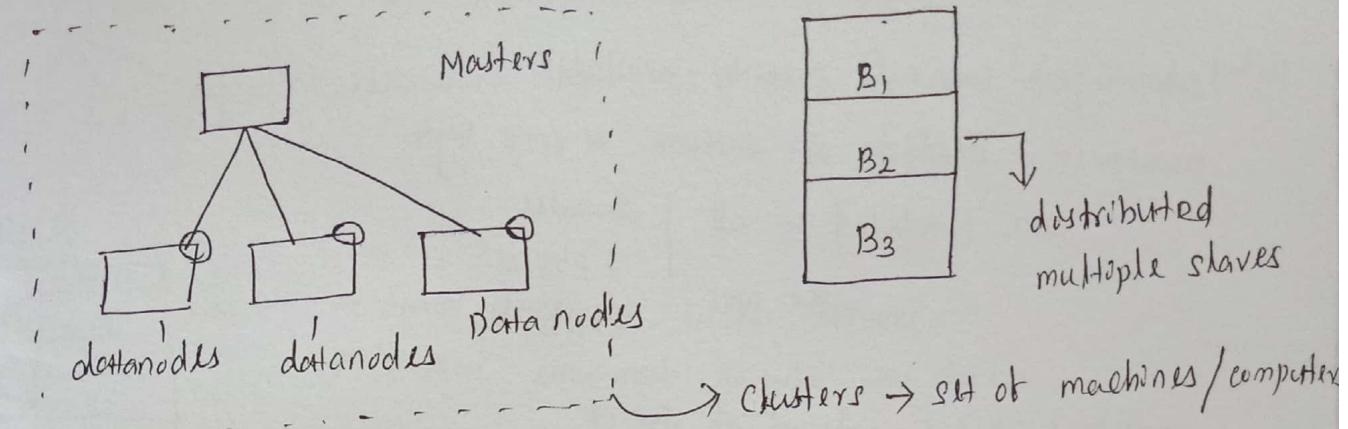


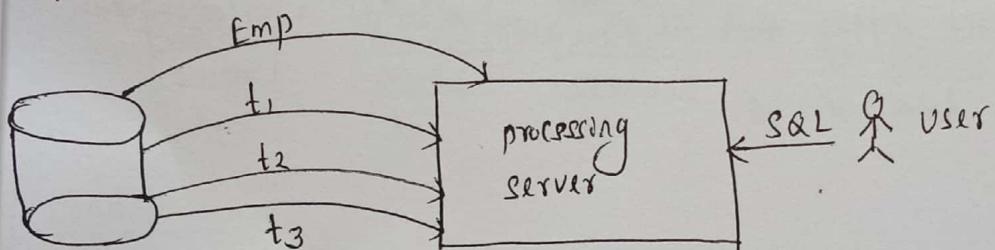
HDFS

HADOOP ARCHITECTURES

05.12.2011
wednesday



* Masters will have the metadata for all the nodes. Metadata will contain filename, size, no of blocks, block ids, block locations, etc.



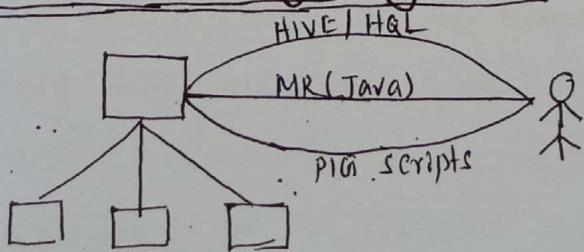
- Huge content of data is travelling over the network.
- Leads to network traffic/congestion.
- poor query performance

* So to solve these issues data locality comes into picture.

Data size	MB, GB, TB
Logic size	B, KB

- ⇒ It will do the logic travel, not the data travel.
- ⇒ Logic will be passed each and every slave and processing happens locally at each slave.
- ⇒ Data + Logic happening in one server is called data locality.

PPM (parallel programming Model)



cheap and commodity HW DNS \Rightarrow slaves

↳ easily available

Advantage - Low cost & easily available

drawback - chances of failure is very high

space contributed by all	20 nodes
the DNS	7 phyts
	7 \times 1021 TB

Heart beat = 3 sec interval of the CF.

* A TCP-IP protocol that each and every DN are connected.

DN failure

DN sends to NN saying that we are down & can do a piece of -

To overcome this data loss RF.

Replication factor

Advantage - Data high availability

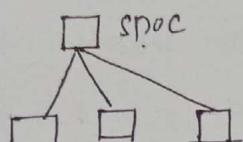
Disadvantage - occupies more spaces.

Backup cluster to take the data (DR), disaster recovery

wait time:- it will wait for a while & take the replica.

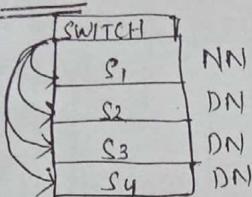
hadoop fs -setrep 2 file.txt

↳ file level replication.



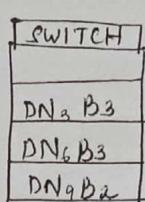
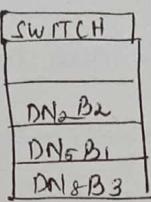
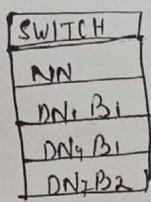
- All the information is coming to the Name node first.
- Replication factor = 3 → copy of same data in 3 nodes within the clusters.

RACK



- Actual placement for the servers.
- physical equipment to hold the servers.

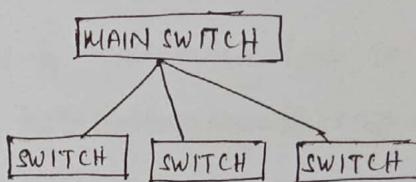
9 DNS, 1 NN, 3 RECS, 3 RRF, 150MB, BS=64MB



R1

R2

R3



conditions

- No 2 copies of the same block will be stored in same data node. Always a replica should be stored in a different data node.
- No 3 copies of the same block will be stored same rack. At least 1 copy of the block should be stored in remote (different) rack.

Identification Unique

Blocks - Block ids

→ Block algorithm for placement

Nodes - IP address

→ Rack Awareness

or
host names

→ Hadoop will follow the Awareness by

Rack - Rack ids

following the block placement algorithm.

Datacentre - Datacenter

- Name Node will have all the metadata.
- config files will have all the information for the master and slaves.
- block size should be configured at the starting. It can't be changed later, same for the replication factor.
- Space wastage from 0 to 63.999... MB

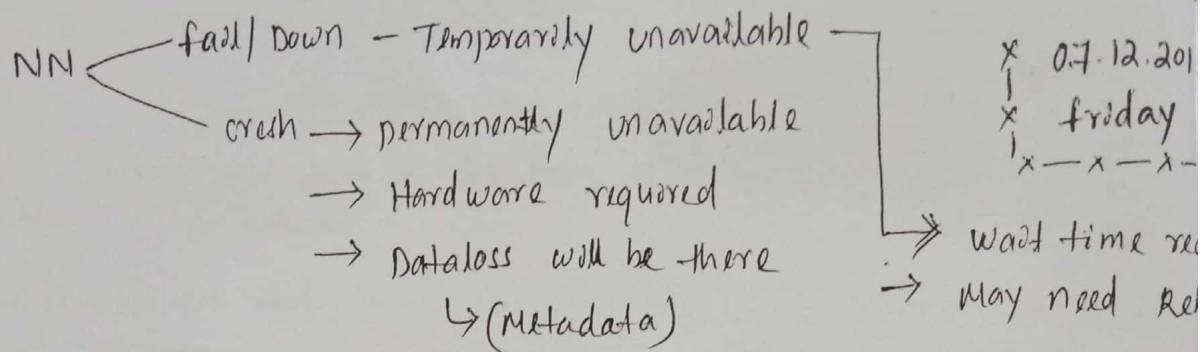
$B_1 - R_1/DN_1, R_1/DN_4, R_2/DN_5$

- If one rack failure B_1 , then

$B_2 - R_2/DN_2, R_1/DN_7, R_3/DN_9$

will be different location for as well to recover.

$B_3 - R_3/DN_3, R_3/DN_6, R_2/DN_8$



→ In first generation of hadoop 1 NN.

→ In 2nd generation of hadoop 2 NN.

1 NN crash → Metadata loss

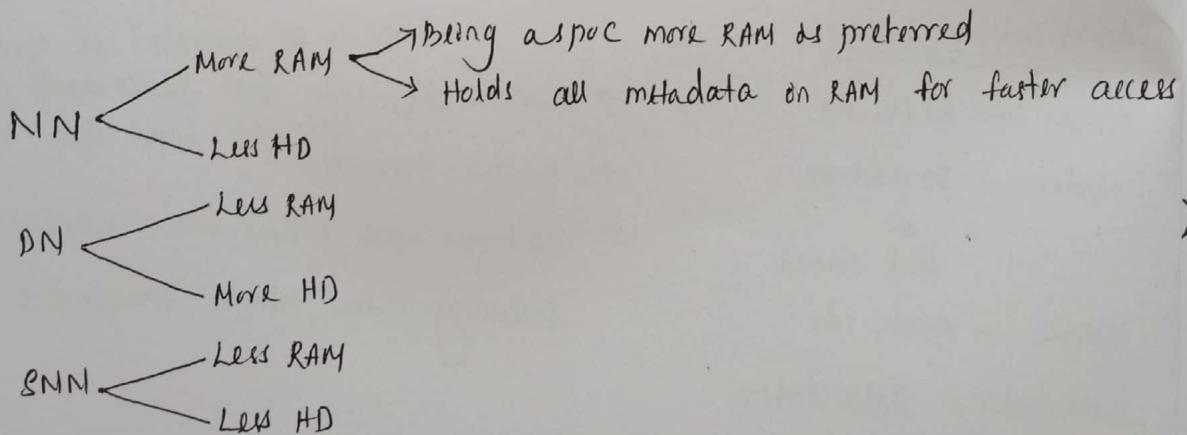
→ So new cluster will be configured.

2 NN → It is not hot backup of Name node

→ It takes metadata backup at some definite intervals of +
(check pointing time = 1 hr by default)

→ Metadata stored in 2NN can be used for rebuilding a new node in case of NN crash.

→ It will take only the delta information.



NN configuration

→ 16 core in one machine - Hot core system.

→ 1x2 TB HDD, 16x16 GB RAM

→ Total capacity : 2TB HDD, 256 GB RAM

Data node

1x2 TB, 4x16 GB

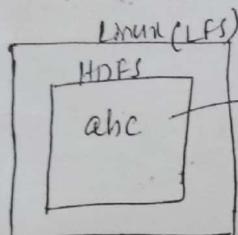
= 2 TB HDD, 64 GB

Datanodes - 32 TB - 30 x 4
= 120 TB

① \$ start-all.sh

② \$ jps - Java related processing

3004 NameNode



LFS
↓
HDFS

- ① HDFS
 - ls
 - cat
 - rm
 - copy from local
 - move from local - cut + paste
 - put
 - get
 - cp
 - mv
-] HDFS → HDFS

Hadoop fs -mkdir abc

\$ hadoop fsck /user/training/abc/emp.csv -files -blocks -locations -racks

Localhost: 50070

- Browse the file system

HDFS - commands

\$ jps

- JPS
- Datanode
- Jobtracker
- Secondary Namenode
- Tasktracker
- Namenode

\$ hadoop fsck → Check the status of the file like size, dirs, files, blocks

\$ hadoop dfs -cat /user/training/abc/emp.csv etc
→ To read the content of the file.

\$ hadoop fs -ls /user/training/abc/emp.csv

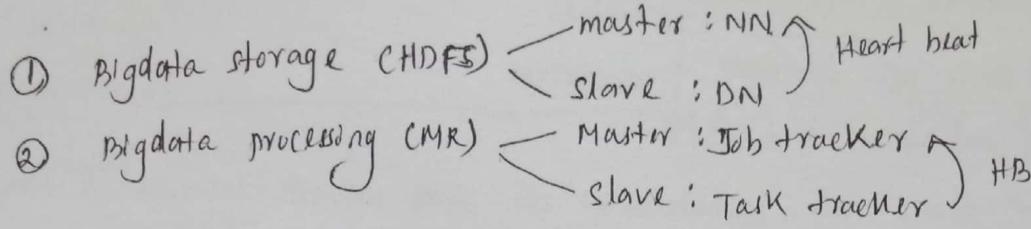
→ To see the content for the specified path, how many files are found, access level

\$ hadoop fs -cp /user/training/abc/emp.csv /user/training/abc/emp1.csv

] Source Destination]

- The content from emp.csv will be directly transformed to emp1.csv
- \$ hadoop fs -cat /user/training/abc/emp1.csv.
- \$ hadoop fs -rm /user/training/abc/emp1.csv /user/training/abc/emp1
- Removing the files from the specified location.
- \$ hadoop fs -cat /user/training/abc/emp1.csv
- The file will not be present more in the specified location.
- \$ hadoop fs -put /user/hive/warehouse/abc.db/emp1/emp1.csv /user
- \$ hadoop fs -cp /user/hive/warehouse/abc.db ----- /abc
- * HDFS to HDFS will be done by the -cp command, -put will not work for HDFS to HDFS.
- \$ hadoop dfs -put /home/training/f1.txt /user/training/abc
- \$ hadoop dfs -cat /home/training/emp1.csv
- \$ hadoop dfs -get /user/training/abc/emp1.csv /home/training
- file already exist
- \$ hadoop dfs -copyFromLocal /home/training/f1.txt /user/training/abc
- \$ hadoop dfs -moveFromLocal /home/training/f1.txt /user/training/abc
- \$ hadoop dfs -ls /user/training/abc
- * file will be present in this location.
- \$ hadoop fs -mv /user/hive/warehouse/abc.db/emp1/emp1.csv /user/training
- \$ hadoop fs -ls /user/training/abc
- emp1.csv
- emp1.csv
- f1.txt
- f2.txt

10.12.2018
Monday



NN will do 2 tasks \rightarrow check file is there / not \rightarrow file already exist
 \rightarrow check storage capability - then RF factor ...

storage jobs \rightarrow hadoop fs -

processing jobs \rightarrow hadoop jar -

\rightarrow Job Tracker will connect with Task Tracker

Job ~~JT~~ \rightarrow TT₁ + TT₂ + TT₃

TT₂ + TT₃ + TT₄

if the DN failed $\frac{1}{3}$ of data not completed, so it will find for the next node
for the same task.

JT URL \rightarrow Username - Ramush

\rightarrow (web user interface)

Hadoop config files

1) HDFS-site.xml (BS, HB, RF ...)

2) core-site.xml

101 (NN)

3) Masters

103 (NN)

3) mapped-site.xml

102 (JT)

5) Slaves

104 (DN+TT)

105 "

106 "

6 servers

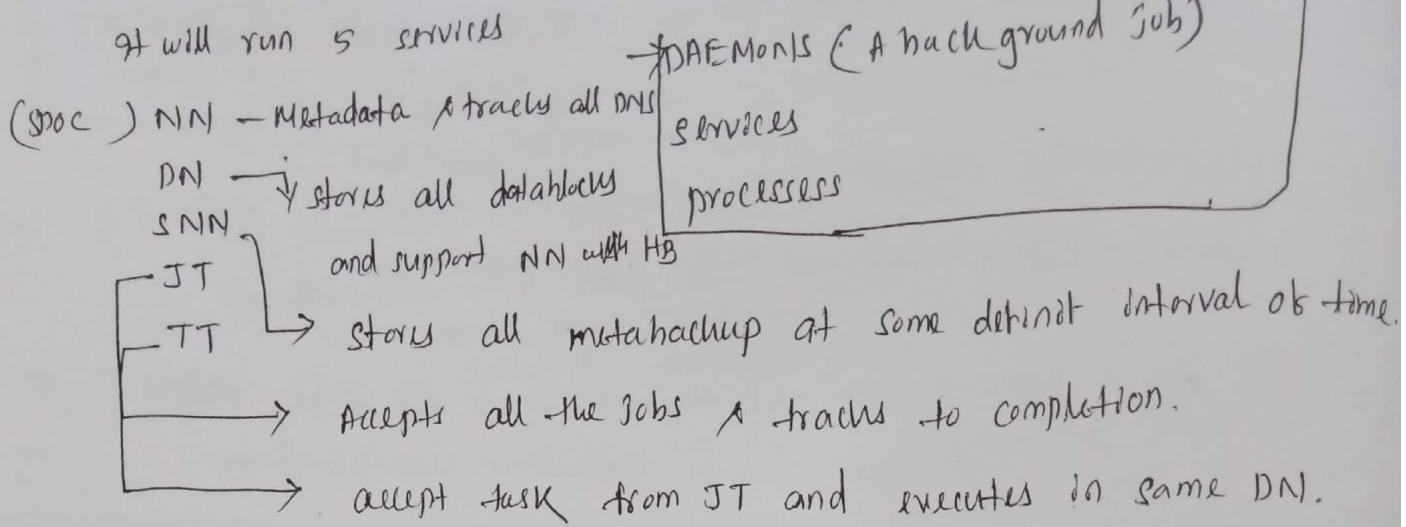
192.168.133.101	(NN)
102	(JT)
103	(NN)
104	(DN1)
105	(DN2)
106	(DN3)



Java + Hadoop both
needed

* Hadoop config folders will be there across the machines
all the machines will have.

* No need to start the servers at each node, in one machine you can do as
all the nodes have the config files.



Types of installation

- 1) ~~WIN~~ Linux
- 2) WIN + Linux (Dual Boot)
- 3) WIN + VM + Linux

CELOS - community enterprise operating system

\$ ifconfig \Rightarrow provide ip address.

\rightarrow \$ stop-all.sh \Rightarrow as JT is not working

localhost:50030 \rightarrow job details \rightarrow GUI for Mapreduce

Ment \rightarrow \$ exit; \rightarrow then logout from vmware

clientnodes/ dummy node

HDFS writes

HDFS reads

HDFS deletes

scaling \leftarrow vs
HS

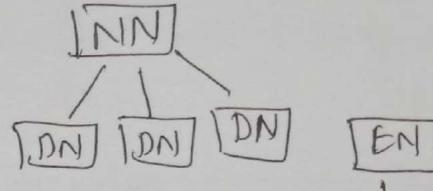
Load balancing

Namenode access

- \rightarrow Name node is busy for other task
- \rightarrow Name node having the sensitive data
- \rightarrow so no access given to developers

{ 11.12.2018
* Tuesday
* root root

→ Edge node
 Client node
 Dummy node
 Gateway node
 Connecting node
 Application node



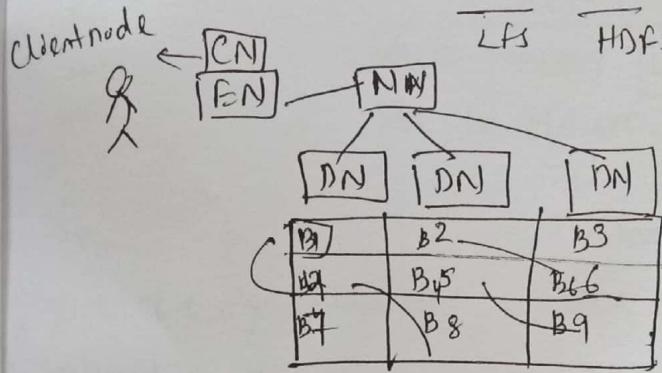
not used for storing data

* all the nodes & EN are interconnected

* can be accessed with ssh command

HDFS writes

hadoop fs -copyFromLocal file.txt abc → to client node



* EN will split a file & store in a location and as the other block to keep the copy of the file with them.

* After data replication all the information/ metadata will be stored in NN.

* CN will ask for available resources from the main node.

DN₈ acknowledged to DN₄ then

DN₄ report to DN₁ then

DN₁ report to client node the blocks are replicated.

* Block placement algorithm

* NN will provide excluding dead and full datanodes, to client node

* CN have the config files.

* the DN confirmation to CN then CN report to NN to update the metadata.

HDFS Reads

→ hadoop fs -cat abc/file.txt

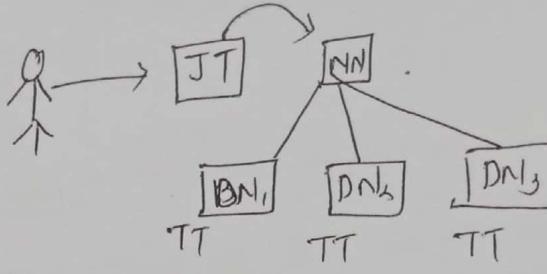
* it will give the nearest node

HDFS deletes

→ hadoop fs -rm abc/file.txt

* All the blocks will be removed.

processing
 PIG (script)
 MR (Java)
 Hive (SQL)



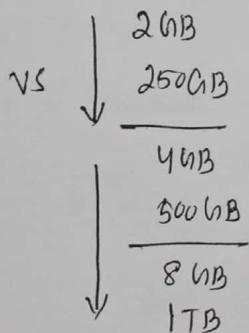
- 1) soft delete \Rightarrow hadoop fs -rm abc/fi.txt - (soft delete)
 2) Hard delete \Rightarrow hadoop fs -skipTrash abc/fi.txt - (Hard delete)
 \Rightarrow hadoop fs -cp abc/fi.txt HDFS HDFS
 \Rightarrow hadoop fs -mv abc/fi.txt abc

* Blocks will be generated in ascending order.

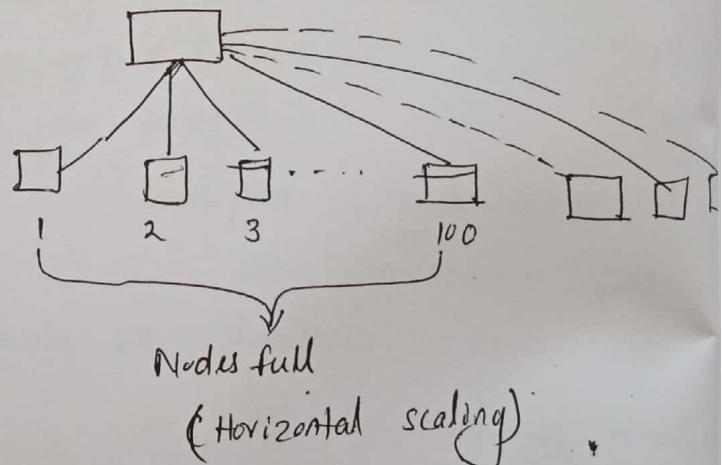
Scaling

12.12.2018
 Wednesday
 ~ ~ ~

\rightarrow increasing / decreasing the capacity



(Vertical scaling)



Load Imbalance

cluster resources under utilization.

Satemode (Read-only mode)

\rightarrow hadoop dfsadmin -satemode enter (ON)

If not satemode then data will come to the newly added node

\rightarrow hadoop dfsadmin -satemode leave (OFF)

\rightarrow hadoop dfsadmin -satemode get (status)

balance.sh

- some of the blocks from the present nodes moved to newly nodes.
- How many are newly added
- how much work should be copied to new nodes.

stopall.sh

→ some GUI tools will be used to change the node configurations and clustering will be handled

→ NN will have new data.

200

201

202

203

204

205

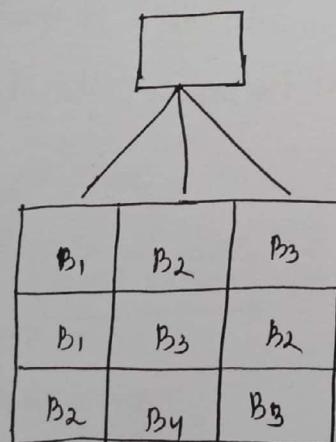
→ another way **Load balancing** will be done.

During cluster start-up BLR → Block Level Report

every 10th heartbeat as BLR

How many data blocks particular node is holding. $3 \times 10 = 30$ sec

When the cluster starts it will be in **safe/readonly mode**, till one heartbeat it should get.



$$3 \times B_1 - 1 = 2$$

$$3 \times B_2 = 3$$

$$3 \times B_3 = 3$$

RF < 3

→ All the blocks lost in their particular node called **"under replication blocks"**

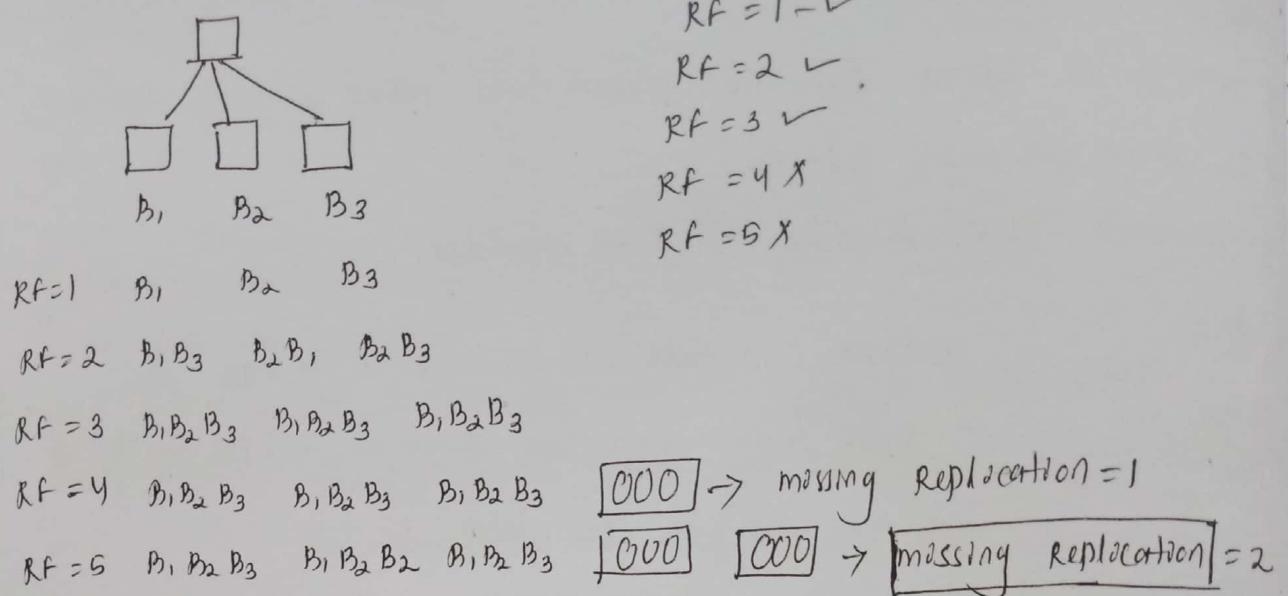
→ When the node come to normal condition then

$$3 \times B_1 + 1 = 4 \text{ nodes}$$

→ All the blocks present in the reconnected node are called as **"OR blocks"**

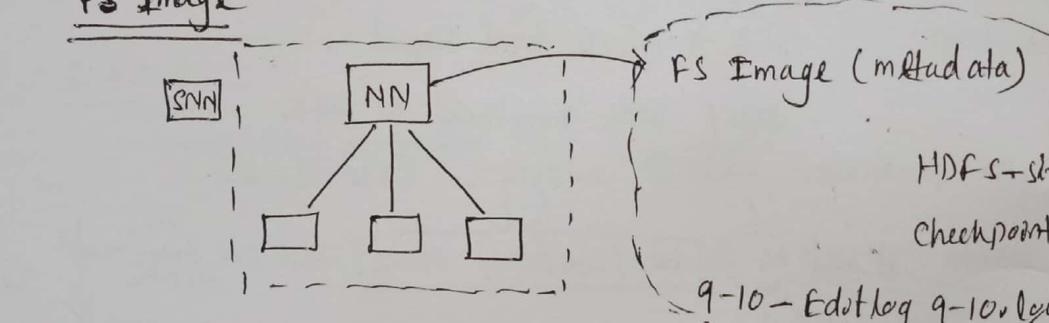
over replication blocks

RF. 7.3



→ RF for a cluster is always less than no. of DNS or atmost it can be equal no. of DNS, else it will be lead to "MR".

FS Image



$HDFS-site.xml$ $\times \rightarrow$ At f
Checkpointing = hr. req

9-10 - Editlog 9-10.log
10-11 - Editlog 10-11.log

→ at 10 AM SNN downloads fsImage + Editlog 9-10

= new fsImage

→ uploads then acopy of new fsImage to NN.

→ at 11 AM SNN downloads Editlog 10-11.log

- merges with new fsImage & creates one more $\boxed{\text{new fs image}}$

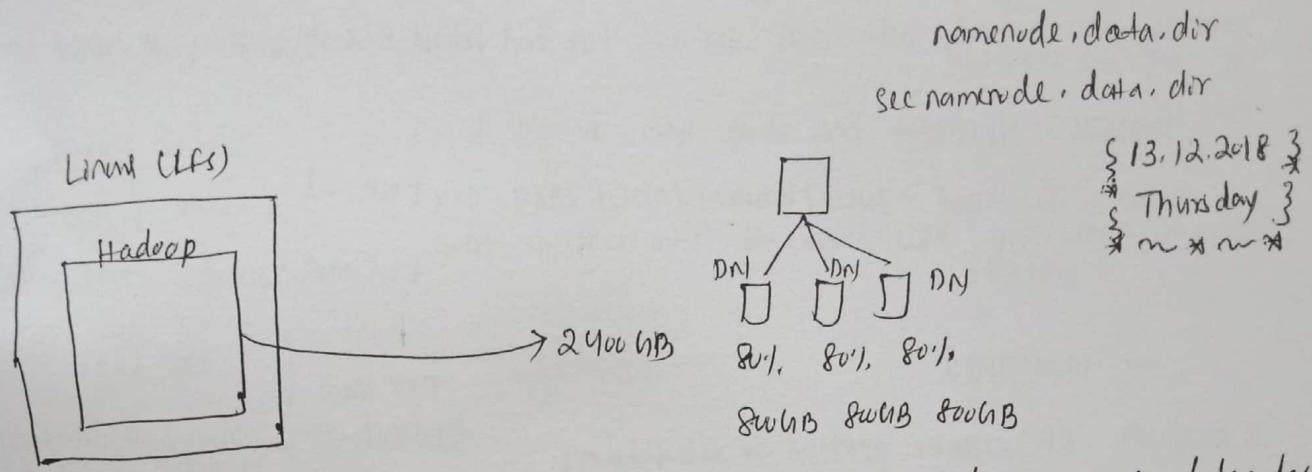
→ This is the recurring process of fsImage.

