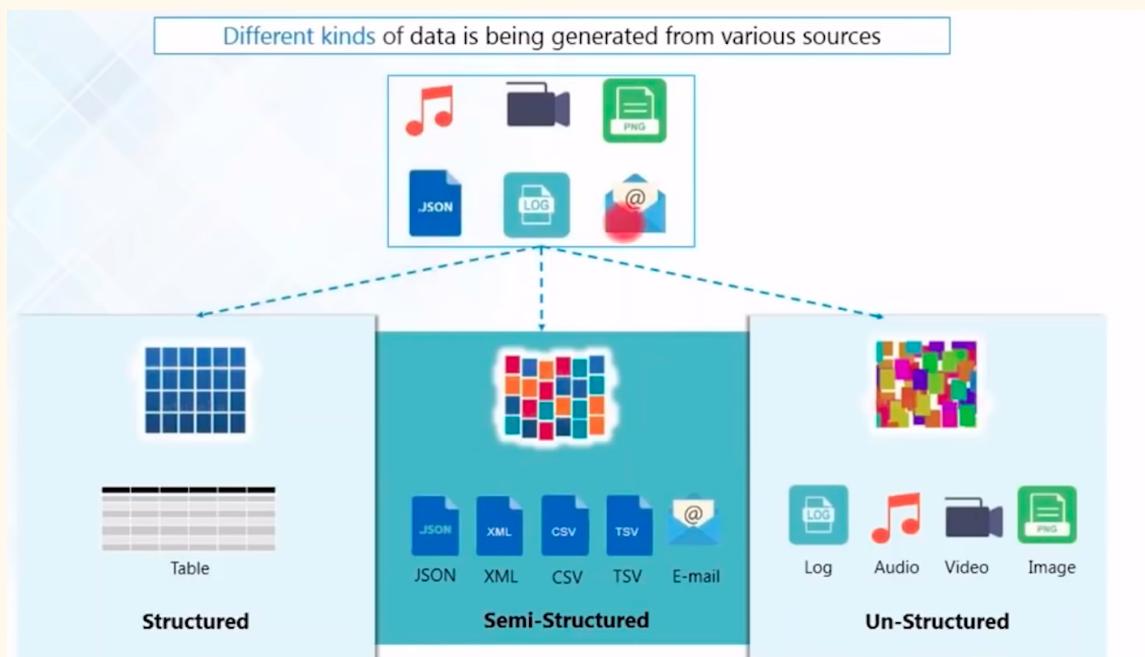
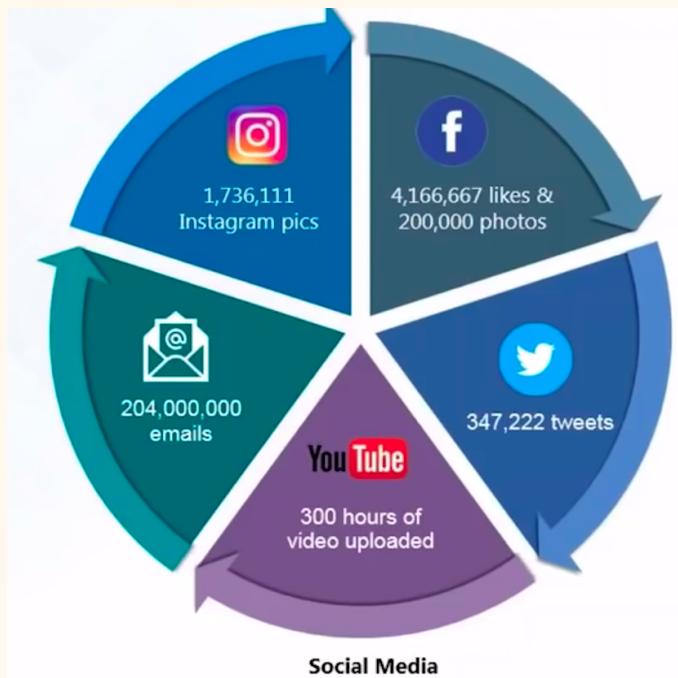


Big Data

What is Big Data?



DATA Analytics -

Everyone is creating almost 40 Exabytes data per month

Careers in Big Data -

- Data Analyst
- Data Scientist
- Big Data Architect
- Data Engineer
- Hadoop Admin
- Hadoop Developer
- Spark Developers

Streaming Components



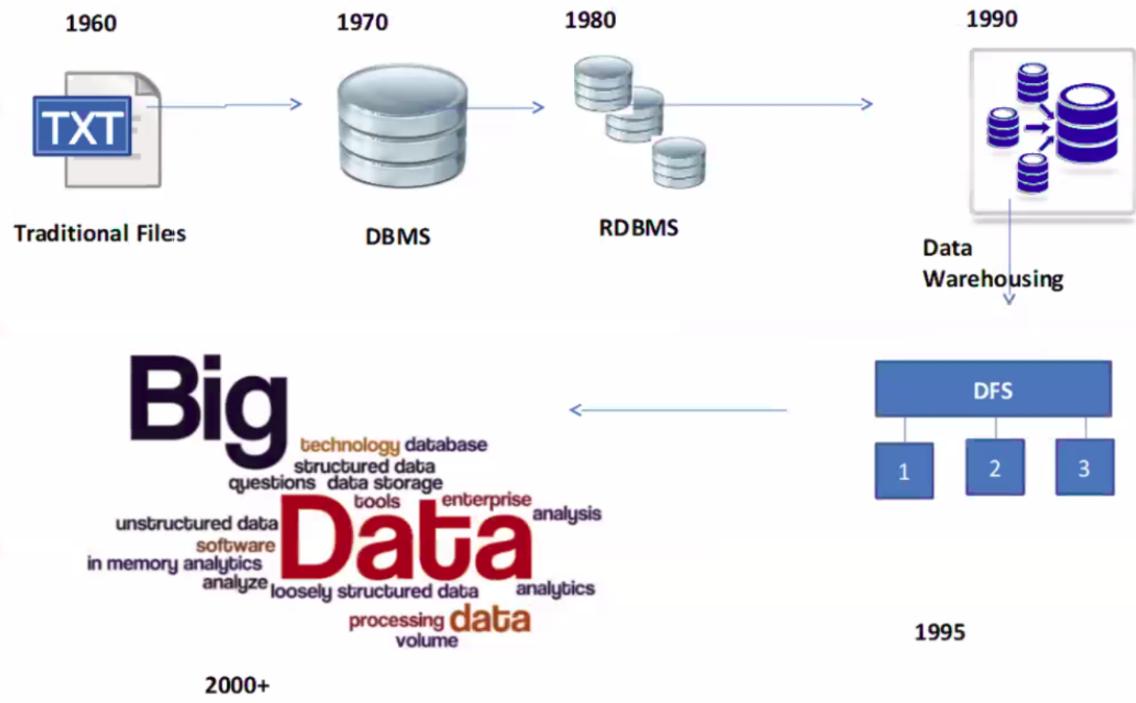
Big data Evaluation -

- 1960 - txt files
- 1970 - spreadsheets/DBMS
- 1980 - RDBMS but Licensed
- 1990 - Data warehouse (many RDBMS) - Damn costly
- 1995 - DFS found by Doug Cutting but failed due to processing wasn't good
- 2003 - Google released GFS paper same as DFS == GFS
- 2004 - Google Released Mapreduce paper → Distributed Processing → Doug Cutting read the paper and was Shocked
- 2005 - DFS enhanced and use Mapreduce
- 2006 - DFS + Mapreduce = Hadoop (name on son's elephant Doll) , HDFS and Mapreduce

- 2008 - Emerged a lot -- Problem -- Hey Mapreduce code is damn hard to understand
JAVA → Doug Cutting said --> do not worry I have a solution to simplify MapReduce
→ Learn SQL → Trigger SQL → I will ensure my TOOL converts SQL query to
Mapreduce JAVA
- I found a tool, it looks like SQL but not --- If you trigger any queries -- it converts into a
Mapreduce JAVA
- In 2008 - Hive and Pig came into the market
 - Mark Zuckerberg --- Facebook using --- HIVE
 - Yahoo--using -- PIG (Not Popular)
- 2009 - **CLOUDERA** - Problem was --for installation and Maintenance -need Laptop
 - Give me Laptop But Installation, maintenance, and Upgrade - are Paid services
- 2011 - **HortonWorks** - Give me laptop and Installation → Free, Maintenance and
Upgrade are Paid
- 2014 - **SPARK** - Birth of a super Power
 - Superpower, Process data very faster, Free of cost, Will perform SQL, Will
support Streaming, Will support Machine learning, Can run ON Hadoop
- 2016 - **Cloud** - I will give u **laptops**, I will **install, upgrade, and Maintenance**
 - AWS EMR - Single Click of Button Hadoop Laptops will be ready within 10 Min by
AWS
 - Azure HDInsights
 - Google data Proc (Google cloud platform - GCP)
 - Alibaba cloud
- 2017 - Cloudera Bought Horton works
- 2022 -- We live in the Cloud Big Data market = Cloud-based Hadoop Big Data analytics

Cloud Hadoop Spark ==> Crazy Technology today

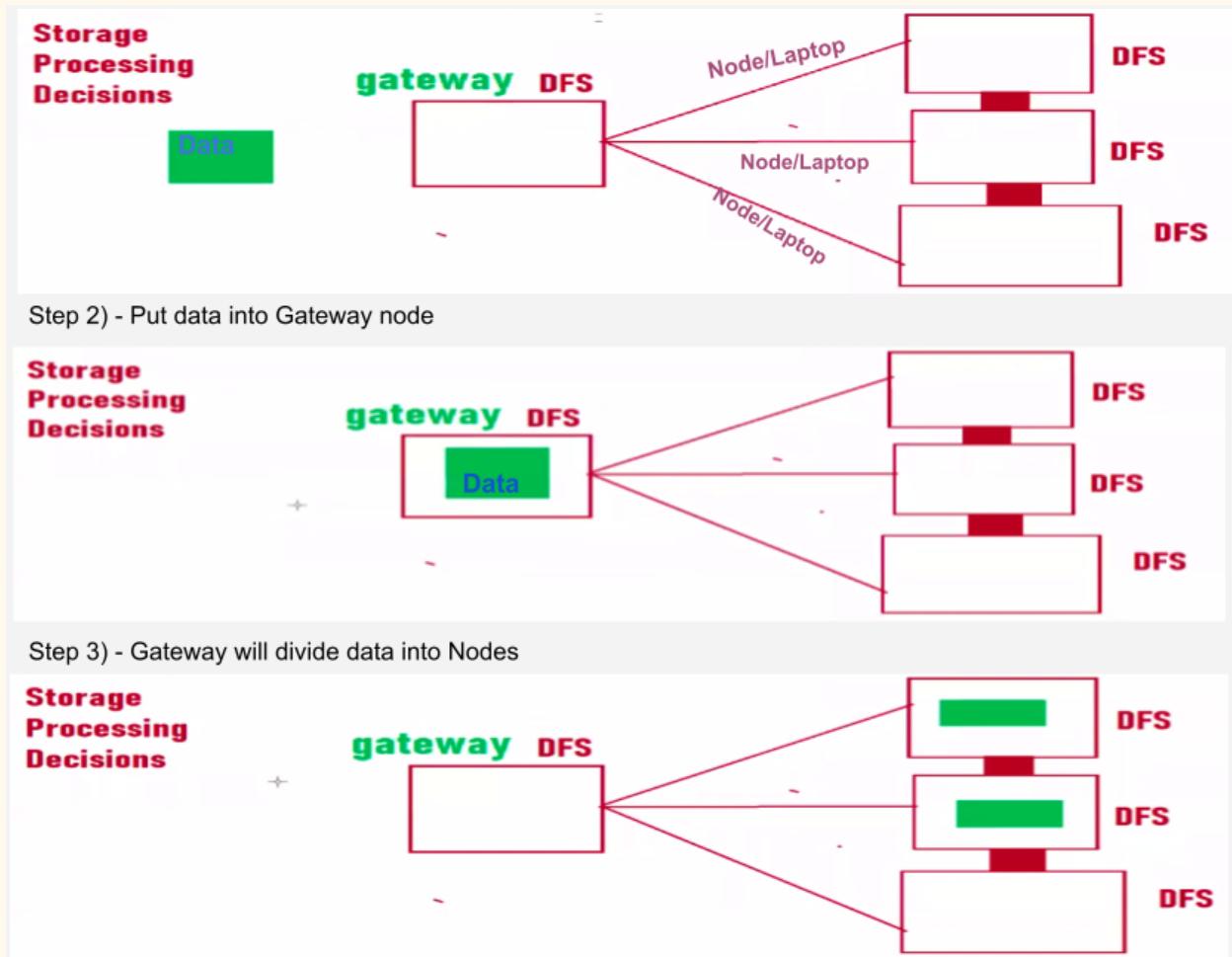
Evolution of Data



What is DFS -

Introduced by **Doug Cutting** in 1995.

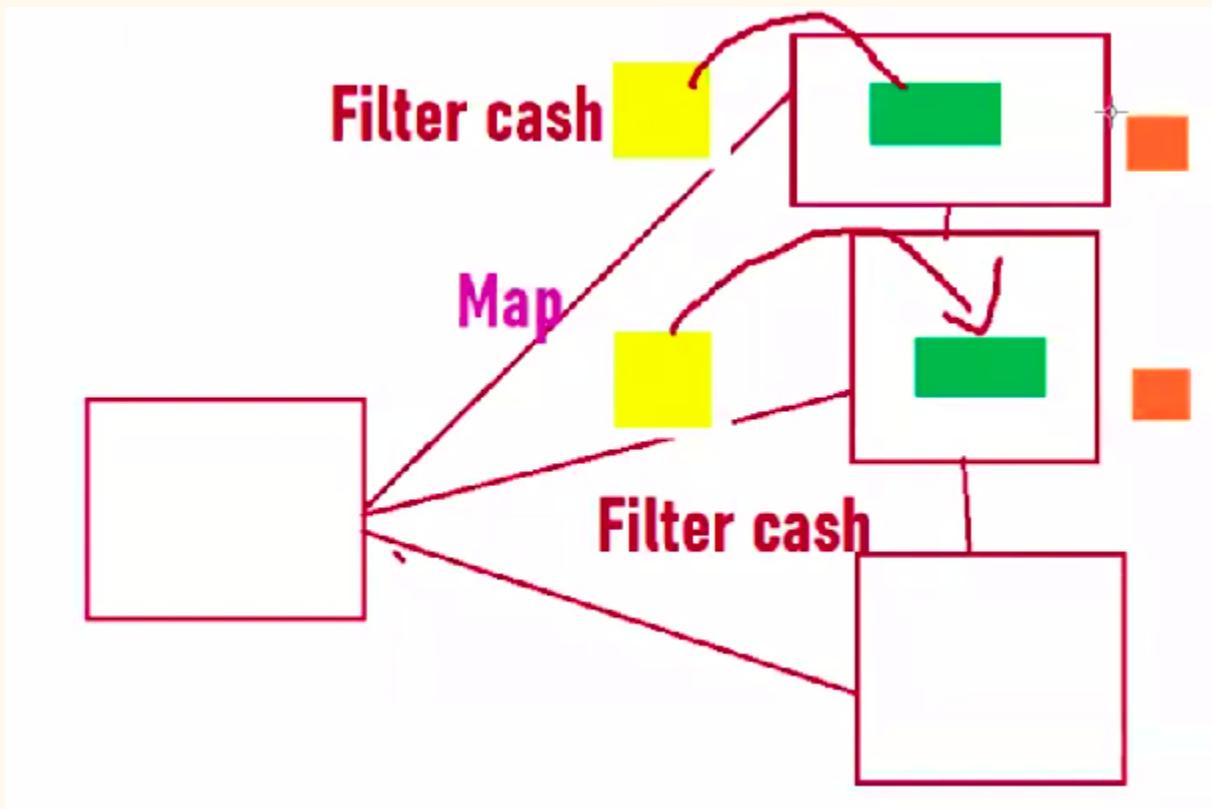
- License-free
- We don't need RDBMS servers
- How DFS works - for example there are 4 laptops - 1 is Gateway and other 3 are interconnected to each like below image -
 - Gateway - Once you give data, data gets stored on the Gateway node first then after running a few commands it is distributed in other 3 nodes.



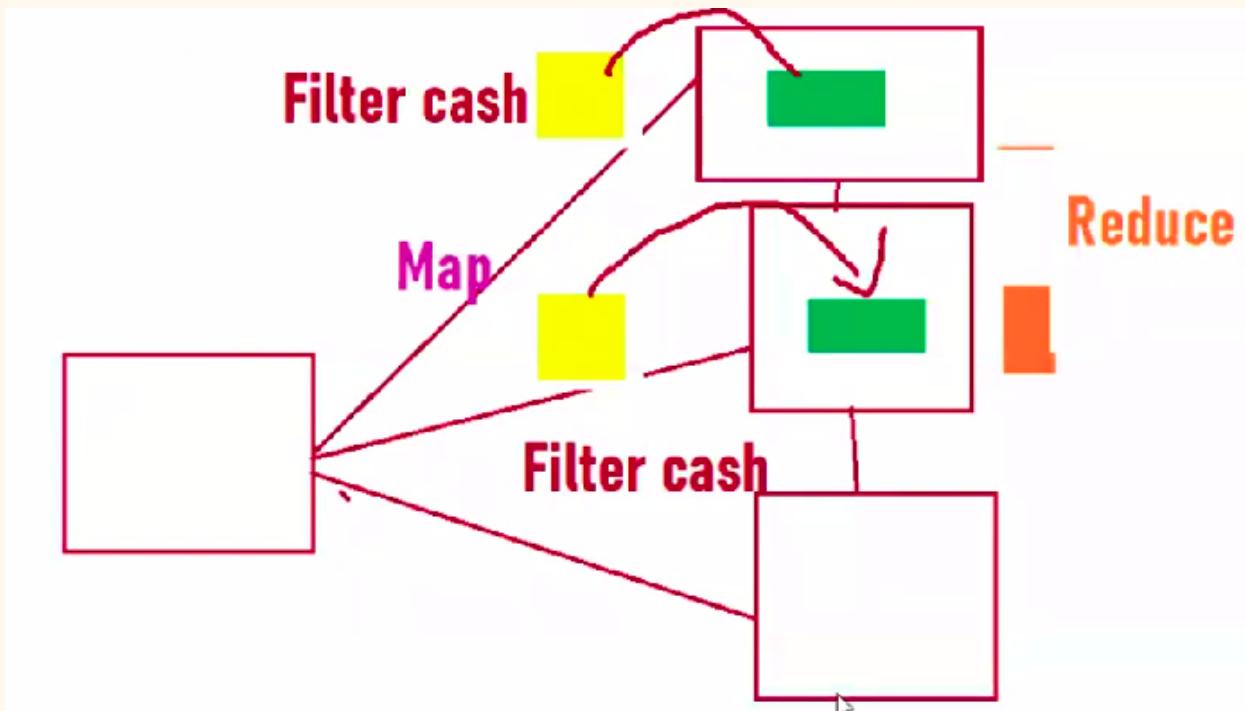
- **Why DFS failed** - As per Doug cutting if data increases then have to increase Gateway Node size as well to process the data. Which was really an issue, as putting a lot of money for Gateway Node

How is GFS different from DFS -

GFS(Google file system), they are using the same mechanism to store data as DFS but for processing data GFS uses Map-Reduce. It means data won't come to process instead process will go to data and Map it after collecting data from all nodes data gets reduced and returned.



Data is getting reduce -

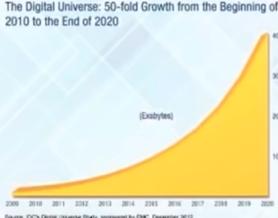
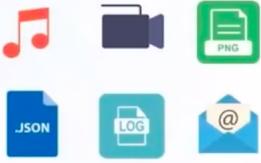


Why Big data?

- 80% of data is unstructured which is challenging to analyze.
- Structured formats limitation in handling large quantities.

5 V's - (now its 6 V's)

1. Volume
2. Variety
3. Value
4. Velocity
5. Veracity - Data can't be perfect in a realistic environment, it will have inconsistencies, errors, and noise. For instance, collecting data for marketing surveys but there can't be data from fake accounts.
6. Variability - Data flow is inconsistent with period peak. The same tweets can have different meanings based on the content.

 <p>The Digital Universe: 50-fold Growth from the Beginning of 2010 to the End of 2020 Source: IDC's Digital Universe Study, sponsored by EMC, December 2012</p>	 <p>Different kinds of data is being generated from various sources</p>	 <p>Data is being generated at an alarming rate</p>																				
 <p>Mechanism to bring the correct meaning out of the data</p>	<table border="1" data-bbox="643 1332 1002 1459"><thead><tr><th>Min</th><th>Max</th><th>Mean</th><th>SD</th></tr></thead><tbody><tr><td>4.3</td><td>7</td><td>5.84</td><td>0.83</td></tr><tr><td>2.0</td><td>4.4</td><td>3.05</td><td>50000000</td></tr><tr><td>15000</td><td>7.9</td><td>1.20</td><td>0.43</td></tr><tr><td>0.1</td><td>2.5</td><td>?</td><td>0.76</td></tr></tbody></table> <p>Uncertainty and inconsistencies in the data</p>	Min	Max	Mean	SD	4.3	7	5.84	0.83	2.0	4.4	3.05	50000000	15000	7.9	1.20	0.43	0.1	2.5	?	0.76	<p>• • •</p> <p>V's associated with Big Data may grow with time</p>
Min	Max	Mean	SD																			
4.3	7	5.84	0.83																			
2.0	4.4	3.05	50000000																			
15000	7.9	1.20	0.43																			
0.1	2.5	?	0.76																			
<p>Value</p>	<p>Variety</p>	<p>Velocity</p>																				

Challenges of Big data -

1. Cost
2. Handle a variety of Big data - structured and unstructured
3. Scalability
4. Store the sheer size of data

5. Process the huge data - There is no sense in just storing lots of data if we can't process it.

“That's when Hadoop came to rescue.”

