfile;

↳ schema is created.

↳ folder is created in hdfs

↳ /user/hive/warehouse/hw.db/  
myuser/

- Hive uses derby to store data base.

- Derby db is very light weight.

hive > select \* from myuser;  
↳ nothing will come bcs table  
is empty.

> load data local ~~temp~~ input  
"/etc/passwd" into table myuser;  
↳ loads the dataset into myuser  
table.

> select uid from myuser;

> select ~~\*~~ from myuser where  
shell = "/bin/bash";

↳ converts the query into mapper  
& reducer program & gives the op

> select count(\*) from myuser  
where shell = "/bin/bash";

> insert ~~onefile directory~~ <sup>onefile directory</sup> into myuser  
select count(\*) -

↳ stores the output in /hol dir.

> desc myuser;  
↳ it describes the table.

> desc formatted myuser;

Hive show

hive> create database;

Creates the db in hdfs

↳ /user/hive/warehouse

# /opt/apache-hive/conf

↳ we can change the default dir  
find ~ <workspace>

hive> create database if not exists tw;

> describe database tw;

↳ gives all the details of tw db.

hive> create database if not exists twtbl;

comment "for team1" with  
dbproperties ('creator'='vimal', 'date'='2021-01-27');

> use twtbl;

> show tables;

↳ 2 types of Tables in hive!

Internal

(Managed type)

- hive owns power over data

External

- If the table is deleted  
the raw data won't  
be removed.

# vim db.csv

1, cricket, India, 1111

2, xyz, US, 2222



- If we have multiple values in a field.

1, cricket:adv:cooking, india, 111

↳ we have to change all delimiter  
ex: space or colon(:).

→ 2, cooking:xyz:abc, us, 2222

→ 3, studying:creating:abc, india, 3333

if not exists

> create table `table1` (id int, hobbies  
array<string>, country string, mob int)

↳ raw formed <sup>delimited</sup> fields terminated by ;

↳ lines terminated by "\n", collection  
items terminated by ":" → stored as  
textfile;

at the  
last

> show tables;

> describe `table1`;

> read data local input '/hivesrc/date/  
ab.csv' into table `table1`;

↳ Data is loaded onto the table

> select \* from `table1`;

↳ Prints the table.

If want to print the column name

then go to the `conf` file

and change `show columns` from  
false to true

External table  
↳ Create Table

> Create external table table2 (id \_\_\_\_\_  
—) ; — stored as textfile;

> Load data

> ~~Load~~ drop table table2;

↳ The table & Schema is removed by  
the raw data is still intact.

↳ Now if we recreate the table  
all the data will come up without  
loading it.

# hadoop fs -mkdir db  
# hadoop fs -put db.csv /db

> Create database if not exists hdb;

> Create external table ~~as~~ hdb(\_\_\_\_\_) —  
— stored as textfile location '/db';

↳ Table created

> describe table hdb;

↳ External table

↳ mapped to db from hdfs

↳ ~~it~~ automatically picks data  
from hdfs.

> Select country, mob from hdb;

> Create table newtable (~~country~~ string,  
~~phone~~ int) store as textfile;  
> select \* from newtable;

> Insert into table newtable select  
country, mob from table3B where  
country = 'india'

↳ first we run and the o/p is  
stored in newtable.

ELT : Extract Load & Transform  
↳ Hadoop uses ELT operation.