Editlog (transaction log)

- Add a file
 - delete a file
 - move a file
 - change in metadata
- Transaction in NN.

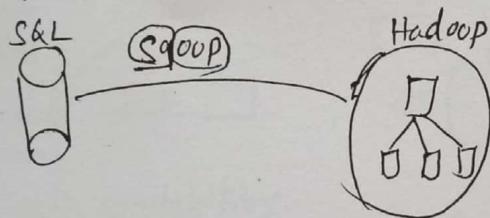
Checkpointing

→ process of merging edit logs to the fsImage & creating a new fsImage.
The node which is doing this activity is called as "checkpointing Node".



file will be maintained distributed across the datanodes.

single node cluster (pseudo distributed mode) ← training practice
fully distributed mode (multinode cluster) ← dev, prod, poc, test



8 hadoop fs -rmr abc

deleted hdfs://localhost:54310/user/training/abc

8 rmr abc

8 hadoop fs -ls abc

1024 B = 1 KB

1024 KB = 1 MB

1024 MB = 1 GB

1024 GB = 1 TB

1024 TB = 1 PB

1024 PB = 1 EB (Exa)

1024 EB = 1 ZB (Yotta)

1024 ZB = 1 YB (Zetta)

1024 YB = 1 PB (Bronto)

dfs] & hadoop fs -copyToLocal abc/emp1.csv

↳ home directory dfs

& cat emp1.csv → file will be displayed.

& hadoop fs -cat /user/training/abc/1987.csv | head -10

| head -10 → bottom 10 lines

& wc -l 1987.csv

1048576 1987.csv

& hadoop fs -cat /user/training/abc/1987.csv | wc -l

→ 1048576

| wc -l

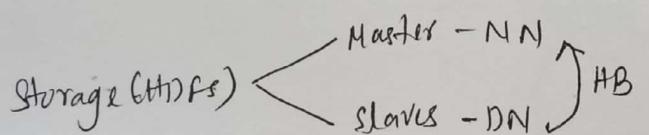
→ 102664080

& recursive RF will be applied to sub-folders. SetRep → To change replication

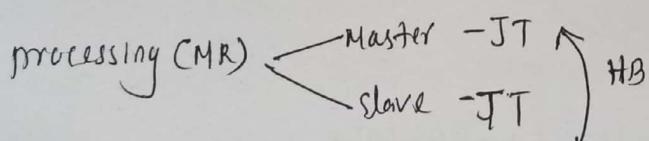
-R → Recursively apply the same RF, sensitive data and subfolder.

-W → waittime where we indicate the R.F. before copying the files

hadoop fs -^Ddfs, replication = 5 -copyFromLocal emp1.csv abc



{ 14.12.2018
* friday
** m



→ All storage jobs directly land in NN.

→ All the processing jobs directly land in JT & indirectly lands in NN.

Drawbacks of 1.x Hadoop

→ NN is SPOF → JT is SPOF

→ No horizontal scaling of NN. - No load sharing of metadata in different NNs.

→ JT is SPOF, in case of failure all jobs are restarted again.

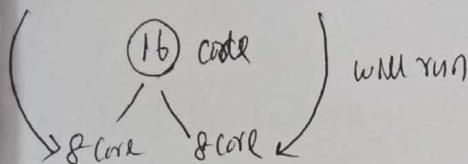
→ It is suitable for batch processing of huge amount of data.

→ It supports 4000 nodes for the cluster.

Job tracker is over burdened. It has a single component. Job tracker to perform many activities like resource management, job scheduling, job monitoring, maintaining all the history of jobs.

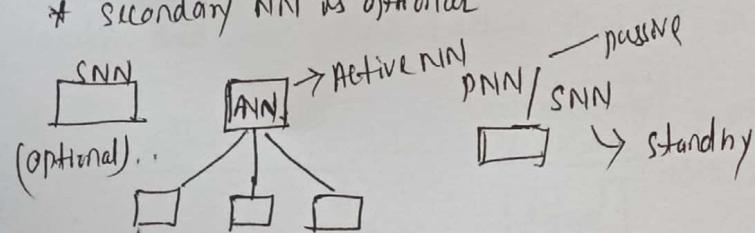
It runs only Map/Reduce jobs but not non-MR Application like spark, storm.

map | reduce
Phase-1 | Phase 2
mapred.java reduce.java
→ fixed no of map slots and fixed no of reduce slots
→ It takes more it can't use more than 8 core.



2.2 Hadoop generation

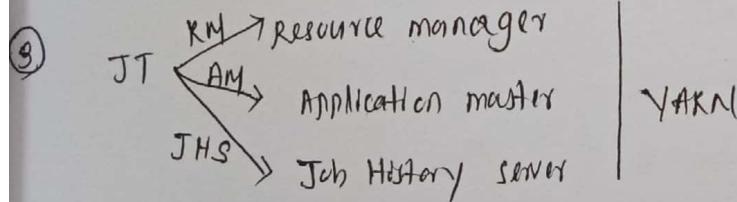
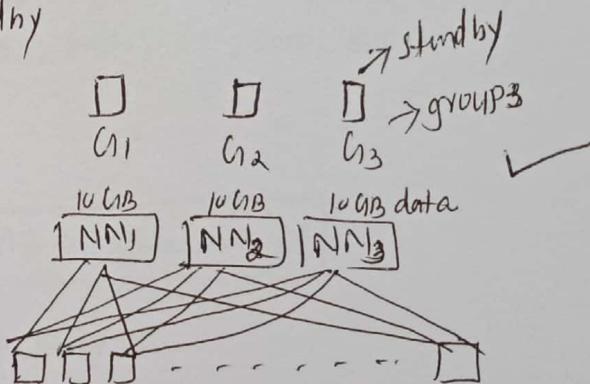
* Secondary NN is optional



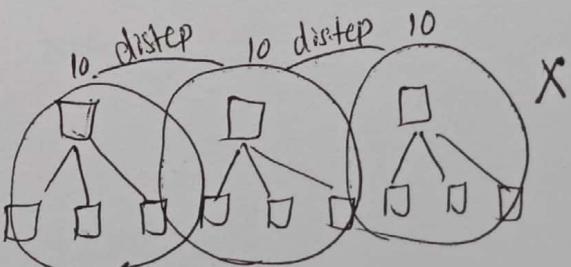
① NN High availability

② User segregation for multiple NameNode.

→ Name node horizontal scaling.



It anyone fails, it has high availability.

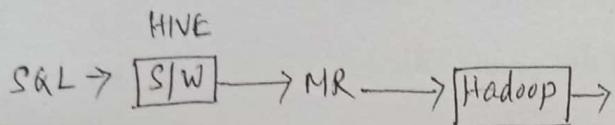
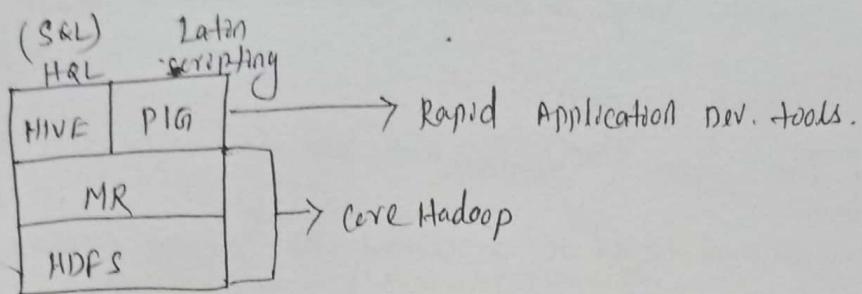


→ Under utilization of cluster resources → max 10 available

→ max 30 resources will be available at any time.

distep \rightarrow from one cluster to another cluster.

HIVE

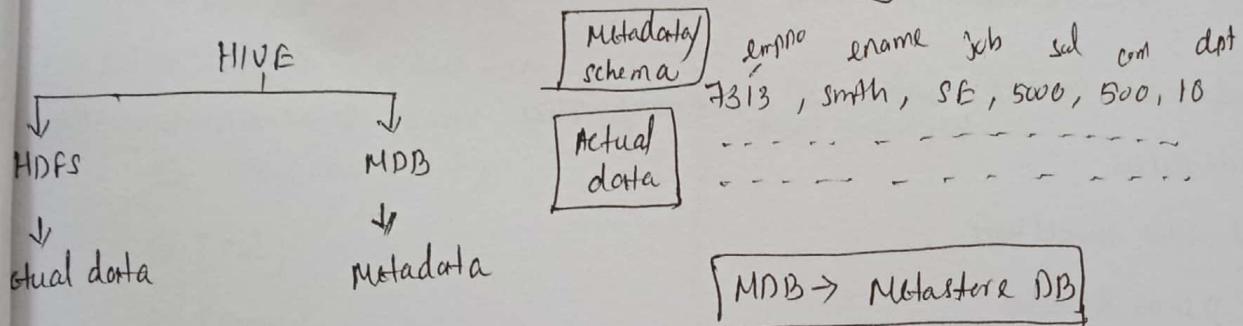


Developed by facebook.

It gives SQL kind of interface on top of data stored in HDFS.

Represents data in the form of tables.

Mainly used for OLAP (Online Analytical Processing)



* Apache Derby DB

* session or user oriented.

* To overcome above problem → MySQL or other DB

* Traditional DB is suitable for well organized metadata.

→ Create database dvs-nov18

- ① An empty dir is created in HDFS under default path
(/user/hive/warehouse)
dvs-nov18.db - directory
- ② Metadata is stored in MySQL

dvs-nov18.db → Database

emp → table belongs to default database without under any db.

hive> describe database dvs-nov18;

→ coming from MySQL. ← Metastore

hive> set hive.warehouse.dir; → setting the current directory.

~~hive>~~ → hive.warehouse.dir = /user/hive/warehouse

→ \$ hadoop fs -ls /user/hive/warehouse/dvs-nov18.db

drwxr-xr-x - - - - - dvs-nov18.db/emp →

→ select * from emp;

dvs-nov18.db

all the data coming from hadoop.

↓
emp (directory)

→ hive-site.xml

↳ All the default properties are set.

↳ emp.csv (file)

emp1.csv

emp2.csv

→ Hadoop will provide write once, read many. Thus version is not supporting the write access.

→ All files should have

i) same delimiter

ii) same no of columns

iii) same column sequences

→ It will 1st search the metadata then display the datasets.

① Managed or internal → table is dropped both metadata & actual

② External → table is dropped, only metadata is deleted (hdfs deleted)

③ Temporary table → metadata is deleted & data in hdfs is untouched

HIGH LEVEL DATATYPES

primitive/scalar

tinyint	1B	char - can hold char values with variable precision
smallint	2B	
int	4B	
bigint	8B	varchar - can hold char + int + several char with variable precision.
float	4B	
double	8B	
string	24B	
date	dd-mm-yyyy	

timestamp dd-mm-yyyy hh:mm:ss

boolean (true/false)

map
key value, key - value

Sal # 7000, comm # 600, 20

c2 map<string, int>

c2 [Sal]

c3 [comm]

} Map datatype

→ Sid, sname, Address
BTM & Bangalore # 560010

struct → collection of heterogeneous datatypes / dissimilar datatypes.

address struct<street: string, city: string, pin: int>

Select * from abc where

Select sid, address.city, address.pin from abc.

In query - while creating use below query.

collection items terminated by ' ; '

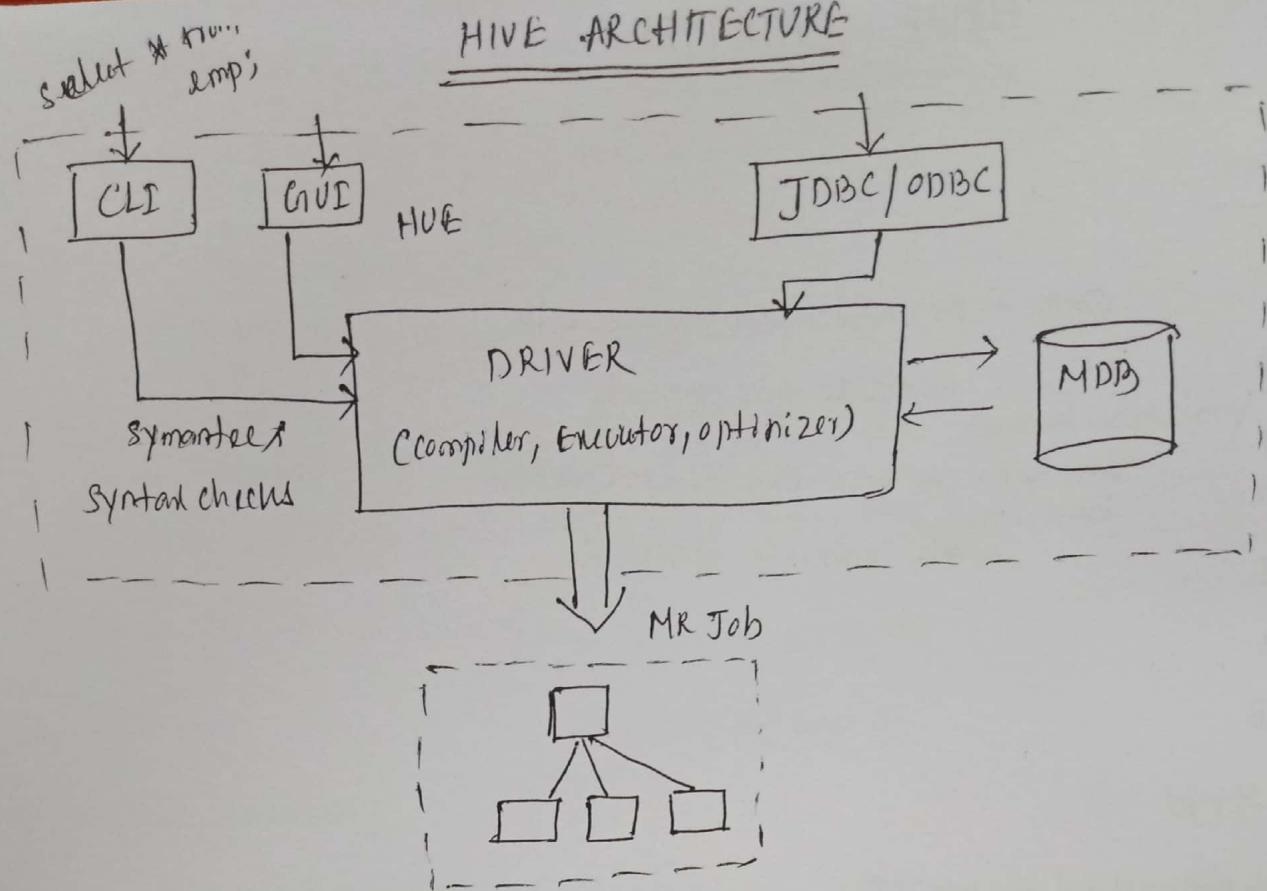
BTM & Bangalore # 560010 — X Not possible with more delimiters
Same: ..

complex

array

map

struct



HUE → GUI (like) → Hadoop User Experience → user friendly

JDBC/ODBC → from frontend like .net/ Java to run the MR Job, but most not used.

NOTE → If hadoop is not running the metastore DB commands will only not all the commands for getting actual data.

* Always table mapped to the Hadoop directory.

* N number of files can be there in the directory. So write once concept came into picture and more data can be stored.

17.12.2011
Monday

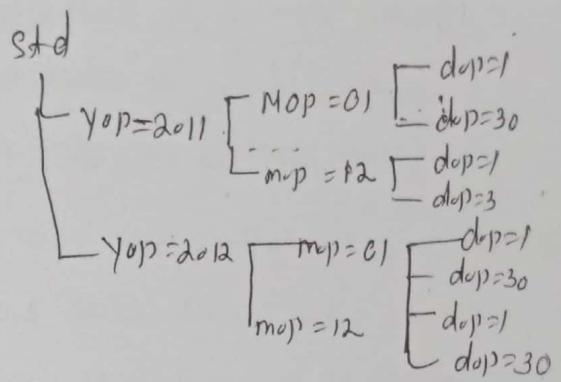
Mixed partitions

Std, sname, yop, mop, dop

/std/dp/ /year=2011 /month=2/ day=20

cat stud.txt

Ravi, 123, 31.12.2011



querying metadata don't need services start/run, but data retrieval will not work.

* select * from std where yop=2011;

360 folders searching

* select * from std where yop=2011 and mop=1 and dop=1;

359 subfolders will be avoided.

Hive Scripts

create database & {hivecont:db3;

use & {hivecont:db3;

create table & {hivecont:db3} (sname string, sid int).

JOINS

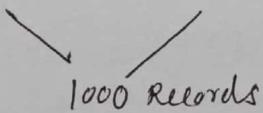
1) inner or Equi

2) Left outer \rightarrow all the rows from left, return NULL value.

3) Right outer \rightarrow all the records from right ^{matching from left,} returning NULL value.

4) full outer

5) cross \rightarrow 100 records \times 10 records



Create view v1 as select * from emp;

> Show tables.

VIEWS

Create view V_1 as select * from emp;

Create view V_2 as select c_1, c_2, \dots, c_{50} from emp;

Create view V_3 as select c_1, c_2, \dots, c_{10} from V_2 ;

519.12
* Thru
* m

* stores in metastore DB.

Create view V_4 as select $A.A, A.C_2, B.C_{12}, B.C_{13}$ from emp A JOIN dept B
on (A.deptno = B.deptno);

Location with partition table

1) Data in HDFS

2) Data should be well partitioned.

std

Year=2011

std-2011.csv

Year=2012

std-2012.csv

Year=2013

std-2013.csv

-mduo duv-hdfs / std / year=2011

=2012

=2013

hadoop fs -put std-2011.csv duv-hdfs / std / year=2011

Year=2012

Year=2013

Hadoop

1) Create table std-et(sname, ...)

Load data loc - "std-2011.csv" into std -
 slab directory created in HDFS after
 copied.

Metadata of partitions are in

After creating location with partition, in MDB,

You may not able to fetch the metadata details.

HIVE> alter table std-et add partition (year=2011);

> (year=2012);

(year=2013);

hive> ~~create~~ table std; ~~et;~~

↳ metastore check.

↑ The way we are fetching and adding data as partitions.

DFS	MDB
3	0 → 3
0	3 → 0
2	3 → 2
3	2 → 3

You can't delete/update/insert applied on ~~hive~~ structures.

hive> alter table std drop partition (year=2014)

→ metastore db data & actual data from DFS will be dropped.

hive> insert into table std1 partition (year=2017) select * from std where year=2017;
(static)

hive> insert into table std1 partition (year) select * from std where year=2017;
(dynamic)

std	std1
year=2011	
2017 TB	<u>empty</u>
2018 TB	

→ hive> alter table std1 exchange partition (year=2015)
with table std;

hive> alter table std1 exchange partition (year=2016) with

table std;

Now

std1
2015
2016

std-
et

hive> Alter table std partition (year=2012) ENABLE NO-DROP;

hive> alter table std-
et drop partition (year=2012);

hive> alter table std

INDEXING

20.12.
Third
27/12/2023

hive > select * from emp where empno = 1111;
→ will scan all the files in the directory.

emp (dir)
emp1.csv
emp2.csv
;
emp100.csv

hive > create index emp_index on table emp(empno)

as 'COMPACT' with 'decreased REBUILD';

↳ unique columns. ↳ partition/delay

when ever you are adding data need to rebuild the index.

Index table is created - dbname - tablename - index name.

empno	Location	Byte offset (Physical address of the record)
7859	1001/hive/worchrw/_dvs-nv18/dp/emp/	[0]
7369	emp1.csv;	[32]
1111	- - - - -	[62]

> create index emp_index on table emp(empno) as "COMPACT" with decreased re
> show index on emp;

> show tables;



dvs-nv18-emp-emp_index

> select * from ..

→ no data, as no rebuilt;

> alter index emp_index on emp rebuilt.

> select * from dvs-nv18-emp-emp_index;

> select * from emp where empno = 7782;

→ it will go to index table, get the address and find 7782.

→ Auto ascending implements.

→ so at very first time it will be taking more time to build the index.

→ When a main table dropped automatically index table is dropped.

BUCKETING

→ for faster joins (join optimization)

→ Table sampling.

clustering (deptno)

→ we have to consider amount of joining with table.

→ If we have deptno, our joins will be efficient.

→ 3 buckets = 3 files.

hive> create table emp-buc (.....);

emp-buc

000000-0

000001-0

000002-0

$$\text{Hash function} = \left\lfloor \frac{\text{deptno}}{\text{no of bucket}} \right\rfloor = \left\lfloor \frac{10}{3} \right\rfloor = 1$$
$$\left\lfloor \frac{20}{3} \right\rfloor = 2$$
$$\left\lfloor \frac{11}{3} \right\rfloor = 2$$
$$\left\lfloor \frac{30}{3} \right\rfloor = 0$$

→ hadoop fs -ls ... du-nov18/emp-buc

hadoop fs -cat ... /emp-buc/00001-0

→ for non-integers it will consider the ascii value.

select * from emp-buc tablesample (bucket 1 out of 3);
union all
(bucket 3 out of 3)] both 183
data will come.

Table sampling

hive> select * from emp-buc tablesample (10 percent);

select * from emp limit 100;

1) Buckets inside partitioning ✓ - possible - files inside folders

2) partitioning inside buckets ✗ - Not possible - folders inside files.

partition inside bucket ✗

Bucket inside partition. ✓

* Friday
sun sun

partition	order-id	order-status	order-date	cust-id
assignments	1	PENDING	08-08-2018	11111

orders

ord-status = PENDING

↳ 000000-0

ord-status = INPROGR

↳ 000000-0

order

ord-year = 2016

↳ 0000000

ord-year = 2015

↳ 0000000

Assignment-

abc

- 1) abc
- type = summer position based on type
 - ↳ 000000-0
 - type = winter
 - ↳ 000000-0

4) OLP
country = US ↳ year = 2000

country = UK ↳ year = 2000

5) OLP
year = 2000 ↳ type = summer

year = 2004 ↳ type = winter

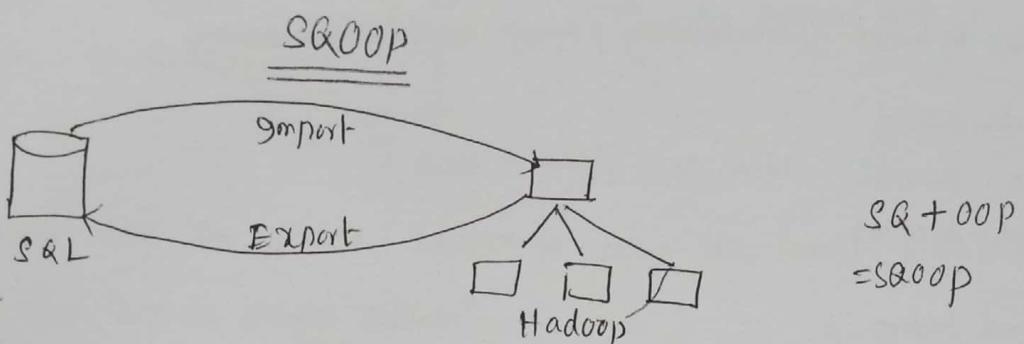
- 3) OLP
- year = 2000 ↳ 000000-0
 - year = 2004 ↳ 000000-0

6) OLP
country = US ↳ type = summer

- abc
- country = US ↳
 - 000000-0
 - 000000-0
 - 000000-0
 - country = UK ↳
 - 000000-0
 - 000000-0
 - 000000-0

* cluster by state - bucket
* sort by city

partitioned by (Country) clustered by (state) into 4 buckets sorted by (city).
 DML]
 UDFs] MR
 OT]



→ Data transfer / Integration tool from RDBMS → Hadoop or

MySQL
 DB → DVS-DB
 TBL → EMP

HDFS
 Sqoop-Nav

Hadoop → RDBMS only.

Import

- 1) from RDBMS (table) to HDFS file
- 2) from RDBMS (table) to hive tables

Export

- 1) from HDFS (folder) to RDBMS (table)
- 2) from hive (table) to RDBMS (table)

sqoop-import > --connect jdbc:mysql://localhost / DVS-DB \

> --username hadoop \

> --password hadoop \

> --table emp \

> --target-dir sqoop-Nav18 \

> -m 1;

can be oracle,
 any address.

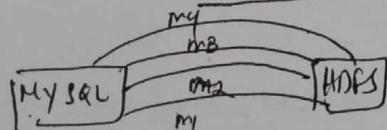
No. of mappers

for any sqoop job -

Any sqoop job running with map=100 \.

reduce=0 \.

→ m stands for mappers.



- 1) Depends on sourceDB no. of connections allowed.
- 2) It should not lead to small files.

mysql fs -ls sqoop-novi8

reduce 0%.

map 100%.

hadoop fs -cat /user/training/sqoop-novi8/part-m-00000

part-m-00000

part-m-000002

part-m-000003

mysql -u hadoop -p

mysql> use dvb_db;

> select * from emp;

- 1) restricting rows — where clause
- 2) restricting columns — columns clause
- 3) changing field delimiter.

sqoop-import --connect jdbc..... --username hadoop --password hadoop --target-dir sqoop-novi8 --map-column-mutations "empno,ename,job" --where "deptno=20" --fields-delimited-by '1';

cat > sqoop-param.txt

--comment

--username hadoop

--password password hadoop

sqoop-import --options-file sqoop-param.txt --table emp --target-dir sqoop-novi8 --map-column-mutations "empno,ename,job" --where "deptno=20" --fields-delimited-by '1';

* --username hadoop -p --table emp

→ prompt for password will come.

* password-file /etc/sqoop/cont/passwords/mysq1-pass.txt

* sqoop-novi8-m 5;

* If $m=1$ - it is sequential input

* If $m>1$ - it is parallel input.

* If $m=2$

part - m - 00000

* It will check the primary key

part - m - 000001

1

2

$$1000 - 1000/2 = 500 \text{ block}$$

If no primary key for a table we have to either do import with $m=1$ (or)

use split-by "some unique column".

-- split-by emp_no;

→ Hadoop fs -ls /sqoop4-nov18;

→ If no primary then use split

→ cat /home/training/dvs/.temp.sqoop/wc-1

HIVE IMPORT

MySQL

DB: retail-db

TB: orders

sqoop-import --options-file sqoop-novam.txt --table orders --hive-import

--hive-database dvs-nov18

→ It skipped default DB

--hive-table orders

→ It skipped same name as mysql

-m 10

→ Default mapper 4

--split-by 'order-id';

→ Distribution according to primary key or
no PK then give composite key

* If $m=1$

no need for split-by

→ split by empid/orderid/.....

→ order is not important, but all the
parameters should be satisfied.

→ max mappers depends on the max

DB connection sessions per user.

hue - hadoop user environment

su - switch user

password :-

JDS

mysql -u root -p

password - cloudera

mysql > show databases

hive > show databases;

hive > create database dvs-nov18;

hive > show databases;

cloudera@quickstart: ~] \$ sqoop-import --connect jdbc:mysql://localhost/~~hive~~ --username root password cloudera --table orders --hive-import --hive-database dvs-nov18 --hive-table orders.

HDFS
user / cloudera
orders

--delete --target-dir

HDFS-HIVE loc
/user/hive/warehouse/dvs-nov18.db/orders

hive > show tables

-- orders (sqoop created)

hadoop - root@quickstart: ~] hadoop fs -ls /user/hive/warehouse/dvs-nov18/

part-m-00000
00001
00002
00003

hadoop fs -ls ... orders/part-m-00001/wc-l - count

part-m-00002/wc-l -
part-m-00003/wc-l -

NOTE -> delimiter should be defined while creating the table struct.

Types of Loads

- ▷ initial / full / one-time / first load
- ▷ incremental / delta load
- Last value method
 - Last modified date method.

25/10/2018

Thursday

MySQL (orders)

ord-id	ord-date	ord-status	cust-id
1			
2			
3			
100			
1001			
1002			
:			
2000			
2001			
2002			
:			
2050			

HIVE (orders)

ord-id	ord-date	ord-status	cust-id
1			
1000			
1001			
1002			
:			
2000			
2001			
2002			
:			
2050			

-- hive-overwrite

Last value method

sqoop-import --cs

-- hive-import

-- hive-database

-- hive-table orders

-- incremental append

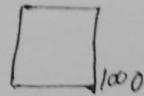
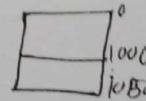
rx = hive -e "select * from (select * from "

-- check-column order-id

-- last-value 1000

last value from the hive table

-m 10;



source → destination

Last modified date method

-- incremental append

-- check-column ord-date

-- last-value '24-12-2018 23:59:59'

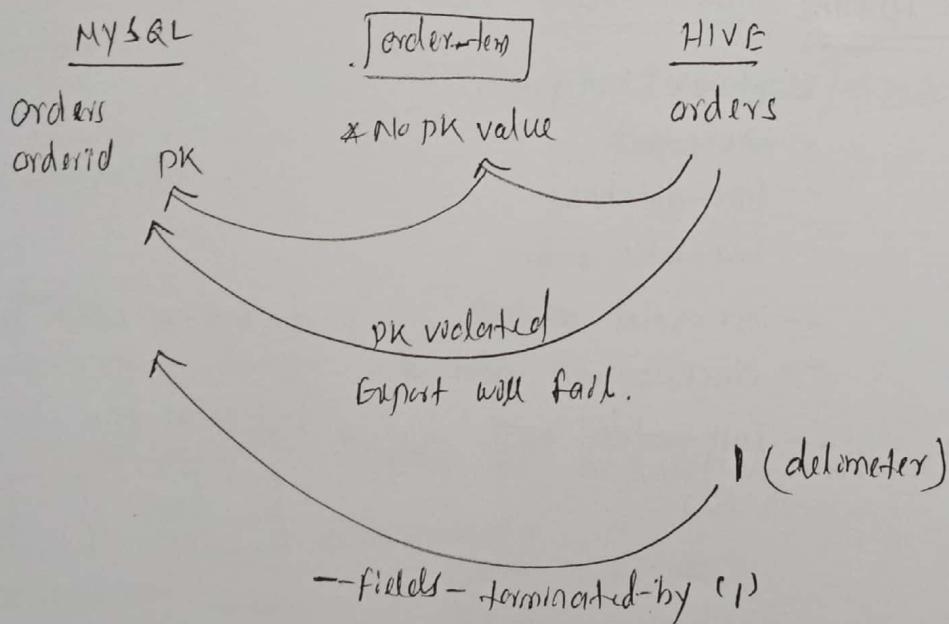
-m 10;

ord-id	ord-dt	ord-status	act-id
1	23-12	pending	1111
2	23-12	inprogress	1112
1	24-12	complete	1111
2	24-12	on-hold	1112
⋮			
1050			

Sqoop Export

```
Sqoop-export --connect username root --password cloudera --table emp
--export-dir '/user/hive/warehouse/dfs.db/emp-bk' --input
string '\n' --input-null-non-string '\n';
```

- MySQL is not part of hadoop ecosystem
- emp table u have to create in mysql then do import.



```
mysql> create table orders1 as select * from orders where ordid > 2100
Cloudera> hive import from sqh → Hive
hive> select * from orders1
```

```
mysql> insert into orders1 select * from orders where ordid between 10
hadoop> hive import from mysql → Hive → lower bound
hive> select * from orders1;
```

upper bound

Sqoop will find out.

sqoop job --create dvs-nov-incr --import --connect -- --last-value 1050;

1) If new table? last-value = 0

2) If existing table? last-value = max (ord-id) from hive

(a) Sqoop metastore - to store job name, details.