Characteristics of Big data -

1. Hard to store - (Solution) HDFS
2. Hard to Process - (Solution) MapReduce, Apache spark

Big Data layers -

1. Storage
2. Processing/Analysis
3. Testing
4. Data science/ML/AI
5. Automation/Scheduling

Big data Use cases -

Use case 1 - Execution 1

WallMart -- Got a use case

I have problem to Solve

I have two tables in a system ---

cashtable -- lot data for cash customers

creditable -- lot data for credit customers

I have target table (Empty) -- hivetar

Can i get a help to process above two tables and get top 10 customers in the target table

I seek Help

Doug Cutting - Dont worry I have solution
walmart -- what is it ?

DC ----- Hadoop (HDFS & MR)

WallMart --- Can you show your Hadoop?

Goahead I have question ? How come data will reach your system

DC --- We have super ingestions tool something ---**SQOOP**

WallMart ---- ok you dump the data .. who will process in system ?

DC -- Mapreduce to file first 10 customers --- I am going to use tool know as **HIVE** which uses Mapreduce

- Ingestion Tool - **SQOOP**
- Process Tool - **HIVE**

WallMart --- Go ahead ---

DC ---- Hey walmart -- I have used sqoop to bring data - I have used hive to process and top 10 customers I have found

WALLMART --- Where are those top 10 customers -- Is it in your hadoop system

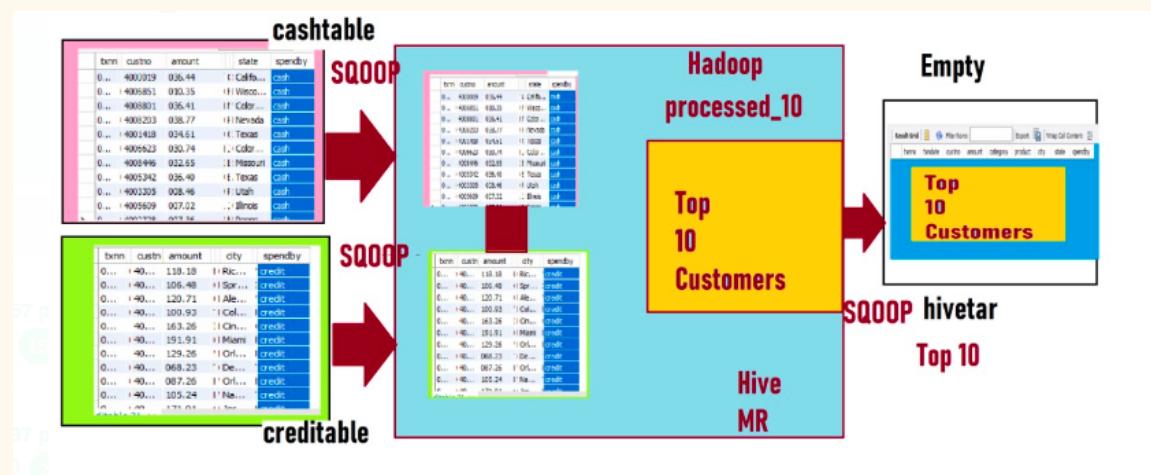
DC --- Yes IT IS Sitting in our hadoop system-- Dont worry

I will use the same **sqoop Tool** to put the top 10 customer to your target system
 Its Done

WallMart ----- KUDOS

WALLMART -- **Doug Cutting** you did well. I will give business To you

DC used followed below method -



Hey This is Jeff Bezos. (AMAZON)

I have another powerfull system —

My system complets the work very faster. Same scenario I can assure I will use damn powerful tools

WALLMART --- What are your Tools

Jeff Bezos --- I am also going to use Hadoop Only. But I going to use it differently

WALLMART --- Tell me your Design like DOUG CUTTING

Jeff -- Give me target Table

Wallmart -- created --- sparktar

Jeff -- I will use my own system S3 AWS

WallMart what is the ingestion you gonna use what is the processing you gonna use "?"

Jeff --- Both are powerfull to me

Ingestion ----- NIFI

Processing --- SUPER POWER --- Spark (Hadoop)

Wallmart --- can you show your system

Jeff --- Sure

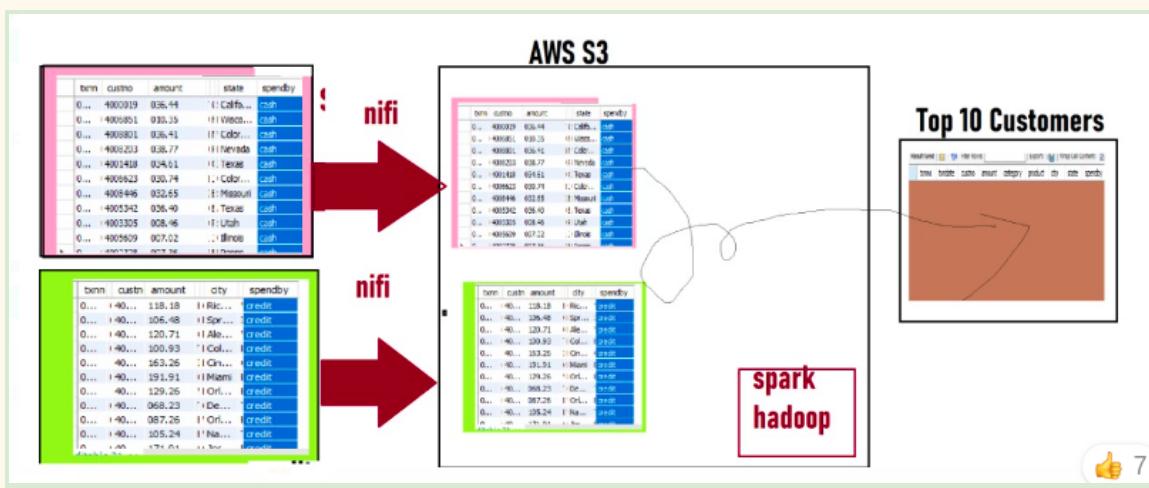
Wallmart --- Go ahead ---

Jeff -- within 3-4 Min-- He said its Completed

WallMart --- Where is the data -- is it in your system

Jeff --- My super directly written your target system

Wallmart --- Perfect I am giving the contract to you



Software needs to install for Big Data learning -

- Windows -

- Putty/Mobaxterm
 - VirtualBox → install Cloudera inside VirtualBox
 - 7zip
- WIndows/Linux/Mac -
 - VirtualBox -> install Cloudera inside VirtualBox

Programming topics need to learn for Big data -

- Architecture
- Access modifier
- OOPs
- Exceptional Handling
- String Handling
- File Handling
- Collections
- Functional Programming
- JDBC connectivity

Vim editor (Edit a file)-

- **Vi <fileName>**
 - Enter ‘i’ to enter in “INSERT” mode
 - To exit from the file press **ESC**
 - then enter “**:wq!**” to exit from vi editor
-

Hadoop

Stored -- Processed- Decision Making

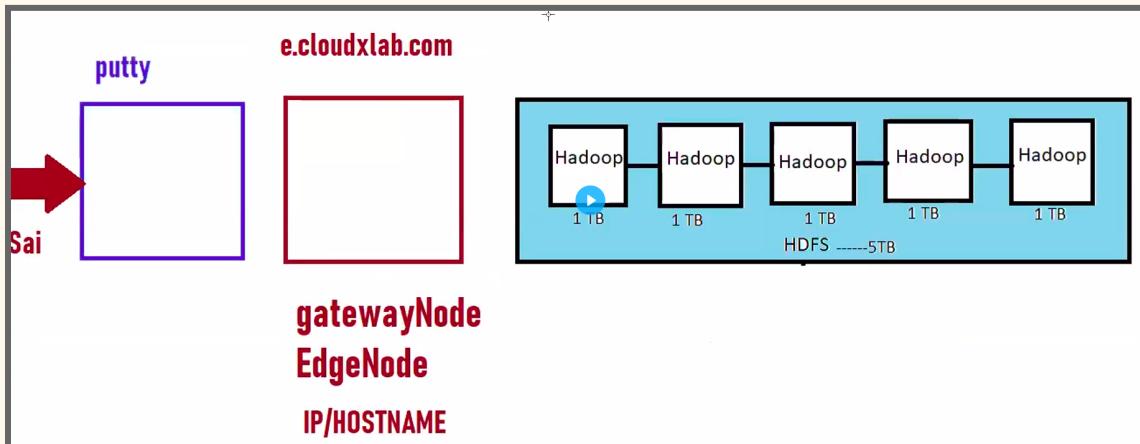
Hadoop is a framework that uses distributed storage and parallel processing to store and manage big data. It is the software most used by data analysts to handle big data, and its market size continues to grow. There are three components of Hadoop:

1. Hadoop HDFS - [Hadoop Distributed File System \(HDFS\)](#) is the storage unit.
2. Hadoop MapReduce - [Hadoop MapReduce](#) is the processing unit.
3. Hadoop YARN - [Yet Another Resource Negotiator \(YARN\)](#) is a resource management unit.

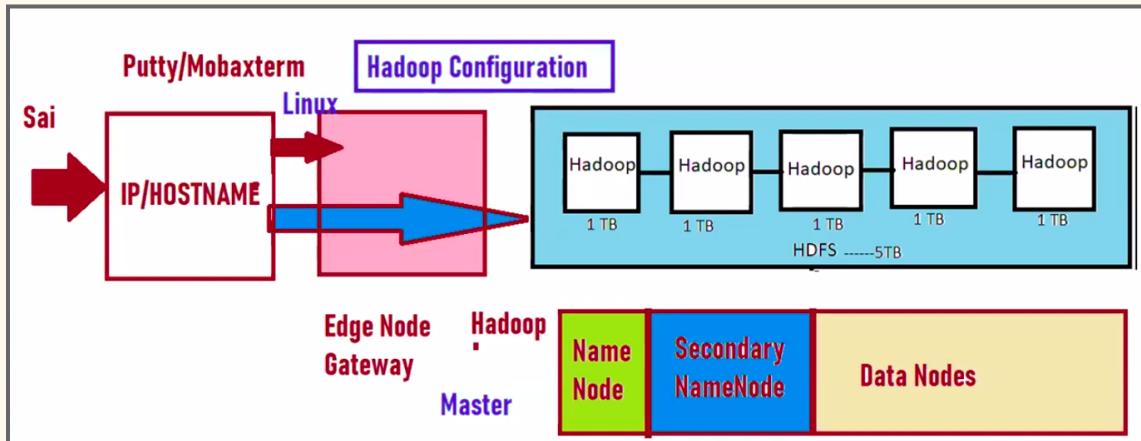
Cluster - Group of Nodes/laptops (including gateway node)

To login into Gateway -

1. Tool to use for login into Gateway - **Putty**, **Mobaxterm** (for Linux, Mac users can also use **ssh@hostname** in terminal)
2. you should have a hostname or IP address and username/password.
3. If want to talk to **Gateway** then use - **Linux** command
4. If want to talk to **HDFS/ Hadoop-Nodes** then use - **Hadoop** Language
5. **For example we have 6 node cluster in these 1 node will be gateway node/Edge node to communicate to the other nodes**, like image (1)



In other 5 nodes, 1 node will be Name node, 1 will be secondary node and remain 3 will be data node, like below image -



Mail - Use Case

Manager --- Sai,, Cluster is ready --- 6 Node Cluster .. Can you see whether you have any data in the GateWay Node
Sai --- Sure Boss - I see whether we have any data in the Gateway Laptop. But can you help how to check it.
Manager --- Sai,every laptop has a ip address or HOSTNAME . Use the gateway hostname-- e.cloudxlab.com
Sai ---- Thanks for the hostname But how can I login to that gateway
Manager --- You have a specific Tools to login to the remote gateway node --- **putty or Mobaxterm**
Sai --- Dear Manager. I have installed Putty I have given hostname also but it is asking login as →username/password
Manager -- Sai use below username/passwrod → twezeyo29867 / BG227QNZ
Sai --- Manager. I got Logged in
Manager --- Its very basic that you have log in - Do the work whether you have any data sitting in the gateway node
Sai --- Can you help me how to check the data
Manager --- Sai, To communicate with Gateway node -- You should know the communication Language
Sai --- Can I communicate in english

Manager --- GATEWAY nodes will not understand English... GATEWAY LAPTOPS ARE INSTALLED WITH LINUX OPERATING SYSTEM --- TO COMMUNICATE WITH gateway you should have LINUX LANGUAGE
Using linux can you list and check whether you have data

Sai --- I listed the command using ls I got a file name 200MB.zip
Manager -- can you tell me what is size of that file
Sai --- Dear manager - file name is 200mb.zip file size also 200MB
Manager--- Thanks. Can you also check whether hdfs is also having the data ?
Sai ----- Then Shall I login to any of the HDFS laptop and check?
Manager --- Idiot , You do not have access to the HDFS laptop even I will not have if You have to communicate with the HDFS .. The only path is you have go through EDGE NODE
Sai ---- ok.. Shall I communicate in the linux language

Manager ---- If you talk in Linux.. Edge Node will consider that you are talk to Edge Node
. If you want to talk to HDFS -- You have to talk in Hadoop Language
If you start talking in hadoop language with Edge Node
Then edge Node will understand that commands are not for edge node It will automatically sends those commands inside the cluster

hadoop fs -ls

Sai -- yes boss.. I spoke in hadoop language -- I got the output --- hadoopdir
Manager -- Thanks

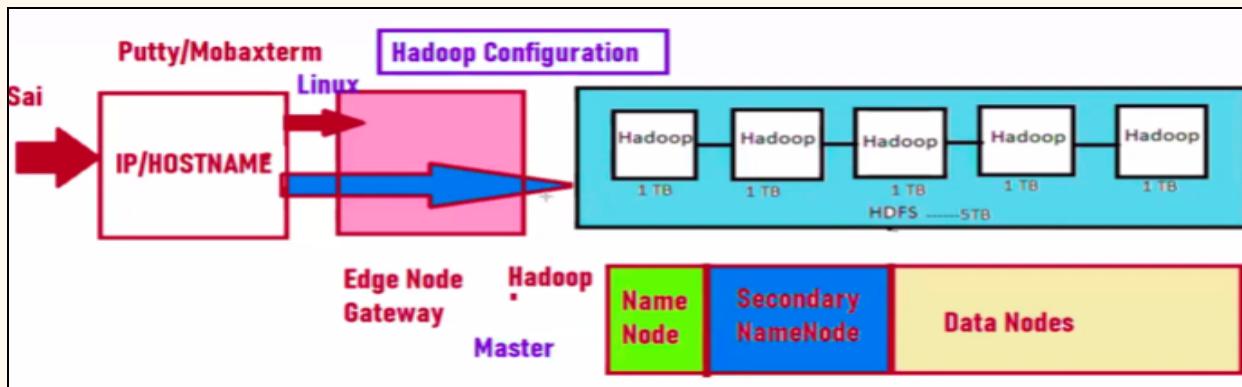
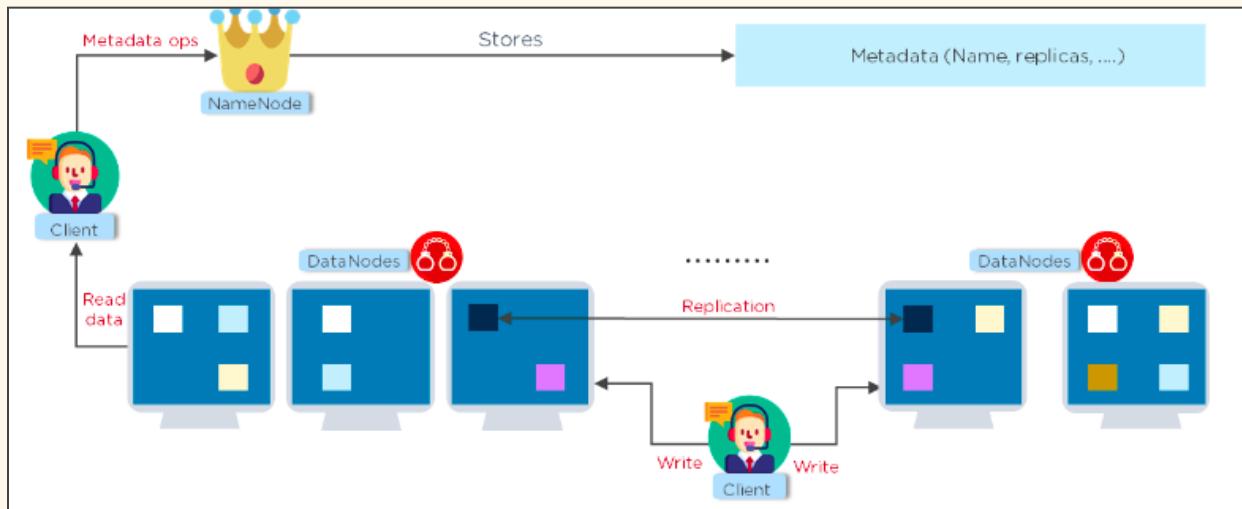
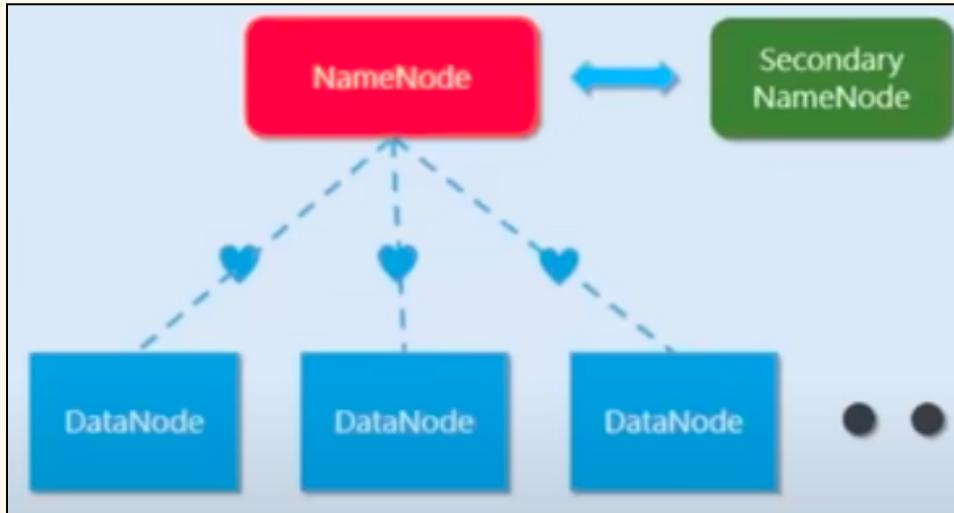
Components of Hadoop

1. **Name node** - Central file system, contains metadata of Data nodes like(which datanode is alive and which datanode has data, etc.)

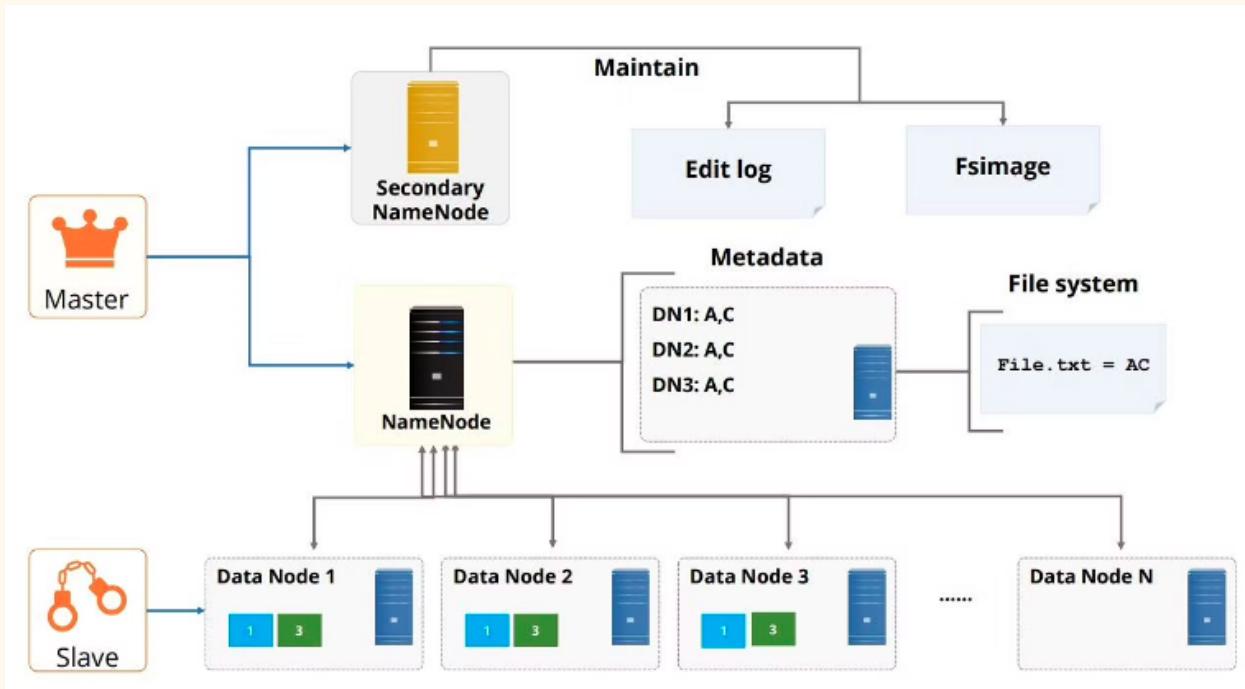
Rules of Namenode -

- a. Send a heartbeat in every 3 sec to namenode.