↳ sqoop job-list;

dvs-nov-incr

cloudera > sqoop job --exec dvs-nov-incr

No new records detected since last import;

mysql > insert extra 200 records → 1300 records

cloudera > sqoop job --exec dvs-nov-incr

hive > select * from orders

1300 → 1300 records

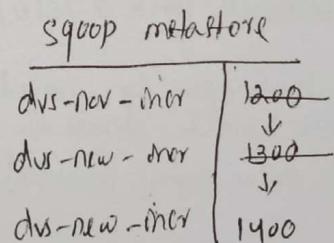
cloudera > sqoop job --show dvs-nov-incr

cloudera > sqoop job --delete dvs-nov-incr

order-dttsle ↓

> 00-00-0000 00:00:00

> 00-00-1900 00:00:00



Sqoop codegen

* To see the corresponding Java class files etc.

eval

Sqoop eval --connect

Hive Metastore

mysql > show databases;

mysql > select * from hive-metastore;

mysql > use hive-metastore;

```
hive> create database dvs-nov18;  
mysql> show databases;  
mysql> select * from DBS;  
-----  
dvs-nov18
```

MySQL tables → DBS → Database names

TBLS → Table names

columns-V2 → column names & their datatypes.

```
mysql> select * from columns-V2;
```

tableid	Column name	Type name	Integer index
16	cat-id	:	:
16	ord-date	:	:
16	:	:	:
16	:	:	:

HUE (Hadoop User Environment)

default path → /user/cloudera

hadoop fs -mkdir dvs

```
select order-status, count(*) from orders group by order-status;
```

→ SU -

→ SU - cloudera

→ jps

vi sample.sh

{

esc

:wq

:chmod 777 sample.sh

./sample.sh

PIG

26.12.2018
Wednesday

- High level data flow language
- susion oriented
- can work without hadoop
- no java required
- 200-300 lines of java code \cong 10-15 lines of pig script.
- for transformation.

A = load 'emp.csv' ... ;

* Bags of data to hold data.

B = filter A by deptno == 20;

C = order B by sal desc;

D = limit C 2;

▷ one by one flow of data A \rightarrow B \rightarrow C \rightarrow D

▷ pig >

grunt >

edit & c

& pig

grunt > then the actions will not work

→ pig stores nothing

→ pig can be stand alone on top of hadoop in simple linux environment.

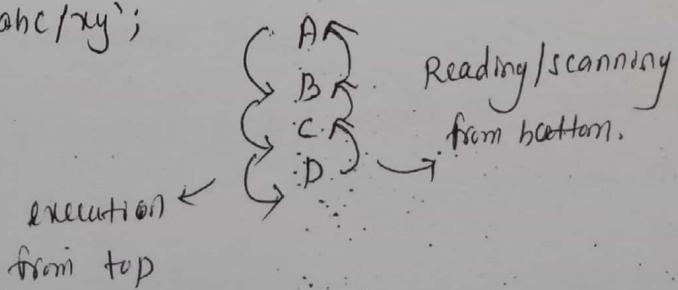


→ Lazy evaluation

→ pig doesn't store data permanently, give data then process, tasks done.

→ grunt > stored into local/xyz;

grunt > dump D;



DAG \rightarrow Direct Acyclic Graph



* By reading from top Pig prepares.

* At any time u can go to the command and

\$ pig -inctype local

grunt > pwd;

file:///home/cloudera

Anaconda

grunt > A = load "chme1training/dvtemp.csv" using PigStorage(',') as (empno:int, ename:chararray, job:chararray, sal:double, comm:float)

grunt > describe A;

grunt > B = filter A by deptno=20;

grunt > describe B;

* After = small

grunt > C = order B by sal desc;

* PigStorage - D.S - 0

grunt > D = limit C 2;

grunt > ctrl + L

\$ pig -x local

grunt > describe A;

] \rightarrow session mounted no data.

-- error ---

grunt > D = limit C 2;

-- only check syntax.

grunt > dump D;

] \rightarrow Lazy evaluation / advantages of pig
actual evaluation when give store / dump

Modes of execution of PIG

Local mode (uses LFS for i/p and o/p)

PIG - execute local (or) PIG -> local

MapReduce mode (uses HDFS for i/p and o/p)

PIG - execute mapreduce (or) PIG -> mapreduce (or) PIG (default)

PIG properties

Default mode is mapreduce in real-time, but local in Apache setup.

home/training/PIG/PIG.properties

grunt > stored into '/home/training/dvs/PIG-nov'

CD '/home/training/dvs/PIG-nov'

\$ ls

--- part-r-00000

\$ cat part-r-00000

logno	name			
1				
2				
3				
4				
5				

* default delimiter is tab

* loading also

grunt > stored into '/home/training/dvs/PIG-nov'

using PIGStorage(',');

\$ cat part-r-00000

728, Qab, analyst, 3000.0, 20

;

foreach --- generate

used to restrict columns from a bag.

used to derive new columns from existing columns.

used to apply column level functions.

grunt > x = foreach A generate empno, empname, job; → it u know the column name

x = foreach A generate \$0, \$1, \$2; → it u know the column positions.

$$2) \frac{A}{6} \quad \frac{B}{8} \quad 6+2$$

$$\text{Bonus} = \text{salary} * 0.005$$

$$\text{Incentive} = \text{salary} * 0.006$$

grunt > foreach A generate *, salary * 0.005 as bonus, comm * 0.006 as

warning for comm. as float

grunt > describe X;

3) grunt > y = foreach A generates ename, SUBSTRING(ename, 1, 3);

warning Implied - cast to double ...

grunt >

↓
starting point → length
to capture

CASE SENSITIVE

→ All table names are c.s.

→ All column names are c.s.

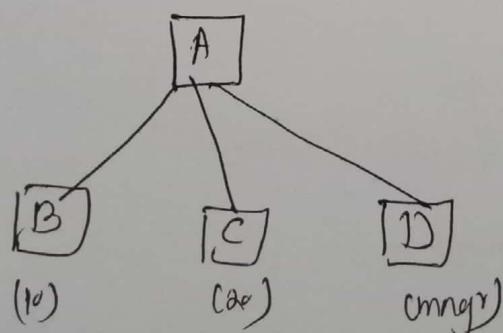
→ function like SUM, MIN, AVG, COUNT, SUBSTRING, LENGTH, TOKENIZE, FLATT

→ keywords like store, load, dump, using, as, join, cross, union, filter, or, limit etc. are not c.s.

9) B = filter A by salary > 2500;

9) B = foreach A generate substring(ename, 1, 3); X

B = foreach A generate SUBSTRING(ename, 1, 3); ✓



B = filter A by depth

C = " " " " depth

D = " " " " job

SPLIT

SPLIT A into B for deptno=10;
deptno = 10;
Dn = "mgr";

→ scans only 1 but the filter will scan 3 times, so best to avoid multiple scans to improve performance.

JOINS

g> B = load 'home/cloudera/DVS/dept.csv' using PigStorage(',') as (deptno:int, dname:chararray, loc:chararray);

g> describe A;

g> describe B;

→ EQUI (INNER)

→ LEFT OUTER

→ RIGHT OUTER

→ FULL OUTER

→ CROSS

g> dump A;

---: tuples
---:

g> dump B;

g> C = join A by deptno, B by deptno;

g> dump C;

g> D = foreach C generate 80, 81, 86, 87;

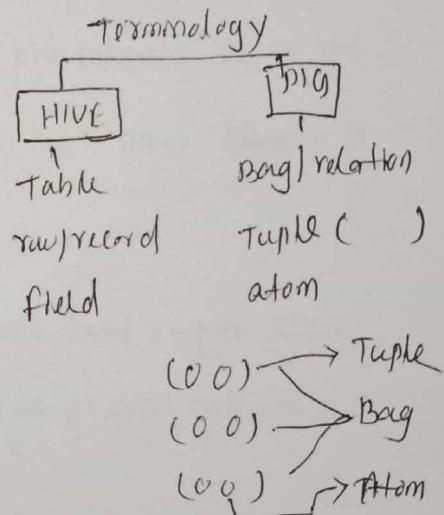
g> describe C;

g> D = foreach C generate A::empno, A::ename, B::deptno, B::dname;

g> dump D;

g> C = join A by deptno right outer, B by deptno;

g> dump C;



g) $C = \text{cross } A, B;$

PIG SCRIPT

(claudia@) DVS > cat sample1.pig

Max Pig
ll & hgl

SPLIT A into B at deptno==10, C at deptno==20, D at deptno==30;

store B into '/home/training/dvs/dept10-data' using PigStorage(',');

" C " " " /dept10-data" " " (");

" D " "/home/claudia/dvs/mngr-data";

> pig -f sample1.pig

→ default map reduce

> pig -x local -f sample1.pig

→ data in Local so change to local

DVS > cat sample2.pig

A = ~~load~~ LOAD ('8a') using PigStorage(',') as (conas int, ename:chararray)

DVS > pig -x local -param a=/home/claudia/dvs/emp.csv -b=/home/claudia/dvs/b-out -c=/home/claudia/dvs/c-out -f sample2.pig.

x - execute

p - param

m - param-file

DVS > cat pig param.txt

DVS > pig -x local (-param-file) pigparam.txt -f sample2.pig

pig -x local -m pigparam.txt -f sample2.pig

difference between Local & Mapreduce mode

Group and aggregate function

deptno	min(sal)	max(sal)
10		
20		
30		

g> B=group A by deptno;

g> C=foreach B generate group, MIN(A.sal), MAX(A.sal);

g> D=foreach B generate &0, MIN(A.sal), MAX(A.sal);

Grouping

▷ Based on one-column

▷ Based on multiple-column

▷ all data in one single group.

B=group A by (deptno, job);

dump B;

▷ A=emp.csv

→ min, max,

B=group A all;

g> dump B;

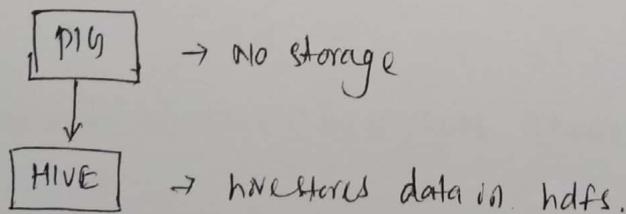
g> describe B;

g> C=foreach B generate MIN(A.sal), MAX(A.sal), AVG(A.sal), SUM(A.sal), COUNT(&1);

g> dump C;

→ &0 = ~~group~~ all
&1 = { } --- }

- 28.1
Date: 2023-01-12
- 1) No schema ($A = \text{load 'emp.csv'}$);
 - 2) full schema ($A = \dots$);
 - 3) semischema ($A = \text{load 'emp.csv' using PigStorage(',') as (emp, ename, ...)}$);
- 1) data will come with comma separated \rightarrow tab separated file
- $\triangleright B = \text{load '/home/claudia/DVS/emp.tsv'}$;
- * dump always represent in form of comma delimiter.
 - * storage will be tab separated.
- 3) $A = \text{load ('/home')... using PigStorage(',') as (empno:int, ename:chararray, job, sal, comm, deptno);}$
- $\triangleright \text{describe } A;$
- $\triangleright \text{ename:chararray, job: bytearray, sal: bytearray...}$
- $\triangleright B = \text{filter } A \text{ by sal desc;}$
- * `bytarray` starts operation with ASCII values, so the ordering will not come properly.
 - * `transformation` / `arithmetic` / `mathematical formula` will not work with `bytarray`.



- * Any Reporting tool talk with HIVE, but not with pig, as file can be read but table in HIVE can be ~~not~~ accessed by any tools.

Unstructured data

$\triangleright \text{cat test.txt}$

O/P (Word occurrence)

Ths, 2

Ms, 1

DVS, 1

grunt > A = load 'hdfs/cloudera/DVS/Test.txt' as (line: chararray);

> describe A;

A: {line: chararray}

→ test the default delimiter with

> dump A;

space.

TOKENIZE

TOKENIZE (col-name, 'delimiter')

Default as 'space'. Apart from space, we need to specify delimiter.

NOTE Excel - Transpose

SQL - PIVOT

* Converting row into the columns

HIVE - EXPLODE

PIG - FLATTEN

({ethn}, {is}, {day})

→ Nested tuples.

Steps

8a

1) A = load 'hdfs/cloudera/DVS/Test.txt' as (line: chararray);

2) B = foreach A generate TOKENIZE(line) as Token;

3) C = foreach B generate FLATTEN (token) as flat;

4) D = group C by flat;

5) E = foreach D generate 80 as word, 81 as word-occurrence;

6) F = store E into '8b' using PigStorage(',');

7) G = ~~foreach~~ order F by 81 desc;

8) H = filter G by 80 != 'is';

9) I = limit G 1;

rating < -1

grant > explain sample.pig

> run sample.pig

* real-time implementation for pig
lets as not strong data for
retrieval.

> Cat emp.csv imp.csv > empp.csv

> A = load empp.csv

1> distinct records

2> distinct jobs

3> B = distinct A;

4> C = foreach B generate job;

5> D = distinct C;

6> dump D;

or > group by job ("Analyst", 1, 3))

("Manager", 2, 1))

NOTE → distinct is better than grouping for getting data.

→ distinct can be applied on the bag not a column or subqueries.

⇒ A = load "emp.csv" --;
B = filter A by deptno = 20;
C = order B by sal desc;
D = limit 10;

⇒ D = limit (order (filter A by deptno = 20) by sal desc) 10;

* innermost query gets executed first followed by outer queries.

⇒ JSON DATA

cat sample.json

grunt > A = load '---' /sample.json'

(food: chararray)

JSON loader → JSON to CSV.

JSON storage → CSV to JSON

A = load exp CSV

convert from CSV to JSON.

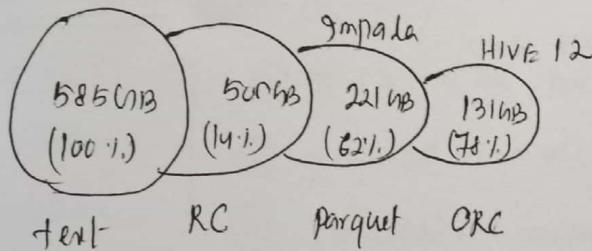
{ "empno": "7369", "ename": "smith", "job": "clerk", ... }

Store A into ' ' using JSONStorage

RC - Record columnar format

ORC - optimized Record columnar format

{ 31.12.2018 }
Monday
max max

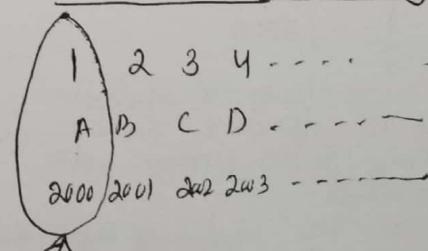


Encoded with RC / ORC / Parquet / ORC

Raw based storage

id	sname	year
1	A	2000
2	B	2001
3	C	2002

columnar based storage



HDFS storage → 1, A, 2000

2, B, 2001

3, C, 2002

4, D, 2003

RC file format

DVS > JPS
SU-
JPS

HIVE

→ create table emp ()

row format delimited fields terminated by ',';

H> load data local inpath '/home/claudera/DVS/emp.csv' into table emp;

Hadoop fs -ls ./user/hive/warehouse/ds-nv18dtb/emp

Hadoop fs -cat /emp/emp.csv

H> create table emp-rc (

terminated by ',' stored as rcfile;

H> insert into emp-rc select * from emp;

1> How data stored in hachend

2> How data retrieval happens

RB

	Yop	filu
→	5000	
→	6000	
	7000	

→ It will traverse row by row with each records in row based

RCBIL

X A

X 2000

✓ 5000

→ Aggregation data retrieval is faster.

Hadoop fs -cat /-----. /emp-rc /00000-0

---- only understood by HIVE.

Static partition

- from CSV → Table
- LFS format \Leftrightarrow hadoop format
- Hive doesn't read the data only take and put in the required directory
- Hive uses schema on read not write.

Traditional DB Schema on "read"
↳ Schema on "write"

cat emp.tsv

null null null null null

hadoop fs -du emp-rc;

DML for HIVE

Insert
Update
Delete

- ① It should be bucketed table
- ② It should be ORC file format
- ③ Transaction manager locked

emp

00000-0

00000-1

...

00000-10

small no of large files \checkmark designed for

large no of small files \rightarrow

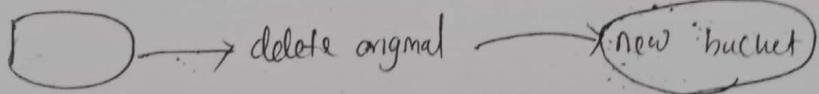
- metadata increases

- NN carries all the metadata.

- 1KB metadata get generated for 1 block of data in Data node.

* Metadata should be less

Update



- copy of file will be taken, modify, delete, & replace.

dynamic partition

- buckets
- file format

Temporary table insertion.

- LFS format \neq hadoop format
- Some manipulation is done

while storing, you have choice

for some intelligence to do with temporary table.

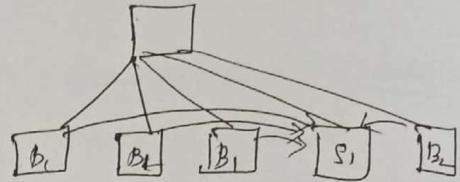
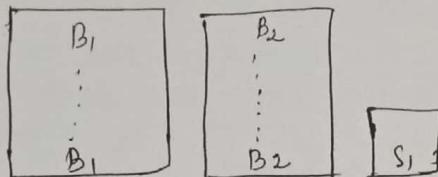
- No history maintained only updation/deletion done.

HIVE OPTIMIZATION TECHNIQUES:-

1) partition / Index/ Buckets

2) Map Joins (Mapper side joins) complete join happening in mapperside.

set hive.auto.convert.join=true;



→ Replicated over the buffer/cache memory at each data node.

→ To reduce the data transferred over the Network.

→ Session oriented set hive.auto.... it changes in hive-site.xml and will be for ever for the whole cluster.

3) Vectorization (instead of processing record by record, it will process bunch records at once)

Default → 1024

HIVE > set hive.vectorized.execution.enabled=true.

4) fileformats (ORC, PARQUET)

5) controlling the no. of Reducers

hive> set hive.mapred.reduce.tasks=4;

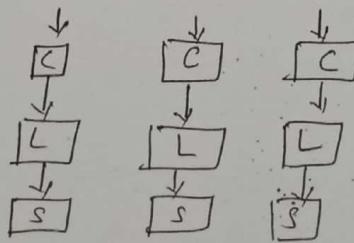
6) Parallel execution

hive> set hive.exec.parallel=true

sample.hql

```
create
create
load
load
select
select
```

Direct Acyclic Graph



NOTE → By default Reducer M sqoop =4

- 8) Data De-normalized data (saves lot of time in terms of joins)
 - Joins will be faster.
- 9) final output compression (compressing final output in HDFS)
 - SET hive.exec.compress.output = true;
 - SET hive.mapred.output.compression.type = block;
 - SET mapred.output.compression.codec = org.apache.hadoop.io.compress.SnappyCodec
- 10) Intermediate compression
 - hive.exec.compress.intermediate = true;
- 11) CBO (cost based optimization)
 - set hive.cbo.enable = true;
 - set hive.compute.query.using.stats = true;
- 12) Using TBZ Execution engine
 - set hive.execution.engine; → MR
TEZ
SPARK

TBZ - only one read and one write

- get the data from the disk, perform operations and give output, instead of having 5-6 disk read/write like MR.

PJO OPTIMIZATION TECHNIQUES

- 1) push up filters to avoid unwanted data and add for each to remove unneeded column.
- 2) if there is a chance, avoid filter and prefer split by.
- 3) Limit optimizer - Limit as early as possible.
- 4) pushDownJoins - Apply join as late as possible.
- 5) specialized joins like replicated, skewed & merge joins.

9) Error handling
pig -x local -f sample.pig

-f or "stop on failure"

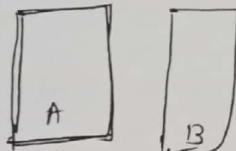
10) prefer DISTINCT over group by / generate

A = foreach A generate $\{ \}$;
B = distinct A;
C = foreach B generate 80;

spatialized joins in pig

1) Replicated - Map side joins

2) Merge Joins



→ When 2 datasets are bigger, more memory.

A₁ = order A by deptno

B₁ = order B by deptno

→ C = Join A₁ by deptno, B₁ by deptno using

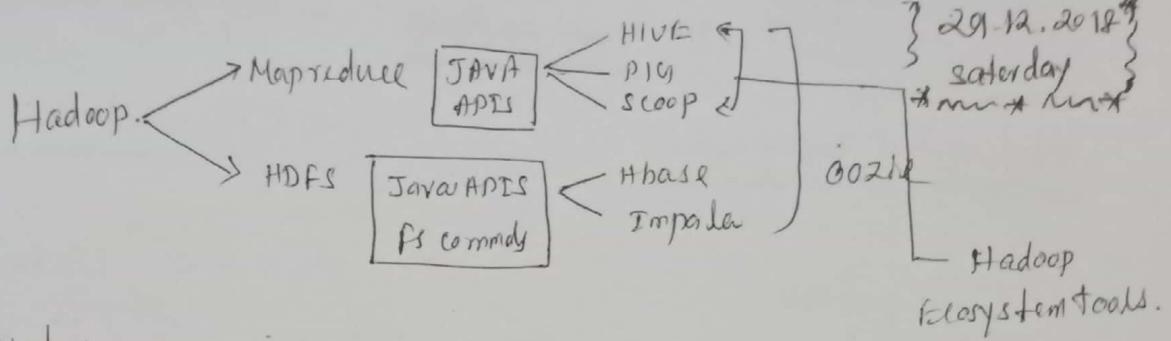
3) SKewed →

R₁ R₂ R₃ R₄

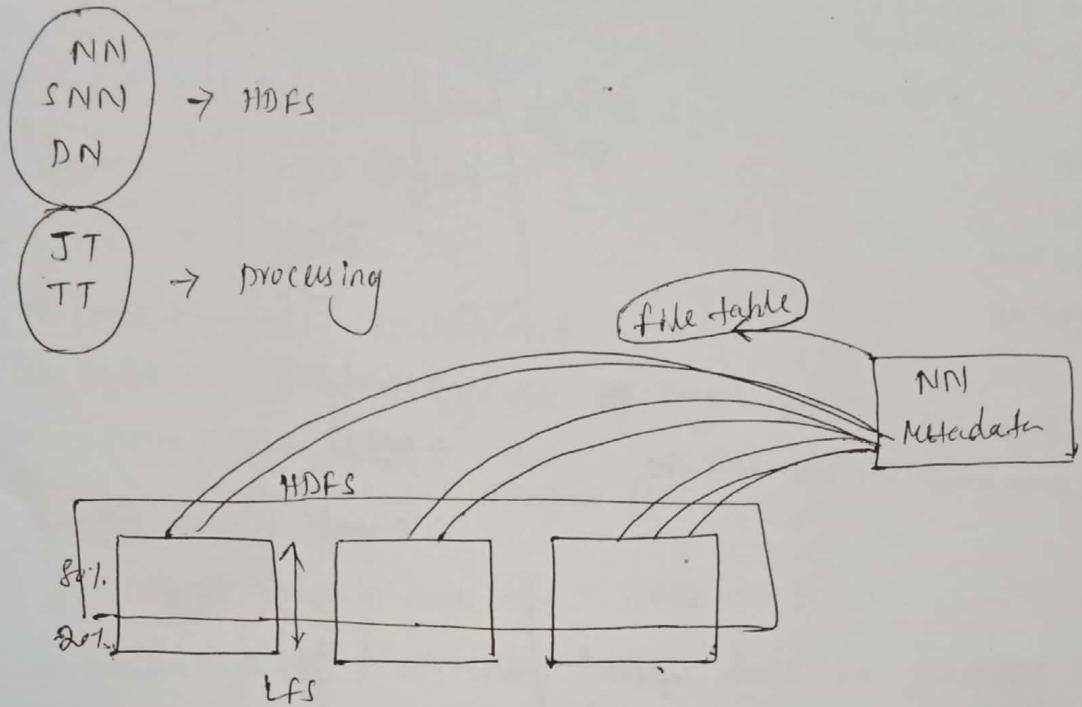
5 5 10 10 min

→ Eventy distribute the Reducer task to the free Reducers.

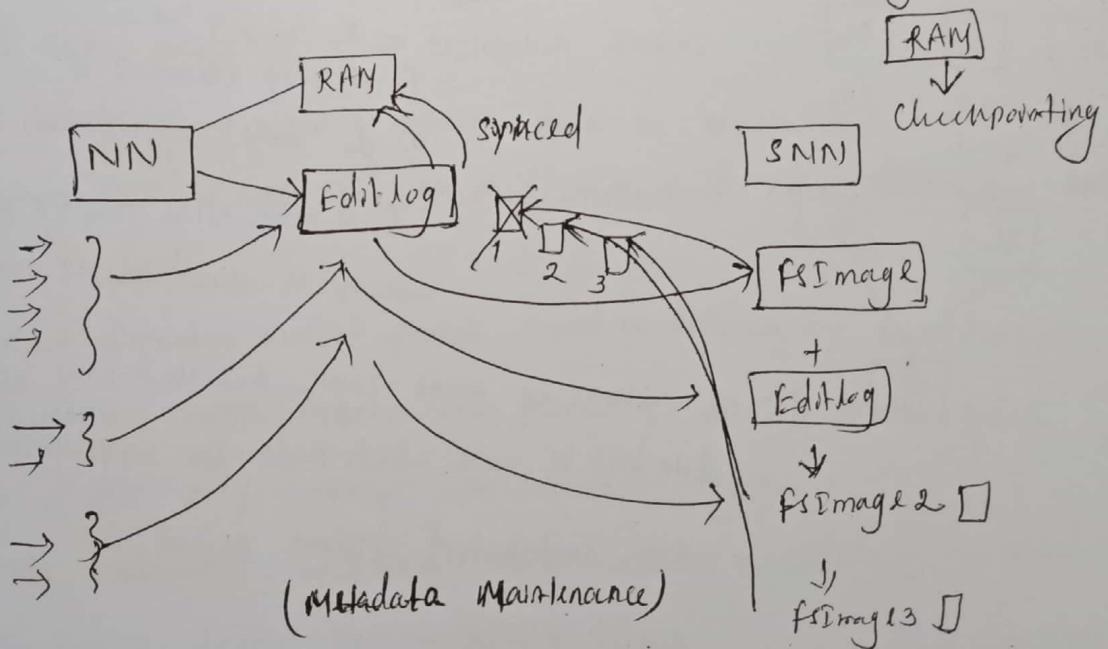
MAP REDUCE



Services for Hadoop

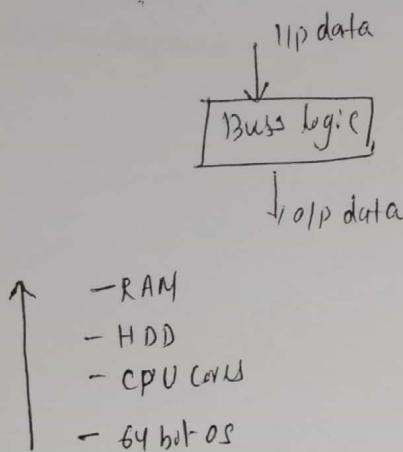


each LFS shares their space for the HDFS filesystem storage.

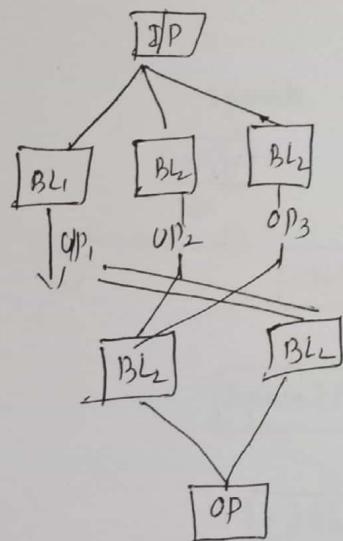


* only 2 images will be there.