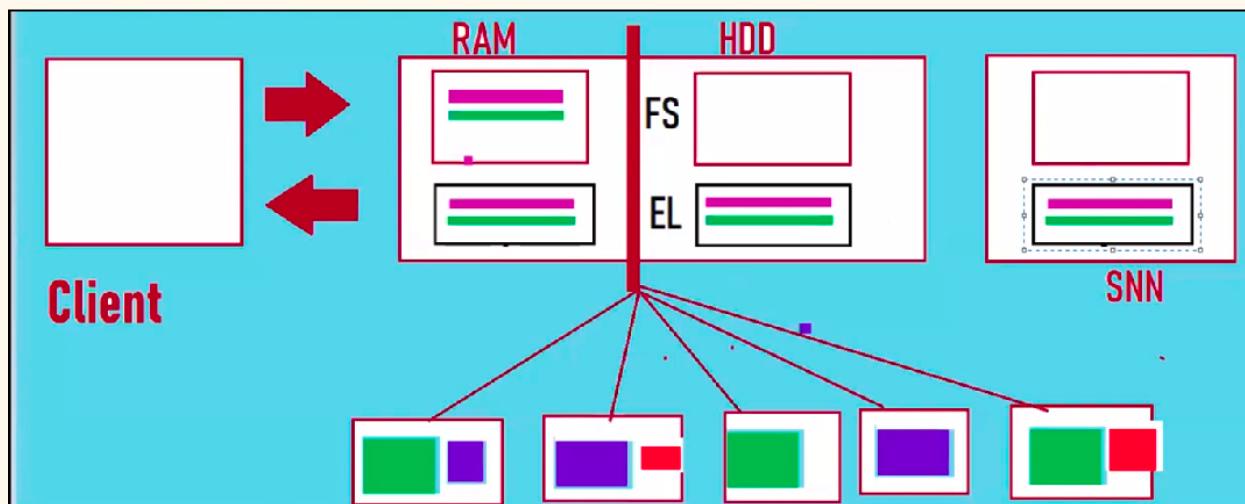
- b. If you don't send a heartbeat for 10 min I will give grace time and still you still don't send, I will declare you a dead HDFS node.
 - c. If you have any work send a copy of the work to 2 other nodes. (There will be 3 copies of Each data)
- 2. **Secondary Name node** - Data backup of Name Node
- 3. **Job Tracker** - Centralized job scheduler
- 4. **Data Nodes** - cluster/machine where data get stored and processed, it is also called **slave** nodes. Every slave node keeps sending a heartbeat signal to the **Name node** every 3 seconds to state that it's alive.
- 5. **Task Trackers** - a software service that monitors the state of the job tracker



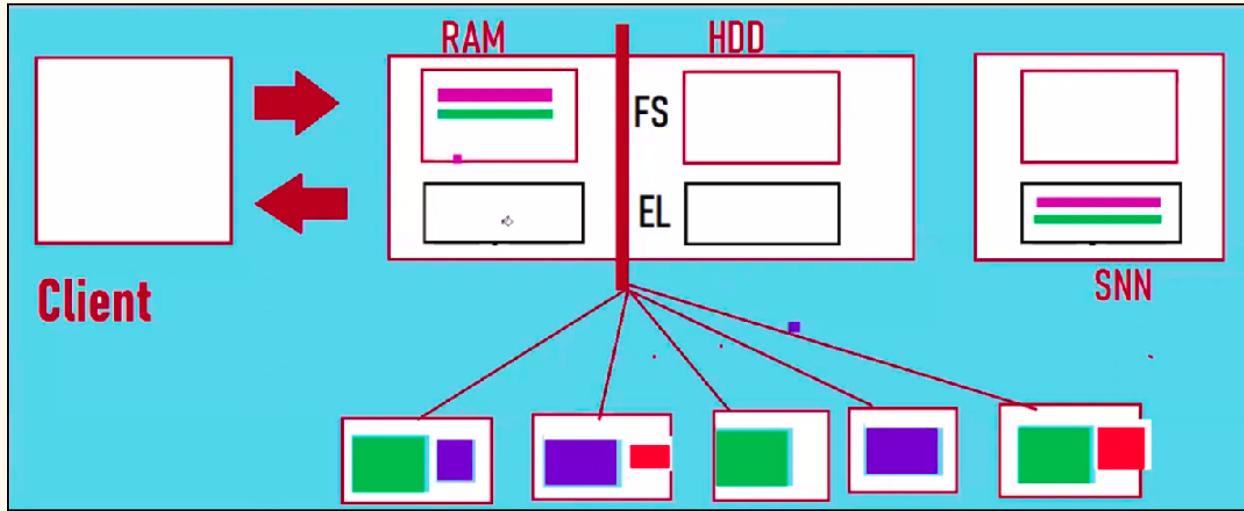
Hadoop 1.0 Architecture -



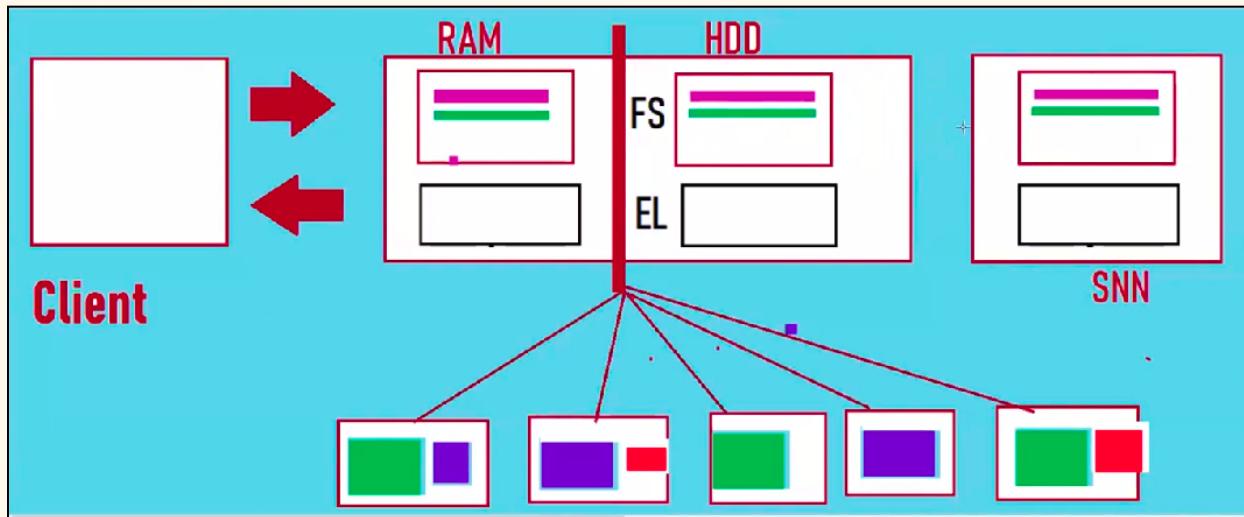
1. Once data translation happens, Save metadata in three places of Namenode -
 - a. FS image - RAM
 - b. Edit Logs(EL) - RAM
 - c. FL - HDD (hard disk)
2. Then After 1 hour,
 - a. if this is first transaction then Namenode sends the FS Image-HDD -> SNN(secondary Name node)-FS image
 - b. Latest EL (RAM) —> SNN-EL



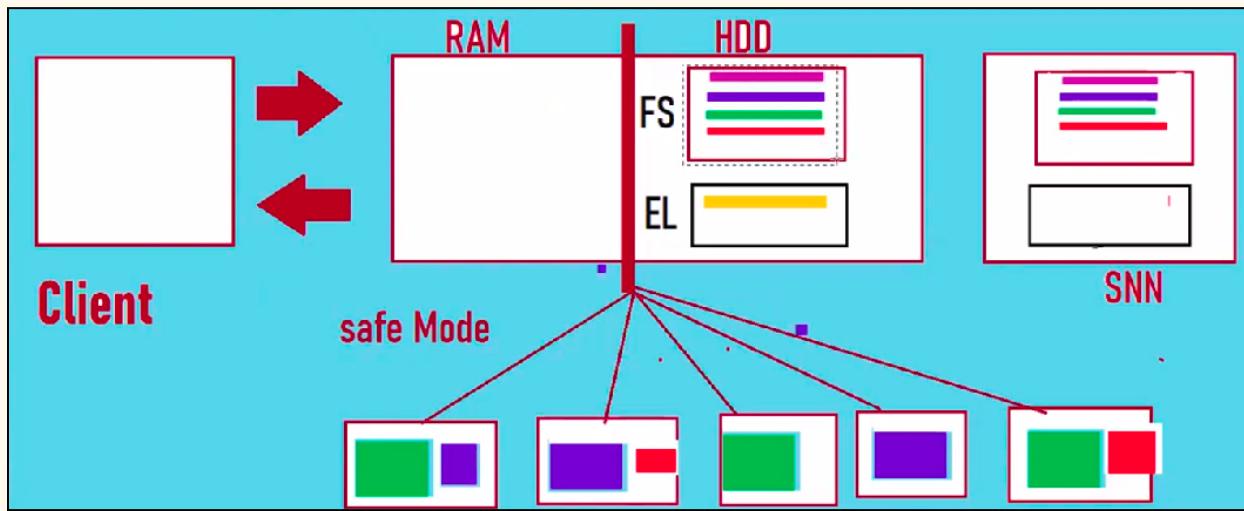
3. And Clear EL - RAM and EL- HDD of Namenode



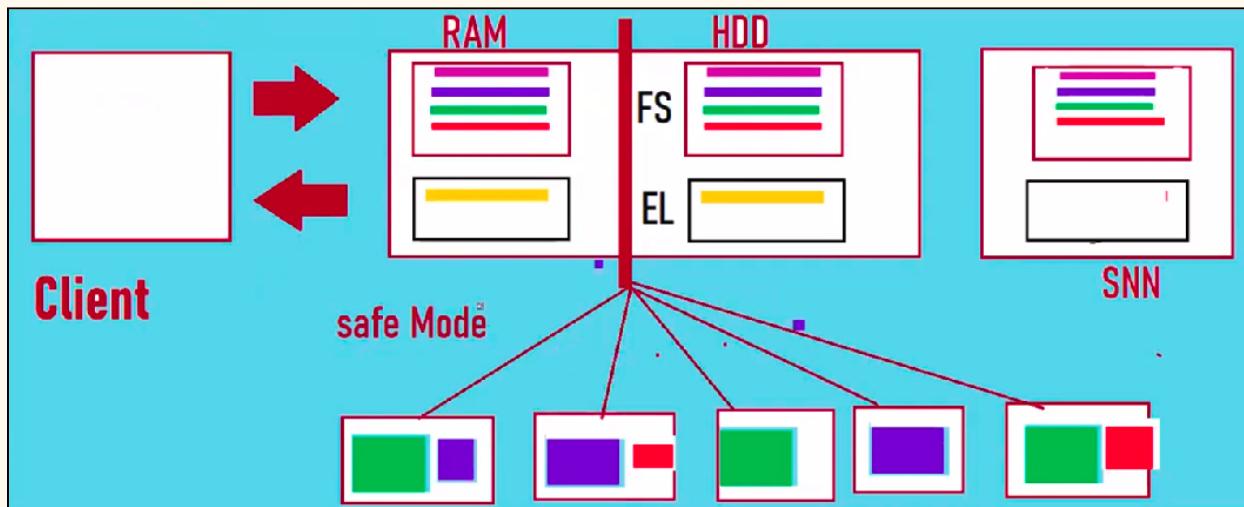
4. Then Checkpointing(Merge) SNN- EL data to SNN-FS image and then send SNN-FS image to HDD-FS image.



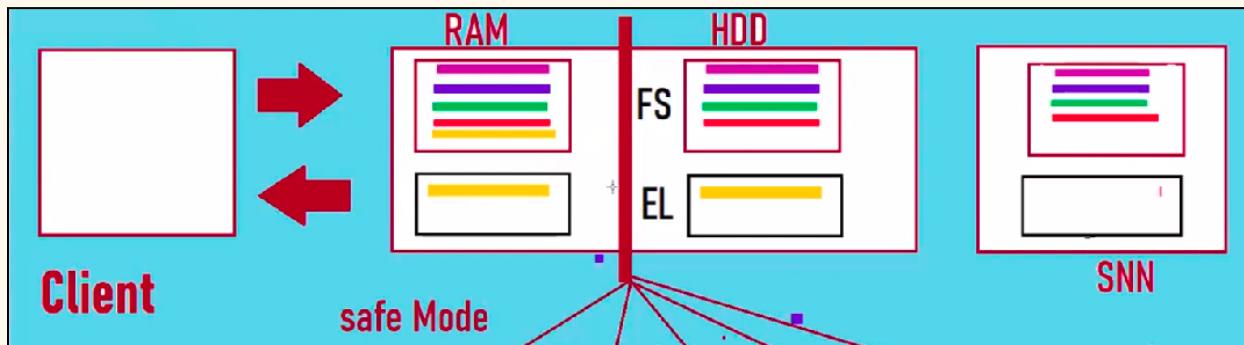
5. What if Namenode- RAM crash, then Namenode immediately goes to **SAFE MODE**,



6. Then Namenode will take FS images and EL from HDD(Namenode).



7. And will Checkpointing(Merge) Data of EL- RAM (namenode) to FS image- RAM(namenode)



FS image -

- Why FS image can't be committed to the hard disk every time there is new data - because FS image is huge in size and it will slow down the system, also if that time Client asks for metadata from Name node it will not respond properly.
- FS image in Hard Disk should be copied during restart time.

Important Notes for Hadoop 1.0 -

- Each transaction will be divided into Blocks, and its size is 128 MB(in Hadoop 1.0 Block size - is 64 MB). Suppose the transaction is 200 MB then it will divide into 2 blocks → 128 MB and 72 MB
- These configurations can be edited and set as per requirement for Edge Node -

Configurable↓

↓

Block size -- 128mb↓

Heartbeat --- 3↓

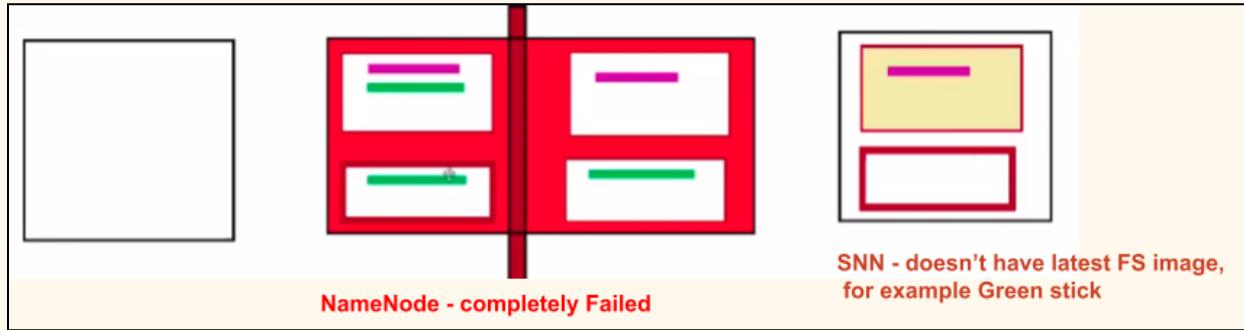
stale time --- 10 Min↓

replicationFactor - 3↓

- If the transaction is not replicating data to other data nodes then the name node will throw notification failure.

Defects of Hadoop 1.0

- Data Loss - Namenode failure was a threat. Suppose after the New data transaction namenode completely fails, then there is no data recovery in that case as we can't depend on SNN because it doesn't have the latest FS image.



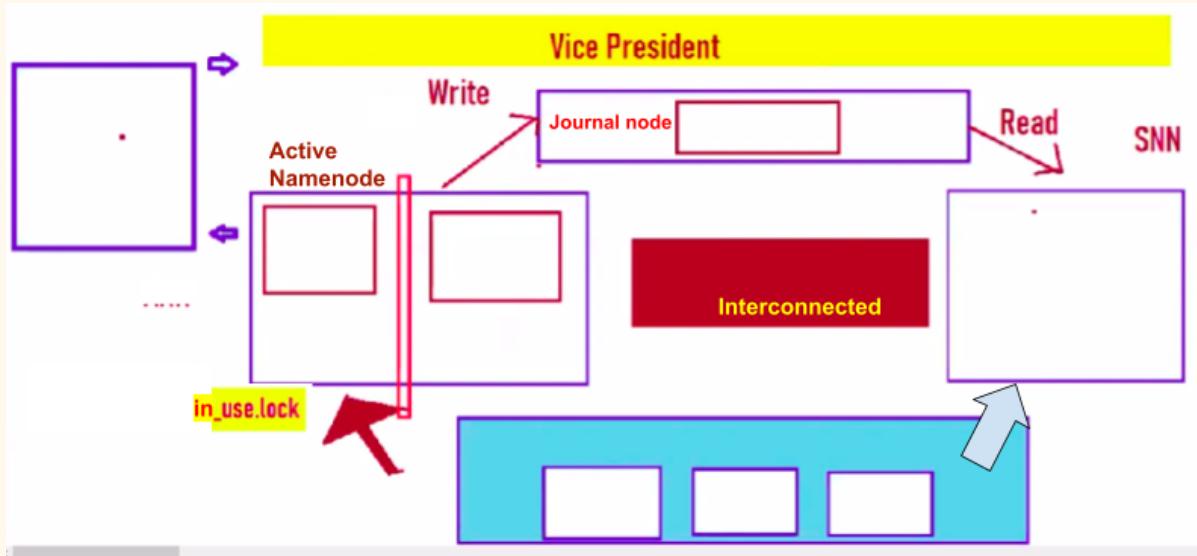
Hadoop 2.0 Architecture

2 New nodes were added -

1. Journal Node
2. Zookeeper

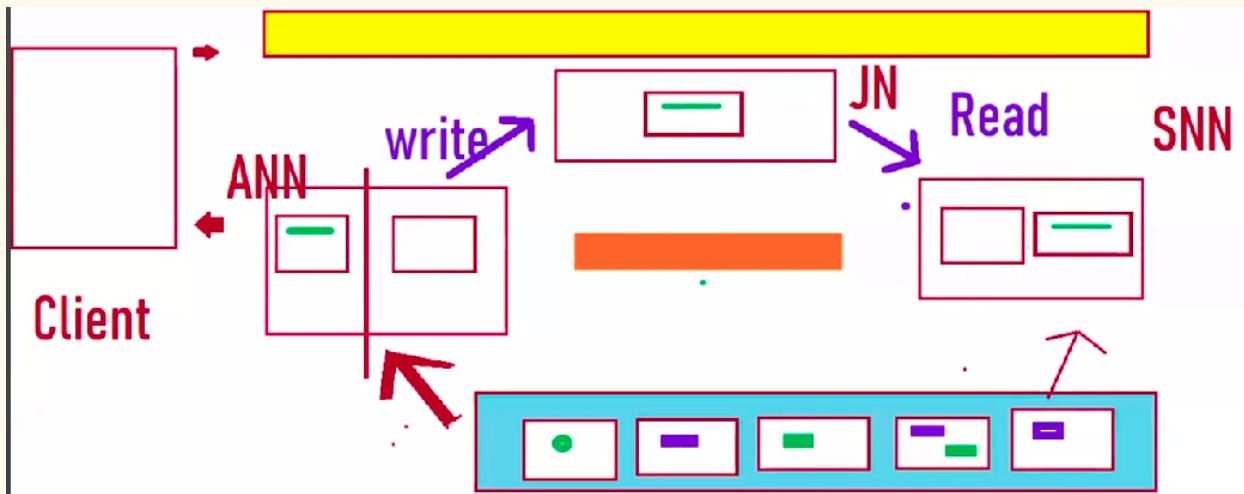
Rules of Zookeeper-

1. Secondary namenode renamed as **Standby node**.
2. Name node(Active name node) will write Edit logs to Journal Node.
3. Standby name node, During **checkpoint** have to read form **Journal node**.
4. If required(in case of name node failure), the **Standby namenode** has to become **Active name node**.
5. Datanodes have to send a heartbeat to both Namenode and Standby Namenode. But metadata will only be sent to Namenode.
6. Whichever is the Active node - it will have file - **in_use.lock** as a representation that this is the current **Active name Node**.
7. ANN and Standby nodes are interconnected.
8. After 1 hour -
 - i. Read **EL - Journal node** and write in **EL-Standby**.
 - ii. If it is the first transaction then copy the **FS-image** of **HDD-ANN** and save it in the **Standby node**.
 - iii. And also clear **EL-Journal node**
 - iv. checkpoint **EL-Standby node** → **HDD-ANN**(active node).

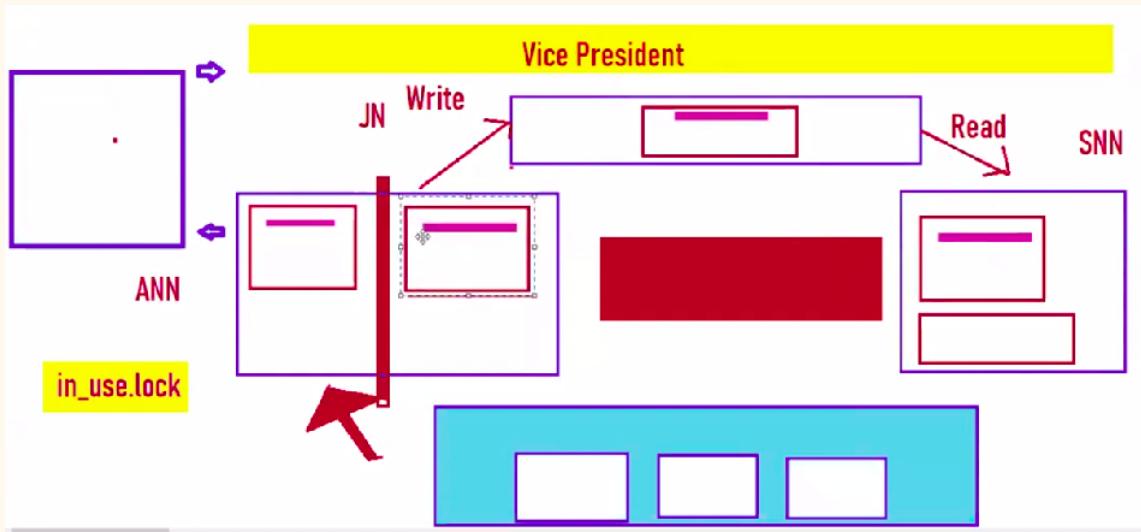


Flow -

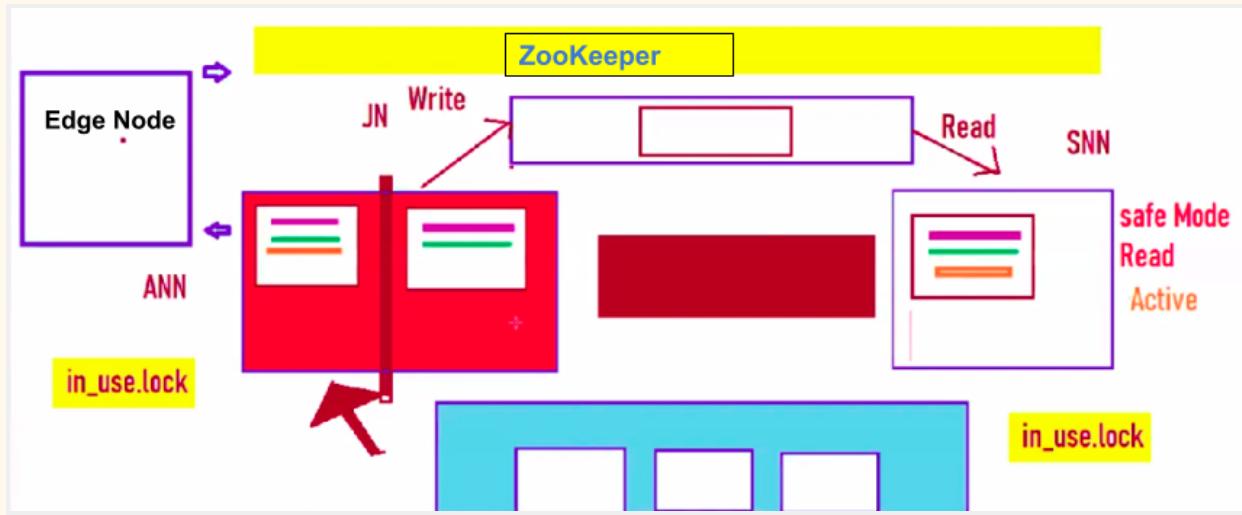
- First Transaction - Write **FS image** in ANN(RAM)and **EL** in **journal node** and after 1 hour **SNN** will **read EL** from **journal node** and copy **old HDD** form **ANN**(if this is first transaction).



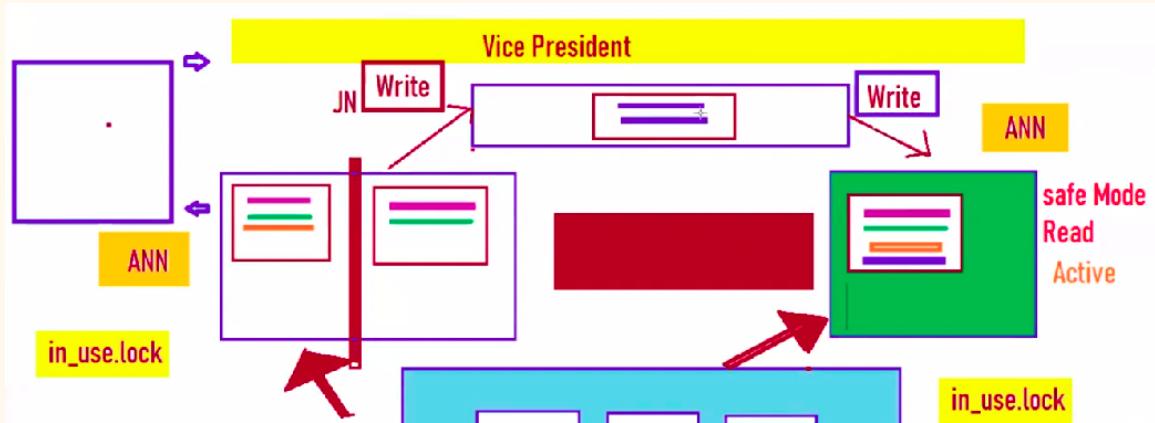
- After one hour there will be checkpointing where SNN merge EL to FS-SNN and copy FS-SNN to FS-ANN(RAM)



3. If ANN goes to complete failure - Then the Standby node will go to **SAFE MODE** and a new file wil be created for **Standby node - in_use.lock** (representing this is the current active node).



4. What if after some time or hours Dead Node(Name node) becomes active - In that case, Both the Dead(name node) and current Active node(standby) will start to write to the **Journal node**, which means now there will be 2 times EL in journal node, which is an issue. This scenario is called **Split Brain Scenario**.



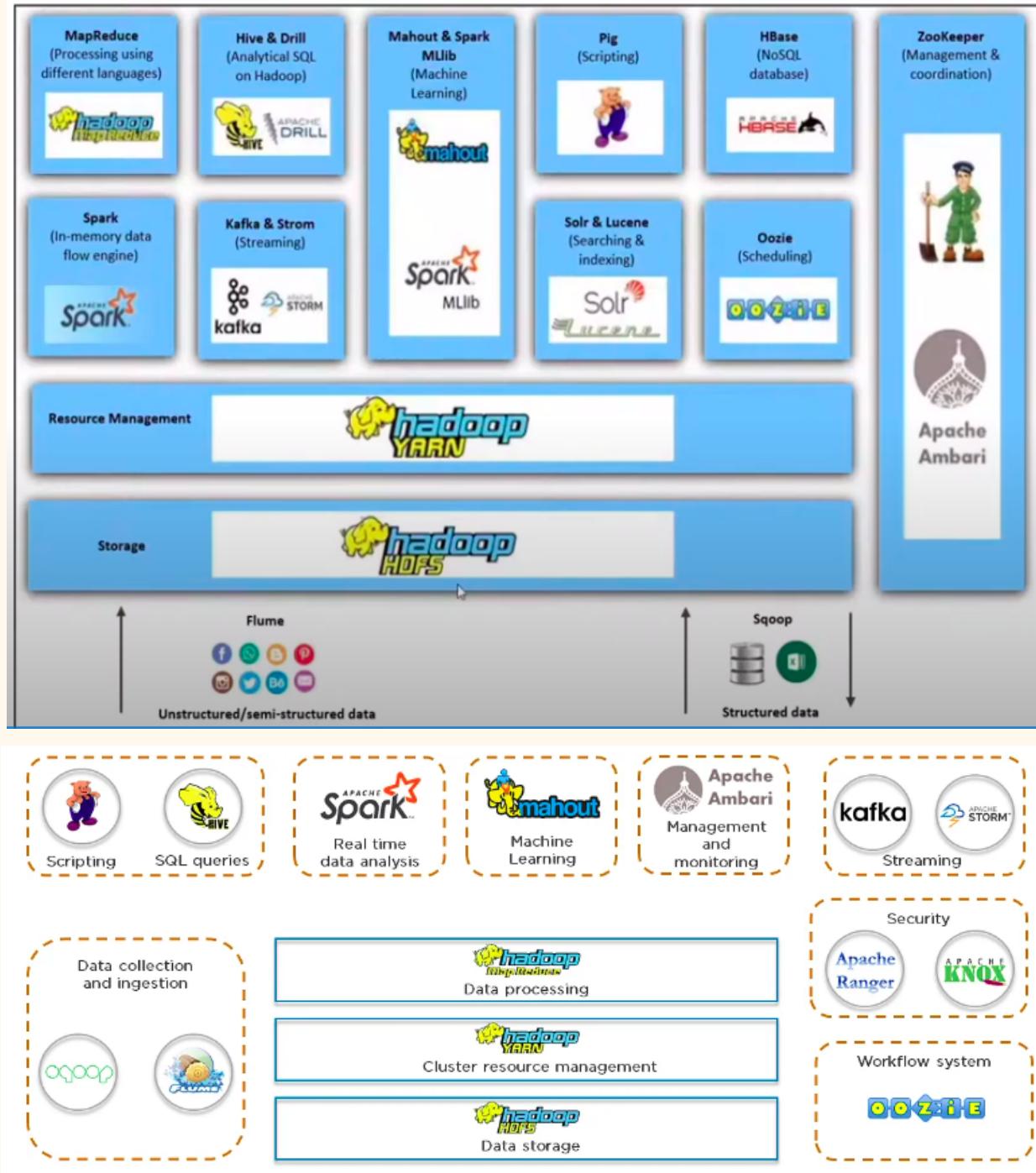
Solution - Doug cutting created a Program **Fencing Java Program** where **Zookeeper** sends messages to Dead node, where dead mode will stop writing to other nodes. And **Dead node** is now **Standby node**.

Hadoop Commands (just like Linux)-

To talk to Hadoop →prefix Linux command with '**hadoop fs -**'

- hadoop fs -ls
- hadoop fs -touchz <fileName> // create a new empty file
- hadoop fs -cat <fileName> //print file content on console
- hadoop fs -put <edgeNodePath> <HDFSPath> // copy file from edge to HDFS
- hadoop fs -get <HDFSPath> <edgeNodePath> // copy file from HDFS to Edge
- hadoop fs -moveFromLocal <edgeNodePath> <HDFSPath>
- hadoop fs -moveToLocal <HDFSPath> <edgeNodePath>
- hadoop fs -appendToFile <sourceFilePath> <targetHDFSPath>
- hadoop fs -rmdir <path> //delete folder
- hadoop fs -rm -r <path> //delete folder and file inside it
- hadoop fs -rm <path> // delete file
- hdfs dfsadmin -safemode get //→find out if you're in safe mode
- hadoop dfsadmin -safemode leave //→turn off safe mode
- hdfs dfsadmin -report //→find out how much disk space is used, free, under--replicated, etc .
-

Hadoop Ecosystem



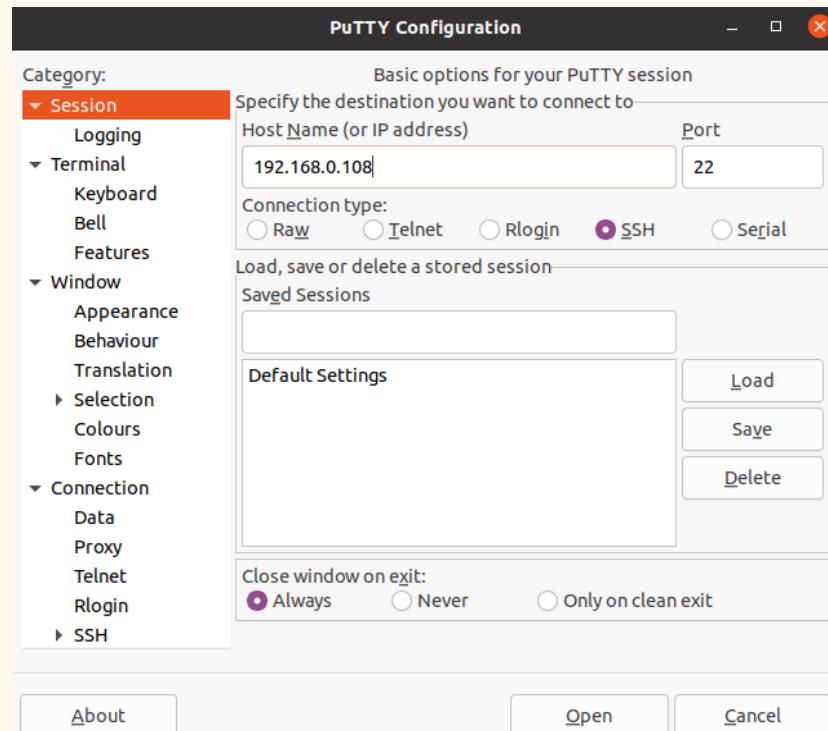
Cloudera

Installation →

- Prerequisites - [Install virtual box](#) [sudo apt install virtualbox]
- Article -
<https://www.simplilearn.com/tutorials/big-data-tutorial/cloudera-quickstart-vm>

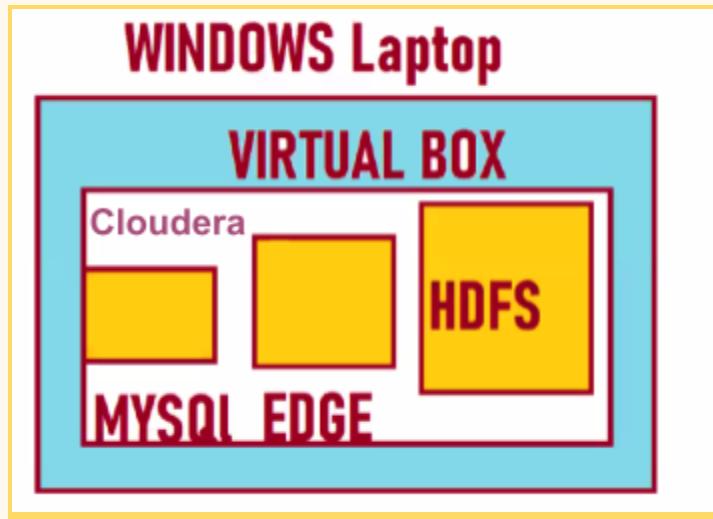
How to Login into Edge node/Gateway node

1. Login direct into Cloudera by opening **VirtualBox** with Cloudera installation
2. By **putty**, **moxaterm**, or **ssh** as per OS -
 - **Windows** - Open putty and provide **host IP** and connect



- Linux/Mac -
 - Either install putty and connect **OR** run below command
 - Go to Terminal → ssh <user_name>@<remote_host>
 - For example - **ssh cloudera@hostIdAddress**

3. Run Command to leave safe mode - **hadoop dfsadmin --safemode leave**

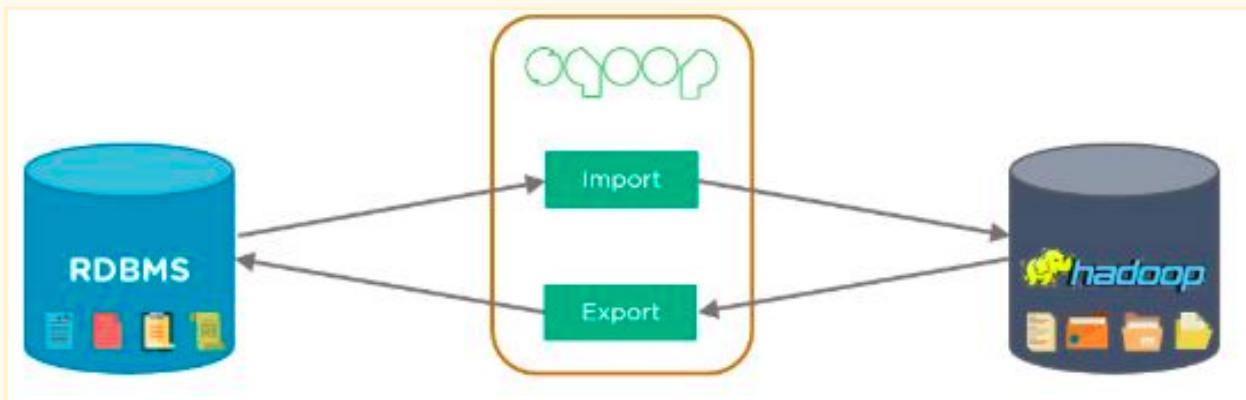


SQOOP -

2009 cloud support - Kathleen Ting

Tool used to transfer bulk data between HDFS and RDBMS

Sqoop word came from SQL+HADOOP=SQOOP



To process data using Hadoop, the **data** first needs to be loaded into Hadoop clusters from several sources. However, it turned out that the process of loading data from several heterogeneous sources was extremely challenging.

Issues Big data developers faced with Hadoop -

1. RDBMS is tough to Handle
2. Damn Slow to import data from RDBS to Hadoop
3. Portion(partial) Data Import is not possible
4. Incremental Data import is not possible
5. Change of data is Hard
6. Serialized data has no support
7. Export no much Support etc

Doug Cutting → That's not my problem

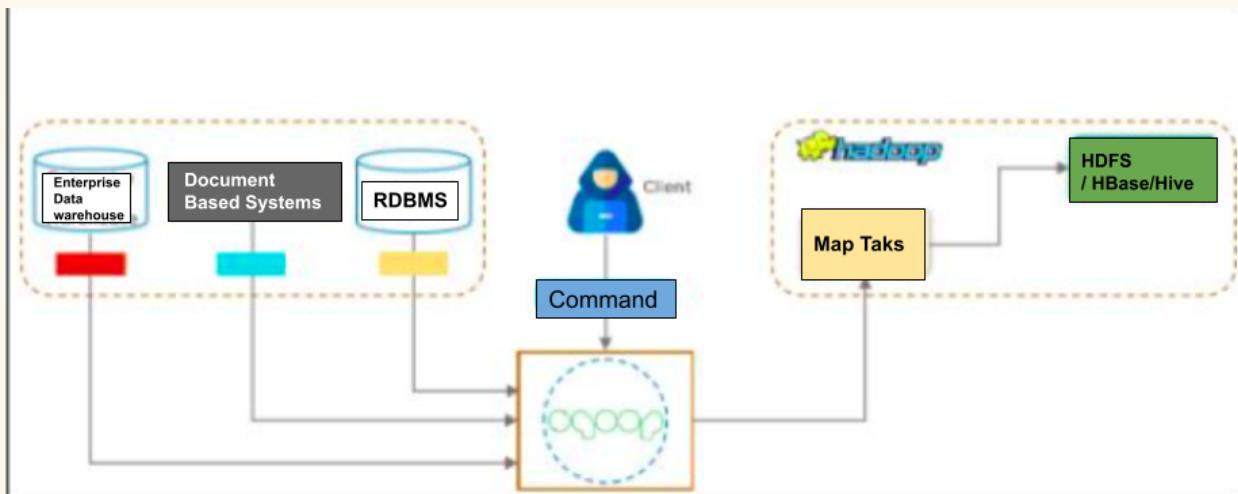
2008 → Data ingestion has become big problem

2009 → **Kathleen** → I have a solution to Integrate SQL and Hadoop

- 1) RDBMS integration is Made Simple
- 2) My tool is damn faster (Multi threading)
- 3) Yes you can import certain portion of data happily
- 4) Obviously incremental data import is possible
- 5) Change of Data =: absolutely
- 6) My tool is made of Java (Serialization Support is good)
- 7) Exportss is possible now.

SQOOP Architecture -

Sqoop fetches data from different databases. Here, we have an enterprise data warehouse, document-based systems, and a relational database. We have a connector for each of these; connectors help to work with a range of accessible databases.



SQOOP Command -

- **sqoop import**
--connect jdbc:mysql://<hostname>:<portName>/<DBname>
--username <userNmae>
--password <password>

```
--table <sourceTableName>
--m 1
--target-dir <targetDirectoryPath>
• sqoop job --create // Create Sqoop job
• sqoop job --list // To check all job list
• sqoop job --exec <jobName> // To run specific job
• sqoop job --show <jobName> // To check saved job details
• sqoop job --delete <jobName> // To delete job
```

Prerequisites -

Login into remote Gateway either from SSH or virtualBox cloudera [How to Login into Edge node/Gateway node](#)

For big data classes -

- **hostname** → localhost
- **portnumber** → 3306
- **database** → zeyodbd
- **username** → root
- **password** --> Aditya908
- **table** → ztab
- **targetlocation** → /user/cloudera/dimport

For Big data practice class - how to run SQOOP

- We are going to use cloudera in VirtualBox - Open Cloudera and login into MySQL to check if SQL is working

mysql --uroot --pcloudera and now run the below command in SQL

```
=====
Mysql comannds
=====

• mysql -uroot -pcloudera
• create database zdb;
• use zdb;
• create table zeyotab(id int,name varchar(100),city varchar(100));
• insert into zeyotab values(1,'Sai','Chennai');
• insert into zeyotab values(2,'Zeyo','Hyderabad');
• insert into zeyotab values(3,'Analytics','Hyderabad');
• insert into zeyotab values(4,'Sai','Chennai');
• insert into zeyotab values(5,'Hema','Hyderabad');
• insert into zeyotab values(6,'Vassu','Chennai');
• select * from zeyotab;
• quit
```

- Open another Tab and Then run below command -

```
sqoop import --connect jdbc:mysql://localhost:3306/zeyodbd --username root --password Aditya908 --table ztab --m 1  
--target-dir /user/cloudera/dimport
```

```
=====
Edge Node
=====

• sqoop import --connect jdbc:mysql://localhost/zdb --username root --password cloudera --table zeyotab --m 1 --delete-target-dir --target-dir  
/user/cloudera/datadir

• hadoop fs -ls /user/cloudera/datadir

Two files will be created like below - only part-m-0000 contains data

[[cloudera@quickstart ~]$ hadoop fs -ls /user/cloudera/inimport
Found 3 items
-rw-r--r-- 1 cloudera cloudera          0 2022-06-26 05:22 /user/cloudera/inimport/_SUCCESS
-rw-r--r-- 1 cloudera cloudera        98 2022-06-26 05:22 /user/cloudera/inimport/part-m-00000
• hadoop fs -cat /user/cloudera/datadir/part-m-00000
```

Partial import data - We have do 2 way

1. --where <condition>

```
sqoop import --connect jdbc:mysql://localhost/zdb --username root --password  
cloudera  
--table zeyotab --where "city='chennai' and id>1"  
--m 1 --delete-target-dir --target-dir /user/cloudera/datadir
```

2. --query - for example let join two tables.

- Important thing for command - this is **where \\$CONDITIONS** compulsory to put wherever using **--query** command

```
sqoop import --connect jdbc:mysql://localhost/zdb1 --username root --password cloudera  
--query "select a.* , b.product from zeyotab a join zeyoprod b on a.id=b.id where \$CONDITIONS"  
--m 1 --delete-target-dir --target-dir /user/cloudera/joindata
```

- even if there is any other condition like id>3, put it as -
where id>3 and \\$CONDITIONS

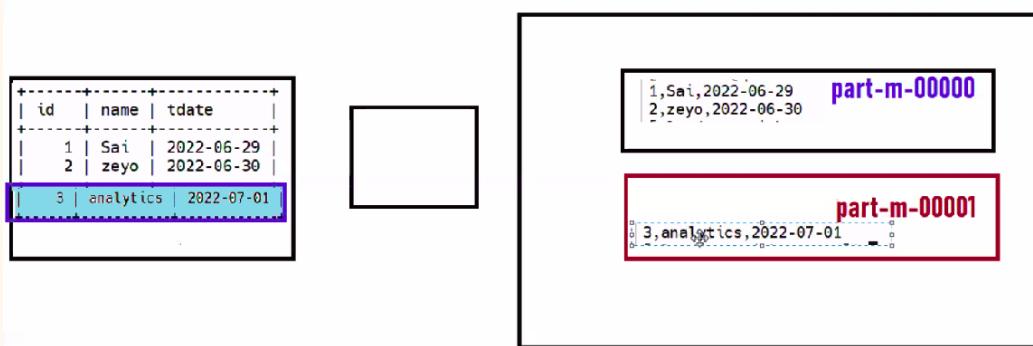
Increatmental import data -

```
sqoop import --connect jdbc:mysql://localhost/zdb2 --username root --password cloudera  
--table zeyotab --m 1 --target-dir /user/cloudera/inimport  
--incremental append --check-column id --last-value 6
```

Argument	Description
--check-column (colName)	Specifies the column to be examined when determining which rows to import.
--incremental (mode)	Specifies how Sqoop determines which rows are new. Legal values for mode include append and lastmodified .
--last-value (value)	Specifies the maximum value of the check column from the previous import.