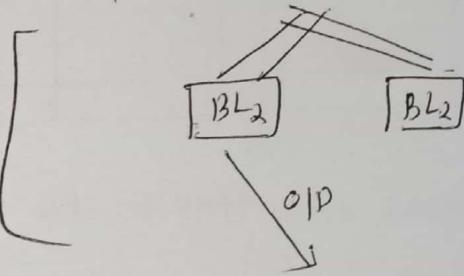
* return NN the fsImage3 will be go to RAM.

MAP REDUCESingle system computing model.

- vertical scaling
- As the data grows the performance decreases.



reduce
phase



- 1) do only selection
- 2) do only aggregation → sum, group, SUM, MAX, MIN, group by all aggregation will be in 2nd phase

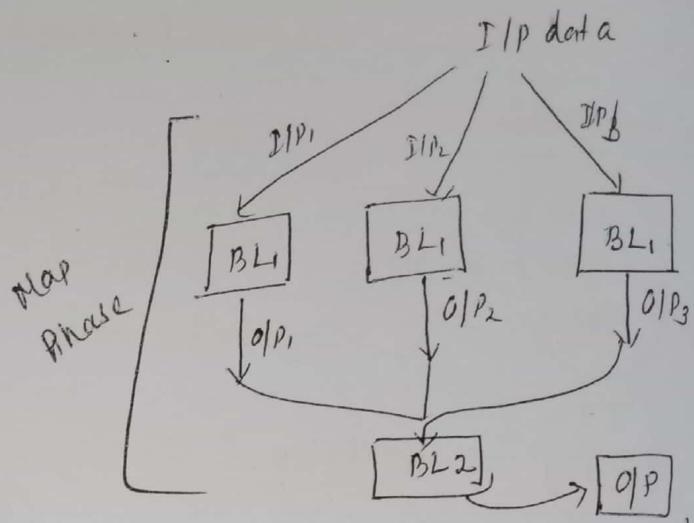
* we have to write logic in such a way that selected data will be go to which node for aggregation function

* This distributed system called as

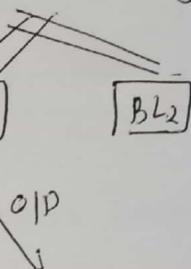
Approach DSCMAP-REDUCE

- 1) Map-reduce
- 2) MPP - Massively parallel process
- 3) CORBA

Distributed system - Computing model



- select prod-id, from selectable graph
- partial aggregation
- 2 level business logic

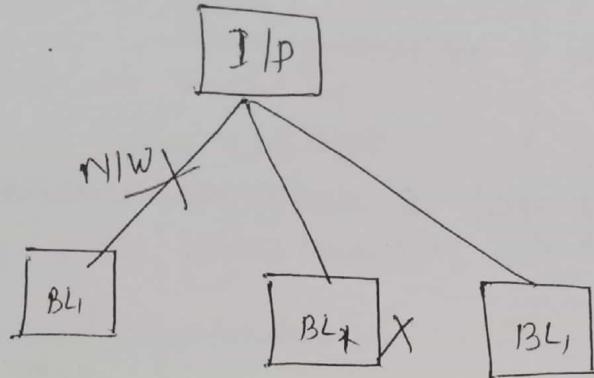


product id
subunit

- 1) do only selection
- 2) do only aggregation → sum, group, SUM, MAX, MIN, group by all aggregation will be in 2nd phase

Map and reduce phase can be numbers depending query.

3



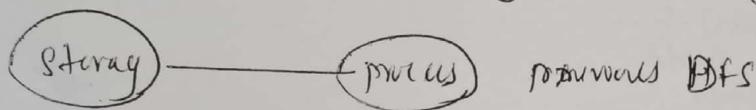
* If N/W will fail, how to transfer data to systems.

> After copying the node failed

> How data distributing.

> How to ensure phase 1 is completed.

> How to ensure final output go the ~~sys~~ single system.



* data and process same node - HDFS

* No need to go whereabouts all the time for processing.

* processing sent to node where data.

* High availability \rightarrow replication factor

* Reliability \rightarrow without failure data in required node.

* same API is used for structured, unstructured same Map reduce API will be used.

* Input formats output format should be followed by Map reduce.

* Job tracker keeps ensure the processing the job at the nodes.

* As a user only writing the business logic

* The framework will take care of all the activities, so in market Hadoop is highly acceptable.

* write once and read many times, concept.

Overview of MapReduce

- 1) Distributing computing model - programming model.
- 2) Written Java
- 3) It provides MR APIs (JAVA APIs) to develop MR applications.
 - ↓
 - * classes
 - * interfaces
 - * abstract classes, * static classes.
- 4) Available as `hadoop-core-v.jar`.
- 5) By default the MR applications are developed in JAVA.
 - ↓
 - can be in C/C++/scripting (shell, python, ...)
- Native support $\xrightarrow{\text{C/C++}}$ hadoop pipes
- other language \Rightarrow hadoop streams
- 6) Data representation in Mapreduce should be always as a key & value pair.
 - ↓
 - aggregation, partitioning will be easy.
- 7) MR applications generally have 2 phases
 - 1) Map phase \rightarrow select part of query
 - 2) Reduce phase \rightarrow group by and aggregate functions.
- 8) MR Application / Job Algorithm = [task] = [Map tasks] + [Reduce tasks]
- 9) MR always read/wdata from HDFS (filesystem) at Hadoop.C.HDFS
- 10) The Map logics and reducer logic will take / give input / output as key & value.
- 11) Though we use Java to develop MR application, the data i.e. key value pair is represented as hadoop datatypes, because hadoop uses its own serialization to represent key & value objects.
- 12) Map logic will be written in a mapper class and reduce logic will be written in a reducer class.
- 13) Set Job properties & Running \rightarrow Driver class

Mapper class() \rightarrow map()

Reducer class \rightarrow reduce()

Driver class \rightarrow main()

using a Job object

*.jar

> hadoop for *.jar Driver class name top-path o.p-path

Map - structured / semi/ unstructured

Hive - structured

Pig - structured / unstructured

MapReduce Architecture

1. Master-slave Architecture

Job Tracker (JT)

Hadoop 2.x

Task Tracker (TT)

\rightarrow Resource Manager (RM)

\rightarrow Node manager (NM)

Job tracker:

responsible for resource management

responsible for job scheduling

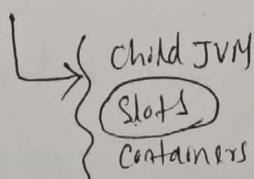
responsible for job execution \rightarrow job coordination, failure management etc.

maintaining job history \rightarrow

Task tracker

executes tasks assigned by Job tracker.

The task can be Map task or the Reducer task.



\uparrow JT \uparrow

\uparrow TT \uparrow
 \rightarrow Task status or
progress information

Resource manager

\rightarrow Resource management

\rightarrow Job scheduling

\rightarrow Job execution is delegated to Application Master (AM).

\rightarrow only one job can be taken care by the Application Master.

\rightarrow Job history will be maintained / handled by job history server. (JHS)

JT & TT are daemon processes, while RM & JHS are daemon processes of back end process running.

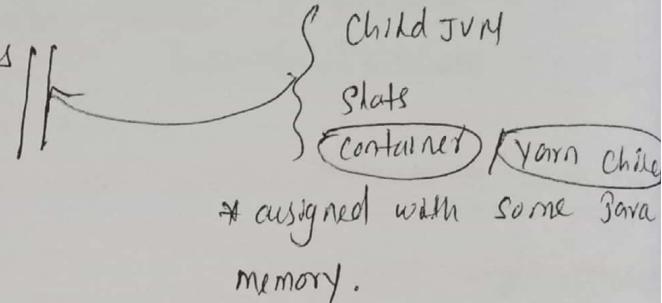
AM & JHS are servers not servers

* one task one JT/RM

in 1.x JT have heavy tasks of memory as many to do

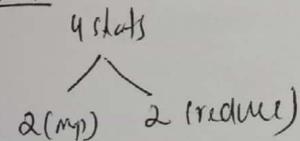
Node manager

- Execute the tasks - either MR tasks
- provides container for AM.



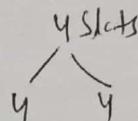
RM instructs the NM to allocate the memory for the task and application manager as assigned for the job execution.

1.x



→ reservation for the slots.

2.x



→ No reservable for slots

MyMRApp.jar → MyMapper. class

MyReducer. class

MyDriver. class

To get the To run the job \Rightarrow hadoop jar <jarfile-path>

<driver-name> <ip path> <op path>

\Rightarrow hadoop jar /home/hari/MyMRApp.jar org.ous.mr. mydriver
LFS
fully qualified class name
HDFS

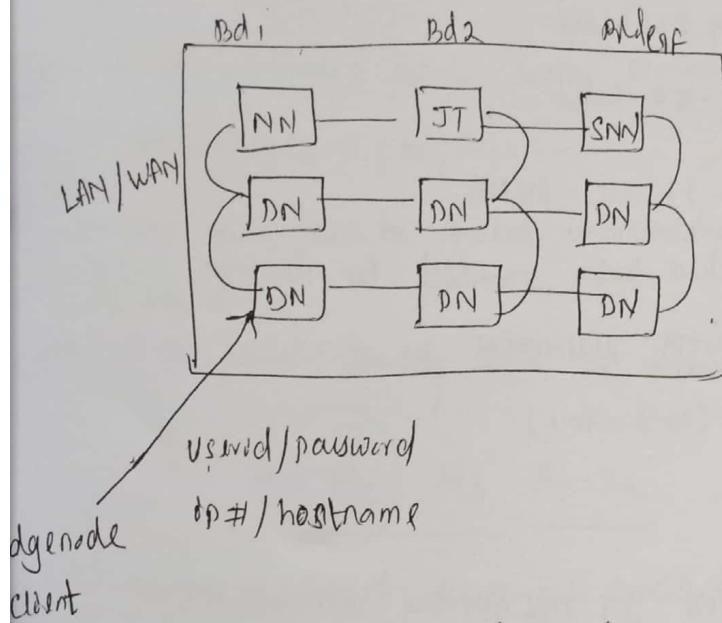
* According to memory slots are decided. 512MB RAM for one container.

hadoop jar → JT

hadoop fs → NN

dfs
dfsadmin → NN with admin

→ shell scripts where we can pass the parameters as fs, dfs, admin
* in hm directory it can be found.



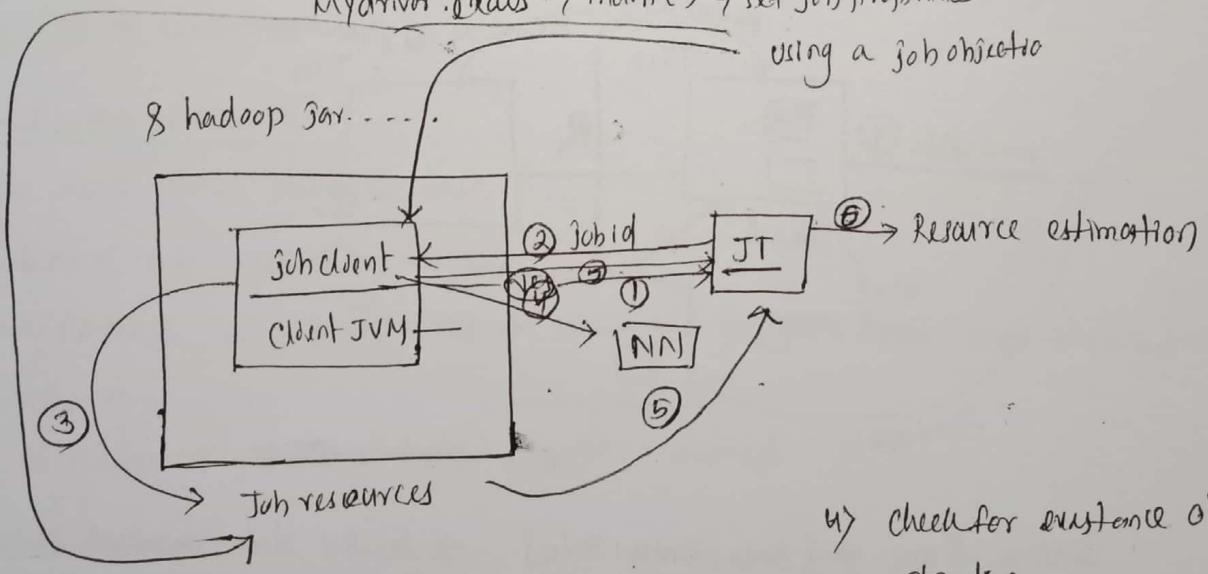
* ssh protocol, they can talk to each other.

* one system access should be done for a developer

* from edge node you can run jar files.

Mydriver.java → main() → set job properties

using a job object



4) check for existence of a/p dir.

→ job client will request JT

→ JT provide jobid to reference future.

→ job client prepare job resources.

{ - jobid info
- job configuration file (job specific properties)

- { - Input splits into - blocks nos.
- make a copy or jar file \rightarrow MyMRApp.jar
- make a copy of miscellaneous files/supporting files.

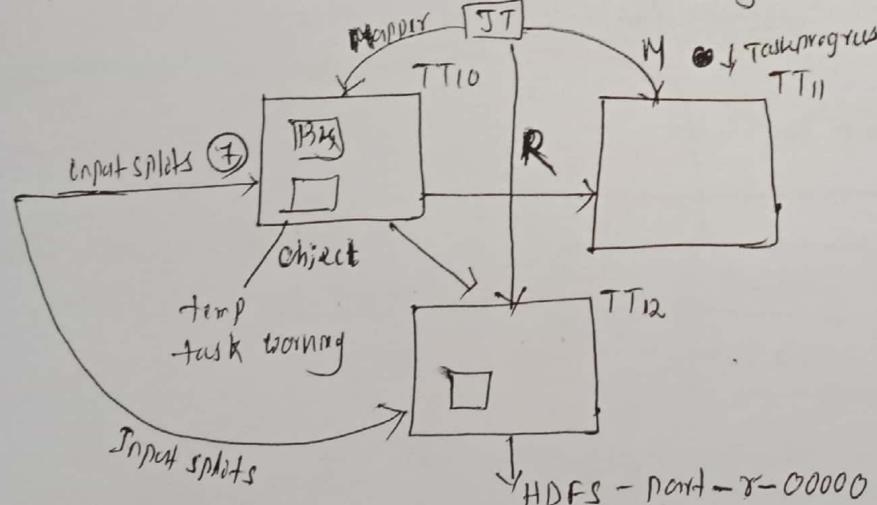
- \rightarrow copy to HDFS \rightarrow cluster(nodes) shared location.
- \rightarrow check for instance it will submit to JT, by sharing the job resources.
- \rightarrow Resource estimation? estimating # M-tasks
and/or # R-tasks

$$\rightarrow \# \text{ Mappers} = \# \text{ Input splits} - 1 \text{ block} - 64 \text{ MB}$$

$\# \text{ Reducers} =$

- ① mapped. reduce-tasks provided by developer. e.g. = EXPLICIT
- ② The framework estimates the #reducers, with an algorithm. (10% - 90%) IMPLICIT

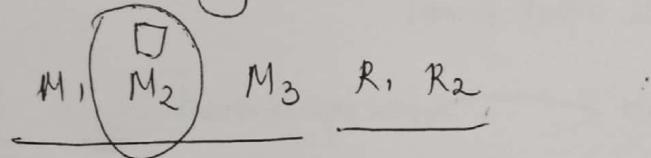
- \rightarrow Heartbeat & status information JT will get the information.



$$\text{file} = \begin{cases} \text{BK1} \rightarrow \{ \text{DN}_{10}, \text{DN}_{15}, \text{DN}_{13} \} & \rightarrow \text{locally data available for the file} \\ \text{BK2} \rightarrow \{ \text{DN}_{11}, \text{DN}_{21}, \text{DN}_{22} \} \end{cases}$$

\rightarrow assign the map task.

\rightarrow If DN - not free for BK, it will assign to adjacent node and then copy and assign the map and process & delete temporary data.

- Mapper definition will be read
 - Create the object in the TTIO
 - Once the map task done it will assign reduce task.
 - The reduce task distributed data availability/locality.
 - After getting all information reduce will do aggregation and gives output to HDFS.
- ⇒ If the job failed due to wrong status
- $DN_{10} \rightarrow DN_{15} \rightarrow DN_{16} \dots$
- If any task can be failed continuously 4-time, the entire job will be marked as failed.
- ⇒ Other running task are killed.
- 
- If one is failed whole task cancelled.
- ⇒ All the resources will be released.
- 2) Speculative execution -
- To avoid process straggler task.
 - Task is not responding to map.
 - Who finishes 1st it will take the result and make the other node to cancel the task.
 - At a time 1 speculative task can be started.
- 3) Hard disk failure
- will be handled on other node, it will take all the information from the HDFS.

3 2.01.2
Wednesday
3/11/2016

\$ hadoop jar ~~mapred~~ driver ip o/p

↓
Main()

↓
at present job will not process.

M = # input-splits

R = mapred. reduce. tasks ← dev

F/W itself will do the estimation.

JT will find TT is free and they have free slots to get the task, Then assign the map to TT. Some of the reducer works parallel to mappers. Reducer will wait till mapper job done.

10% - 90% - number of reducers as per mappers.

input split resides at the file input format.

Mappers → input splits ← file input format

R ← mapred. reduce. tasks ← Jobobj. setnumreducers (n)

1. X

2. X (YARN) - Job tracker heavy in X. / restriction for number of M or R task
- MR is rewritten.

Master - RM, JHS, AM/APP

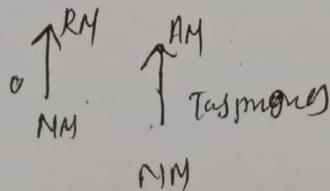
Slave - NM (Map logic, Reduce logic, AM logic)

Yarn jar · ~~jar~~ driver ip o/p

↓

Yet Another Resource Negotiator

AM estimates the mapper and reducer and request to Resouce m



JHS

AM

- everything done AM handle history on JHS.
- cleans all memory, garbage collector.
- HDFS stored location to hold all the input splits.
- Temporary working directory present in AM will get the input splits and calculate how many mappers and reducers needed and inform to resource manager to get the free container/JVM in node manager.

- 1) AM Task
- 2) Map Task
- 3) Reducer Task.

TT \rightarrow JVM

\downarrow M₁

TT \rightarrow JVM

M₁ \rightarrow JVM
 \swarrow child

JVM
 \swarrow child

if one fails another will continue to work.

- + once the HDFS written, the AM will maintain history, the signal goes to resource manager, then resource manager will inform node manager to release the memory and resources for the application manager.

Hadoop 1.x

spark library

\downarrow local

client } resource manager
 cluster

\downarrow

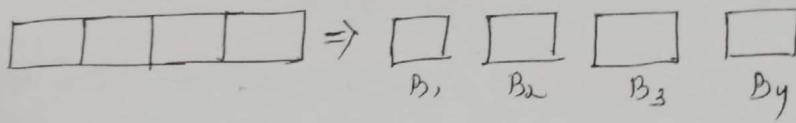
\rightarrow MOSC
 \Rightarrow TT

Hadoop 2.x \rightarrow As spark already did the resource negotiation, in hadoop 2.x

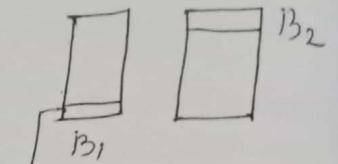
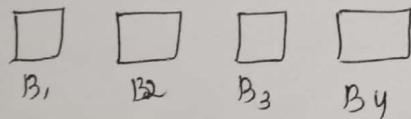
named as yet another resource negotiator.

Writing

pos-2018.csv



Processing



→ Array out of band
60 records to return

where block have no rec

input splits → Logical blocks

IS1 = B1 + B2* → define data boundaries

IS2 = B2 → unit of data on which a map can be run log C

IS3 = B3 + B4* → # M = # IS.

IS4 = B4

Properties

Input split size → Min

→ Max

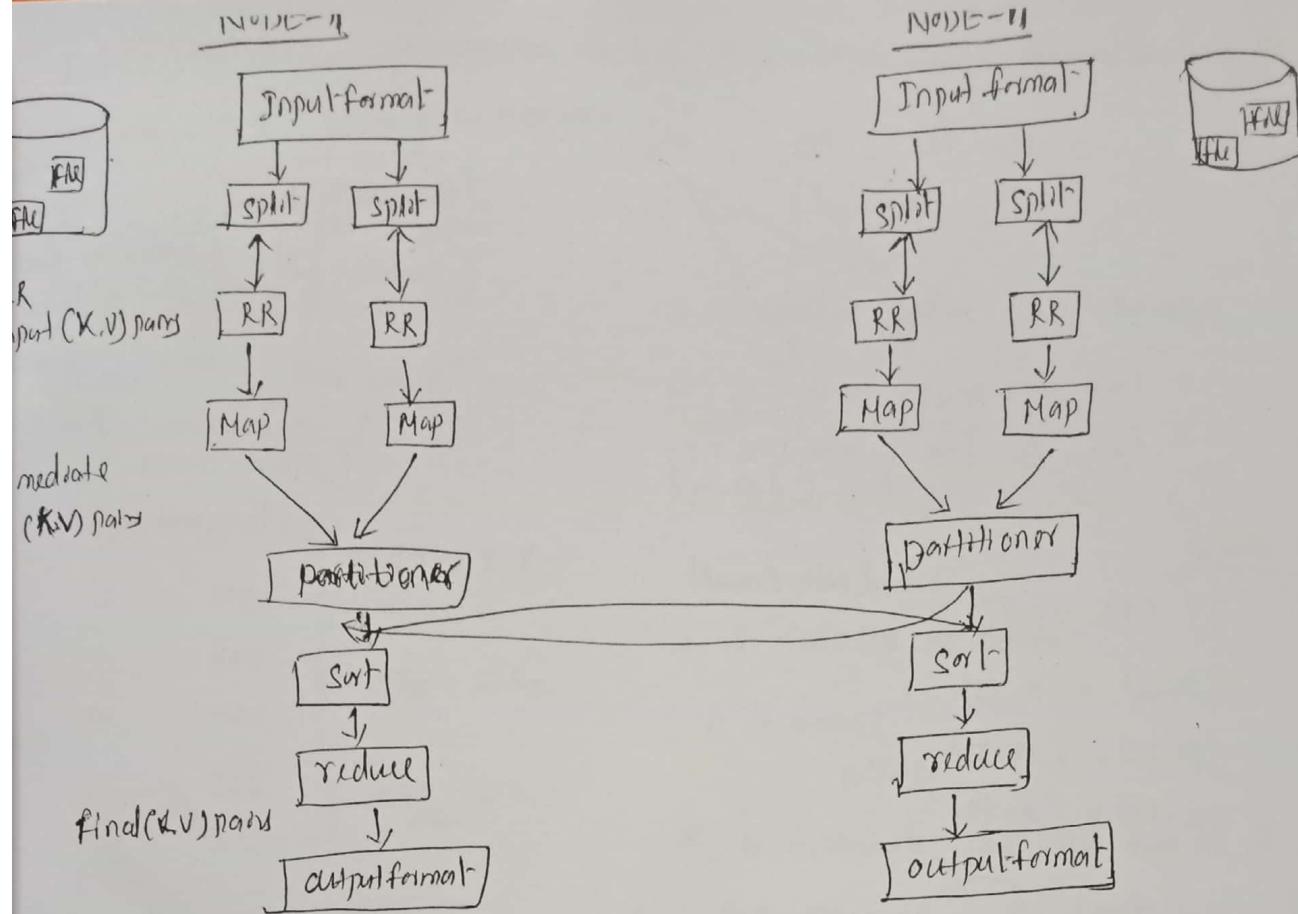
Input split size = block size

NOTE → Hadoop read is faster than writing.

→ Writing is sequential, reading is parallel.

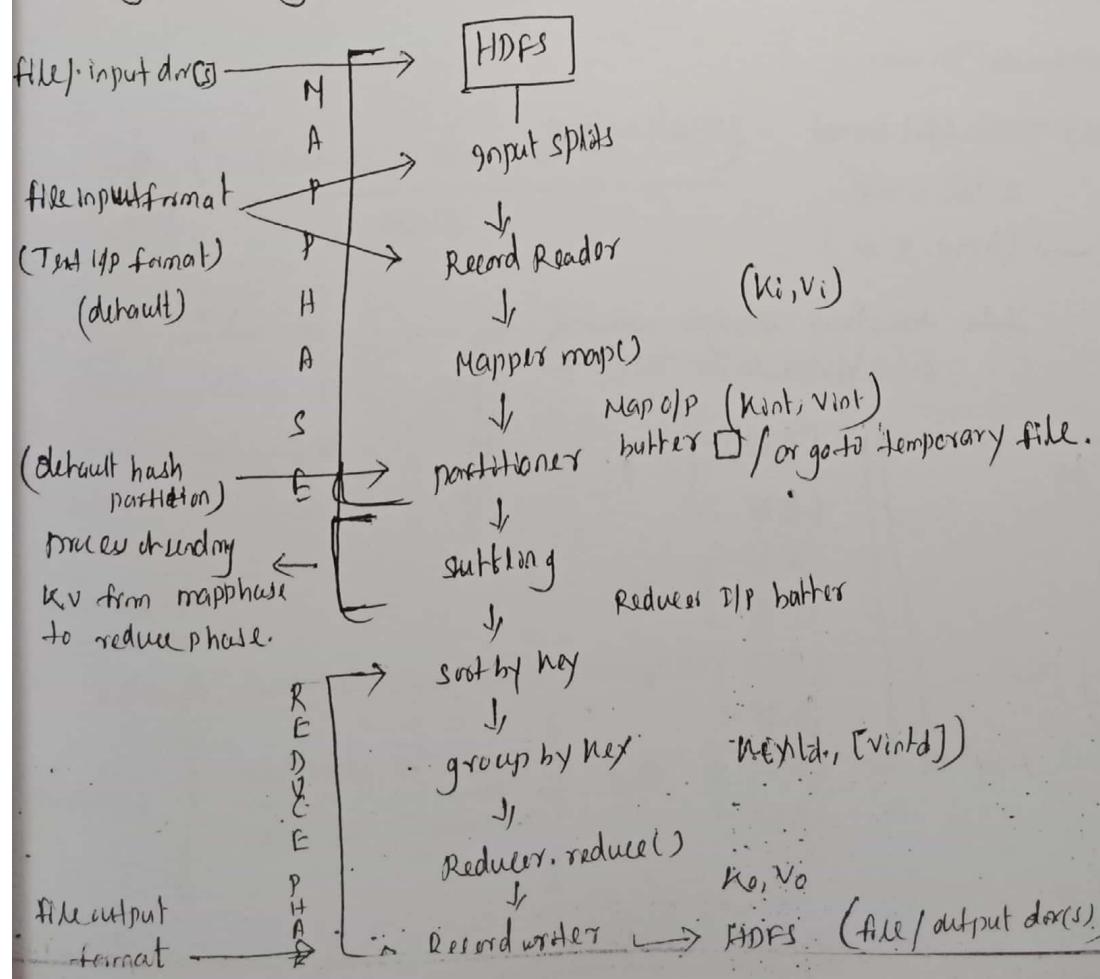
→ Pipeline writing / sequential writing.

- 1) How to see the block sizes of a particular file after copying to hadoop.
- 2) Getting the files from the materials folder.

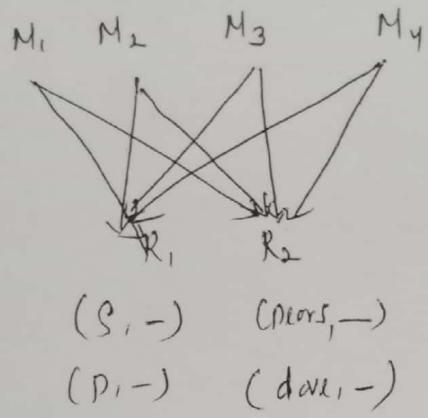


partitioner will decide which pair will go to which reducers.

Sorting, grouping done by the framework, then data will go to outputformat.



→ partitions will ensure same KV gets the corresponding Reducer.



(key, value)
(S, -) (P, -) (d, -) comes to R₁
(P, -) (d, -) comes to R₂

output dir
HDFS
(driver class)
(job properties)

Output

(S, C) = (S, TS) → total sale amount

(P, C) = (P, TS)

(P, C) = (P, TS)

(d, C) = (d, TS)

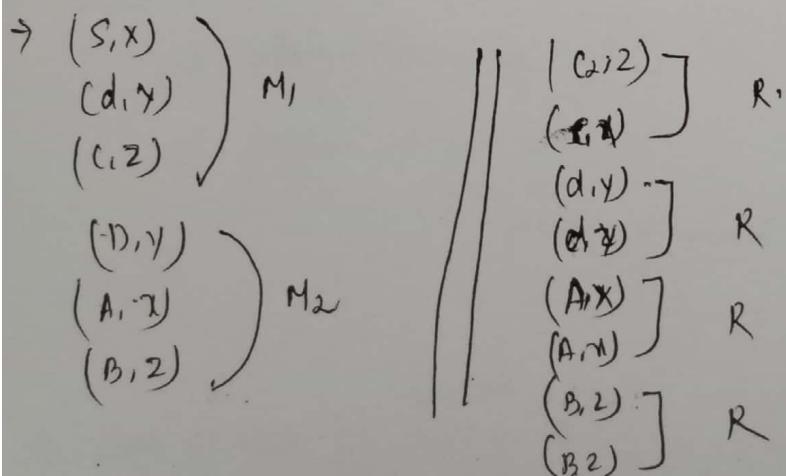
→ semi-structured
structured } Entity extraction
structured } → it will try to convert as much
as possible to structured data.

→ describe formatted table

(text input format)

Hive lineage key to output format - serialization.

→ shuffling → (partition, value)
data transform through network



Reducer $> 0 \rightarrow$ Sorting in mappers
Reducer $< 0 \rightarrow$ Sorting in Reducers

<u>Date</u>	<u>Transd</u>	<u>prodid</u>	<u>amount</u>
9/12/2016	Two	A	100

01/21/2016 1001 B 200

9/21/2016 1002 E 300

→ Line number will be the offset value.