★ Rules

- Pick a column that should have only integer-like - numbers, date/Time
- That column shouldn't have duplicates otherwise increment process won't work for duplicates
- A new part file will be created for every increment data, Process -



```

=====
Mysql commands
=====
• mysql -uroot -pcloudera
• create database zdb2;
• use zdb2;
• create table zeyotab(id int ,name varchar(100),city varchar(100));
• insert into zeyotab values(1,'Sai','Chennai');
• insert into zeyotab values(2,'Zeyo','Chennai');
• insert into zeyotab values(3,'Analytics','Hyderabad');
• insert into zeyotab values(4,'Sai','Chennai');
• insert into zeyotab values(5,'Hema','Hyderabad');
• insert into zeyotab values(6,'Vassu','Chennai');
• select * from zeyotab;

=====
Edge Node
=====
• sqoop import --connect jdbc:mysql://localhost/zdb2 --username root --password cloudera --table zeyotab --m 1
--delete-target-dir --target-dir /user/cloudera/inimport
• hadoop fs -ls /user/cloudera/inimport
• hadoop fs -cat /user/cloudera/inimport/part-m-00000

=====
Again login mysql
=====
• mysql -uroot -pcloudera
• use zdb2;
• insert into zeyotab values(7,'Hema','Hyderabad');
• insert into zeyotab values(8,'Vassu','Chennai');
• select * from zeyotab;
• quit

=====
Edge Node( for increament new part-m-00001 file will be created and will only contain increased data)
=====
• sqoop import --connect jdbc:mysql://localhost/zdb2 --username root --password cloudera --table zeyotab --m 1
--target-dir /user/cloudera/inimport --incremental append --check-column id --last-value 6
• hadoop fs -ls /user/cloudera/inimport
• hadoop fs -cat /user/cloudera/inimport/part-m-00001
• hadoop fs -cat /user/cloudera/inimport/*

```

Updated/Modified import data

Sai Video - <https://youtu.be/CGiJf7tAyIk>

```

sqoop import --connect jdbc:mysql://localhost/zdb2 --username root
--password cloudera --table zeyotab --m 1 --target-dir /user/cloudera/inimport
--incremental lastmodified --check-column dataTime --last-value 6 --merge-key
id

```

Argument	Description
--check-column (colName)	Specifies the column to be examined when determining which rows to import.

--incremental (mode)	Specifies how Sqoop determines which rows are new. Legal values for mode include append and lastmodified .
--last-value (value)	Specifies the maximum value of the check column from the previous import.
--merge-key (colName)	This column shouldn't have duplicates

If providing these above 4 details, there will run 2 Map-reduce jobs -

1. To bring Data from RDBMS
2. To override the data to HDFS

One Case - For example, there are multiple part-m-0000 files in target dir, and then running the update/modified command, there will be only 1 part-m-00000 file and it will have all the data

Using the options file to Pass Arguments

For example, the following Sqoop invocation for import can be specified alternatively as shown below:

```
sqoop import --connect jdbc:mysql://localhost/dbName --username foo --password cloudera --table TEST
```

TO

```
sqoop --options-file /users/homer/work/import.txt --table TEST --target-dir /user/cloudera/dataArgFile
```

- ★ where the options file `/users/cloudera/work/import.txt` contains the following: Make sure to provide a new line for every argument and value like below, also we can add more arguments in the file if want -

```
import
--connect
jdbc:mysql://localhost/db
--username
root
--password
cloudera
```

SQOOP-All-Tables-Import-

sqoop import-all-tables (generic-args) (import-args)

SQOOP Job -

Question - Do we need to save the last value every time, it is hard to keep remembering the last value every time as data goes in bulk in a day many times?

Sai → kathleen Do i need keep remembering the last value everytime. as its very Tedious

Kathleen → Not required. There is a concept of Automation Jobs, Take your sqoop incremental import and come to me

Sai →>

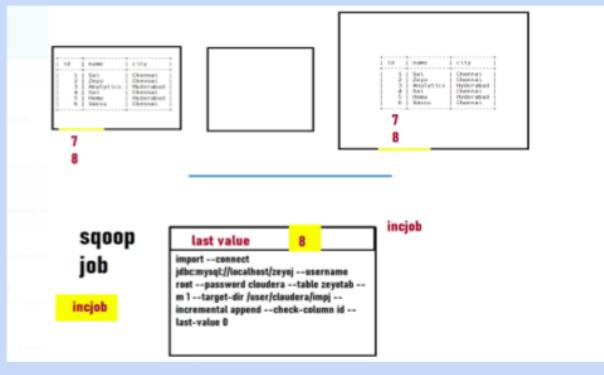
```
sqoop import --connect jdbc:mysql://localhost/zeyoj --username root --password cloudera --table zeyotab --m 1 --target-dir /user/cloudera/impj --incremental append --check-column id --last-value 0
```

Kathleen → Automation job will have below role -

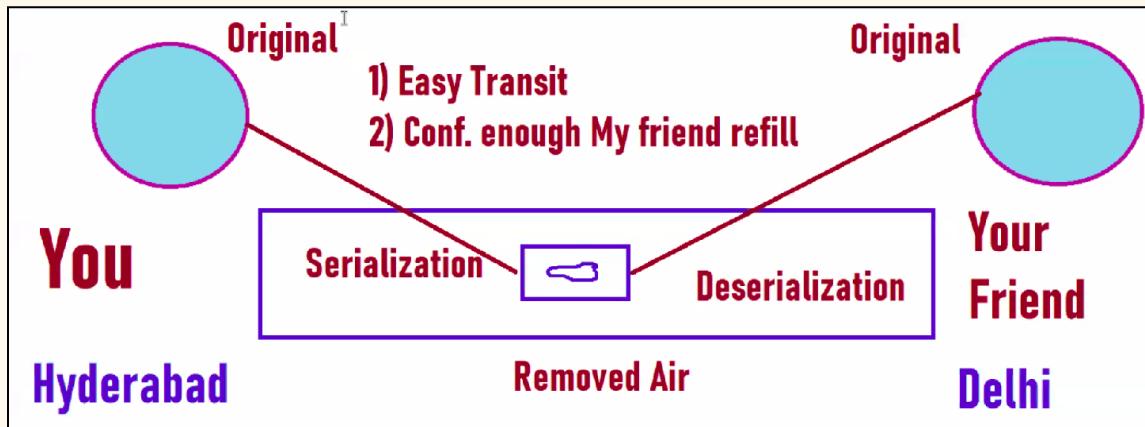
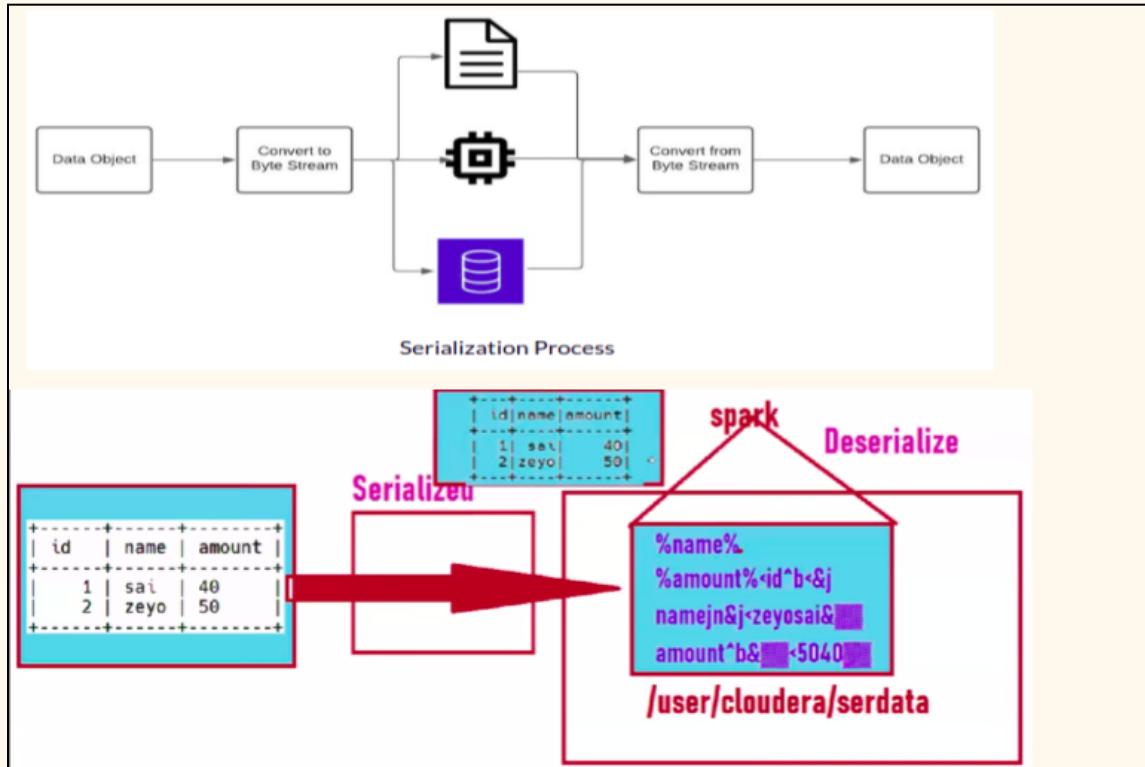
1. Send the Incremental Import Command to a job and assign name to it , like → **incrementJob**

Steps

1. Create a job using below command -
sqoop job --create incrementJob
2. Assign incremental command -
sqoop job --create incrementJob -- import --connect jdbc:mysql://localhost/zeyoj --username root --password cloudera --table zeyotab --m 1 --target-dir /user/cloudera/impj --incremental append --check-column id --last-value 0
3. Check whether it got created -
sqoop job --list
4. Execute the Job and check target -
sqoop job --exec incrementJob
5. Insert some more records in sql.
6. Execute the Job and check target -
sqoop job --exec incrementJob



Serialization File Types



Below are serialization file formats.

1. Text files

--as-textfile

Features -

- Huge in size
- Readable format
- If you query on TextFile is huge
- Very bad option for Big Data

2. Sequence data files - 2008

--as-sequencefile

```
sqoop import --connect jdbc:mysql://localhost/serdb --username root --password cloudera --table avtab --m 1 --delete-target-dir --target-dir /user/cloudera/sequencedir  
--as-sequencefile  
hadoop fs -cat /user/cloudera/sequencedir/*
```

Features -

- Java File format - it is made of java and since Mapreduce is also Made Java it is very friendly
- Not much used Today as Mapreduce is not used Nowadays
- Sequence is also Huge in Size
- It's not good for Queries

3. Avro Data file - (2009) Add below command to serialization -

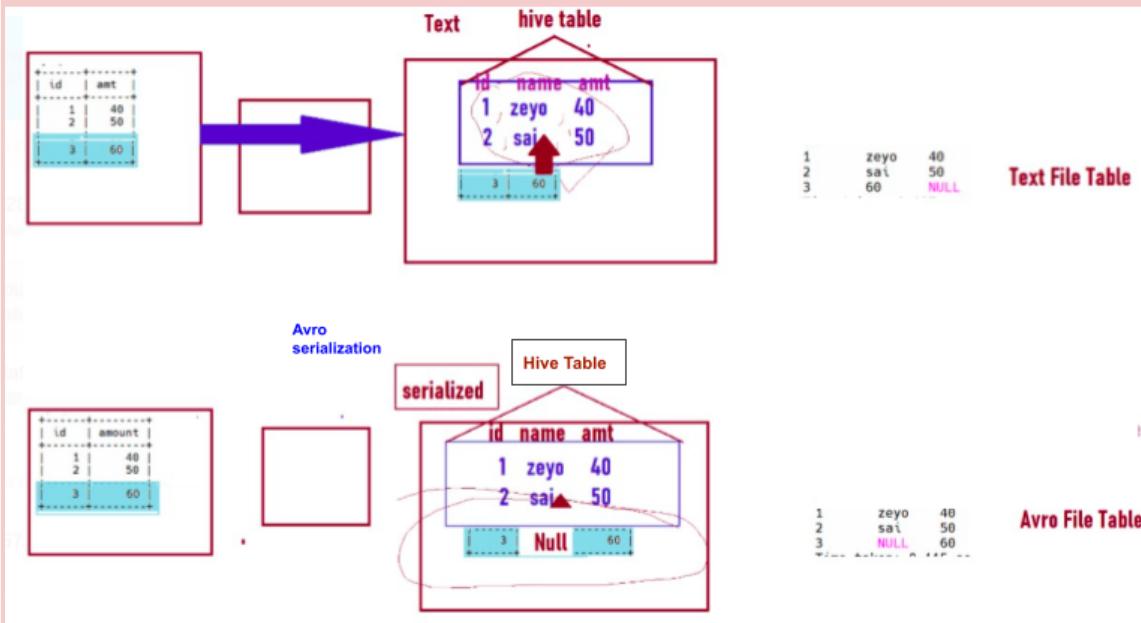
--as-avrodatafile

```
sqoop import --connect jdbc:mysql://localhost/serdb --username root --password cloudera --table avtab --m 1 --delete-target-dir --target-dir  
/user/cloudera/avrodir --as-avrodatafile  
hadoop fs -cat /user/cloudera/avrodir/*
```

Features -

- It can handle the Schema evolution, like delete a column, add new column--- It can map data to right columns
- 40-50% Compression
- Row format storage

- Sai, I have table created in mysql, and Done very normal Import (text file) in HDFS.
- I checked the data its fine with textfile, To process the data I have created hive table also on top
- Queried the data, Data is also Good
- But When the schema evolution (Schema changes -- Column changes like delete name Column just like in below example)**
- Hive or mapreduce could not tackle that



```
=====
mysql
=====
• create database avrodb;
• use avrodb;
• create table avro_tab(id varchar(100),name varchar(100),amount varchar(100));
• insert into avro_tab values(1,'zeyo',40);
• insert into avro_tab values(2,'sai',50);
=====
Edge Node
=====
sqoop import --connect jdbc:mysql://localhost/avrodb --username root --password cloudera --table avro_tab --m 1 --target-dir /user/cloudera/avrodir --as-avrodatafile
=====
hive
=====
• hive
• CREATE EXTERNAL TABLE ahivetab ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.avro.AvroSerDe' STORED as AVRO LOCATION '/user/cloudera/avrodata'
TBLPROPERTIES ('avro.schema.url'='/user/cloudera/avro_tab.avsc');
• select * from phivetab;
1      zeyo    40
2      sai     50
=====
mysql
=====
• alter table avro_tab drop column name;
• insert into avro_tab values(3,50);
=====
Edge Node
=====
Now run increment command
• sqoop import --connect jdbc:mysql://localhost/avrodb --username root --password cloudera --table avro_tab --m 1 --target-dir /user/cloudera/avrodir --as-avrodatafile
--increment append --check-column id --last value 2
=====
hive
=====
• select * from phivetab;
```

3	NULL	40
4	NULL	60
1	zeyo	40
2	sai	40

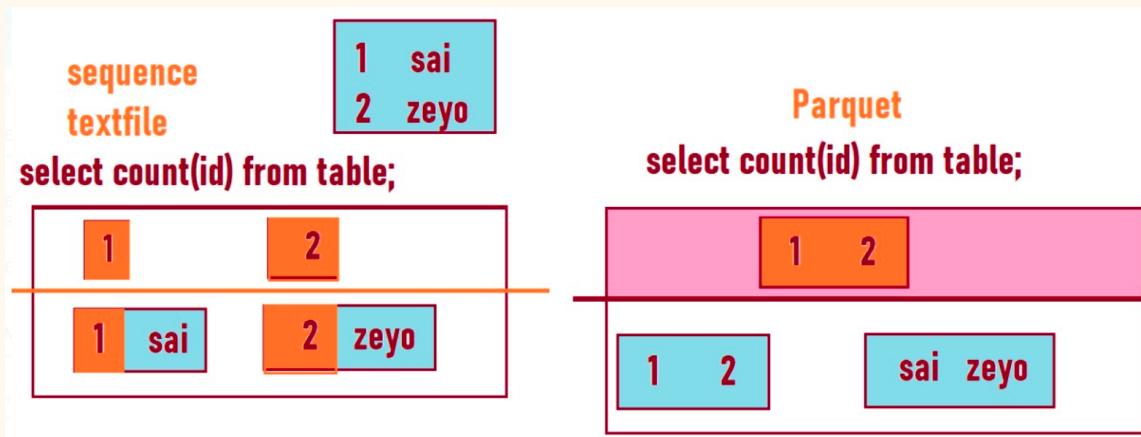
4. PARQUET - Came in 2012 by Doug Cutting, it is a very powerful file format yet

--as-parquetfile

```
sqoop import --connect jdbc:mysql://localhost/serdb --username root --password cloudera --table avtab --m 1 --delete-target-dir --target-dir /user/cloudera/parquetdir --as-parquetfile  
hadoop fs -cat /user/cloudera/parquetdir/*
```

Features -

- i. Reduce compression from 60-80%
- ii. It saves the file data in **columnar file format**
- iii. **Faster query file format**
- iv. Parquet uses Predicate pushdown mechanism



```

=====
mysql
=====
• create database pardb;
• use avrodb;
• create table par_tab(id varchar(100),name varchar(100),amount varchar(100));
• insert into par_tab values(1,'zeyo',40);
• insert into par_tab values(2,'sai',50);

=====
Edge Node
=====
sqoop import --connect jdbc:mysql://localhost/pardb --username root --password cloudera --table par_tab --m 1 --target-dir /user/cloudera/pdir --as-parquetfile
=====
hive
=====
• hive
• CREATE TABLE phivetab (id string, name STRING,amount string) STORED AS PARQUET location '/user/cloudera/pdir';
• select * from phivetab;
+---+---+---+
| 1 | zeyo | 40 |
| 2 | sai | 40 |
+---+---+---+
=====
mysql
=====
• alter table par_tab drop column name;
• insert into par_tab values(3, 50);
=====
Edge Node
=====
Now run increment command
• sqoop import --connect jdbc:mysql://localhost/pardb --username root --password cloudera --table par_tab --m 1 --target-dir /user/cloudera/pdir
--as-parquetfile --increment append --check-column id --last value 2

=====
hive
=====
• select * from phivetab;
+---+---+---+
| 3 | NULL | 40 |
| 4 | NULL | 60 |
| 1 | zeyo | 40 |
| 2 | sai | 40 |
+---+---+---+

```

5. ORC (Optimized Row Columnar)-

We can not do SQOOP import for ORC

Features -

- Compression rate is 90-95%
- Query takes time because de-compression takes time
- Historical Data
- Columnar File
- Faster but not much as Parquet file

Tasks/Interview questions

#1 - Find a way to import only Chennai records, But I need only id and name column
sqoop import --connect jdbc:mysql://localhost/zdb --username root --password cloudera --table zeyotab --where "city='chennai' and id>1"
--columns id,name
--m 1 --delete-target-dir --target-dir /user/cloudera/datadir

#2 - Now I have created Automation jobs but i have do enter Password every time, it is too much manual work?
To reduce this work, create [Password file](#) and provide path of file in sqoop command

Create file either in Edge node or HDFS node - (*but one condition there should be no new line character in file*)
echo -n cloudera>spfile

SQOOP command -

```
sqoop job --create increamentJob -- import --connect jdbc:mysql://localhost/zeyotab6 --username root --password-file  
file:///home/cloudera/spfile  
--table zeyod --m 1 --target-dir /user/cloudera/djimport --incremental append --check-column id --last-value 0
```

For HDFS node provide file like this - [hdfs:/user/cloudera/spfile](#)

#3 - Where does all jobs data is saved like last value , upper value and everthing -
There is file name [metastore.db.script](#) under [.sqoop](#) directory

```
1. ls .sqoop [cloudera@quickstart ~]$ ls .sqoop  
metastore.db.properties metastore.db.script  
[cloudera@quickstart ~]$
```

2. [cat metastore.db.script](#)

#4 - How can we import data from particular row or column? What is the destination types allowed in Sqoop import command?

- sqoop import --connect jdbc:mysql://db.one.com/corp --table INTELLIPAAT_EMP --where "start_date>'2016-07-20'"
- Sqoop import --connect jdbc:mysql://db.test.com/corp --query "SELECT * FROM intellipaat_emp LIMIT 20"
- sqoop import --connect jdbc:mysql://localhost/database --username root --password aaaaa --columns "name,emp_id,jobtitle"

#5 - I am getting connection failure exception during connecting to Mysql through Sqoop, what is the root cause and fix for this error scenario?

This will happen when there is lack of permissions to access our Mysql database over the network

Go to MySql and run below command -

```
mysql> GRANT ALL PRIVILEGES ON *.* TO ''@'localhost';
```

#6 I am having around 500 tables in a database. I want to import all the tables from the database except the tables named Table 498, Table 323, and Table 199. How can we do this without having to import the tables one by one?

```
sqoop import-all-tables --connect --username --password --exclude-tables Table498, Table 323, Table 199
```

Multi-Mappers -

Multi-mapper is nothing but multi-threading. (Arg in sqoop ⇒ [--m 1](#)).

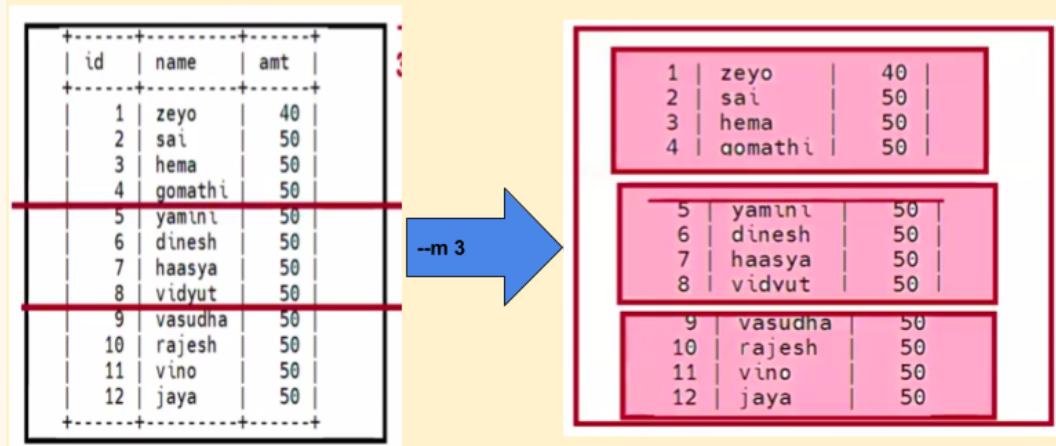
- We use multi-mapper to make imports faster
- If you want to increase mapper value([--m 2](#)) then you have to provide column value by adding an argument in command [--split-by <columnName>](#) or there should be a primary key in the table. If not provide any of these values then there will be an error like the below -

```
22/07/10 09:11:00 ERROR tool.ImportTool: Import failed: No primary key could be found for table avro.tab. Please specify one  
with --split-by or perform a sequential import with '-m 1'.
```

- By default, sqoop export uses 4 ([--m 4](#))threads or number of mappers to export the data

- For each mapper one separate part is created.