(0, 2017-01-01, 1000, A, 100)

(1, 2016-01-01, 1000, B2, 200)

(2, 2016-01-01, 1002, 1, 300)

$$\begin{array}{r} 100 \\ 100 \end{array} \quad \boxed{J} \quad 65/2 = R,$$

$$B \quad 200 \quad] \quad 66/2 = R_0$$

B 200 7 661 - P * 2 redlers are there

$$B_{200} \int w_2 = \infty$$

$$C \quad 300 \quad J \quad 67/2 = R,$$

$\mathbf{n} \rightarrow \mathbf{map}()$

antimes

→ to distribute the load more no of Reducers.

Example

-file Input format-

✓	2018-01-01	1234	Santosh	1	26
✓	2018-01-01	1234	debtal	2	40

we will decide how to read

* one record per one line

* converts to key, value

* 0..121 character as repeat symbols

offset value

- 0 $\hookrightarrow 0 \xrightarrow{120} \text{v0}$
- 1 $\hookrightarrow 121 \xrightarrow{120} \text{v1}$
- 2 $\hookrightarrow 242 \longrightarrow \text{v2}$
- 3 $\hookrightarrow \text{v3}$

* offset value calculator.

Partitioner

default partitioner \Rightarrow Hash partitioner

$$\text{Modulus } \left\{ \frac{\text{hashcode of Key}}{\text{num of Reducers}} \right\} = 0 \text{ to } (n-1) \Rightarrow \text{Reducer Index}$$

partitioner



int getPartition()



Keys $n=1$

10	$ 10/1 = 0$	$ 11/1 = 0$	$ 12/1 = 0$	$ 13/1 = 0$	$ 14/1 = 0$	$ 15/1 = 0$
----	--------------	--------------	--------------	--------------	--------------	--------------

 $n=1$ $r_0 \rightarrow 11 - \text{part-}r_0$ R₀ (one reducer)

Keys $n=2$

10	$ 10/2 = 0$	$ 11/2 = 1$	$ 12/2 = 0$	$ 13/2 = 1$	$ 14/2 = 0$	$ 15/2 = 1$
----	--------------	--------------	--------------	--------------	--------------	--------------

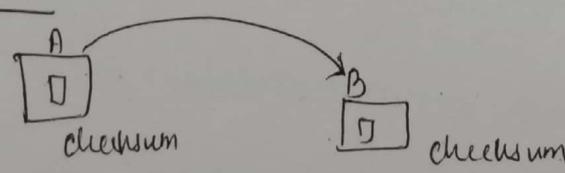
 $n=2$ $r_0 = 6$ $r_1 = 5$ part- $r-00000$ part- $r-00001$ R₁

Keys $n=3$

10	$ 10/3 = 1$	$ 11/3 = 2$	$ 12/3 = 0$	$ 13/3 = 1$	$ 14/3 = 1$	$ 15/3 = 2$
----	--------------	--------------	--------------	--------------	--------------	--------------

 $n=3$ $r_0 = 4$ part- $r-00000$ $r_1 = 4$ part- $r-00001$ $r_2 = 3$ part- $r-00002$

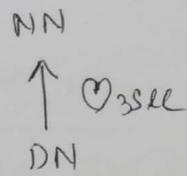
checksum



* If checksum same at 2nd
file is corrupted.

Map reduce flow example HDFS-Excel

How stored data in HDFS is reliable.



- Every 3sec ↑ block report information
- Checksum of each block is also mentioned.
- file : blocks : checksum
- After 30 sec it will match the checksum for corresponding blocks.
- If down it will replicate in another node & delete the corrupted datablock.
- NameNode continuously doing the checksum matching.

HIVE



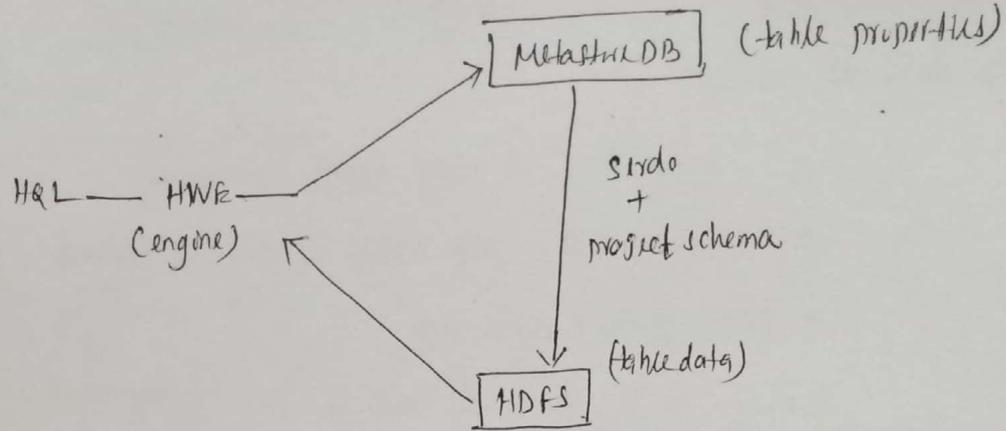
Metastore DB

tbl properties

- Table names - location
- Column names - datatypes
- Input output format
- Partition info
- Bucketing info
- Function definitions.

HQL → Hive
(Engine)

1. syntax check
2. resolve the query
3. logical plan
4. logical generate MR code
5. generate & execute HDFS
physical plan
6. output HDFS console



Skeleton code for Mapper, Reducer, Driver class

* 04.01.2018
* Friday

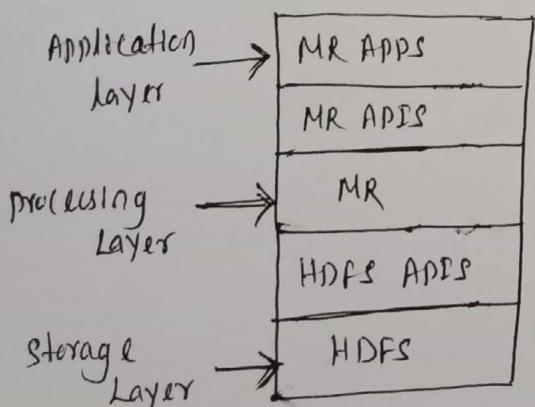
3 clauses for MR Application

- 1) Mapper() → Map()
 - 2) Reducer → Reduce()
 - 3) Driver → main()
- } compile

*.java → *.class → *.jar

- Eclipse 'export' option
- Maven build
- gradle build
- Ant build

MR APIs to develop MP Apps



Code for Mapper class

public class MyMapper extends Mapper<K1, V1, K2, V2>

{
 @Override

void map(K1, V1, context)

{

{

Mapper (Ni, Vi - - -)

API class from MR APIs

*.jar



hadoop-core.jar



import org.apache.hadoop.mapreduce.Mapper;

Skeleton code for Reducer class

public class MyReducer extends Reducer<Kintd, Vintd, Kv, V0>

{ ②) override

void reduce (Kintd, Iterable<Vintd>, context obj)

{

----- Actual reduce logic -----

contentobj. write (Kv, V0);

{

{

Skeleton code for Driver class

public class MyDriver

{

① public static void main (String[] args)

{

②) // create job object

Job myjob = new Job();

③) // set your job properties through job object

myjob.setJobName ("MyfirstMR");

myjob.setMapperClass (MyMapper.class);

myjob.setReducerClass (MyReducer.class);

④ // submit and run the job

myjob.waitForCompletion(true);

3
3

1) List All 3 site.xml file

+

2) Adding job specific properties

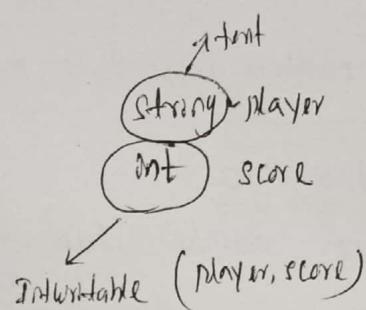
=

Job configuration.

Ex) Data modeling

↓

Data understanding.



JAVA

int
long
float
double
String

⋮

HADOOP

IntWritable
LongWritable
FloatWritable
DoubleWritable
Text

⋮

Wrapper classes

class implements

WritableComparable

* Hadoop have it's own serialization which is different from Java.

* Hadoop serialization is implemented through "Writable" interface.

*
1) Read fields()
2) write()
→ serialization
→ Deserialization.

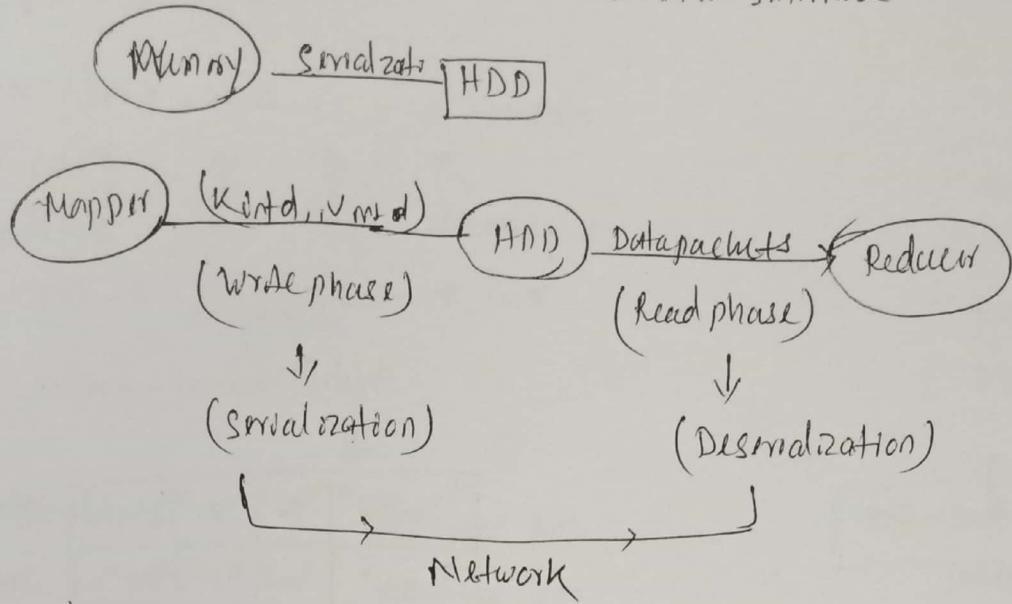
|| 2 methods of "Writable" class

* "Writable Comparable" Interface has addition method -

1) compareTo()
2) readFields()
3) write()
+ isSet()
+ get()

Writtable compatible class extends from Writable Interface.

21



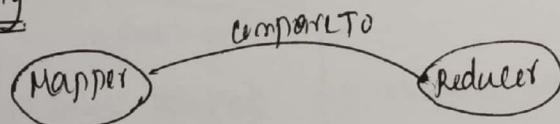
Same content can be serialized in different ways - 16 bit, 32 bit - - - - -.

Java serialization not optimized / slow to handle Hadoop serialization / to represent data keep memory Java processing with writable interface before sending to HDD.

{
 IntWritable x = new IntWritable();
 x.set(10);
 x.set("DVS");
} // object writing and reading will
 take more memory.

Or
 x = new IntWritable(10);

Shuffling



To compare and sort the keys **comparator** needed.

set() \Rightarrow int \rightarrow IntWritable \longrightarrow boxing

get() \Rightarrow IntWritable \rightarrow int \longrightarrow unboxing

Sachin - 98

Kohli - 100

Sachin - 120

Kohli - 150

[K, [V]]

↓
aggregation

Mapper < K, V, Writable, VIntWritable >

Writable ↑
Text ↑
Text ↑
IntWritable ↑
↓
toString()
split(" ")

↓			
"suchin"	"98"	"120"	"150"
"kohli"	"100"	"130"	"150"

- (std::vector)

Content obj. write (KIntd, VIntd)

* Map() will be done one by one for all the records.

public class MyReducer extends Reducer<KIntd, VIntd, K, V>

↑ ↑ ↑ ↑
Text IntWritable Text IntWritable

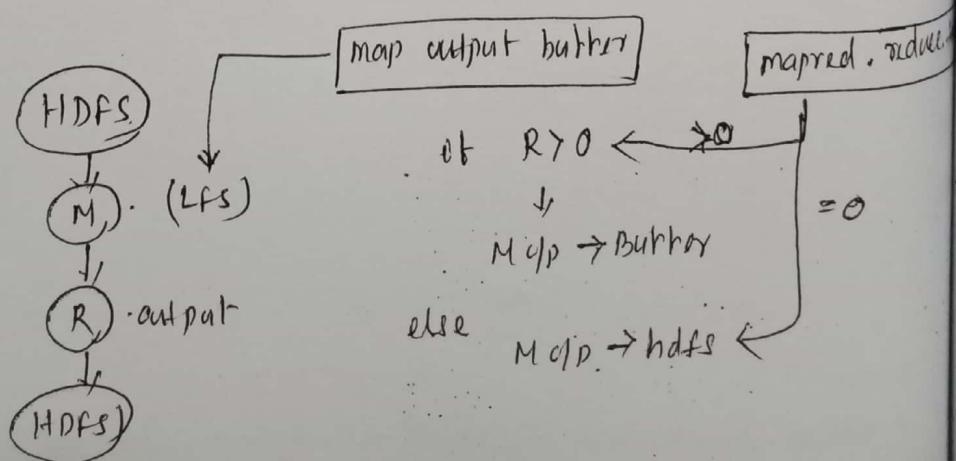
② override

void reduce (KIntd, Iterable<VIntd>, ReducerContext obj) {

↳ [K, [V]]
----- per key one one reducer
[reduceContentObj.write (K, V)]
} reduceContentObj.write (KIntd, new Int (sum)); }
foreach value in list
{ sum = sum + value.get(); }

Content obj. -

Content inner class



which hadoop
 ed lib
 cd hadoop
 ls -al
 cd
 find . | grep -i hadoop -ire -jar

at 8.25 - hadoop-common.jar

U) Client

MR Old API

MR New API

MR 1.X API

MR 2.X API

NOTE - Don't miss the APIs

ds MapReduce base class

→ extends Mapper (abstract class)

org.apache.hadoop.mapred

→ org.apache.hadoop.mapreduce

InputFormat → input split logic
 → Record Reader logic

TextInputFormat → data boundaries are based on '\n'
 → take line offset as key and
 line content as value;

outputformat → TextOutputFormat

↓
 RecordWriter
 logic

↓
 str(k) \n str(v) "\n"

add
 set

i.p o.p
 [0] [1] - main(args[1])

→ fileInputFormat.addInputFormat(mylib, new Path(args[0]));

→ fully qualified class name - class with package name.

→ import * option in eclipse dev phase

→ build.xml in QA & production

```

24
N
pwd
cd ~ / desktop / 
pwd
ls -lrt > CD - MR - Apps /
jar - tf TotalScore . jar
TotalScoreMapper . class
Reducer . class
Driver . class
META - INF / MANIFEST . MF

```

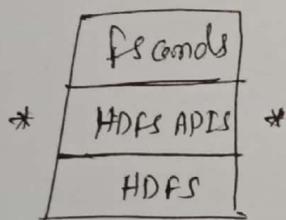
→ Reducer task will start simultaneously, else the task will be go to the initial loop and it may not get the reducer ~~map~~ slots.

C > wc -l cricketscores.txt

103

Mapreduce . Job % the url to track the Job : <http://localhost:8088> .
copy in HUE

Jobhistory server the details will come.



filesystem ps = filesystem . get (new configuration)

fs . delete (" path " , true) ;

→ // delete output directory if exist

- code . m4

{ filesystem hdfs = filesystem . get (new configuration ()) ;
hdfs . delete (new Path (args [1]) , true) ;

// submit the job

→ filesystem hdfs = filesystem . get (new configuration ()) ;
hdfs . delete (new Path (args [1]) , true) ;

Reducer can be set directly

inputsplit size setting will change the map part.

myJob.setNumReduceTasks(4);

hadoop jar abc.jar Driver i.p o.p → Directory

① directory: All files in the dir

→ Do not consider subdirectories

② directory mention file name

③ directly filenames with wild chars

dfs/data/*.csv

/*2018*.csv

④ multiple directories → inputfile.addpath()

MR

↓

R2

↓

MR3

Map only jobs → where no reducers are there

→ simple transformation without aggregation

→ sorting and ordering.

Identity Reducer → when reducer class commented, but no of reducer mentioned

→ sorted output

→ after map → shuffling → sort by, group by done.

→ input => output (same output)

→ Default Reducer = Identity Reducer.

→ No mapper - No Reducer

1	standard	13	9	2
1	Kohls	1	2	3
3100	Lyon	4	5	5
key				value

018 } jobobj. getNumReduceTasks()

mapred. reduce. tasks

$r=0 \rightarrow$ No reduce phase

\rightarrow The o/p will be at Map phase o/p.

\rightarrow \times partitioner, shuffling, sort by key, group by key \rightarrow

$r>0 \rightarrow$ Identity Reducer

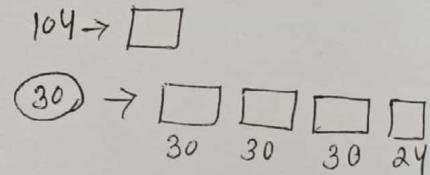
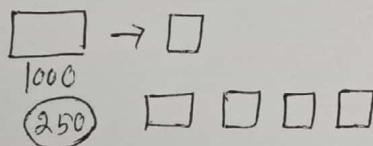
\star No custom Reducers \rightarrow The o/p will be at Map phase o/p but in a sorted order.

No custom \rightarrow Identity mapper

Mapper \rightarrow The o/p will be Record Reader o/p.

\Rightarrow # Mappers = # Input splits

① NLine Input format \equiv Text Input format + No. of lines per split



② By setting the max. input split size

`<Input format>.setMaxInputSplitSize()`

$$80\text{ MB} = \boxed{\quad} + \boxed{\quad}$$

64 MB 16 MB
↑ ↑

\star good for large volume of data, else for small split also the allocated with 512 MB, so not for small splits.

\star overhead for creating container will be negligible for large file \rightarrow block size

\Rightarrow Key Value Input format

`(Sachin, 0, 32, 100...)`

$(\frac{0}{K}, \frac{\text{Sachin}, 0, 32, \dots}{\text{value}}) \rightarrow$ TextInputFileFormat

$(\frac{\text{"Sachin"}, \text{"102, 30, \dots}}{\text{key}} \frac{\text{value}}{\text{value}})$ tab as default delimiter

DB Inputformat

→ sqoop replaces DBInputformat

TableInputformat (HBASE)

→ faster lookups of data

→ unstructured data processing

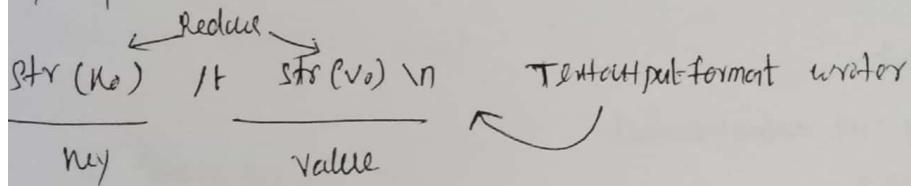
StreamInputformat

→ JSON, XML format

→

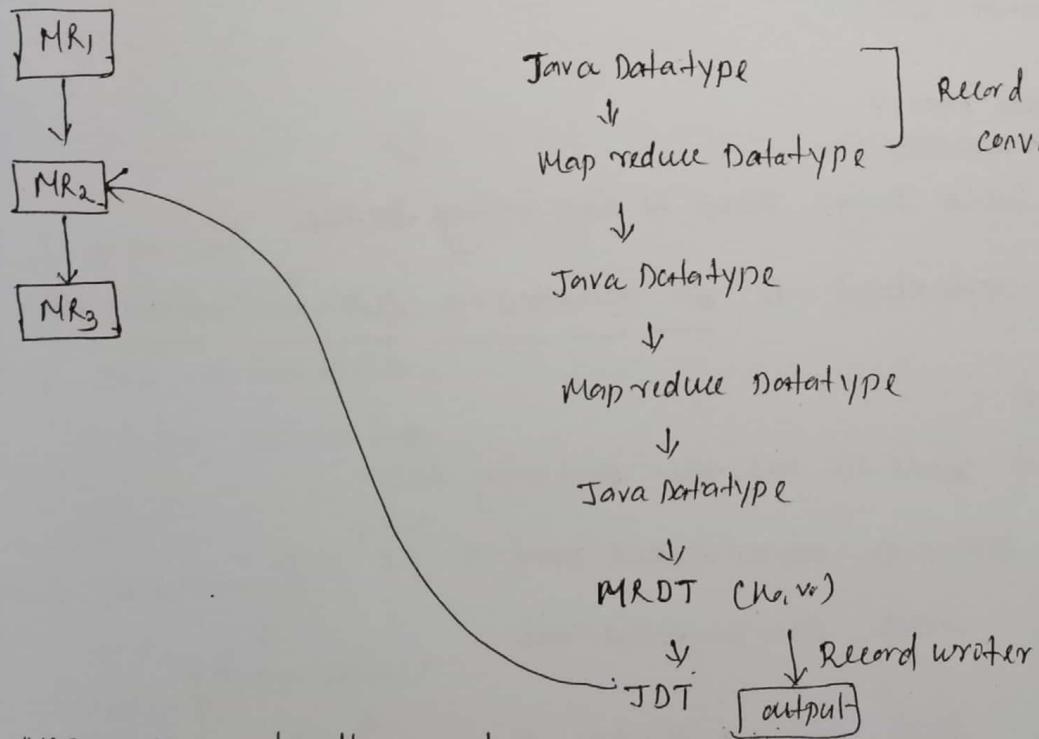
SequencefileInputformat

→ sequence file outputformat



SequenceFOF → Ko # Vo *

* Internal key value delimiter and delimiters

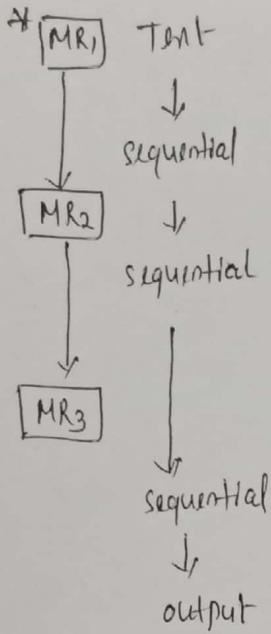


MRDT will be directly input

to then MR1 MR2. No need to convert. MRDT → JDT and

JDT → MRDT

* Performance will be increased.



- * outside HDFS ecosystem sequential file format
- can't be read
- * Hive, wharf etc can read the file and converts to human readable format.
- * convert Text output format to sequential output format.

HIVE

input → Text Input Format

output → Hive Ignore key Text output format

YARN

* Yarn jar Totalscores.jar

// Yarn logs - applicationId job ... > job3.log

Yarn application - list

Generic options parser

Replication factor change setup as after writing the file.

hadoop fs -copyFromLocal -D dfs.replication=5 <lfspath> <hdfs-path>

Java program

→ internal client reads the XML files from config files

→ hadoop jar MRWIR -D mapred.reduce.tasks=0 i.p o.p
-conf ABC-properties.xml

- ||
- * copy from local & put
 - * Hive Ignore key Text output format
 - * hadoop fs

<configuration>

<property>

<name> dfs.replication </name>

<value> 3 </value>

<description>

</property>

<configuration>

{ -p mapred.reduce.tasks=10 } X will not work as it don't know the -p.

public class MyDriver

{

 p.s.v.m (strong ej args)

}

① create job obj

② set job properties through job object

③ Run the job through job object

3-site.xml

+
job specific properties

=
Jobconf file

⑥

3



① 3-site.xml

② command line

③ Job → Jobconf file

run the job

→ ToolRunner.run()

↓
① Job.parse()

config obj with cmd file properties or set

getconf()

⑤

Y()

run

→ All the job task will be done by Y().

public class MyJobDriver

① extends configured

② implements Tool

{
 p.s.v.m() ③

 {
 ToolRunner.run()
 }

④ run()
 {

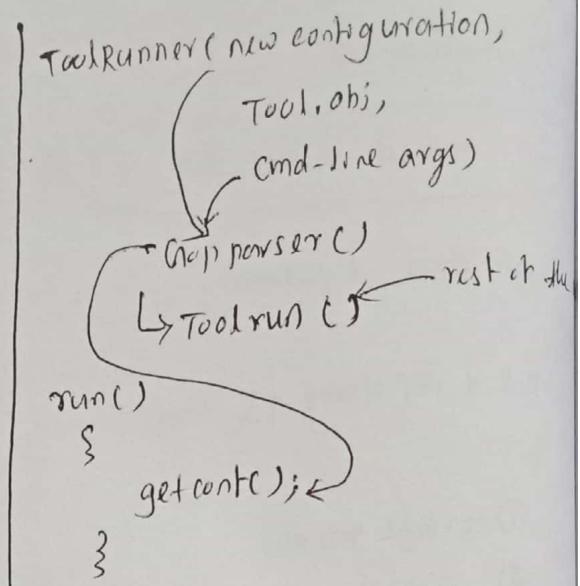
 // create job obj with CJob configuration

 ⑤ Job myJob = new Job(getConf());

 ⑥ // set job static properties

 ⑦ // run the job.

}



⇒ Crop parser will get the cmd-line args and then read the configuration with new configuration and replaces the cmd-line values and create the new configuration file.

→ But it will not read the input and output file values. Those values have to set in the Toolrun();

→ So in run() with getConf() we will get all the configuration set values and with input/output format set the files names as arguments with static properties.

→ Then run the job. CJob creates the configuration object.

- D
- conf
- fs
- jt
- JobJars
- files
- archives

Util, ToolRunner, run()

1) invoke Crop
2) invoke ToolRun

ToolRun

.. ↑

Interface

Toolrunner.run(new configuration, tool obj, and
 || myJob class obj)

args $[-D \text{ x1=y1} \ -D \text{ x2=y2}]$ [i.p o.p]
 rest of arguments

Tool.Run()

{
 Job myJob = new Job (getConf());
 }

Yarn for Total scores. Jar org/bigdata/MapReduceDriver -D mapred.reducer.tasks=4
 /dvs/data -----.

return (myJob.waitForCompletion)

NOTE →

- LFS we can check and run at local
- -fs, -jt local read to local and wrote to local.

Yarn

-fs local -jt local /home/cloudera ... cricketscores.txt /tmp/lop1/

hdfs://nn:port/

file:///

gs:// - cloud

→ Warning - "local" is a deprecated filesystem name. use "file:///x" instead.

mapreduce.task.coortmb=100

mapTaskW & mapTaskB

→ cat /tmp/lop1/

Local execution is faster than HDFS process, as the whole process needs to be involved as assigning container, namenode process, resourcemanagement etc etc.

-libjars
 - files
 - archives

submit → job-id created / Application-ID
 it will take the jar file and ships to working node, and
 addded to local task tracker
 → to temporary directory
 → doing join operation in MapReduce.

Distributed cache → -files

-archives → compress file distribution

→ working node will uncompress the files and do their job.

→ default value - 1 GB

* Node manager will allocate memory, it will take the distributed cache files and process them.

* Not used for Non-Jar concept | Non-Java platform.

JOINS IN MAPREDUCE

* Mean- for structured data table.

emp-tbl

A

dept-tbl

B

select A.* , B.* from emp A join dept B on A.dept-id = B.dept-id

Map - 1 map logic - n mappers() → Map side joins

Reduce - 1 Reduce logic - n Reducers() → Reducer side joins

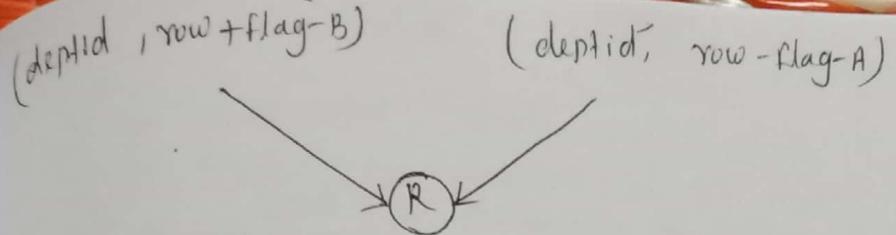
MultiInputFormat

MultipleInputFormat.addFormat ("fileA", "Mapper A", TextInputFormat.class)
-path , class

multipleInputFormat.addFormat (file-B-path , MapperB.class , keyvalueInputFormat)

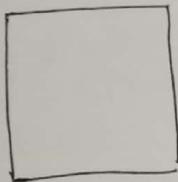
* Hive operation doing joins do the multipleInput format in HDFS backend to do the operation.

(dept-id, rowflagA) → to know which records are coming from which file
Kintd vintd

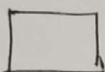


- 1) $(\text{deptid}, [\text{row+flag-A}, \text{row+flag-B}])$
- 2) $(\text{deptid}, [\text{row+flag-A}])$
- 3) $(\text{deptid}, [\text{row+flag-B}])$

Explain:-



cricketdata.csv



playercountry.csv

$(\text{player}, \text{Totalscore}) \rightarrow (\text{player}, \text{country name}, \text{totalscore})$

As the files is very small we can't keep in HDFS to involve a lot of processing. So we can keep in the distributed cache to process our job.