Suppose there are 3 mapper \Rightarrow `--m 3`
Below is the table which have 12 record, now each mapper will have 4 record in their part file



- Mappers divide rows between multi-mappers using MIN and MAX values.

Performance Tuning By Mappers

In General, performance tuning in Sqoop can be achieved by:

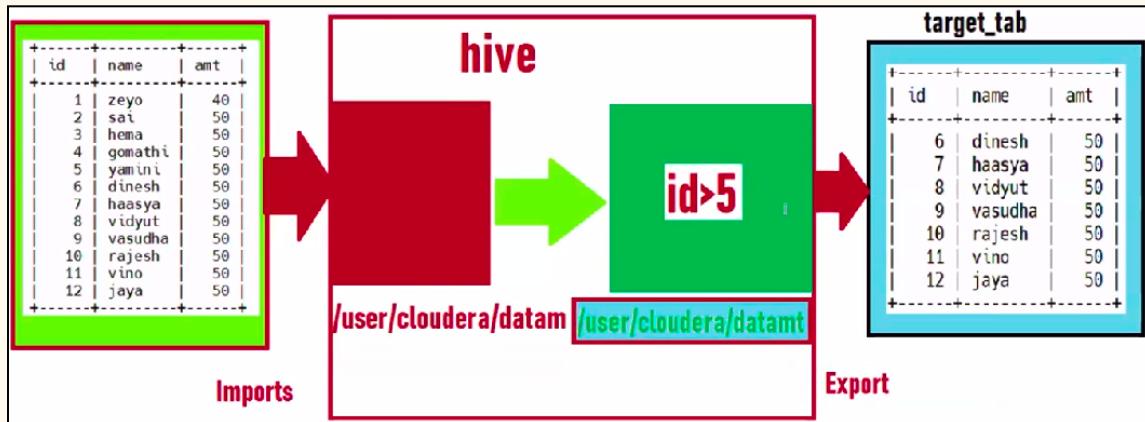
- Controlling Parallelism
 - Using Mappers
- Controlling Data Transfer Process
 - Batching
 - Direct

Export in SQOOP -

Using SQOOP Export data to MySQL from HDFS

```
SQOOP EXPORT --connect jdbc://localhost/<dbName> --username root --password cloudera --m 1 \
--table <targetTableName> --staging-tab <stagingTable> --export-dir <dirPath>
```

***Note** - Here providing staging table(--stage-tab) role is “**Data will export to Stage table first then it will go to Target_table**”, so that if the process couldn’t be completed due to any reason, we can easily delete data from the staged table.



[Using Sqoop move HDFS data to Hive](#)

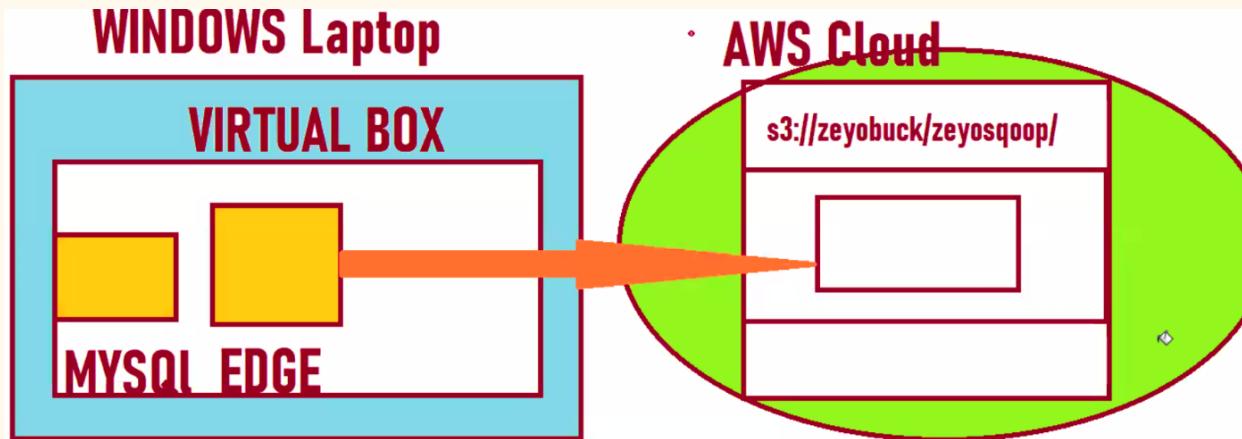
[Using Sqoop move MySql data to Hive](#)

[Locations in Hive](#)

Cloud -

Cloud is the usage of services. Few are the below services -

- **EC2** (Elastic compute cloud) - It is for getting a virtual laptop
- **RDS** (Relational database system) - It is for storing RDBMS data
- **S3** (Simple Storage service/Scalable Storage) - It is storage to store any data.
- **EMR**(Elastic Map Reduce) - For rent Hadoop system, EMR is group of nodes with hadoop install in it.



BOSS --- Sai You have to do a course certification
You need to write an exam but to write an exam you need windows SSD laptop

But sai does not have it he has very basic laptop. lets rent a laptop. Sai is searching is there any one can give WINDOWS SSD laptops for Rental

Sai found three ways to rent Laptop -

1. AWS - **JEFF**
2. Google Cloud Platform - Sunder
3. Azure - Satya Nadela

Sai Decided to go with Jeff.

Sai --- Hey Jeff can you courier the Laptop

Jeff ---- No No no I will not do that . I wil have the SSD laptop with me -- you have to connect remotely

Sai ----- Sure I will do it But how can I connect to that Laptop

Jeff ---- You have webportal just like Flipkart,Amazon - i have my own portal for rental business. Create your own account and add windows laptop to Cart page and check it out

I assure you, your laptop will be available in just 2 minutes.
Once done. Pay me for the minutes you used.

Sai --- whats your WEBPAGE

Jeff --- My rental Product name is Amazon Web Services - AWS

Sai ---- I will create my own account and reach you out. Hey I have create my own account.Next steps

Jeff --- Login first

Sai --- Sure ,i have selected Windows Laptop with SSD and launched its in Running state State

Jeff ---- Yes that mean you laptop is ready to use but wait for 2 Minutes ---- Steps to Connect ██████████

Sai ---- Thanks jeff I completed my work. I will shut it down I used for 2 hours

Jeff --- Sure just pay me 30rs

Sai --- Dear Jeff. I like your service
But I need another Help -- I own zeyobron.
My students was mysql for practise. Can you help me

JEFF --- Dont worry,I have service known RDS.You can go ahead and create MYSQL server

Cloud integration with SQOOP -

Command - Requires 3 arguments -

```
sqoop import
--Dfs.s3a.access.key=<accessKey>
```

```
--Dfs.s3a.secret.key=<SecretKey>
--Dfs.s3a.endpoint=<AWS-s3-endpoint>
--connect jdbc:mysql://localhost/dd2 --username root --password cloudera --table ztab
--m 1 --target-dir s3a://zeyobuck/zeyosqoop/Datadir
```

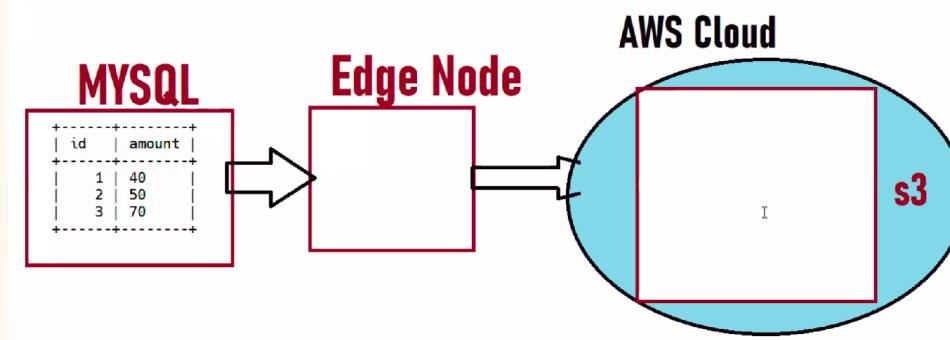
```
=====
mysql
=====

• mysql -uroot -pcloudera
• create database dd2;
• use dd2;
• create table ztab(id int,name varchar(100),tdate date);
• insert into ztab values(1,'Sai', now() - interval 3 day);
• insert into ztab values(2,'zeypo', now() - interval 2 day);
• quit

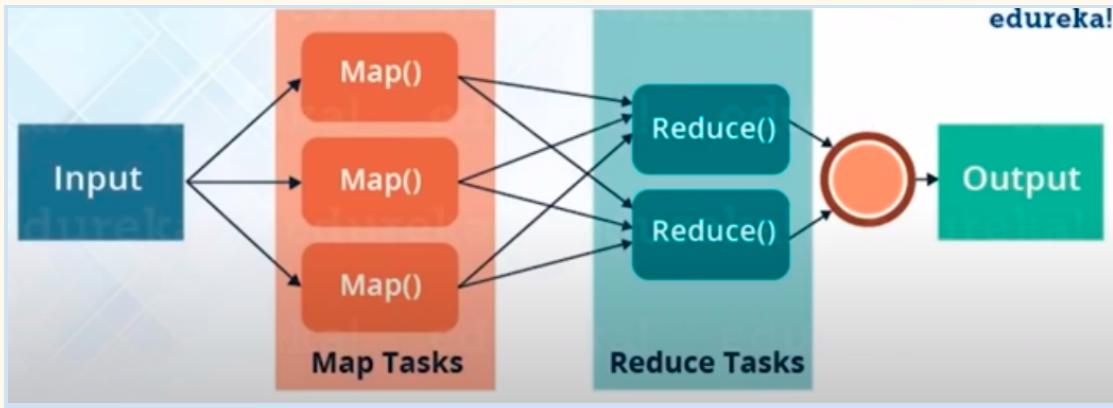
=====
Edge Node
=====

sqoop import
-Dfs.s3a.access.key=AKIAUDT4DPHYY3EHGND3 -Dfs.s3a.secret.key=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
-Dfs.s3a.endpoint=s3.ap-south-1.amazonaws.com
--connect jdbc:mysql://localhost/dd2 --username root --password cloudera --table ztab --m 1 --target-dir s3a://zeyobuck/zeyosqoop/<URNAME>dir

• To verify folder is created in AWS, run below command
○ cd .aws
○ wget https://zeyobuck.s3.ap-south-1.amazonaws.com/credentials
○ cd
○ aws s3 ls s3://zeyobuck/zeyosqoop/
```



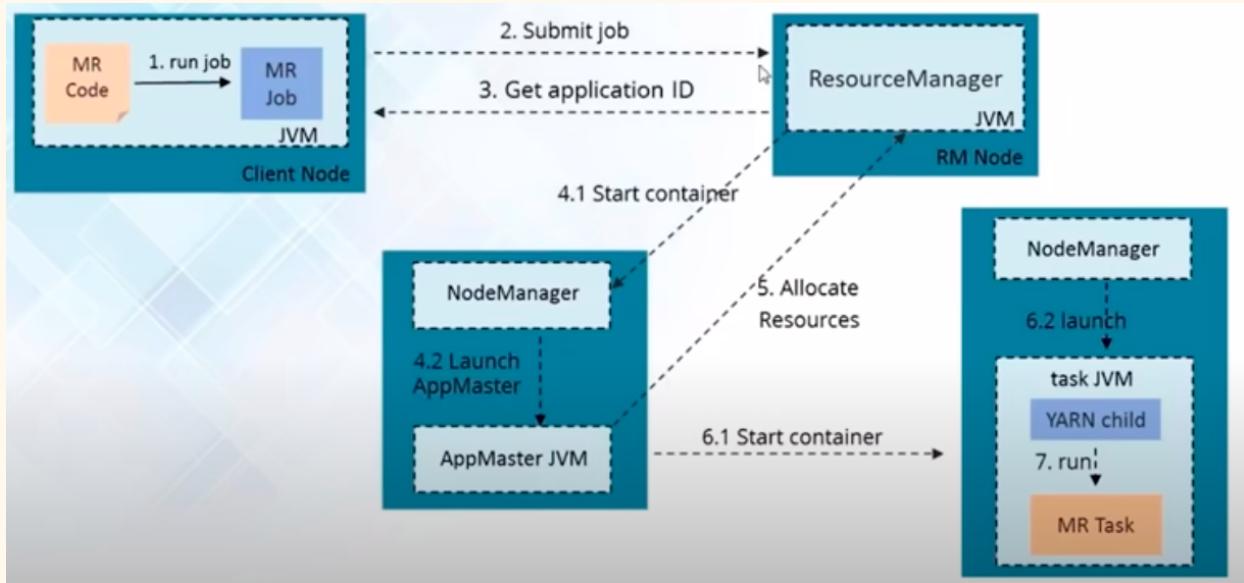
MapReduce



Components of MapReduce

1. **Payload:** The applications implement Map and Reduce functions and form the core of the job
2. **MRUnit:** Unit test framework for MapReduce
3. **Mapper:** Mapper maps the input key/value pairs to the set of intermediate key/value pairs
4. **NameNode:** Node that manages the HDFS is known as namednode
5. **DataNode:** Node where the data is presented before processing takes place
6. **MasterNode:** Node where the job trackers runs and accept the job request from the clients
7. **SlaveNode:** Node where the Map and Reduce program runs
8. **JobTracker:** Schedules jobs and tracks the assigned jobs to the task tracker
9. **TaskTracker:** Tracks the task and updates the status to the job tracker
10. **Job:** A program that is an execution of a Mapper and Reducer across a dataset
11. **Task:** An execution of Mapper and Reducer on a piece of data
12. **Task Attempt:** A particular instance of an attempt to execute a task on a SlaveNode

MapReduce Job workflow -



Abstraction of MapReduce -

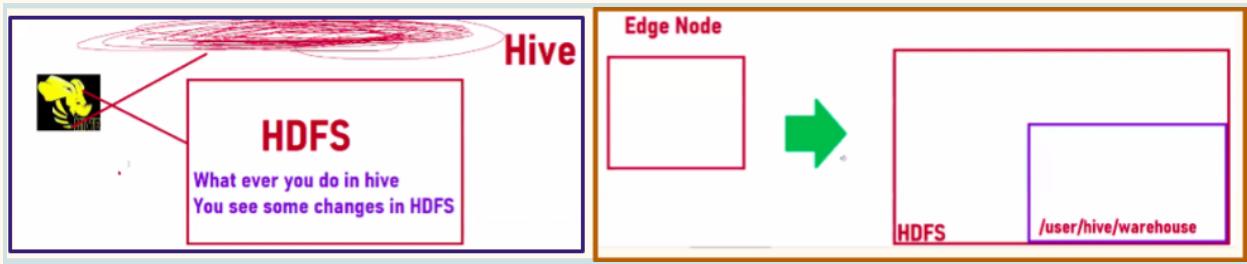
1. Hive - Query engine(uses SQL)
2. Pig - yahoo invented
3. Sqoop
4. Oozie - yahoo invented use for automating/scheduling

HIVE -

Developed by JoyDeep Sen Sarma.

- Hive runs on top of HDFS. In simple, whatever we do in Hive there will be the same changes save in HDFS (dir ⇒ /user/hive/warehouse/)





- Hive table should always be pointed to a location. While creating a table, please ensure to provide a directory Path to the table. If you do not specify any pointing directory, Hive will create a directory in the table's name and get pointed to it.

- Check table Location - **describe formatted <tableName>**

```
# Detailed Table Information
Database:          hivedb
Owner:            cloudera
CreateTime:       Sun Jul 17 02:58:46 PDT 2022
LastAccessTime:   UNKNOWN
Protect Mode:    None
Retention:        0
Location:         hdfs://quickstart.cloudera:8020/user/hive/warehouse/hivedb.db/hivetab
Table Type:      MANAGED_TABLE
```

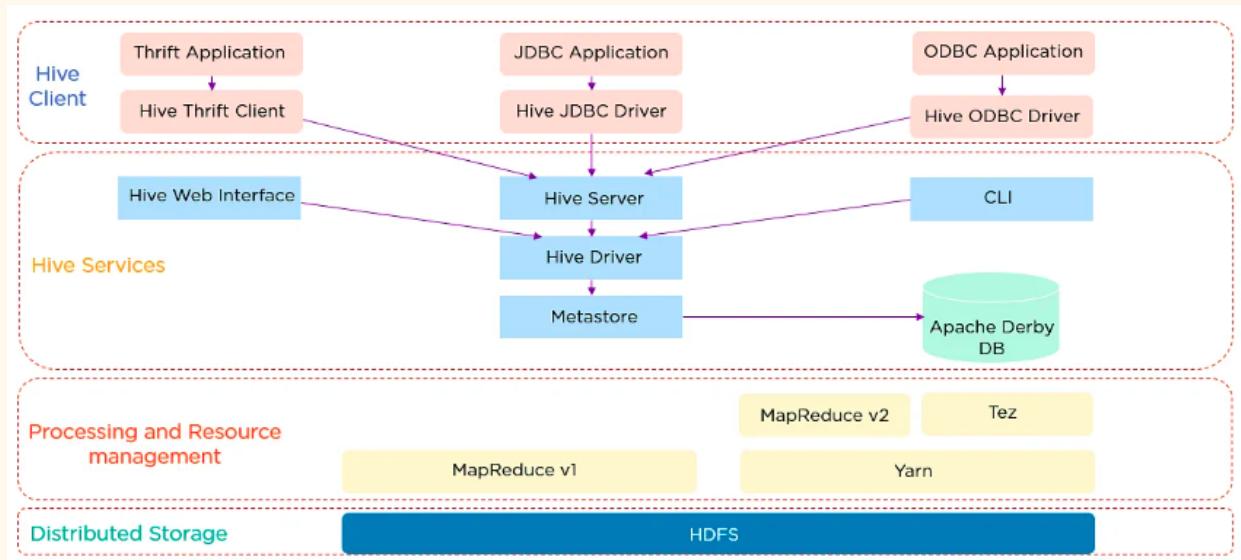
- Location **/user/hive/warehouse** is the default hive location to store all tables and its data if the location is not provided. To provide Location, Read - [Locations in Hive](#)
- Apache Hive is a data warehouse system built on top of Hadoop and is used for analyzing structured and semi-structured data. Hive abstracts the complexity of Hadoop MapReduce. Basically, it provides a mechanism to project structure onto the data and perform queries written in HQL (Hive Query Language) that are similar to SQL statements. Internally, these queries or HQL gets converted to map reduce jobs by the Hive compiler. Therefore, you don't need to worry about writing complex MapReduce programs to process your data using Hadoop. It is targeted towards users who are comfortable with SQL.

Why HIVE come?

- As Map-Reduce is made of JAVA. Earlier Map-Reduce code was difficult to understand because java code wasn't easy at that time.
- To process data in HDFS from one directory to another, you have to follow the below steps -
 - ◆ Should knowledge of JAVA
 - ◆ Should knowledge of Map-Reduce code

- ◆ Do Tune Mappers
 - ◆ Do Tune Reducers
 - ◆ And then export it in JAR
 - ◆ Take that JAR to cluster and deploy it
- Then the Hive tool comes to eliminate the use of Map-Reduce, For HIve, you don't need any Java or Map-Reduce code (although in the backend it uses Map-Reduce).
- Hive looks like SQL but not SQL
- Easy to Work with it
- **Hive internally uses MapReduce to process queries.**
- Hive only replaces JAVA use, not MapReduce.
- Hive can on the top these 3 engines-
- ◆ MapReduce
 - ◆ Tez
 - ◆ Spark

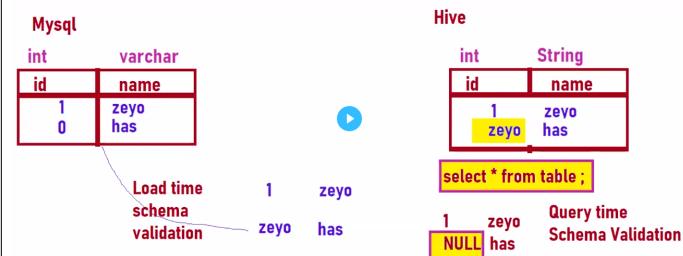
Hive Architecture



HIVE vs MySql

MySQL	HIVE
<p>Data storage for MySQL is DATABASE</p> <p>E T L (Extracted Transform Load) - E \Rightarrow T \Rightarrow L</p> <p>Extracted - Extract de-serialize data</p> <p>Transform - Convert it into Serialized data</p> <p>Load - Then Load it in SQL</p> <p>Run below command to Load data in MySql -</p> <pre>mysql> load data infile 'filePath' into table <targetTableName> fields terminated by ',';</pre>	<p>Data storage for HIVE is HDFS</p> <p>E L T (Extracted Load Transform) - E \Rightarrow L \Rightarrow T</p> <p>In Hive we just need to Extract data and run load command, Transform will take care By HIVE.</p> <p>Load - To Load data in Hive, Run below command to Load data in MySql -</p> <pre>hive> load data local infile '/user/cloudera/data.parquet' into table '<targetTableName>';</pre>
Good for Small Data, for large data it is costly like server, nodes, set up, installation, license	Good and free for large data

Load time schema validation (Schema on write)- once inserting data and changing the data type , data will set to 0 at load time.



Query time schema validation (Schema on Read)- once inserting data and changing the data type , there will be no change in data, But once you query the data it will give you NULL. But originally at backend data have same value that was provided during insert query.