Map side join

DC

Reduce side join

DC

Multiple Input Format

Adding a distributed cache file(s)

1) $\text{map} \Rightarrow -\text{files}$ option \Rightarrow Refer files from LFS

2) in driver class \Rightarrow `DistributedCache.addCacheFile()`

\Rightarrow Refer files from HDFS

Input \rightarrow e.g. \rightarrow hadoop jar ABC.jar Driver -files /dfs/lookup files/ playercountry.csv
input-path \quad output-path

end \rightarrow e.g. \rightarrow `DistributedCache.addCacheFile(new URI("lookupfile-HDFS-path"), jobobj, getcache)`

Refer DC files in mapper/reducer

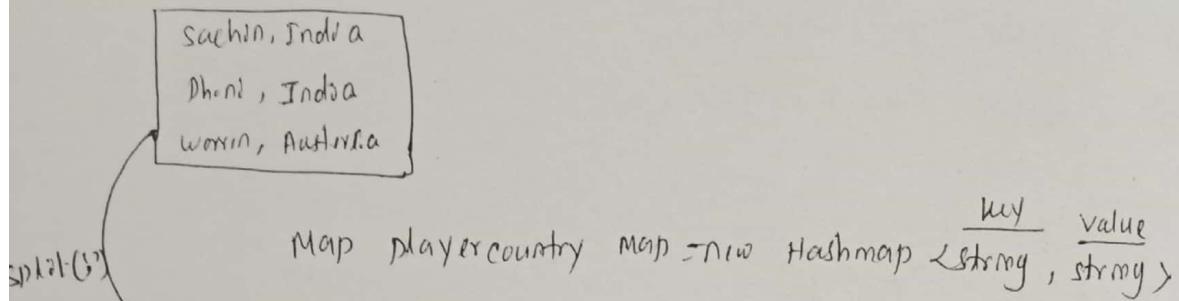
`DistributedCache.getLocalCacheFiles(contextobj, getcache);`

\Leftrightarrow `mapped.cache.files`

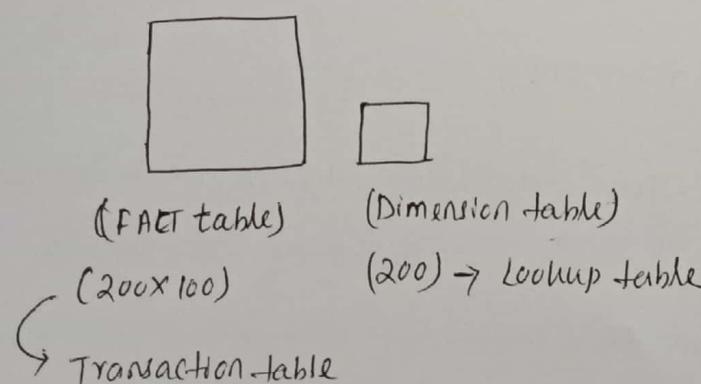
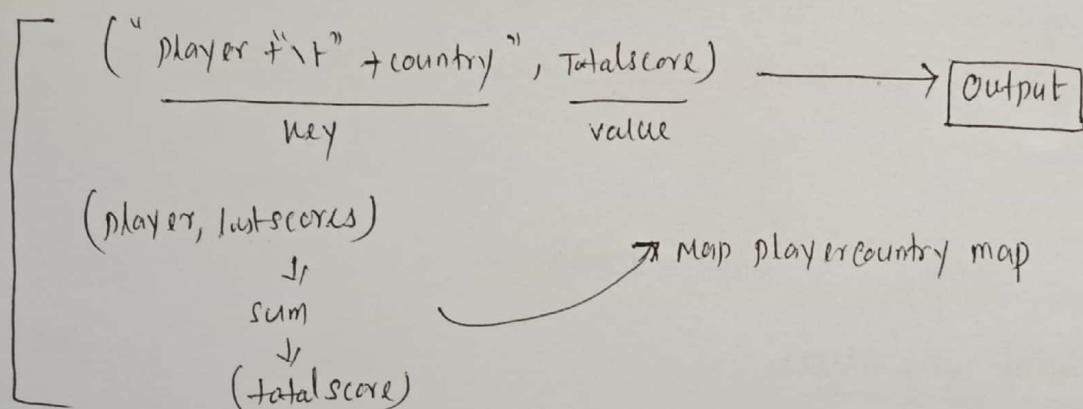
* will get this property through `getLocalCacheFiles`

- 1) get cache file name
- 2) Read file and create a lookup obj \Rightarrow Map obj List Array

playercountry.csv



playercountry Map: put (split[0], split[1]);



9.01.21
 Wednesday
 2023

Distributed cache

→ files are needed are added to the job resources, and pass to the local node manager / worker node, the files will be distributed.

* Create the lookup object, and picked by mappers.

- No of mapper reduces
- Reducing I/O operation, as DC will be faster
- one time transfer of data through network.
- Map side Joins are preferable compared to Reduce side Joins

] Ex:- select * from tableA, table-B;

→ $c = \underbrace{\text{Join A by } x_1}_{\substack{\text{Left side} \\ \text{Large table}}}, \underbrace{\text{Join B by } x_2}_{\substack{\text{Right side} \\ \text{small table}}} \text{ using replicated}$

→ Scan operations will be costlier if big table will be in right table.

Reducer side Joins

Join \rightarrow aggregation

aggregation \rightarrow join \rightarrow acceptable \rightarrow still depends on the situation.

A		
player	score	
(1000)		

B
country
(50)

player, country \rightarrow total score

case-2

1000 \rightarrow Join

records

45 \rightarrow Avg \rightarrow Join

records

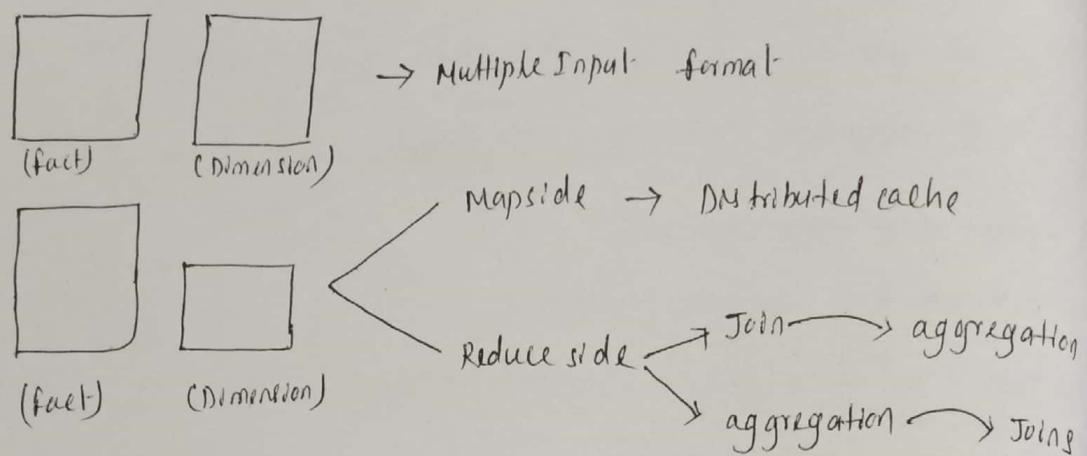
\rightarrow Reduces joining cost operation.

\rightarrow Avg will be done in Reducer.

case-1

* If the aggregation depends on both the table, the join is then aggregation and is acceptable.

→ If both the files are bigger use multipleInput format for reduce joins.



Default methods

Mapper()

setup() →
→ 1st method to call
→ at beginning of the task
→ only once it will be called

Reducer()

← setup

map() →
→ for n-time
→ # records
→ the input splits

→ for m-times
→ # m where m is
the distinct keys
that reducer received

← reduce

cleanup() →
→ Last method to call
→ At the end of the task
→ called only once

← cleanup

- * Distributed cache will be read at setup() method for one time read.
- * HMapobj. will be created in the class
- * Call in the setup() method, so it will be build
- * Map() method it can be used for join operations.

Sort by key
↓
group by key
↓
groupkey, [list value]
↓
Reducer (distinct keys)

- hadoop -p /hadoop-data/
- hadoop fs -put playcountry.txt /hadoop-data
→ SH -ovi
→ jar - ... TotalScoreHadoop.jar
/dfs/data/cricket-data/ /dfs/op1/

COMPRESSION

* 10.01.2018
* Thursday
* am * am *

- space
- faster to transfer through N/W / Harddisk.
- HDFS - compressed file can be stored.
- No overhead for unzip of files in HDFS MapReduce.

cricketscore.txt

↓
cricketscore.txt.gz → Based on file extension it will invoke
↓ HDFS corresponding Decompression.
MR
↓
O/P → compressed O/P

TYPES

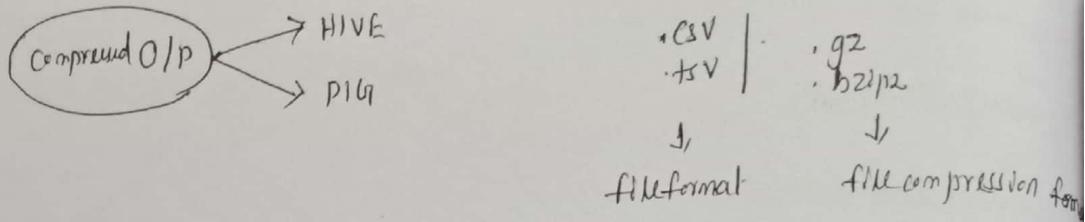
Codecs → compressor, decompressor

- default
- gzip
- snappy
- Lzip
- bzip2

codec → org.apache.hadoop.io.compress.BzipCodec.class
1) through command line
2) in Job definition

hadoop fs -text <compressed file in hdfs loc> | head -10

On top of compressed o/p we can define the hive table.



① C> file cricketscore.txt

ASCII text with CRLF line terminators.

C> gzip cricketscore.txt

desktop> file cricketscore.txt.g2

gzip compressed data

desktop> cat cricketscore.txt.g2

--- symbolic representation

> gunzip cricketscore.txt.g2

② > gzip cricketscore.txt

> hadoop fs -mkdir -p /dvs/data/g2/

> hadoop fs -put cricketscore.s.txt.g2 /dvs/data/g2/

desktop> hadoop jar

> hadoop jar org/hadoop/mr/myDCDriver

/dvs/data/g2/ /dvs/op3/

> hadoop fs -cat /dvs/op3/

↳ player country scores

> hadoop fs -ls /dvs/op3/

③ MyDriver.java

FileOutputFormat.setOutputCompressorClass(myjob, true);

FileOutputFormat.setOutputCompressorClass(myjob, GzipCodec.class);

γ hadoop 3 or Total scores. jar

```
7 hadoop fs -ls /dfs/opu/part-r-00000.g2
```

```
7 hadoop fs -test / dvc / 074 / part-r-00000.gz | head -10
```

player country scores ↪ output

-D mapped-output.compression=true -D mapped-output.compression, code=

org.apache.hadoop.io.compress.GzipCodec /dvs/gz/

cricket scores 'tel. 92

→ Generic options parser.

- download to local the configuration file

- search the compress

- need debt value

- set generic options parser value

10. compression . codee

org.apache.hadoop.io.compress.DefaultCodec

- do - · Zip Codec
- do - · BZip2 Codec
- do - · Snappy Codec

└ mapped.output.compress - false

mapped. output. compression. codec = org.apache.hadoop.io.compress.DefaultCodec

→ Default values

file 1GB — A — 2min — 300MB

B — 30sec — 750MB

Tradeoff between compression ratio & time to compress

▷ depends on time and size availability

2) Command options parser availability long [group].

⑧ create table in hive - load data grip - Read the file.

PARTITIONER

values associated with same key will go to the same reducer.

Hashpartitioner

By Default

$$\text{partitioner} \rightarrow \text{Hashpartitioner} \Rightarrow \text{modulus} \left| \frac{\text{hashcode(key)}}{\# \text{ Reduces}} \right| = 0 \dots (n)$$

Custom partition can be implemented through partitioner class

partitioner <Kint, Vint>

↳ getpartition (, , , no of arguments reducers)

{

- Actual partitioner logic

- Return <int> from 0 to (n-1)

}

Example :- group by Jan year - Jan 19, Jan 18, Jan 17 - - - - .

- Key \rightarrow

YYYY	-	MM	-	dd
------	---	----	---	----

- And modify your partitioner code to get the data.

or

(Hashcode)

Code :- Mypartitioner.java

public class mypartitioner extends

-MyMapper.java

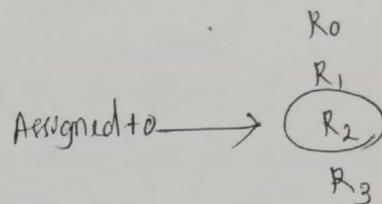
your jar

-> mapred. reduce-tasks = 4

* set \rightarrow myjob. setpartitionerclass (mypartitioner.class);

Return 2 partition will go to the $R_3 \rightarrow$ partition return will be the partition number.

\rightarrow partitioner will tag to partition reducer



Case \rightarrow

$R=2$, partition = 2

* After 1 execution when trying to tag it will show error.

R_0
 R_1

* Illegal partition for D.A. Warner (2)

* No of mappers - 1 time
maptask - 4 time] it will attempt for 4 times

partitioner comes under mapper()

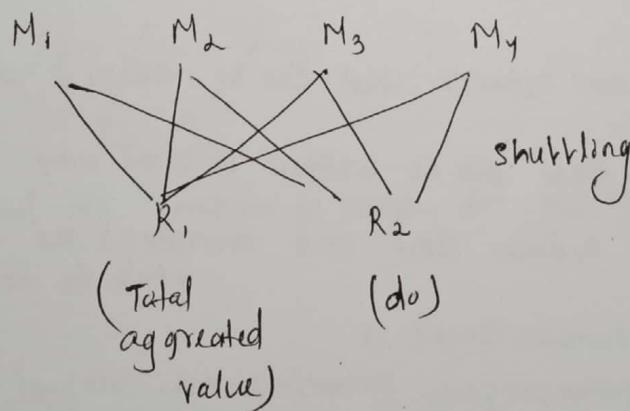
COMBINER LOGIC

11.01.2018
friday

selection or
transformation filter)

aggregation

Case-1

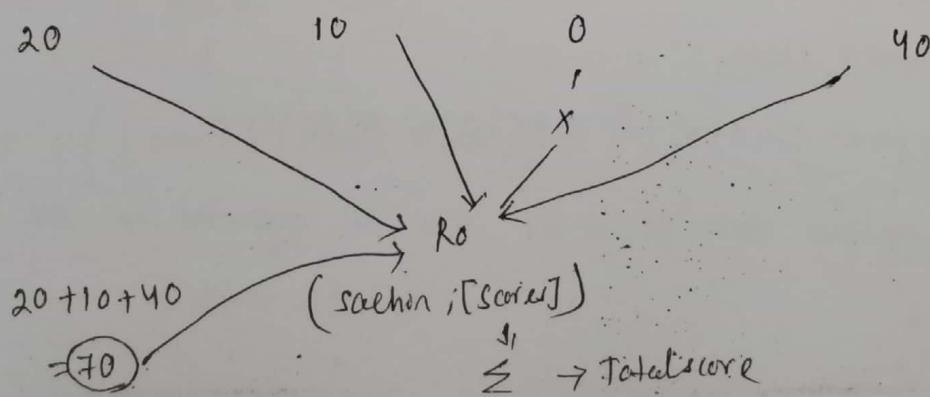


M_1
(sachin, s₁)
20

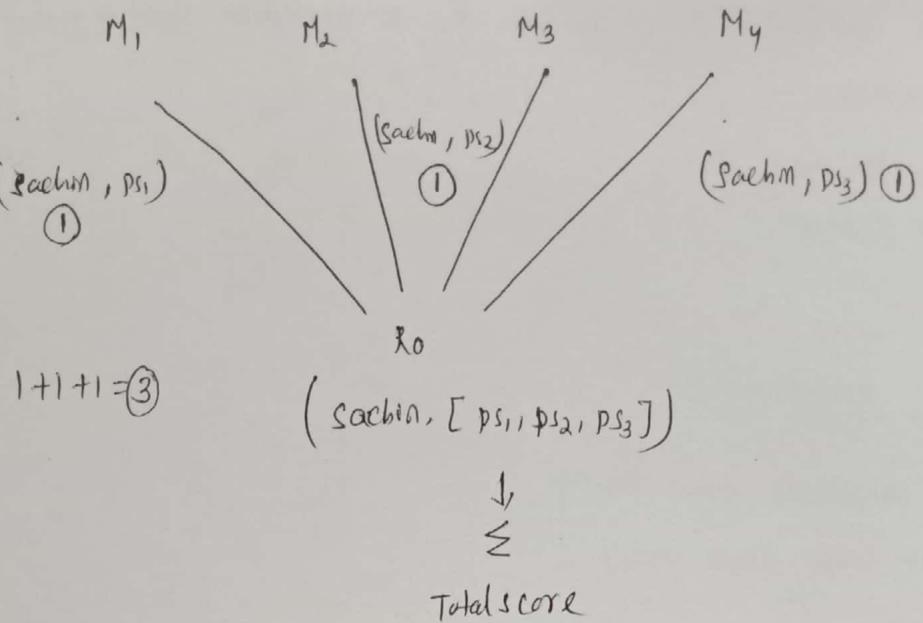
M_2
(sachin, s₂)
10

M_3
0

M_4
(sachin, s₃)
40



case-2



- * while passing data through the network will be faster to send 3 (keys, values) instead of 70 (keys, values).
 - * Combiner logic \equiv reducer logic
 - * Map() \rightarrow Combiner() \rightarrow partition() \rightarrow Reducer()
- (Flow of DATA)

Keypoints

- 1) Combiner is a mini Reducer that can be invoked at Mapper class.
 - 2) The combiner class will be written in same way we write Reducer class, by extending Reducer class and overriding the Reduce() method.
 - 3) Job object. setCombinerClass();
- or
- Mapper. combiner. name = " ";
- 4) If any aggregation follows associative rule then Reducer itself can be set as combiner class.

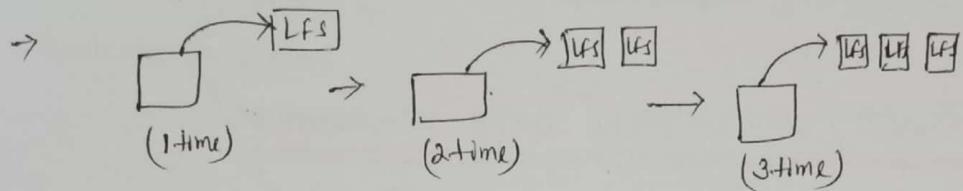
$$A+B+C \equiv (A+B)+C \equiv A+(B+C) \equiv (A+C)+B$$

- 5) Though a combiner class is a set it is not guaranteed to run for all mapper class.

If any mapper's - mapper output buffer spills minimum 3 times then the any combiner will be invoked for that mapper.

→ framework counter will be increased by one, when the buffer spills.

While exceeding the buffer size, it will send to LFS.

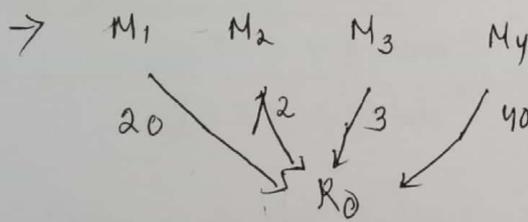


→ for less data no combiner needed.

→ creating the combiners & running for small amount of data itself is the overhead of running the reducer class for all the mapper class.

→ one combiner will be invoked for only one mapper.

→ $M + C \Rightarrow (K_{int}, V_{int})$] for datatype
 or $M \Rightarrow (K_{int}, V_{int})$] value should be always same

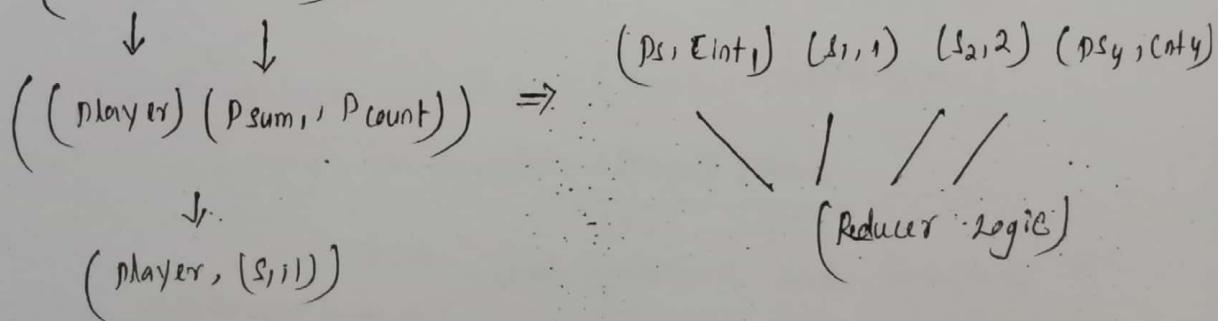


→ for M_1 and M_4 combiner is need, for M_2, M_3 combiner not needed for less volume of data.

How to use combiner in case non-associative aggregation.

$$\text{e.g: avg} = \frac{\sum v_i}{n}$$

Combiner $\rightarrow (K_{int}, V_{int})$



mapred.io.sort.mb = 100MB \rightarrow Map output buffer size.

Job object set (p_1, v_1)

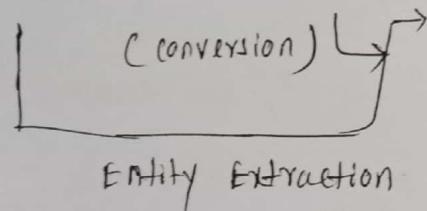
set $(\text{mapred}, \text{io.sort.mb}, 5)$

set $(\text{mapred}, \text{map}, \text{class})$, $\text{org.apache.hadoop.mapred.Mapper}$

SEMISTRUCTURED / UNSTRUCTURED DATA PROCESSING

semistructured data
unstructured data \Rightarrow structured data

Unstructured \rightarrow semistructured \rightarrow structured data



$(\text{senderemail}, \text{date}, \text{amount}) \rightarrow (\text{mobileno} \rightarrow \text{xxx},$
 $\text{accountno} \rightarrow \text{xxx},$
 $\text{emailid} \rightarrow \text{xxx},$
 $\text{date} \rightarrow \text{xxx})$

\downarrow

(key)

\downarrow

(value)

\rightarrow Data probabilizing

\rightarrow Data sufficiency (misused Entity to get the Business requirements)

Ex:- (Task 2)

phone [3] [4] ## I have - - - - -

Battery [-1] ##

(Longwritable) [d]

$\frac{\text{phone}, [3] [4] - - - -}{[0] \quad [1]}$

* Unstructured data processing tools

NLP

Regular expression

→ semi-structured data

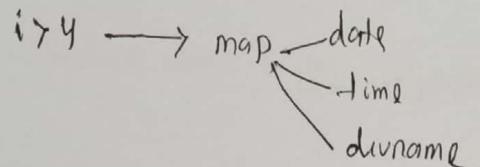
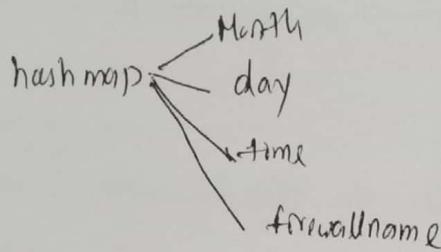
Trans-part

trans-part

1) No Reducer

2) split

3) split [5] → property



Headers → source, destination, ...

map ← At each mapper()

- Hive external table → load data → structured data

UDFs (User Defined Function)

{19-01-2018}
{friday}
* * * * *

HIVE — PIG — MR
① ② ③ ⇒ preference

Hive → Java → Built-in functions → UDF extends evaluate()

Pig → Java → Built-in functions → EvalFunc extends exec() → *.jar
added to
classpath

Ex:- (H, lg) ⇒ Google APIs



→ determine the nearby location for particular object

*.jar
↓
classpath → hive → add jar <jar-file-path>; ①
→ pig → register <jar-file-path>;

→ create temporary function getlocation as ...
org. dvs

have → ⁽²⁾ create function getlocation as org. dvs. vdfs. have. GetLocation
 select vehicleid, getlocation (lt-col, lg-col)
 from gprs-data-tbl;
 → extends vdfs
 → evaluate
 → obj. evaluate()

PL/SQL

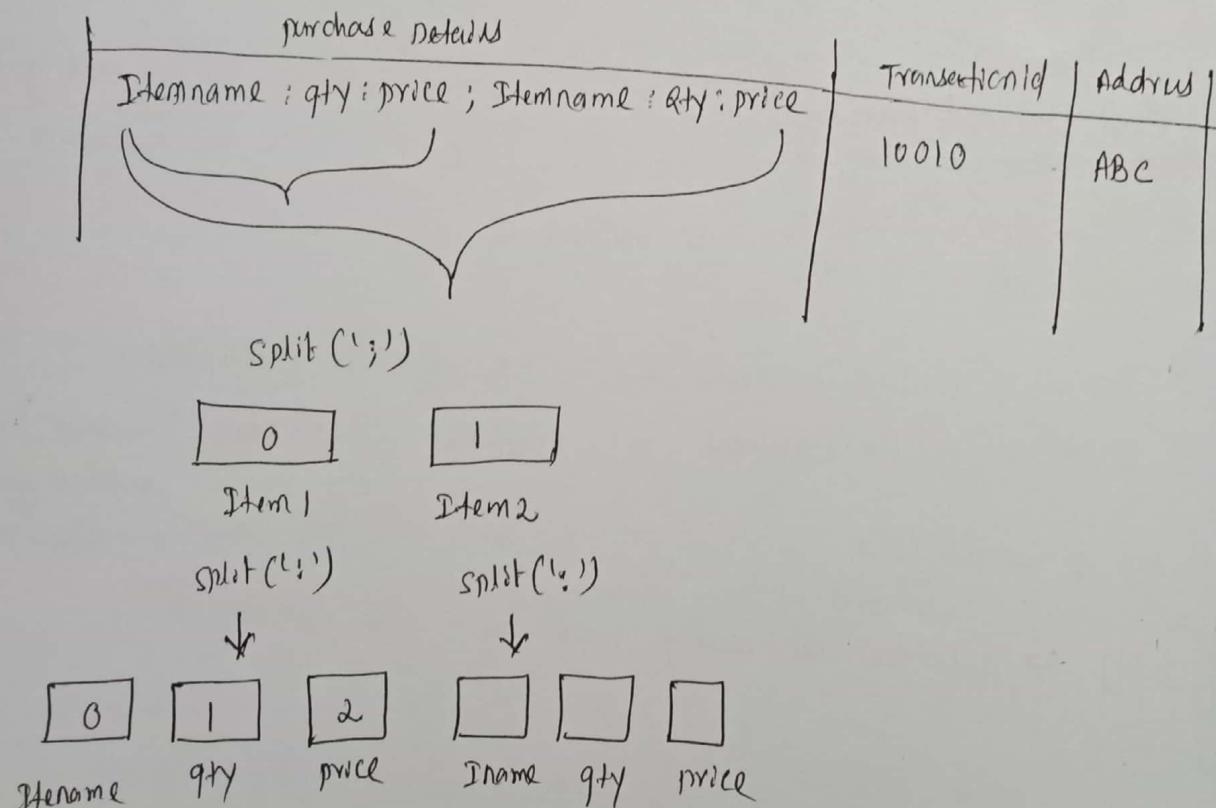
define get-location org. dvs. vdfs. PL/SQL · GetLocation();

A = load . . .

B = for each A generate & o, GetLocation (x1, x2);

* check for
 - keyword
 - function

- | | |
|-----------------------|--|
| <u>Old</u> | <u>New</u> |
| - temporary functions | - functions are available in DB level. |
| - session oriented | - reusable.
- permanent functions. |



Error : —

Log : — Different actions accordingly

Intor : —

HIVE

hive.jar \Rightarrow <hive-install-path> \rightarrow pig.jar \Rightarrow <pig-install-path>

unlike class MyUDF extends UDF { \rightarrow public class MyUDF extends

* evaluate(x_1, x_2, x_3) \downarrow Hadoop/Java datatype

{

- Actual HIVE UDF logic

return <hadoop data type>

}

}

* evaluate(x_1, x_2, x_3, x_4, x_5)

{

}

}

}

PIG

public class MyUDF extends

EvalFunc<Return data type>

\downarrow

Java data type

{

② override pig datatype

exec(Tuple t)

{

- Actual PIG UDF logic

return <Java data type>

}

}

}

Add jar HIVE Env

↳ local path

create .hql for future reference

and also .pig

Example

date string

22/10/2018

date

↓

YYYY-MM-dd

2018-10-22

datatype in Hive

date — YYYY-MM-dd

→ month()

day()

year()

dt1, dt2

timestamp — YYYY-MM-dd HH:mm:ss...

where dt between '2018-01-01' and '2018-01-31'.

* select store_id, sum(sale_amount) from pos_dt
group by month(sale_date)
group by day
group by year

SimpleDateFormat

Raw date → parse() — input pattern
String → format() — output pattern

NOTE - showfunctions() in hive to see all the existing UDFs.

Pig → Register → grunt shell.

Hive → addjar <path to jar file>;

*.jar → classpath
↓
JVM - [-cp]

B = foreach A generate p0, s1;

org.apache.pig.Udfs.formatDate(\$3, "yyyy-MM-dd")

formatDate

↓

obj

↓

exec

↑ tuple

> Define myDateformat org.apache. ... formatDate();

- function name
- Alias name
- short name

> B = -- -- . \$0, \$1, myDateformat(\$3, "yyyy-MM-dd");

/* string fun? date fun?, Regular Expression

We should create a function, because Hive stores function details separately in Metastore DB.

> create [temporary] function <FUNC-NAME> As org.udfs.---formatDateClass;
add jar path-to-jar; ↓
↓
session level only at database level

> select * , y, z, <Udf-fn-name> (parameters) from table ---;
↓ ex
MyformatDate (W, "yyyy-MM-dd")

or// hive-site.xml

<name> hive.aux.jar.paths </name>
<value> jar1-path, jar2-path, --- </value>

Simple Dateformat

↳ parse(dt) ← to read the dt string
↳ format(dt) ← to write to the dt string

→ raw date string ← column
→ dd/MM/yyyy ← arg1 ←
→ yyyy-MM-dd ← arg2 ←

Unix Time	00000 00:00:00	Jan 1st 1970
System Time	→ 0000 0000	
Epoch Time	→ 984512435	

Pig

Pig> pig -x local

grunt> register /home/training/hari/udfs/pig/pig-udf-format-date.jar;

grunt> Load A . . . sample.txt

g> B = foreach A Generate \$0, org.higdata.pig.udfs.formatDateUDF1(\$0);

(25/12/2015, 2015-12-25)

(26/01/2015, 2015-01-26)

! ! ! !