Now after load and query data, get the location of table and check original data at backend

\$ describe formatted <tableName>; ⇒ get Hive file location
\$ hadoop fs -cat /hive/cloudera/<dbName>/<tabName>/*

```
=====
Edge Node
=====
vi schemaCheck 1,zeyo
2,bron
zeyi,bron
bron,4
```

```
=====
HIVE Node
=====
Now load 'schemaCheck' file in SQL
```

- CREATE TABLE scktab(id int, name string) row format delimited fields terminated by ',';
 - load data local inpath '/home/cloudera/sck' into table scktab;
 - SELECT * from scktab;
- For example if have provide string value in id column like in 'sechmaCheck' file, but in table ID column is only accepting integer,
- But once you QUERY- mean SELECT data from table, for these value in HIVE is set to NULL in table -

```
1      zeyo
2      bron
NULL   bron
NULL   4
```

- describe formatted <tableName>; ⇒ get Hive file location

```
# Detailed Table Information
Database:          hivedb
Owner:             cloudera
CreateTime:        Sun Jul 17 11:35:05 PDT 2011
LastAccessTime:    UNKNOWN
Protect Mode:     None
Retention:         0
Location:          hdfs://quickstart.cloudera:8020/user/hive/warehouse/hivedb.db/scktab
Table Type:        MANAGED_TABLE
```

- hadoop fs -cat /hive/cloudera/<dbName>/<tabName>/*

```
[cloudera@quickstart ~]$ hadoop fs -cat /user/hive/warehouse/hivedb.db/scktab/*
1,zeyo
2,bron
zeyi,bron
bron,4
[cloudera@quickstart ~]$
```

=====

Edge Node

=====

```
vi schemaCheck -
1,zeyo
2,bron
zeyi,bron
bron,4
```

=====

SQL Node

=====

Now load 'schemaCheck' file in SQL

- load data local infile '/home/cloudera/sck' into table scktab fields terminated by ',';
 - SELECT * from scktab;
- For example if have provide string value in id column like in 'sechmaCheck' file, but in table ID column is only accepting integer, so for these value in SQL value is set to 0 in table - **it means data changes at LOAD TIME SCHEMA**

SQL only query **STRUCTURED** data

Hive can query **STRUCTURED, Semi-STRUCTURED** data

SQL don't have capability to query large data

Hive have capacity to only query PetaByte data

Type of Table in Hive -

1. MANAGED/Internal Table -

An *internal table* is a table that Hive manages. If you delete an internal table, both the definition in Hive *and* the data directory are deleted inside HDFS.

- ```
CREATE TABLE IF NOT EXISTS Cars(Name STRING, id INT)
 ROW FORMAT DELIMITED
 FIELDS TERMINATED BY ','
 STORED AS TEXTFILE
 LOCATION '/user/<username>/dataDir';

• Drop table Cars;
```

- **Now if you Go to the location and check, directory won't be there anymore**
- Table created for Intermediate data processing, should be Managed type. So we can drop the table once the process is done.

### 2. EXTERNAL Table

An *external table* is a table for which Hive does not manage storage. If you delete an external table in Hive, only the definition in Hive is deleted. The data directory remains inside HDFS.

- ```
CREATE EXTERNAL TABLE IF NOT EXISTS Cars(Name STRING, id INT)
  COMMENT 'Data about cars from a public database'
  ROW FORMAT DELIMITED
    FIELDS TERMINATED BY ','
  STORED AS TEXTFILE
  LOCATION '/user/<username>/dataDir';

• Drop table Cars;
```

- **Now if you Go to the location and check, the directory will be presented there**
- Target table should always be EXTERNAL. Suppose if the table is deleted by mistake, it shouldn't impact on target analytics data.

Question #1 - How to check if table is INTERNAL or EXTERNAL?

- DESCRIBE FORMATTED <tableName>

```
# Detailed Table Information
Database:          hivedb
Owner:             cloudera
CreateTime:        Mon Jul 18 04:30:18
LastAccessTime:    UNKNOWN
Protect Mode:     None
Retention:         0
Location:          hdfs://quickstart.c
Table Type:        MANAGED_TABLE
```

```
# Detailed Table Information
Database:          hivedb
Owner:             cloudera
CreateTime:        Mon Jul 18 05:25:31
LastAccessTime:    UNKNOWN
Protect Mode:     None
Retention:         0
Location:          hdfs://quickstart.c
Table Type:        EXTERNAL TABLE
```

Load data from Edge Node into Hive

1. Create data file (text or .csv, etc) in EDGE node
2. CREATE TABLE table1(id int, name string) row format delimited fields terminated by ',';
3. LOAD DATA LOCAL inpath '/home/cloudera/<filePath>' into table <tableName>;

Load data from HDFS Node into Hive

1. Create data file (text or .csv, etc) in EDGE node
2. Put data in HDFS - hadoop fs -put datacsv.csv /user/cloudera/datafiles
3. CREATE TABLE tab(id int, name string) row format delimited fields terminated by ',';
4. LOAD DATA inpath '/home/cloudera/<filePath>' into table <tableName>;

Note - In step-4 after running the command ⇒ data file will be moved to Hive location (/user/hive/warehouse/<dbName>.db/<tabName>) and CSV or text file will be removed from HDFS location. But once loading a file from EDGE, the data file will be copied instead of completely moving, so the data file remains at both locations EDGE nodes and /user/hive/warehouse/...

```
[cloudera@quickstart ~]$ hadoop fs -ls /user/hive/warehouse/hivedb.db/stab;
Found 2 items
-rwxrwxrwx  1 cloudera supergroup      25 2022-07-18 02:20 /user/hive/warehouse/hivedb.db/stab/datacsv.csv
-rw-r--r--  1 cloudera cloudera       24 2022-07-18 02:44 /user/hive/warehouse/hivedb.db/stab/datacsv copy 1.csv
```

Requirement #1 - Find a way to insert the above source tab data into another Hive table

5. CREATE TABLE targetTab(id int, name string) row format delimited fields terminated by ',';
6. INSERT INTO targetTab select * from tab where <condition>;

Using SQOOP move MYSQL data to Hive

To import data into Hive, below arguments are important -

```
hive-import  
--create-hive-table  
--hive-table HiveDatabase.hiveTableName
```

Sqoop Command Option	Description
--hive-import	Imports tables into Hive using Hive's default delimiters if none are explicitly set.
--hive-overwrite	Overwrites existing data in the Hive table.
--create-hive-table	Creates a hive table during the operation. If this option is set and the Hive table already exists, the job will fail. Set to <code>false</code> by default.
--hive-table <table_name>	Specifies the table name to use when importing data into Hive.

If the database name is not provided in this argument `--hive-table <tableName>` then All data or tables which are imported from MySql will store in the DEFAULT database of Hive. You can find a newly created table under the Default database.

1. Import table into Default database -

Below command will import MySql table to Hive into Default database

- `sqoop import --connect jdbc:mysql://localhost/<dbName> --username root --password cloudera --table <tableName> --m 1 --hive-import`

2. Create table in hive then import SQL table data

```
sqoop import --connect jdbc:mysql://localhost/<dbName> --username root  
--password cloudera --table <SqlTableName> --m 1  
--hive-import --create-hive-table --hive-table hiveDB.<hiveTableName>
```

Note - A data file will be created under '/user/hive/warehouse/' in the name of the Hive table.

```
[cloudera@quickstart ~]$ hadoop fs -ls /user/hive/warehouse/  
Found 3 items  
drwxrwxrwx    - cloudera supergroup          0 2022-07-18 06:56 /user/hive/warehouse/hivedb.db  
drwxrwxrwx    - cloudera supergroup          0 2022-07-18 07:04 /user/hive/warehouse/tabhive2
```

3. We can import partition data, and incremental data as We do in SQOOP, just need to add the below values -

```
--hive-import --create-hive-table --hive-table <hiveTableName>
```

Using SQOOP move HDFS data to Hive

Below command will first ingest Data from Mysql to HDFS ⇒ Then created **process-m-00000** file ⇒ will process to Hive table

```
sqoop import --connect jdbc:mysql://localhost:3306/sqoop --username root --p  
cloudera --m 1 --table customer  
--target-dir /user/cloudera/sqlToHdfsDir  
--fields-terminated-by ","  
--hive-import  
--create-hive-table  
--hive-table sqoop_workspace.customers
```

Locations in Hive

Suppose we are processing data like that -> HDFS ⇒ Source Table ⇒ Target Table

Requirement #1 - Find a way to process only specific data from HDFS to Target table but ensure do not use Source Table(As you did here - [Requirement #1 - Find a way to insert above tab data into another Hive table](#))

1. While creating the table provide the location directory where your data is stored, and once you query the table it will show you data.

```
CREATE EXTERNAL TABLE IF NOT EXISTS Cars(Name STRING, id INT)
ROW FORMAT DELIMITED
FIELDS TERMINATED BY ','
STORED AS TEXTFILE
LOCATION '/user/<username>/dataDir';
```

2. And once you query the table it will show you data. **SELECT * from cars;**

Question #1- What If we increase data at the source file end, will Increased rows show at the Hive end? ⇒ Yes

- Before data file have only 3 rows -

```
[cloudera@quickstart ~]$ hadoop fs -cat /user/cloudera/datafiles/*
4,rubal
5,rubal
6,rubal
```

- Added 2 more rows and new rows are showing at Hive end -

```
hive> select * from ttab;
OK
4      rubal
5      rubal
6      rubal
7      ayushi
8      zeyo
```

Question #2- What If we insert new Row in Hive table using INSERT query, will increased data reflect at Source file Location? ⇒ Yes

- Inserted data in Hive table - **hive> insert into ttab values(9, 'bron');**
- Checked at file location if new row is added there or not -

```
[cloudera@quickstart ~]$ hadoop fs -cat /user/cloudera/datafiles/*
9,bron
4,rubal
5,rubal
6,rubal
7,ayushi
8,zeyo
```

Hive File Formats

1. Load Text Data to Hive

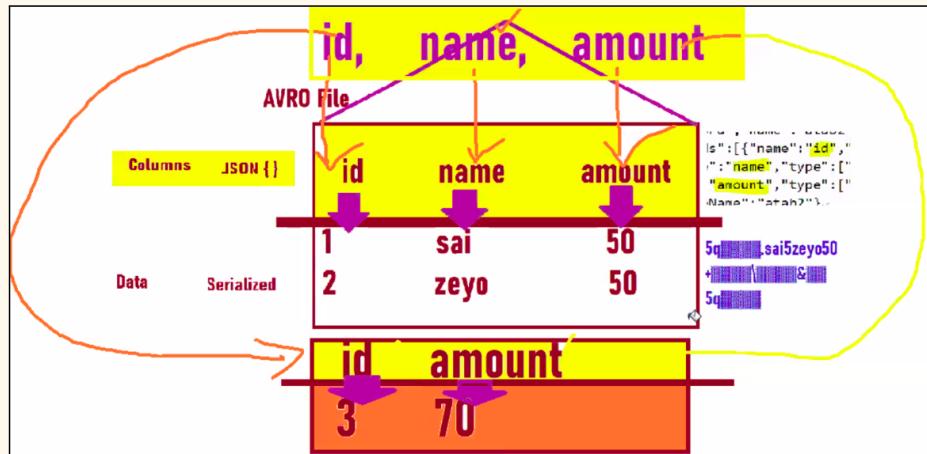
```
create table avroTab(id int, name string) row format delimited fields terminated by ',' LOCATION
'/user/cloudera/project1';
```

2. Load AVRO Data to Hive

```
create table avroTab row format  
SERDE 'org.apache.hadoop.hive.serde2.avro.AvroSerDe' STORED AS AVRO LOCATION '/user/cloudera/project1'  
TBLPROPERTIES ('avro.schema.url'='/user/cloudera/customer.avsc');
```

For more Practice refer [Phase 1 - Project](#)

- AVRO supports schema evolution
- AVRO file has two sections ⇒ JSON schema and serialized data
- Since data has columns coupled, since which columns we call it gets fetched from File
- When you create a hive table for avro data, you cannot specify column names.
Instead you need AVSC file to be available



3. Load Parquet Data to Hive

- sqoop import --connect jdbc:mysql://localhost/pdb1 --username root --password cloudera --table par_tab --m 1 --target-dir /user/cloudera/pdir **--as-parquetfile**
- CREATE TABLE hiveTab (id STRING, name STRING) STORED AS PARQUET location '/user/cloudera/pdir';

4. Load ORC Data to Hive

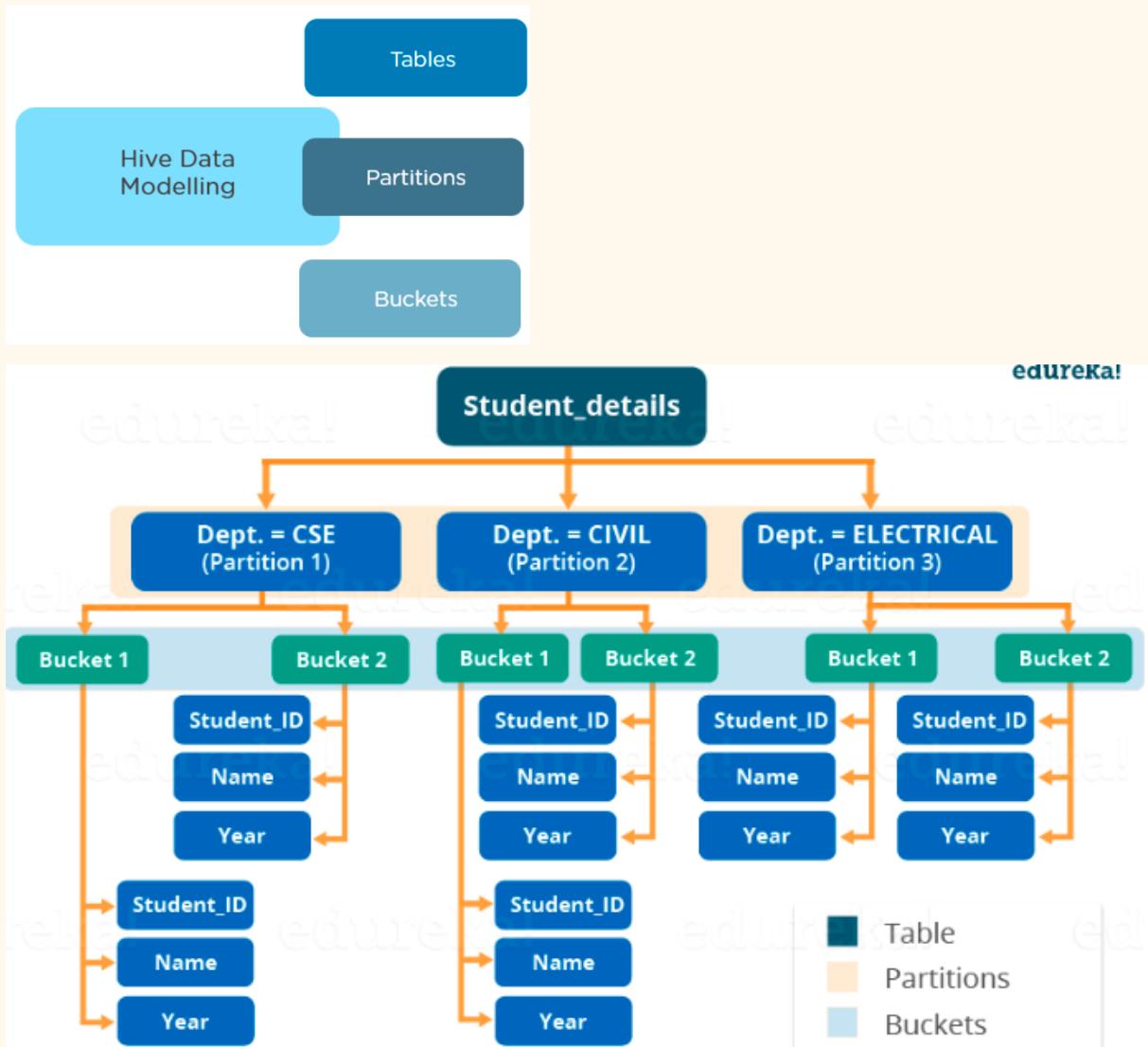
We can not upload data in ORC table using 'load data inpath' (because there is no SQOOP import for ORC). If we want to load data in ORC table we have to use the 'insert into or insert overwrite' command.

```
1. CREATE TABLE IF NOT EXISTS mycars(id INT, Name STRING)
   ROW FORMAT DELIMITED
   FIELDS TERMINATED BY ','
   STORED AS ORC;
2. INSERT INTO orcTab SELECT * FROM targetTab;
----- or -----
INSERT OVERWRITE orcTab SELECT * FROM targetTab;
```

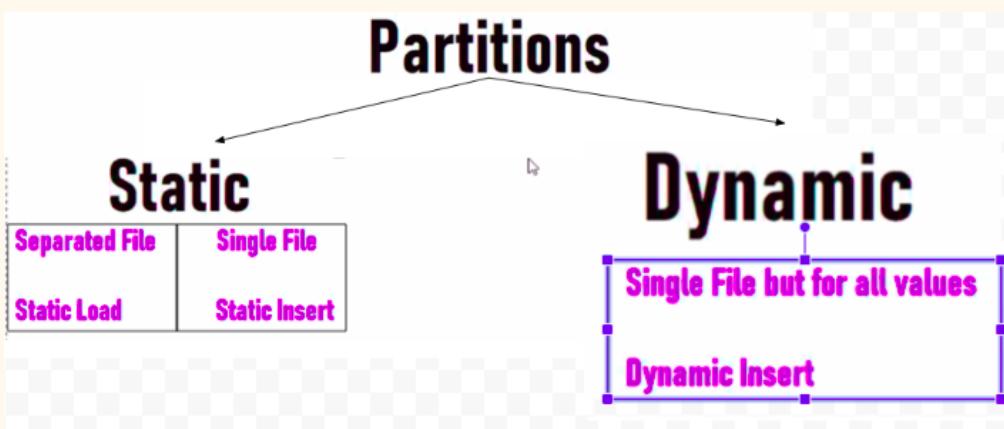
Hive Data Modeling

That was how data flows in the Hive. Let's now take a look at Hive data modeling, which consists of tables, partitions, and buckets:

1. Tables - Tables in Hive are created the same way it is done in RDBMS
2. Partitions - Here, tables are organized into partitions for grouping similar types of data based on partition key
3. Buckets - Data present in partitions can be further divided into buckets for efficient querying



Partitions -



1. Static Partition -

a. Static Load -

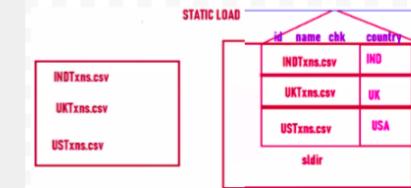
1. CREATE TABLE slPart(id int, name string)
PARTITIONED BY(country string)
ROW FORMAT DELIMITED FIELDS TERMINATED BY ',';
2. load data local inpath '/home/cloudera/partdata/IND.csv' INTO table slPart
PARTITION(country='IND');

#Scenario - We have data separated data in individual files for example country-wise(IND, UK, USA). Create Partitions Table and load each file individual.

IND.csv	UK.csv	US.csv
1,Sai,I	5,Hema,K	9,Jai,S
2,zeyo,I	6,Gomathi,K	10,Swathi,S
3,ze,I	7,Ayushi,K	11,Gupta,S
4,bron,I	8,Ayu,K	

1. Create Static Partition Table -
CREATE TABLE slPart(id int, name string, countryCode string) PARTITIONED BY(country string) \row format delimited fields terminated by ',' location '/user/cloudera/slPartDir';
2. Load each file in table and create individual Partition
load data local inpath '/home/cloudera/partdata/IND.csv' into table slPart PARTITION(country='IND');
load data local inpath '/home/cloudera/partdata/UK.csv' into table slPart PARTITION(country='UK');
load data local inpath '/home/cloudera/partdata/US.csv' into table slPart PARTITION(country='USA');
3. Go to HDFS directory and check partition files are created -

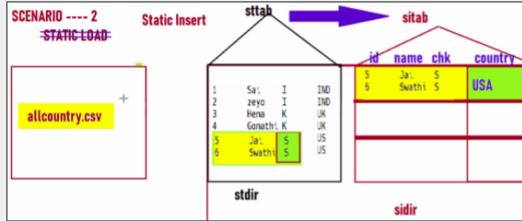
```
[cloudera@quickstart partdata]$ hadoop fs -ls /user/cloudera/slPartDir
Found 3 items
drwxr-xr-x - cloudera cloudera          0 2022-07-28 00:49 /user/cloudera/slPartDir/country=IND
drwxr-xr-x - cloudera cloudera          0 2022-07-28 00:54 /user/cloudera/slPartDir/country=UK
drwxr-xr-x - cloudera cloudera          0 2022-07-28 00:54 /user/cloudera/slPartDir/country=US
```



b. Static Insert -

1. CREATE TABLE siPartTable(id int, name string) PARTITIONED BY(country string) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' ;
2. INSERT INTO siPartTable PARTITION(country='US') SELECT id,name from dataTab WHERE country='US' ;

#Scenario - We have one file and stores data. For example one file contains all country(IND, UK, US) data. Find a way to filter 'US' data and write data into Partition table with Partition 'US'.



1. Create normal Table and load all data from file -
 - create table dataTab(id int, name string, countryCode string, country string) row format delimited fields terminated by ',';
 - load data local inpath '/home/cloudera/partdata/allcountry.csv' into table dataTab;
2. Create Static Partition Table -


```
CREATE TABLE siPart(id int, name string, countryCode string) PARTITIONED BY(country string)
row format delimited fields terminated by ',' location '/user/cloudera/siPartDir';
```
3. Insert data from normal table in Partition table and create individual Partition


```
insert into siPart PARTITION(country=US) SELECT id, name, countryCode from dataTab WHERE country='US';
```
4. Go to HDFS directory and check partition files are created -

```
[cloudera@quickstart partdata]$ hadoop fs -ls /user/cloudera/siPartDir
Found 3 items
drwxr-xr-x - cloudera cloudera          0 2022-07-28 00:49 /user/cloudera/siPartDir/country=IND
drwxr-xr-x - cloudera cloudera          0 2022-07-28 00:54 /user/cloudera/siPartDir/country=UK
drwxr-xr-x - cloudera cloudera          0 2022-07-28 00:54 /user/cloudera/siPartDir/country=US
```

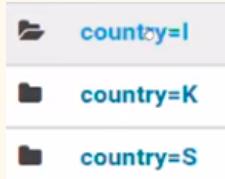
2. Dynamic Partition -

1. CREATE TABLE dynPartTable(id int, name string) PARTITIONED BY(country string) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',';
2. INSERT INTO dynPartTable PARTITION(country) SELECT id, name, country from dataTab;

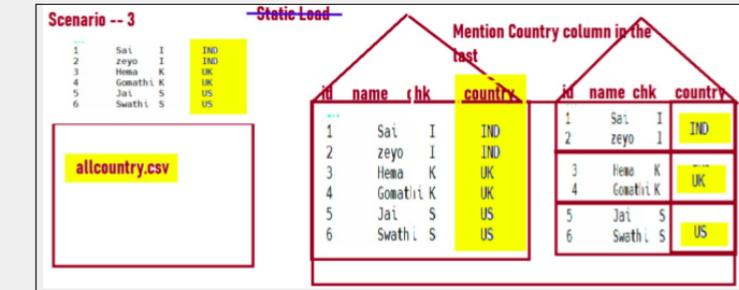
Rules of Dynamic Partition -

1. The last column will be considered for creating dynamic partition files.

Like below, provided partition 'country' but the last column was countryCode -



#Scenario - We have one file and stores data. For example one file contains all country(IND, UK, US) data. Find a way to filter 'US' data and write data into dynamic Partition table considering last column 'country' as reference.