> define function org. bigdata.pig.udfs.formatDateUDF();

pig> which pig

cd contrib/piggybank/java/lib

HIVE

localhost> show databases;

> add JAR /home/training/hari/udfs/hive/hive-udf-format-date.jar

> list jars;

> create temporary function fmt_date AS 'org. bigdata.hive.udfs.

> select *, fmt_date(sale_date) as formatted_date from;

> show function like 'fmt*';

> tab key yes

> describe function sum;

> describe function max;

> describe function substr;

> describe function extended substr;

('facebook', -5) → ebook
5932 ↴

('facebook', 5) → book

> describe function fmt_date;

> describe function extended fmt_date;



Explicitly code to maintain this

② Description (

name = " " ,

value = "-FUNC-takes string as parameter and
return int";

) extended = "example:"

```
import org.apache.hadoop.hive.ql.exec.Description;
```

↓

```
④ Description (name = " " ) get DisplayName()  
↓ {  
factory class } return " "  
}
```

adding the files automatically

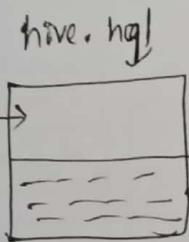
```
<home-dir> | <user> | .hiverc  
| home | .hadoop | profile | .hashrc  
| home | Nagendra | .profile | .hashrc
```

→ all hive related setting can be done for the user.

,hiverc

- use dvfs-db;
- set hive.exec.dynamic.partition=true
- add jar &,jar
- create temp table ...

hive ↴



~]\$ pwd

/home/training/

② ls -l

mc .hive

mc .hivernl .hiverc

cat .hiverc

UDFs - one row → 1 row output

UDAFs → n row → 1 row output

VDTFs → 1 row → n row output

Aggregate

Table Generations

e.g - Ucase, Lcase, substr.

e.g - sum, avg, min, max

e.g - explode \Leftrightarrow flatten (pig)
(hive)

Map	sachin [98, 22, 45, 36]
Array	↓↓
struct	sachin 98
	sachin 22
	sachin 45
	sachin 36

pivoting

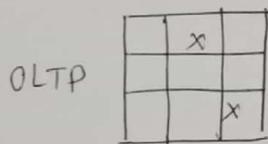
transpose

→ 4 rows.

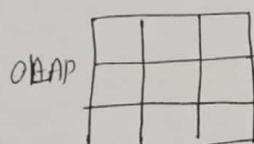
HBASE

OLTP \rightarrow RDBMS \rightarrow supports transactionality

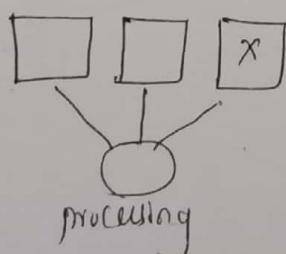
OLAP \rightarrow Data warehouse



- * Live data
- * Transaction purpose



- * past data
- * Analysis purpose
- * hive supports OLAP
- * hive is the data warehousing of hadoop
- * offline/ batch processing can be only done
- * No transactionality
- * Not good for faster lookup of data.
- * doesn't support random read/write operations, as HDFS doesn't support the same.



- * All the blocks needs to be read parallelly then sequentially to get a particular data.

HBase

fast lookups data

CRUD operations supports - create, read, update & delete a record.

Huge volume of data.

can store unstructured data/ semistructured / structured data.

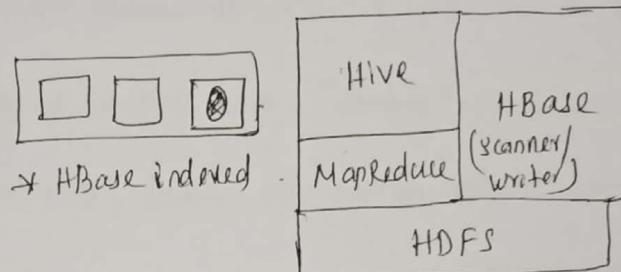
Web tables . Ex:- Table storing a Recommendations

Huge catalog

Google's Big Table.

Database of hadoop

stores read the data , not
meant for Analysis



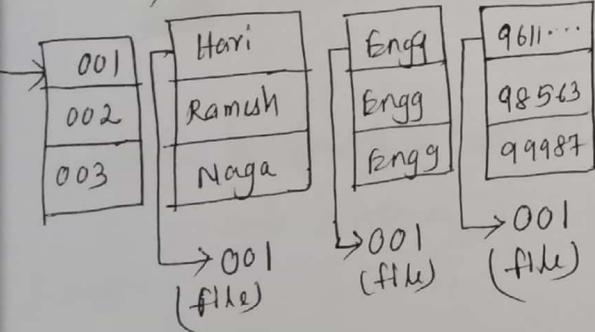
now not supports joins and aggregations.

it stores data as columns/ columnar database.

||
||

(HBase)

→ Row Key



→ + + + + → row
(HIVE)

→ primary key / unique key retrieval

→

001	Harri	Engg	20000
-----	-------	------	-------

(id) (name) (prob) (mobile)

001 → Harri
↓
(Name)

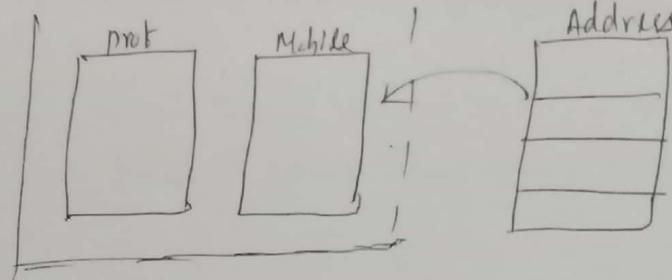
Engg
↓
(prob)

99987
↓
(mobile)

→ HBase stores data as sorted key, value pair.

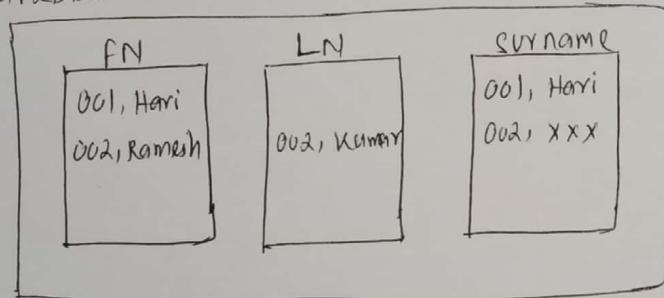
group of columns is called the

Column family



- data insertion will be column by column.
 - Dynamic schema / schema less.
 - At the time of insertion we decide the column family, row key, column name and value.
 - schema of HBase table is column family names.
- Ex:- create 'std-thl', NameDetails, BatchDetails, contact Details, Identity Det
- column family names are virtual / directory in HDFS only.

NameDetails



(NameDetails column family) (HBase Serde)

- Key based access, same as JSON format.
- Sparse Table = No reservation of columns for a record in advance. So we can store unstructured / semi-structured / structured data.

JSON - { ["name": "Hari",
 "id": "001",
 "key": value],
 ["name": "Ramesh",
 "id": "002",
 "prof": "Engg",
 "no": "9861"] }

There is no ordering or limitation of columns so it is dynamic schema, but in hive ordering and no of columns is important.

001	Harvi	"Engg"	96115
002	Ramesh	Prob	9999

(Hive table representation)

(CSV serde)

Sorted key value pairs \Rightarrow Metadata + Actual data

Region 1 { 001
 : }
 1000 } DN10

Region 2 { 1001
:
2000 } DN₁₂

Scanner
get std::file, 105, contactDetails : homeNumber

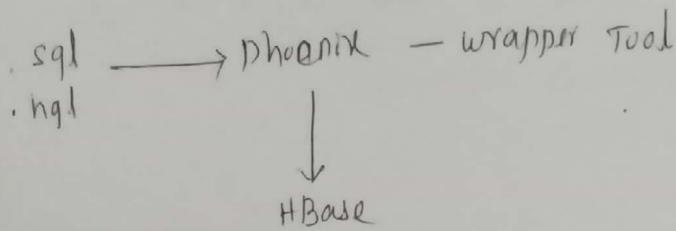
Table name	row key	column family	column
Table	row	column	column

105 belongs to region 1

scanner
put std::list<ios::internal>::name, number, 1234

* It not present create.

* at present overwrite with new version.



HMaster

HRegion Server

ZooKeeper

Memstore (write cache)

Hfiles

WAL (write ahead log)

Block cache (read cache)

HBase scanner

HBase storage handler - Read / write operations

HRegions

• META.

HBase clients

- CLI → HBase shell
- REST API
- JAVA APIs

COMPACTON

- Minor compaction
- Major compaction

TTL

HBase and HIVE Integration.] → **IMPORTANT**

* hadoop cluster should up and running. → Name node

* don't need Map reduce

secondary name node

* Hbase engine to read / write

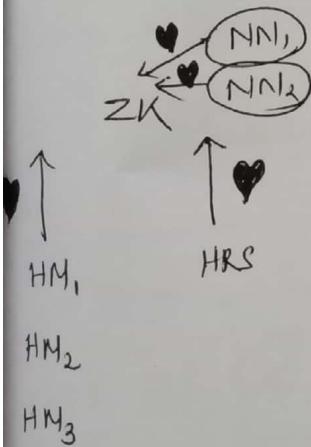
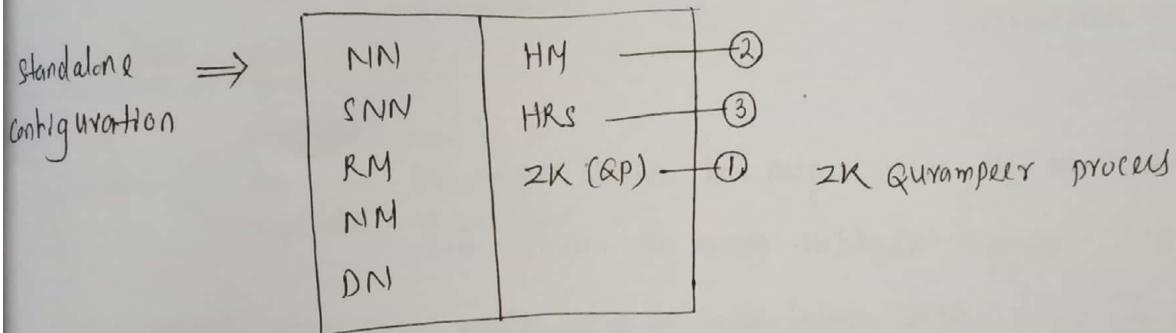
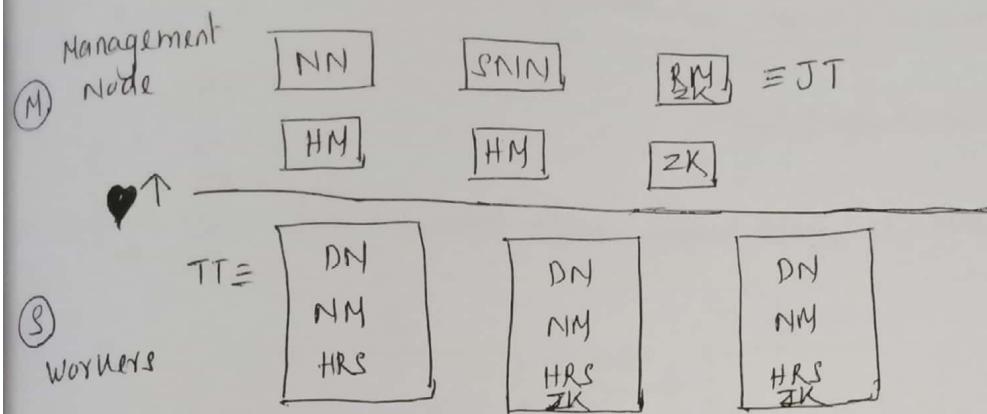
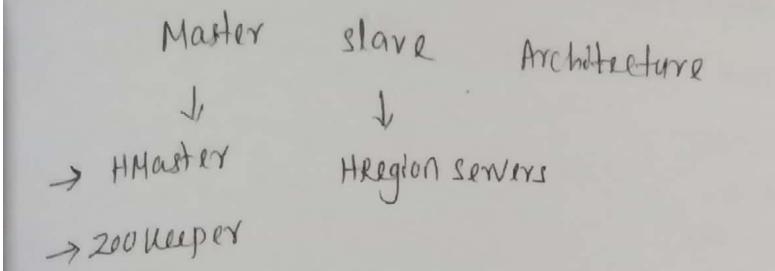
* start-hbase.sh & stop-hbase.sh needs to run.

{ 24/01/2018 ?
* Friday }

→ Read / write operation

→ connecting clients

→ Data storage



- * Minimum 3 HM have to be started.
- * They will go and register to ZK that they are alive with corresponding ipaddress.
- * One will be the lead HM.
- * Signals will be send to all the HM.
- * polling will be happened between HMs and one will be selected.
- * In config file hadoop master will have set up.
- * ZK will do the notification services.
- * If HM down, it will inform other 2 and one of them will behave as HM.
- * No need to do anything manually.

* HRS is down ZK will inform HM that HRS is down.

* Hbase table whose regions allocated to HRM.

* Regions as set of buffers.

→ as HRM assign HRegion to HRS

- activities known as Region assignment.

- ZK ^{Distributed} ~~Resource~~ Resource co-ordinator

* if NN1 is down NN2 will become active by the ZK.

HM

→ responsible for DDL operation

→ Region assignment.

→ ,META.

→ stored in the file, after region and column region is assigned.

→ All DDL changes should be stored in ,META. file.

→ only one region will be created when 1st time Hbase run, with RegionName.

→ 4GB Region size

→ in the middle of region it will be splitted across the key values.

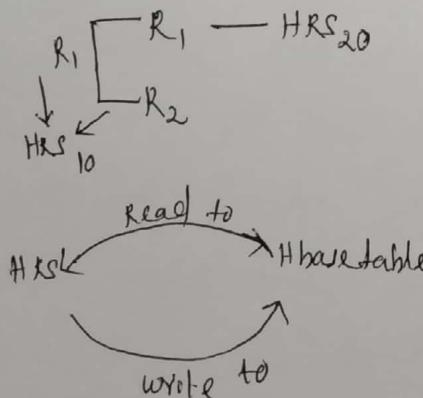


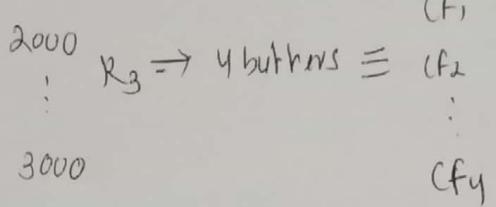
Table \rightarrow CF₁ CF₂ CF₃ CF₄

$$0. \quad R_1 \rightarrow 4 \text{ buffer} = \begin{matrix} CF_1 \\ CF_2 \\ CF_3 \\ CF_4 \end{matrix} \rightarrow \boxed{Hfile}$$

$$1001 \quad R_2 \rightarrow 4 \text{ buffer} = \begin{matrix} CF_1 \\ CF_2 \\ CF_3 \\ CF_4 \end{matrix}$$

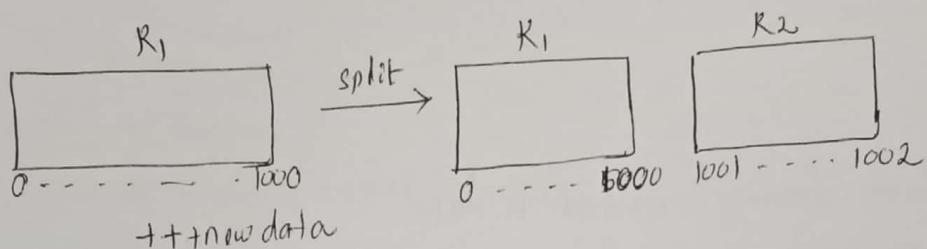
- * When the threshold will happens the data will be saved to HDFS.

- * When data grows, region split will be happens, with required buffer.



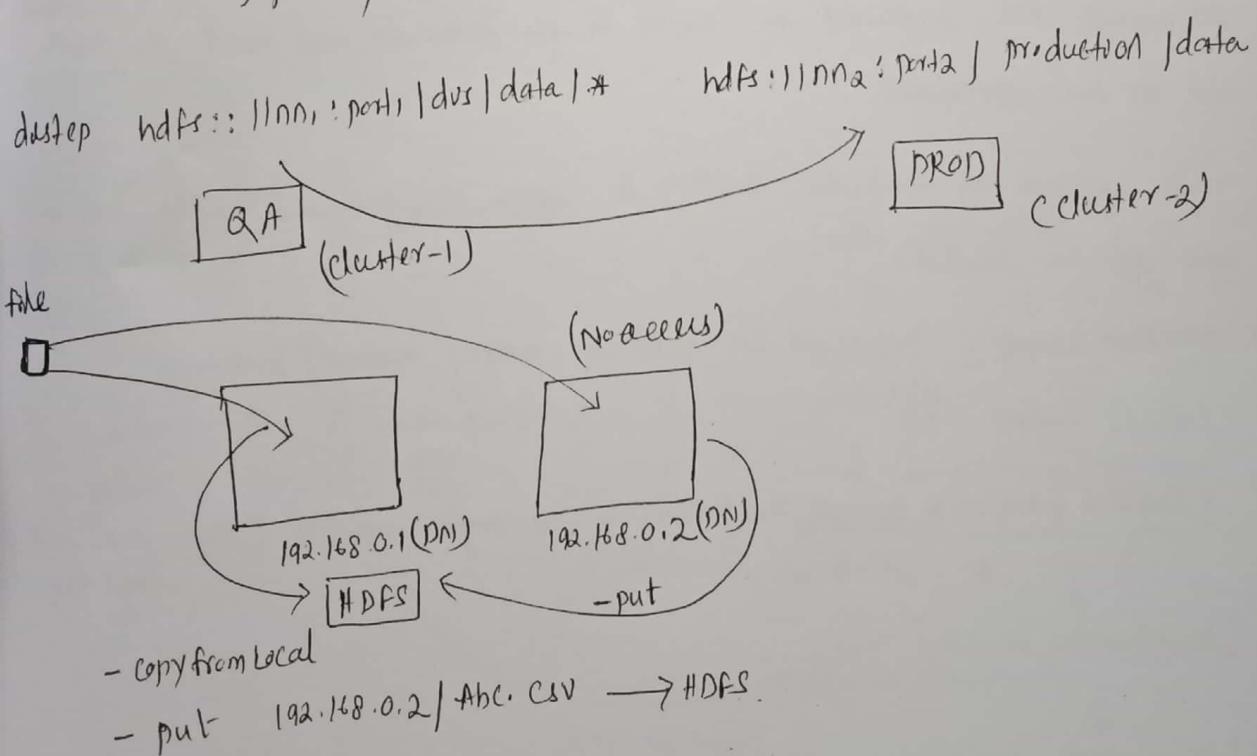
(organisation of data)

- * Regions make data readability easier as reading 1 region as easy than 4 regions.

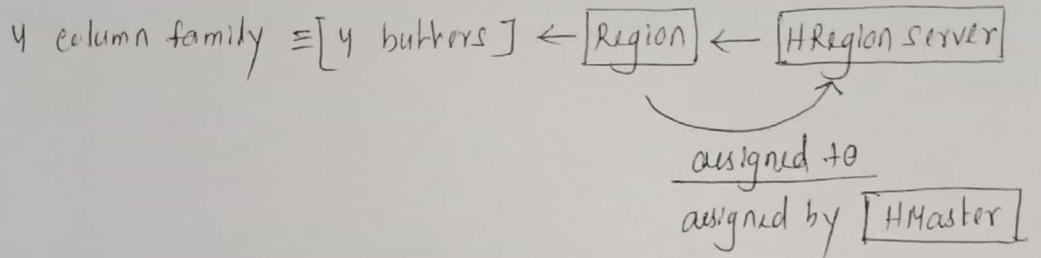


put & copy from local

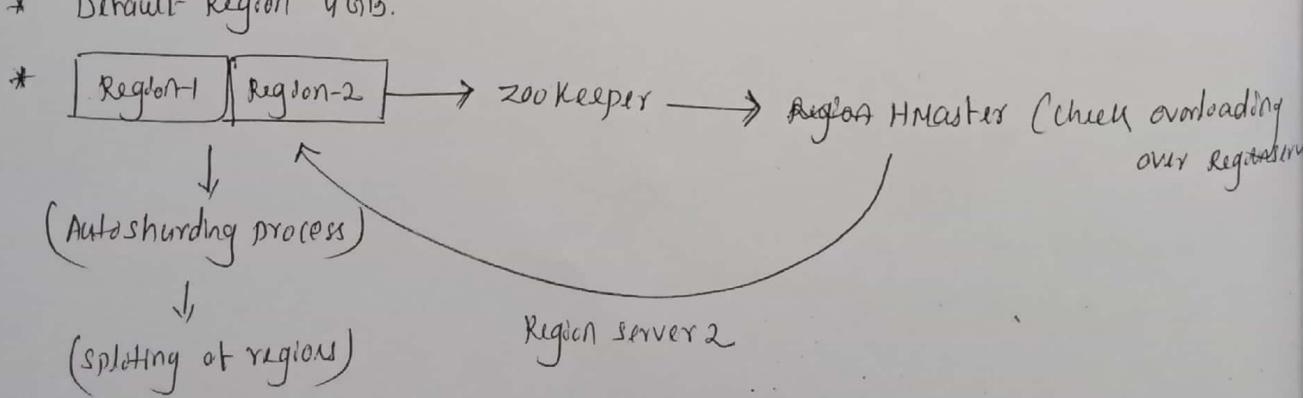
- put → explicitly Remote IP
implicitly → Local host



NOTE → - copy from Local → Local
- put → Different systems / datanode
- distcp → cluster to cluster.



- * An empty directory created in HDFS
 - * .META directory created in HDFS \rightarrow HBase Root directory \rightarrow zookeeper
 - * 1 region created 1st time
 - * RAM of datanode will be dedicated for buffers.
 - * zookeeper will download the .META to the datanode and read the RAM and do the operations.
 - * When request for particular Rowkey, it merge memstore and Hfile then sort and get the details.
 - * Range of Rowkey considered for point of making regions/partitions.



* Region = Memstore + Hfile.size

Region ID : (start key, HRS Name) }
 ↗ value
 ↘ key
 ↘ value

[hr1 hr2 hr3 hr4] - Buffers

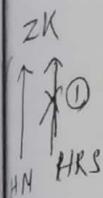
[cf1 cf2 cf3 cf4] - column families

put std::tbl, 100, {cf2 : {cf1, xxx}}

R1 → 0
 HSR10

R2 → 1001
 HSR15

case1 → 500 → ↗ key → R2 key
 compare and assign to the
 region.



→ will check which regions assigned to HRS.

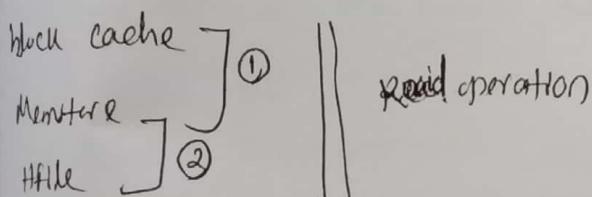
WAL file will have the delta values for the regions.

EN → NN

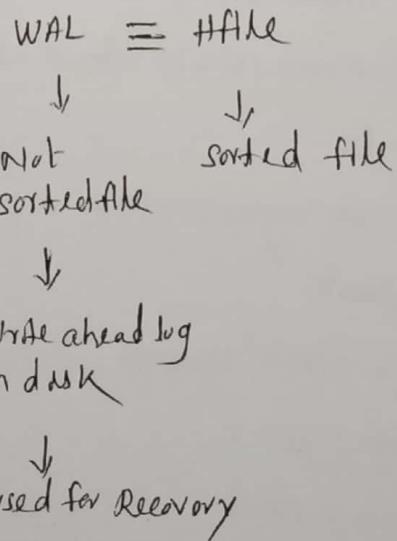
writeJob → JT

Unit → ZK

→ Metacache will be changed as META is changed
 with new timestamp.



Memstore || write operation
 Hfile



Deleted cells marked will be removed from the buffer and put in a large file - hard delete.

Collect files → merge → sort → remove deleted file → one large file.

Read Amplification

Write Amplification → maintenance

time frame → done by the hadoop Admin

{ 26.01.2018
Saturday }

Apache hadoop

[Cloudera
Hortonworks
MapR]

commercial distribution of hadoop

→ All ecosystem tools like hive, pig have to download and install manually.

→ ls -p

- hbase-0.94.10.tar.gz

- cd hbase

- cd lib

- ls

- Jars for hbase

- cd doc

- ls

- hbase conf

- vi hbase-site.xml

- cat hbase-policy.xml

hbase-managed-zk=true-Internal

cont > ls

cont > cat region servers

localhost

hbase > which hadoop

hadoop > ls -p

hadoop > pwd

hadoop > cd conf

cont > ls

(core-site.xml)

cont > cat masters

cont > cat slaves

- pseudo distributed mode - one system

- fully distributed mode

- jps - 3280 JPS

(cont) > hadoop dfsadmin -safemode get
safe mode is on.

(cont) > hadoop dfsadmin -safemode get
safe mode is off.

> start-hbase.sh

- sleep 5

hadoop > jps

4335 HQuorumPeer

bin > hbase shell

hbase(main):001:0> help

hbase(main):002:0>

003:0> status

- 1 servers

: 004:0> version

- 0.94.10

> whoami

- hadoop training

hbase → Log structured filesystem

HDFS → Block structured filesystem

hadoop fs -ls /hbase/

hadoop fs -ls /hbase/student/

hbase > describe 'std10'

TTL - Time to live - the record will mark as deleted when time elapsed.

hbase > count 'aadhar-tbl'

hbase > scan 'aadhar-tbl'

> scan 'std-tbl'

row column + cell

1 column = cf1:bc ...

1 column = cf1: name ...

1 column = cf1: no ...

2 - - - - -

2 - - - - -

2 - - - - -

> put 'std-tbl', 9, 'cf1: name', "Nagendra"

> get 'std-tbl', 9

column cell
cf1: name timestamp = 15..., value = Nag

> get 'std-tbl', 9, 'cf1: loc'

column cell

BLANKS

> get 'std-tbl', 9, { column ⇒ 'cf1: name', versions ⇒ 3 }

column

cf1: name

(f1): name

cell

timestamp = ..., value = Nag

timestamp = ..., value = Nagendra

> put 'std-tbl', 9, 'cf1: name', "Nage"

> put 'std-tbl', 9, 'cf1: name', 'Nag'

CF1: name timestamp = \dots value = 'Nage'
 CF1: name timestamp = \dots value = 'Nage'
 CF1: name timestamp = \dots value = 'Nag'

* 3 versions will be kept
and old 1st version will be soft-deleted

hbase > delete 'std-tbl', 9, 'cf1: name'

hbase > deleted 'std-tbl', 9

hbase > create 'nov-2018-emp-tbl', 'nameDetails', 'parentDetails', 'IdDetails', 'deptDetails', 'dept'

hbase > exists 'nov-2018-emp-tbl'

→ Table 'nov-2018-emp-tbl' exists

hbase > scan 'nov-2018-emp-tbl'

hbase > count 'nov-2018-emp-tbl'

hbase > describe 'nov-2018-emp-tbl'

{ 'nov-2018-emp-tbl', { name => "IdDetails", data block \dots }

\dots { name => "parentDetails", \dots }

{ name => "nameDetails", \dots }

{ name => "DeptDetails", \dots }

→ column family

hbase > create 'nov-2018-emp-tbl', { name => 'nameDetails', versions => 5 }, { parentDetails }

hbase > describe 'nov-2018-emp-tbl'

{ 'nov-2018-emp-tbl', { name => 'nameDetails', VERSIONS => 5, \dots }

Properties can be changed while creating table/attribute

hbase > drop 'nov-2018-emp-tbl'

Table must be described to delete the table → error.

hbase > truncate 'nov-2018-emp-hbl'

- deleting table
- dropping table
- creating table

hbase > describe 'nov-2018-emp-hbl'

HBase & Hive Integration

{ 29.01.2018 }
{ Tuesday }
* must *

HBASE

- Joins and aggregation
- (for unstructured data, may not be common columns) → flexible

→ CRUD

- can use any serialization

Key based access

- No many where conditions and all.

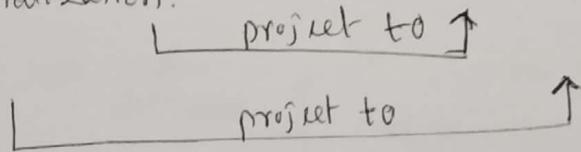
Rigid - it uses its own serialization.

HBase storage handler



HBase server

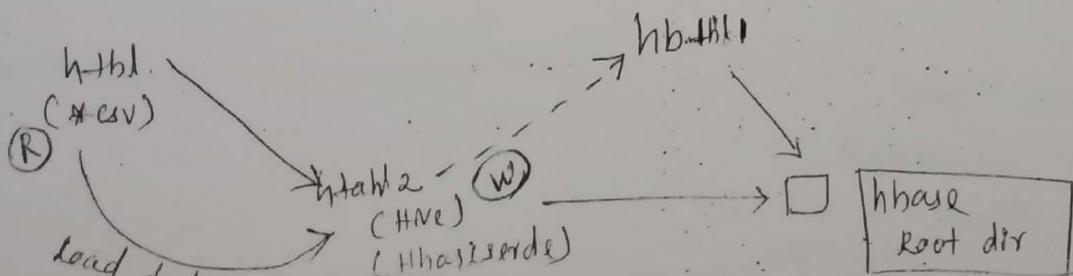
(serialization / Deserialization)

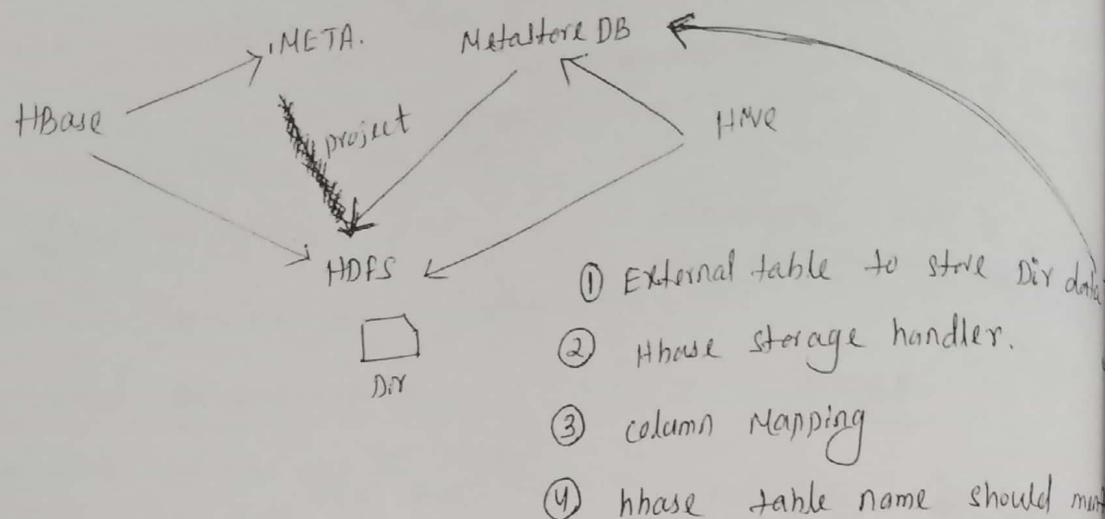
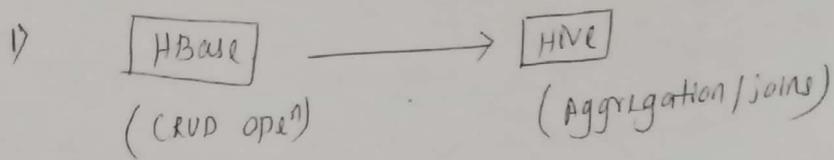


Scenario

1) HBase table $\xrightarrow{\text{view as}} \text{Hive table.}$ → Deserialization
read as

2) Hive table $\xrightarrow{\text{Load into}} \text{HBase table}$ → serialization
Bulk load





Schema on Read \rightarrow Hive uses

Schema on write

- * Blank directory will be created in Hive data warehouse, but read/write will be mapped to the hbase HDFS directory. Later you can drop that directory.

zookeeper-* .jar

hive-hbase-handler-* .jar

guava-* .jar

hbase-* .jar

hive/lib

hbase/lib

hive > add jar /home/...../zookeeper-3.4.3.jar

:

or// hive-site.xml

< name> hive.aux.jars.path</name>

< value> file:///home/.....jar, -..... .jar, - - - .jar </value>

hbase > scan (stdio)

Row column + cell
 | column = cf:lc
 | cf2: name

hive > add jar

hive > list jars

hive > create external table

hive > show tables

hive > select * from htbl_std10_in_hive;

stdid	name	number	loc
1	Ram	123	Hyd
2	Anu	234	null
3	Kum	456	Mys
4	Ashok	555	Bang

hbase > put, ('std10', 9, 'cf:name', "Naga")

hive > select * from htbl_std10_in_hive;

stdid	name	number	loc
1	Ram	123	Hyd
2	Anu	234	NULL
3	Kum	456	NULL
4	Ashok	555	NULL
9	Naga	NULL	NULL

hive > describe htbl_std10_in_hive;

Table parameters:

External TRUE

hbase.table.name std10

storagehandler org.apache.hadoop.hive.hbase.HBaseStorageHandler

storage description params:

hadoop >

/ htbl_std10_in_hive/

hadoop > hadoop fs -ls /hbase/std10/

② Hive to hbase

hive > select * from ds-db.emp-test-tbl;
> create table tmp-hive-thl-emp-to-hbase (---)
stored by 'org.apache.---'
with serdeproperties ('hbase.column.mapping' ---)
TBLPROPERTIES ('Jan-2018-emp-thl');

> show tables

hbase > exists 'Jan-2018-emp-thl'
hbase > scan 'Jan-2018-emp-thl'
hive > insert overwrite table tmp-hive-thl-emp-to-hbase select *
from ds-db.emp-test-thl;
--- maponly job

hive > select * from tmp-hive-thl-emp-to-hbase
--- Data ---
--- Data ---

hbase > scan 'Jan-2018-emp-thl'

⇒ web-table

- recommendation results

Huge Dimensional Table

- customer information.

streaming data injection

- hbase as sink

- Any networking data

- Any web apps data

⇒ Mobile apps

1) CLI → hbase shell → cluster data node

2) REST API → Any application / framework

Any lang supports REST APIs

3) JAVA API → JDBC like

OOZIE

Automation for Hadoop Ecosystem.

Workflow / Job flow

DAG - Directed Acyclic Graph

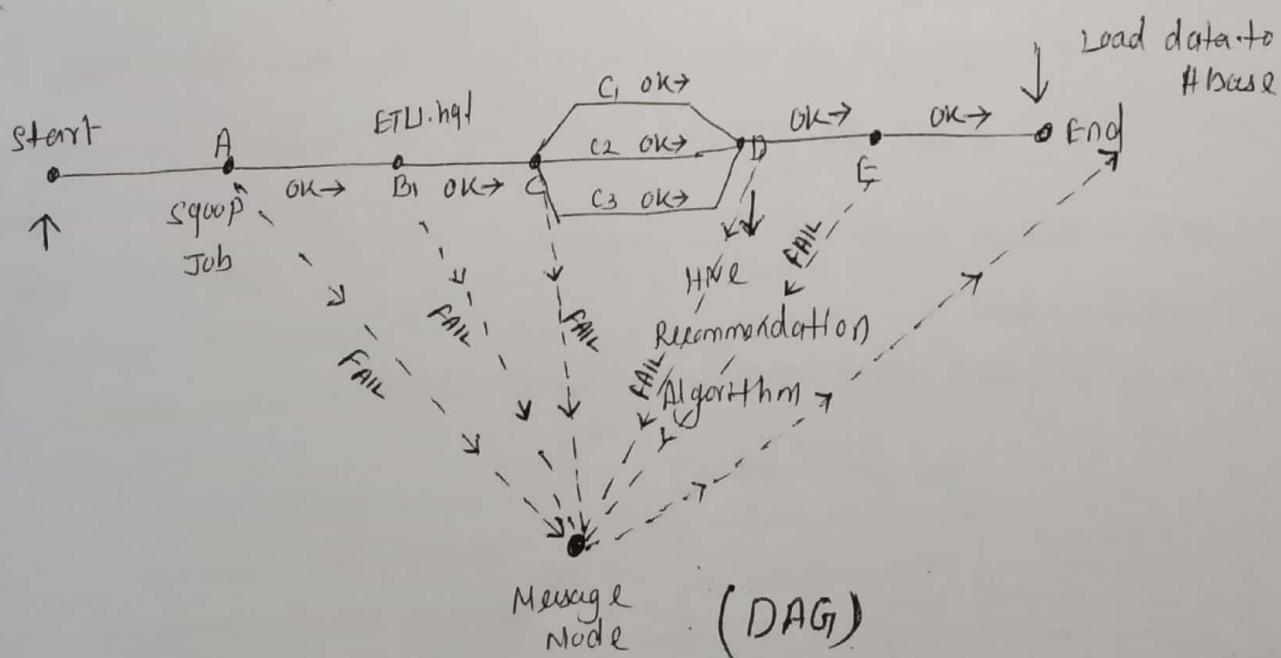
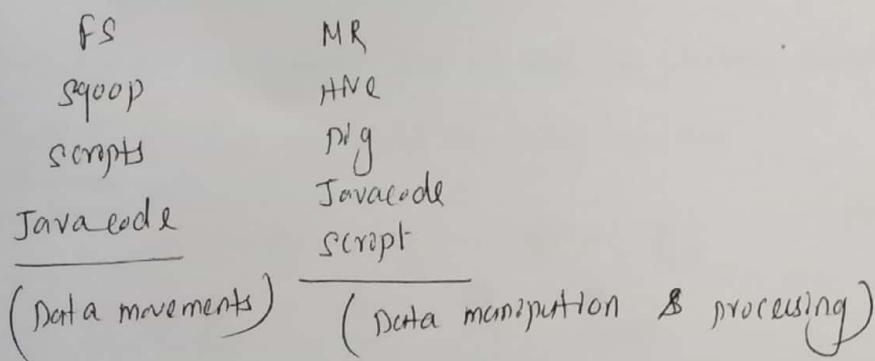
Scheduling → creating oozie co-ordinators

Automation → creating oozie workflows

Scheduling → creating oozie co-ordinators

> Oozie bundles

→ combination of automation and scheduling



workflow.xml → Expression language

HPDL → Hadoop process definition Language

communication between the user and Oozie job definition to understand the workflow with XML language.

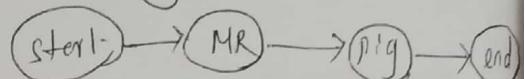
< workflow-app >

```
< start to = "Action-1" >
  < action name = "Action-1" >
    < hive >
      < script name = ... >
        < args ... >
        < /args >
      < /script >
    < /hive >
    < ok to = "Action-m" >
    < fail to = "Action-n" >
  < /action >
< /end >
```

< /workflow-app >

* At each step the action is defined for ok and fail scenario.

* We can see the diagram while running the job itself

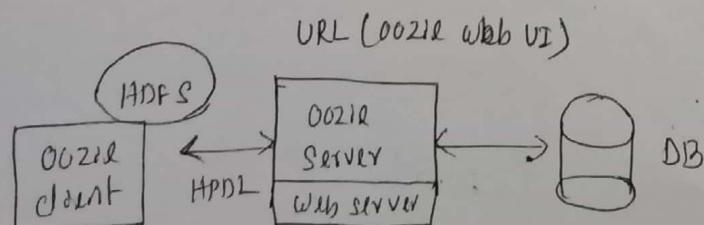


(DAG)

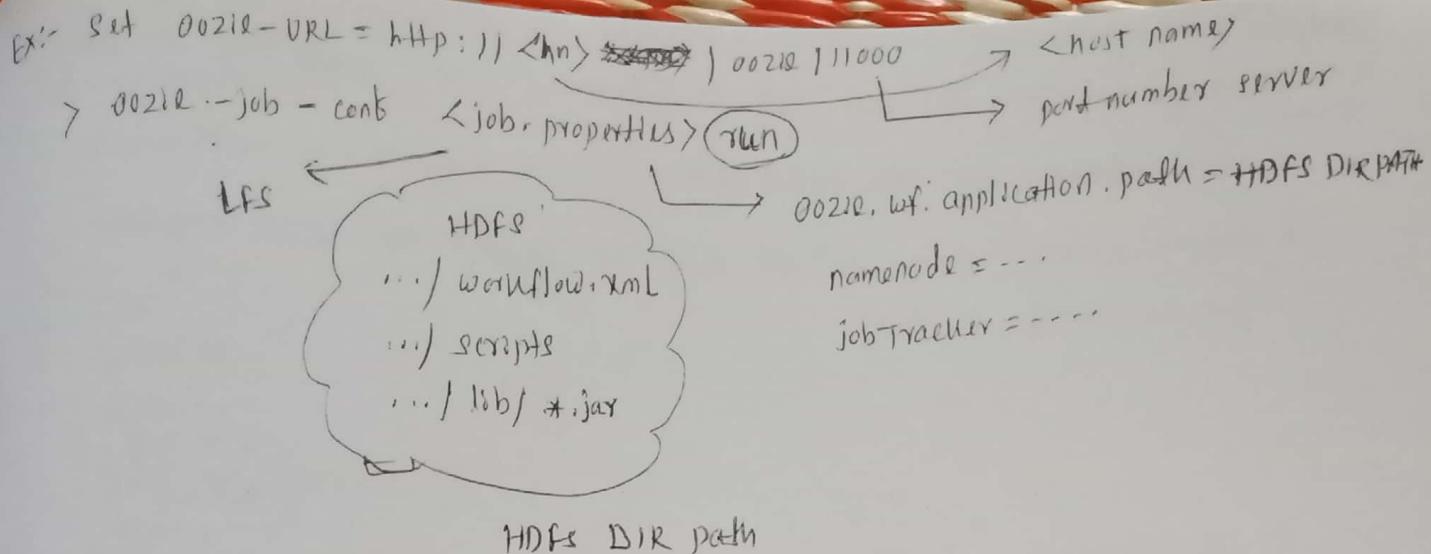
Oozie Architecture

{ 30.01.2018
* wednesday
* mon day

- Client-Server Architecture
- Oozie Server
- Oozie Client → CLI *
 - Java API
 - Rest API
- HDPD - Hadoop process Definition Language
- Written in Java
- Closely integrated with Hadoop and ecosystem tools.



Edge node
Cmd: oozie-job
{ Tomcat
Weblogic
Glassfish } { MySQL
Oracle
DB2 }



> connected to server with 11000 port.

> job.properties file read.

> go to mentioned directory, submitted to DB

> jobid created - oozie.job-id: wf-01234

> process the job.

> Transition information (from one node to another) stored in Database.

> JDBC to DB is stored in oozie.xml file.



> oozie.job - the wf-id → can be access the details for the job id.

> workflow.xml - contains the xml tags containing the HDL to process the data.

<workflow-app>

<start-to="ACTION-1"/>

<action name="ACTION-1">

<pig>

<script> etl.pig </script>

<param> 1 date </param>

<param> \$outdir </param>

</pig>

<ok to="ACTION-2"/>

<fail to="Error-Node"/>

</action>

<action name="..."/>

</action>

<action name="..."/>

</action>

<null name="Error-node"/>

<message>

</message>

<kill>

<email> --- </email>

<end name="end"/>

</workflow-app>

* of will start from the last returned failure/ success

* we can define the step failed node.

* current running status for the node is stored in DB, so at any time you can resume the job.

1) job properties → LFS

↳ Nomnode = hdfs://nm: port /
Job tracker = host-name: port /

0012, wf, application.path = <hdfs-path>

input dir =
output dir =

2) workflow.xml

↳ workflow-app →
<start-to=...>

<action name=...>

<hne> --

<action>

<kill...>

<end>

<workflow-app>

3) ↳ Script & Jar → <hdfs-path>/
mu files

Job /

* Queue name - Dev Testing, Prod

choose one queue under which you want to run the jobs, within the cluster.

* response and configuration tags are optional.

* test or export

host - URL = ...

host - config Job.properties - run

host example > host.job - into 000000 - ...

host example > host.job - into 000000 - ...

Oozie Co-ordinator

01.02.2019
friday
not not

Automation → workflow.xml → oozie-job - cont job.properties - run

Scheduling → Running in a future time

[
Cron
control M
autorsys
Tivoli WF manager

+ oozie co-ordinator tool

→ frequently
→ window

01-02-2019 → Data dependency

} co-ordinator.xml
co-ordinator.properties / job.properties

..... /year/ month/ day/ hr/

2019 | 01 | 31 | 12

- date

2019 | 02 | 01 | 00 |

time

| 00 |

date - v

| 02 |

| 03 |

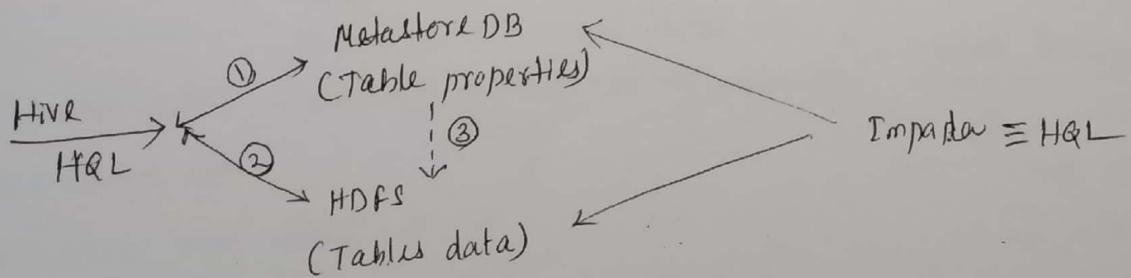
<input-events>

<work-flow>

oozie-job - config co-ordinator.properties start - run

02.02.2019
saturday
not not

IMPALA



1) syntax check
2) query resolution
3) query planning
4) MR code → MR Architecture

5) physical execution

6) result → console

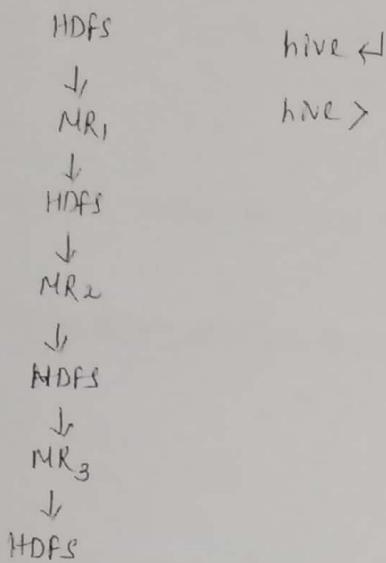
7) result → written HDFS

Metastore DB project details to HDFS to fetch the data.

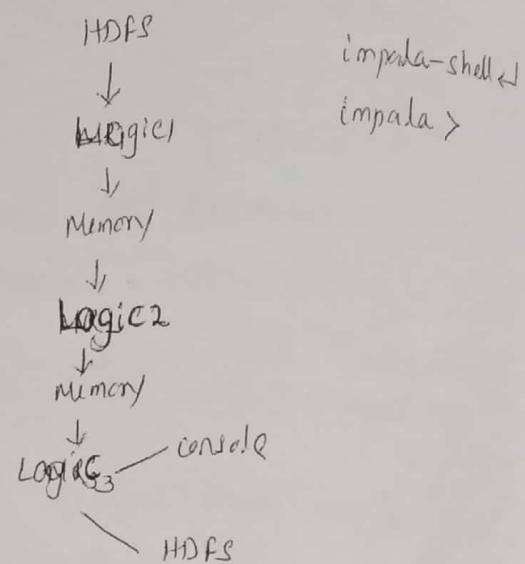
Difference - ① MPP Architecture

Massively parallel processing

HIVE



Impala



- * Table created in hive can be seen in impala.
- * developed by facebook
- * More HDFS process, I/O cycle, serialization needed, so it is slower than Impala.
- * Good for large dataset
- * Long running batch processes.
- * slow compared to Impala for aggregated data.

- * created in Impala can be seen in hive as they are sharing the common MetastoreDB.
- * developed by cloudera
- * intermediate memory are there, so Impala is faster than hive.
- * Good for small dataset, as memory limit
- * good for quick, short queries.
- * ideal for aggregated data, line reporting

H CATALOG

HCatalog → Metastore DB
Interface (HDFS)

→ DDL operations. It supports not DML operations.

→ HeatLoaders

→ pig

A = Load Table <dh-table> using HeatLoader();

→ HeatLoader will get the dh-table → goto metastore DB → Load data to
Redirection A → convert hive datatype to pig datatype.

B = for each A generate trans-id, trans-data.

result = ...

store result into table <dh-table> using Heatstore();

- HCatalog input format

- HCatalog output Format

- HCat store will convert pig datatypes to hive datatypes.

→ 1) writing to HDFS

2) converting datatype

3) metastore DB updated with corresponding tables.

Cloudera@quackerstart:~] heat

] hive -e "show databases;"

] heat -e "show tables;"

] heat -e "create table ---;"

Show create table - used for transformation/automation.

Create table like

hive > show create table emp-tbl;

hive > create table emp-dup-tbl like emp-tbl;

hive > desc emp-dup-tbl;

HQL for running

hive <

BEELINE

hive >

hive based on Thrift server protocol not providing any security. anyone can connect to the edgenode. it will not ask for authentication.

Thrift server 2 protocol

→ hive JDBC url

→ userid

→ password] LDAP authentication

① || → beeline - connect "jdbc:hive://<host:port>: <port> ..."
-u <userid>
-p <password>

→ beeline > --- hql

> -- all hive query u can run

③ || > beeline -file ABC.hql -hivevar X=x,
-hivevar Y=y,

Thrift server 2 protocol developed in sqlite.

beeline <

② || beeline > ! connect <hive-JDBC-URL>
> enter userid
> enter password

* Either 1 or 2 or 3 you can choose to work with beeline.

* -i <filename> - initialization session for hive.

cloudera] which beeline.

/usr/hive/beeline

] beeline--help

* beeline provides the security feature to hive.

NameNode High Availability

1. X

2. X → hadoop MR written again

→ under utilization of resources mapper/reducers

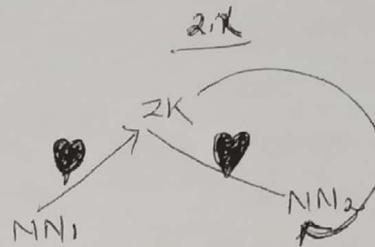
→ container can do map/reduce both.

→ TT supposed to do all the processes, if any process down all the process go down.

→ SNN - replica of fsimage, metastore.

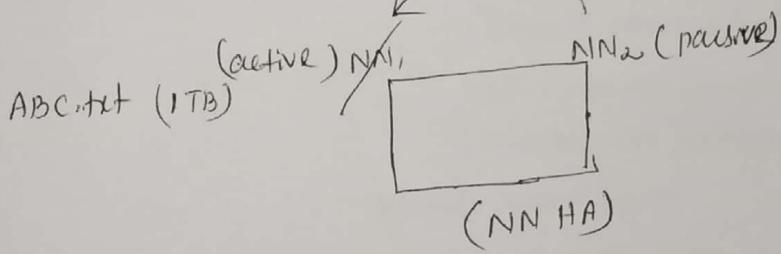
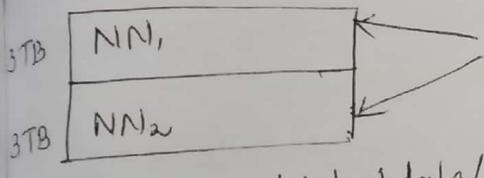
1. X

NN SNN



* when NN1 sending heart beat it will inform NN2 to become active.

NameNode Federation



* headache over hadoop admin.

* More space will be consumed as 6 replica will be done with single file.

* so NN HA is preferable in real time.

hadoop fs -put <dfs> hdfs://nn2:port/dvs/data/

you have to write to NN2 to read the data.

* you have to write to NN2 to read the data.

* NN fed → Having more than one NameNode in the same cluster.

Scheduling

Job scheduling

FIFO: Default

dev - Job1
QA - Job2
prod - Job3

* one after one
* prod have to wait for dev and QA to complete the job
* so accepted in real time.

1) fair scheduling → facebook
2) capacity scheduling → yahoo

Queues → To assign the resources the groups as per the task. Admin suggests que

default - 5%.

batch - 35%.

task - 25%.

I/O - 35%.

1) A cluster

- About

- nodes

- node label

- Application

New

--- Submitted, accepted

--- finished, killed

- scheduler

2)

Application queues



→ queuepool

→ queueidefault - green - means job is running

- queue: batch

- queue: I/O

3) Queue status: running

Used capacity: 0.0%

- Submit the job

Configured capacity: 35%.

- Accepted the job

Configured max capacity:

- Application master allocation

↳ can borrow capacity

- running state

Absolute used capacity:

- finished

Absolute configured capacity:

- success / killed

Absolute configured max capacity:

Used resources: (memory: 355352, vcores: 13)

Configured max application master limit: 2000

Max application master resources:

Used application master resources:

Num container: 13

Max Applications: 3000

Max application per user: 3000

Configured minimum user limit percent: 100%.

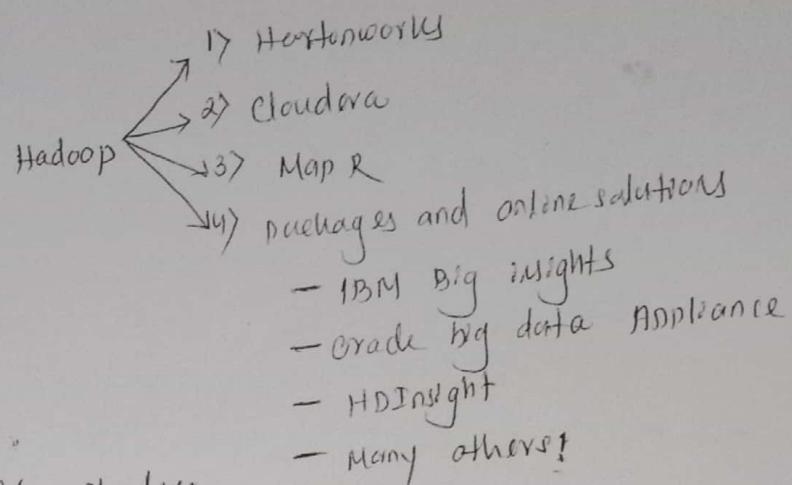
Configured user limit factor: 1

Hadoop Distributions

NASSCOM - Journals

Choosing a hadoop distribution

- company philosophy
- current relationships
- acceptable risk
- specialized functionality

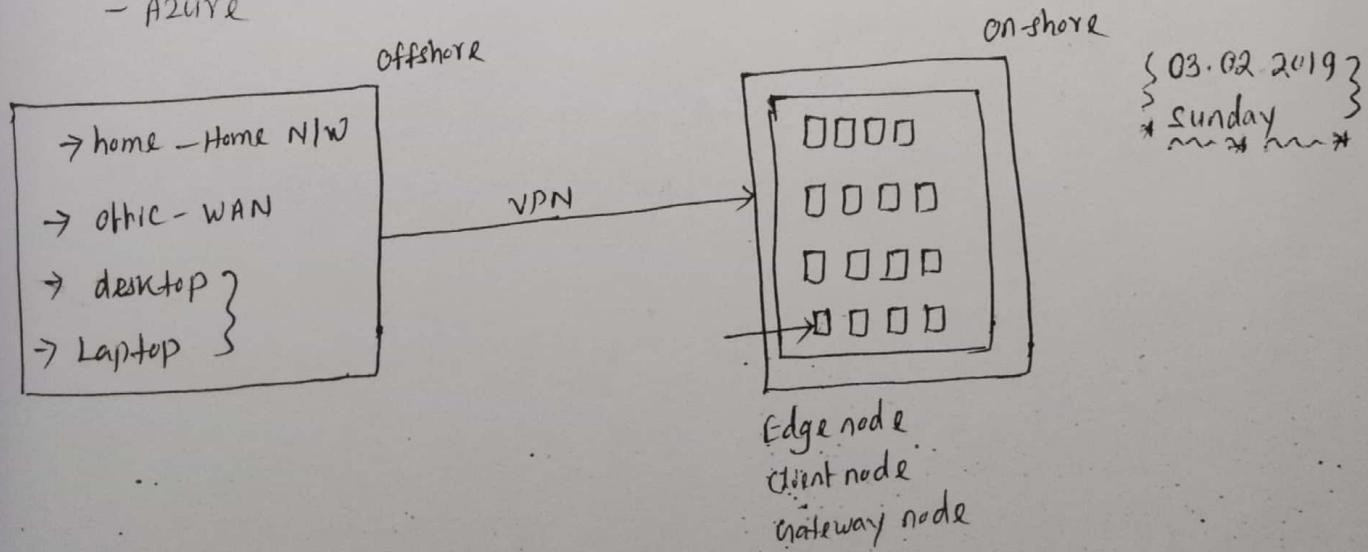


- 1) Apache projects only - Hortonworks
- 2) Apache project + proprietary Add-ons - Cloudera
- 3) Map R - proprietary value adds & redevelopment - own filesystem.
 - ↳ providing high performance.
- Ambari same as Cloudera manager to provide Admin job easier.
 - ↓
 - + HUE, IMPALA

* Hortonwork started YARN: Data operating system.

Cluster Types

- own cluster
- cloud cluster → difficulty while re-migration from cloud
 - AWS
 - GCP
 - AZURE



- 1) VPN - Virtual private Network
- 2) Cisco Anyconnect - client
- 3) RSA APP \Rightarrow pin number \Rightarrow passcode

http://fcs.in/VPN_conn

userid	<input type="text"/>
passcode	<input type="text"/>

option-1 Edge Node : IP#

host name

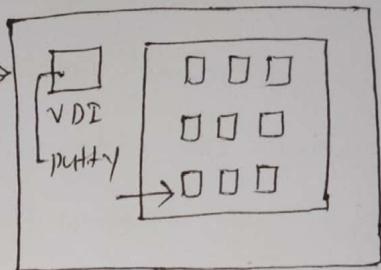
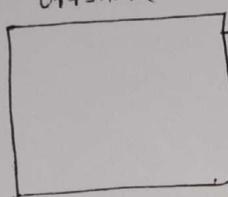
Unix userid : Nagendra
sg80012AB

password : * * * * *

SSH to connect the edgenode
through putty terminal.

option-2 VDI

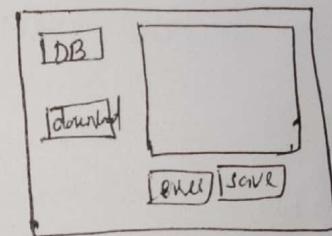
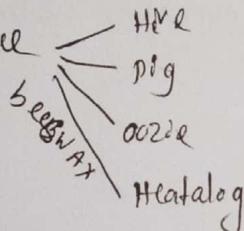
offshore



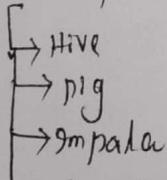
VPN \rightarrow VDI \rightarrow putty \rightarrow EdgeNode

* same sequence for onprem cluster or the cloud cluster.

option-3 VDI \rightarrow HUE Interface



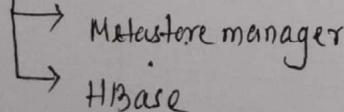
HUE Query Editors



DB query - backend for MySQL

- Manually added through JDBC connection.

Data browsers



Workflows

\hookrightarrow oozie Dashboards - resume, suspend, kill the jobs.

Search

\hookrightarrow solar search