1. Create normal Table and load all data from file -
 - `create table dataTab(id int, name string, countryCode string, country string) row format delimited fields terminated by ',';`
 - `load data local inpath '/home/cloudera/partdata/allcountry.csv' into table dataTab;`
2. Create Dynamic Partition Table taking last column as reference -

`CREATE TABLE dynaPart(id int, name string, countryCode string) PARTITIONED BY(country string)`
 row format delimited fields terminated by ',' location '/user/cloudera/dynamicPartDir';
3. Run command to enable dynamic partition -

`set hive.exec.dynamic.partition.mode=nonstrict;`
4. Insert data from normal table in Partition table and create individual Partition

`INSERT INTO dynamicPart PARTITION(country) SELECT id, name, countryCode, country FROM dataTab;`
5. Go to HDFS directory and check all dynamic partition files are created as per country name -


```
[cloudera@quickstart partdata]$ hadoop fs -ls /user/cloudera/dynamicPartDir
Found 3 items
drwxr-xr-x - cloudera cloudera          0 2022-07-28 03:13 /user/cloudera/dynamicPartDir/country=IND
drwxr-xr-x - cloudera cloudera          0 2022-07-28 03:13 /user/cloudera/dynamicPartDir/country=UK
drwxr-xr-x - cloudera cloudera          0 2022-07-28 03:13 /user/cloudera/dynamicPartDir/country=US
```

Bucketing

Sai Video - <https://www.youtube.com/watch?v=WGATmuJTKGO>

1. Analytics column is highly cardinal
2. Joins are easy
3. Good for Integer columns
4. Reduces runs in this operations
5. Bucketing is not folder creation, it's file creation.
6. It uses existing column to create buckets

1. `CREATE EXTERNAL TABLE part_Buck(no INT, name STRING, city STRING, state STRING)`
`PARTITIONED BY(country STRING) CLUSTERED BY(no) into 5 BUCKETS`
 row format delimited fields terminated by ',' location '/user/cloudera/partBucketData';
2. /* Enable Partition and Bucketing property in Hive - */
`set hive.enforce.bucketing=true;`

```

set hive.exec.dynamic.partition.mode=nonstrict;

3. /* Insert data from Staging table to Partition Bucket target Table */
INSERT INTO part_buck PARTITION (country) SELECT no, name, city, state, country from
stagingtab;

```

Steps -

1. Load data into MySql
`load into infile '/user/cloudera/datafile/data.csv' into table data_tab;`
2. Sqoop Import to HDFS from SQL -
`sqoop import --connect jdbc:mysql://localhost/sqldb --username root --password cloudera --table bucktab --m 1 --target-dir /user/cloudera/stagingData`
3. Create Staging Hive Table and location HDFS dir-
`CREATE TABLE stagingtab(no INT, name STRING, city STRING, state STRING, country string) row format delimited fields terminated by ',' location '/user/cloudera/stagingData';`
4. Create Partition Bucket target Table in Hive
`CREATE EXTERNAL TABLE part_Buck(no INT, name STRING, city STRING, state STRING)
PARTITIONED BY(country STRING) CLUSTERED BY(no) into 3 BUCKETS
row format delimited fields terminated by ',' location '/user/cloudera/partBucketData';`
5. Enable Partition and Bucketing property in Hive -
`set hive.enforce.bucketing=true;
set hive.exec.dynamic.partition.mode=nonstrict;`
6. Insert data from Staging table to Partition Bucket target Table
`INSERT INTO part_buck PARTITION(country) SELECT no, name, city, state, country from stagingtab;`

Partition and bucketing is created like below for 'country' Partition -

```

[cloudera@quickstart partdata]$ hadoop fs -ls /user/cloudera/partBucketData/country=India
Found 3 items
-rw-r--r-x  1 cloudera cloudera      43 2022-07-29 03:32 /user/cloudera/partBucketData/country=India/000000_0
-rw-r--r-x  1 cloudera cloudera      58 2022-07-29 03:32 /user/cloudera/partBucketData/country=India/000001_0
-rw-r--r-x  1 cloudera cloudera      58 2022-07-29 03:32 /user/cloudera/partBucketData/country=India/000002_0
[cloudera@quickstart partdata]$ hadoop fs -cat /user/cloudera/partBucketData/country=India/000000_0
0,Rubal,Kapur,U.P
3,Ayushi,banglore,DELHI
[cloudera@quickstart partdata]$ hadoop fs -cat /user/cloudera/partBucketData/country=India/000001_0
7,Rubal,banglore,T.M
4,Adity,LKO,DELHI
1,sai,banglore,T.N

```

7. Now Go to HDFS location and check Partition Folder and Bucketing files are created -

```

[cloudera@quickstart partdata]$ hadoop fs -ls /user/cloudera/partBucketData/country=India
Found 3 items
-rw-r--r-x  1 cloudera cloudera      43 2022-07-29 03:32 /user/cloudera/partBucketData/country=India/000000_0
-rw-r--r-x  1 cloudera cloudera      58 2022-07-29 03:32 /user/cloudera/partBucketData/country=India/000001_0
-rw-r--r-x  1 cloudera cloudera      58 2022-07-29 03:32 /user/cloudera/partBucketData/country=India/000002_0
[cloudera@quickstart partdata]$ hadoop fs -cat /user/cloudera/partBucketData/country=India/000000_0
0,Rubal,Kapur,U.P
3,Ayushi,banglore,DELHI
[cloudera@quickstart partdata]$ hadoop fs -cat /user/cloudera/partBucketData/country=India/000001_0
7,Rubal,banglore,T.M
4,Adity,LKO,DELHI
1,sai,banglore,T.N

```

How Bucketing decide which part file records will go

Bucketing use Hash Function -

Suppose we have 3 bucketing, now there will be 3 part files created

	Name	Size	User
			cloud
	-		cloud
	000000_0	46 bytes	cloud
	000001_0	17 bytes	cloud
	000002_0	0 bytes	cloud

- Now suppose you have created bucketing for column **id** which contains values like 1,2,3,4,5,6,.....
- Hash function ⇒**

Column value % Number of Bucketing = 00000(remainder) 0

- For example id value is = 1

1 % 3 = 1 (means it will go to **000001_0**)

- For example id value is = 2

2 % 3 = 1 (means it will go to **000002_0**)

- For example id value is = 3

3 % 3 = 0 (means it will go to **000000_0**)

- For example id value is = 5

5 % 3 = 1 (means it will go to **000002_0**)

How to Decide [Bucket Count] in Hive

Youtube -

https://www.youtube.com/watch?v=pHvyhfWK43s&list=PLLa_h7BriLH2kU021Tkxbn5N6LhIp9dpu&index=8

$$\begin{aligned}
 ① \text{Table size} &= 2300 \text{ MB} \\
 ② \text{Block size} &= 128 \text{ KB} \\
 2300 / 128 &= 17.96 \\
 2^{\lceil D \rceil} &\geq 17.6 \\
 2^5 &= 32 \rightarrow \text{Bucket count}
 \end{aligned}$$

Partition Vs Bucketing

Partition	Bucketing
We decide which data have to go in which partition/folder	Here bucketing internally decides using a hash function which record will go in which file.
Partition creates folder and inside that folder we have part-m-00000 file	Bucketing creates part-m-00000 files as per number of buckets
	<ul style="list-style-type: none"> It provides faster query responses like partitioning. In bucketing due to equal volumes of data in each partition, joins at Map side will be quicker.
When to use Partitioning?	When to use Bucketing?
<ul style="list-style-type: none"> When the column with a high search query has low cardinality. For example, if you create a partition by the country name then a maximum of 195 partitions will be made and these number of directories are manageable by the hive. On the other hand, do not create partitions on the columns with very high cardinality. For example- product IDs, timestamp, and 	<ul style="list-style-type: none"> We cannot do partitioning on a column with very high cardinality. Too many partitions will result in multiple Hadoop files which will increase the load on the same node as it has to carry the metadata of each of the partitions. If some map-side joins are involved in your queries, then bucketed tables are a good option. Map side join is a process where two tables are joined using the map function only without any reduce function. I would recommend you to go

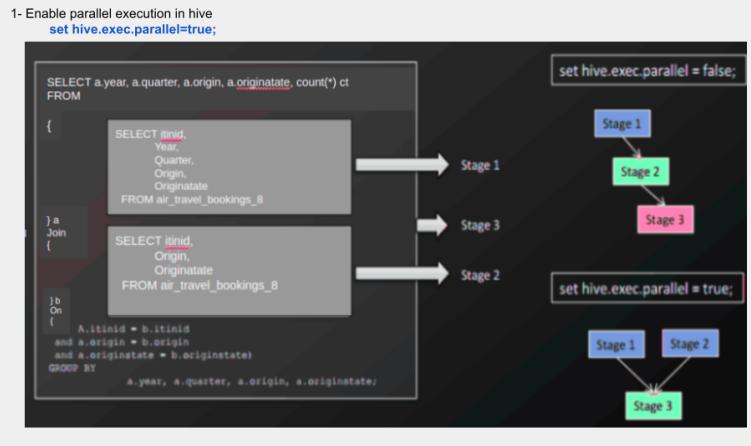
price because it will create millions of directories which will be impossible for the hive to manage.

- It is effective when the data volume in each partition is not very high. For example, if you have the airline data and you want to calculate the total number of flights in a day. In that case, the result will take more time to calculate over the partition "Dubai" as it has one of the busiest airports in the world whereas for a country like "Albania" will return results quicker.

through this article for more understanding about map-side joins.

Hive Performance Tuning

1. Partitioning
2. Bucketing
3. Parquet
4. Parallel Execution



5. Vectorization

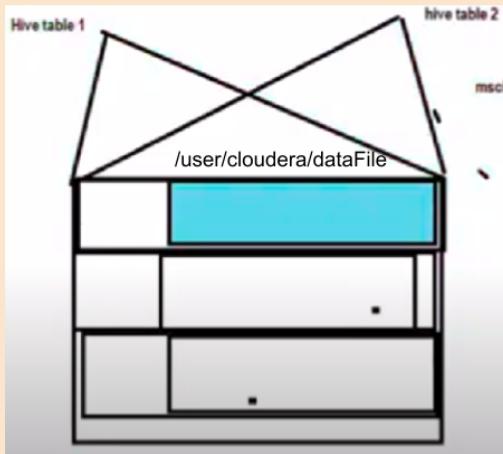
```
set hive.vectorization.execution.enable=true;
```

Usually hive considers only 1 row at a time to process but if you enable vectorization it would consider 1024 rows at a time to process but it kills a lot of memory and core occupancy.

6. Map side join

Task - What will happen if create another partition Bucket Table on the Top of same Location where external Table already created -

Table will be created and its location is same as first one Partition table , Also location have data **but it will show nothing while Query→**



Run below command to Repair this -

- `msck repair table part_buck2;`

```
hive> msck repair table part_buck2;
OK
Partitions not in metastore:    part_buck2:country=INDIA
buck2:country=U.S
Repair: Added partition to metastore part_buck2:country=INDIA
Repair: Added partition to metastore part_buck2:country=India
Repair: Added partition to metastore part_buck2:country=U.K
Repair: Added partition to metastore part_buck2:country=U.S
```

Phase 1 - Project

Phase 1 - Project

Step 1 ⇒ Download winscp (Windows) or Filezilla for (MAC)/Ubuntu

- Windows Download Winscp and install
<https://winscp.net/download/WInSCP-5.19.6-Setup.exe>

Step 2 ⇒ Start filezilla (Login Just like putty) and copy prodata to the /home/cloudera/ ⇒

Host:	192.168.0.101	Username:	cloudera	Password:	*****	Port:	22
-------	---------------	-----------	----------	-----------	-------	-------	----

Step 3 ⇒ Cloudera Folks Goto Mysql and create a table

- create table customer_total(id int(10),username varchar(100),sub_port varchar(100),host varchar(100),date_time varchar(100),hit_count_val_1 varchar(100),hit_count_val_2 varchar(100),hit_count_val_3 varchar(100),timezone varchar(100),method varchar(100),`procedure` varchar(100),value varchar(100),sub_product varchar(100),web_info varchar(100),status_code varchar(100));
- load data infile '/home/cloudera/prodata.txt' into table customer_total fields terminated by ',';
- select * from customer_total;
- create table customer_src(id int(10),username varchar(100),sub_port varchar(100),host varchar(100),date_time varchar(100),hit_count_val_1 varchar(100),hit_count_val_2 varchar(100),hit_count_val_3 varchar(100),timezone varchar(100),method varchar(100),`procedure` varchar(100),value varchar(100),sub_product varchar(100),web_info varchar(100),status_code varchar(100));
- insert into customer_src select * From customer_total where id>0 and id<101;

=====

Edge Node

=====

- mkdir /home/cloudera/avsrcdir
- cd /home/cloudera/avsrcdir
- echo -n cloudera>/home/cloudera/passfile
- sqoop job --delete inpjob
- sqoop job --create inpjob -- import --connect jdbc:mysql://localhost/prodb --username root --password-file file:///home/cloudera/passfile -m 1 --table customer_src --target-dir /user/cloudera/customer_stage_loc --incremental append --check-column id --last-value 0 --as-avrodatafile
- sqoop job --exec inpjob
- hadoop fs -mkdir /user/cloudera/avscdirpro
- hadoop fs -put /home/cloudera/avsrcdir/customer_src.avsc /user/cloudera/avscdirpro

=====

Hive shell

=====

- create table customer_src ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.avro.AvroSerDe' STORED AS AVRO LOCATION '/user/cloudera/customer_stage_loc' TBLPROPERTIES ('avro.schema.url'='/user/cloudera/avscdirpro/customer_src.avsc');
- select * from customer_src; === U will see the data
- create external table customer_target_tab partitioned by (current_day string,year string,month string,day string) ROW FORMAT SERDE 'org.apache.hadoop.hive.serde2.avro.AvroSerDe' STORED AS AVRO LOCATION '/user/cloudera/customer_target_tab' TBLPROPERTIES ('avro.schema.url'='/user/cloudera/avscdirpro/customer_src.avsc');
- select * from customer_target_tab; ===== U will not see the data

SPARK -

Invented by Matai Zaharia in 2010 but initially released by Apache in 2014.

Big Data Processing

- Hadoop introduced a radical new approach based on two key concepts
 - - Distribute the data when it is stored
 - Run computation where the data is
- Spark takes this new approach to the next level
 - Data is distributed in memory

Spark Installation

1. Prerequisites ⇒ Java and Scala
2. Download Java in case it is not installed using the below commands.

```
sudo apt-get install python-software-properties  
sudo apt-add-repository ppa:webupd8team/java  
sudo apt-get update  
sudo apt-get install oracle-java8-installer
```

3. Download the latest Scala version from [Scala Lang Official](#) page. Once installed, set the scala path in the `~/.bashrc` file as shown below.

```
export SCALA_HOME=Path_Where_Scala_File_Is_Located  
export PATH=$SCALA_HOME/bin:$PATH
```

4. Download Spark from the [Apache Spark Downloads](#) page.
5. Extract Spark tar using the below command.

```
tar -xvf spark-2.1.0-bin-hadoop2.7.tgz
```

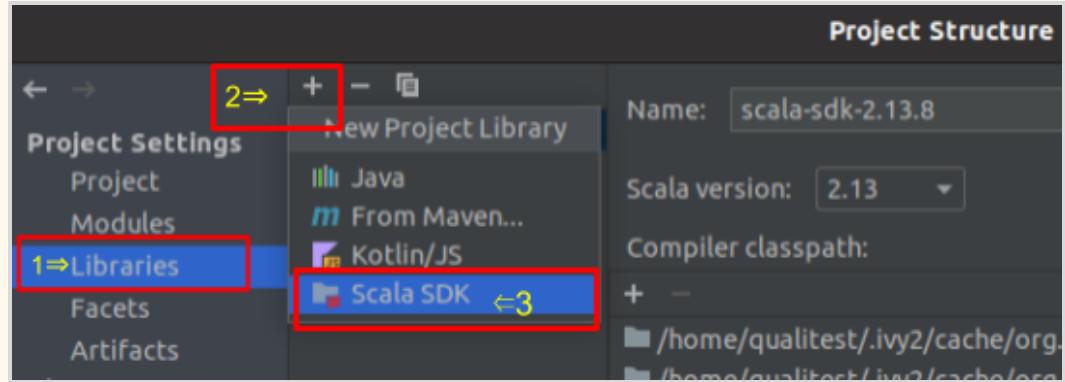
6. Set the Spark_Path in `~/.bashrc` file.

```
export SPARK_HOME=Path_Where_Spark_Is_Installed  
export PATH=$PATH:$SPARK_HOME/bin
```

Spark set up with Scala and Run on IntelliJ

1. Make sure Java is installed on 8/11/17. (windows - Download winUtils and save in dir → `D:/hadoop/bin/winutils.exe` and set environment variable path)

- i. HADOOP_HOME ⇒ D:/hadoop
- ii. Update ‘Path’ variable with add a new path item there ⇒ “%HADOOP_HOME%\bin”
2. Create Maven Project with Java
3. Install Scala Plugin in IntelliJ
4. Rename /src/main/java with /src/main/scala
5. Setup Scala SDK - Go to **Project-structure** in IntelliJ → Go to Libraries → **Add Scala SDK** → Click on APPLY and click OK



6. Add Spark Dependencies → Now either **Add Spark Jar files** ([spark-2.4.7-bin-hadoop2.7.tgz](#) → jars folder) or **Add Spark Maven dependency** in pom.xml

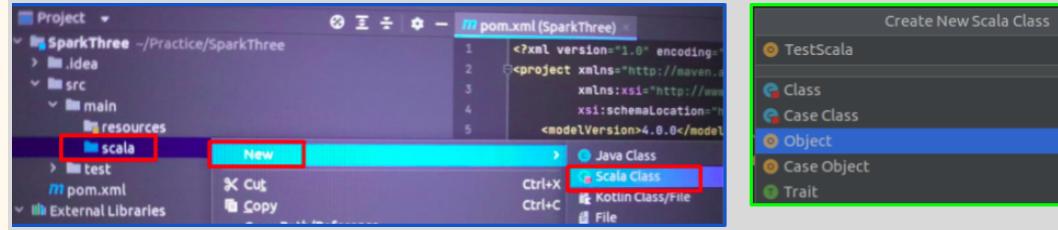
```

<dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-core_2.12</artifactId>
    <version>3.3.0</version>
</dependency>

<dependency>
    <groupId>org.apache.spark</groupId>
    <artifactId>spark-sql_2.12</artifactId>
    <version>3.3.0</version>
    <scope>compile</scope>
</dependency>

```

7. Create first Scala Object under /src/main/scala
 - i. Right Click on /scala/ package and select New
 - ii. Select Scala class
 - iii. Give the object name and select **scala object**. Like the below image -



8. Run Maven Build → `mvn clean install`
9. Now run the Spark application `SparkSessionTest` program.
10. Run jar in Cloudera - `spark-submit --class pack.test`
`SparkThree-1.0-SNAPSHOT.jar`

Why does Spark come into the Market?

1. Hive used to satisfy only a few requirements.
2. Hive runs on MR
3. Hive is quite slow for very very very large data set
4. Hive RUNS on HDFS (No other platform)
5. Hive No support on Streaming
6. Hive no support on Machine Learning
7. Hive cannot integrate with any sources to bring the data
8. Hive is not customizable easily

Matei Zaharia

2010--- I got Super power

Features of My Tool

- Simple To Use**
- Environment Independent (HDFS,LINUX,WINDOWS,S3,AZURE blob)**
- My Tool Can bring the data,process the data, put the data**
- My tool Support SQL**
- My Tool is DAMNNNNNNNNNNNNNN faster**

My Tool Support Streaming

My tools Support Machine Learning

Good Cloud Support

It uses only one formalue to get integrated with any tools

Advantages of Spark

- faster
- in-memory processing
- lazy computations
- environment independent
- MLV streaming
- Cloud Support
- Customization Serialization
- Unified Formulae
- GraphDB support
- SQL

Why is SPARK Faster?

- Spark uses in-memory processing while MR(MapReduce) uses I/O processing)
- Lazy in nature - Once Action is triggered then only Transactions will run
 - Transaction 1 <==== Transaction 2 <==== Action
- The lightning-fast performance due to its In-Memory Processing Capability brought Spark its place amongst the top-level Apache Projects.
-
-

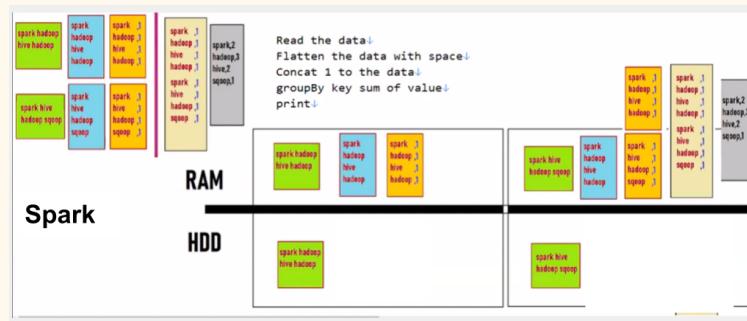
Brief Introduction to Apache Spark

- Apache SparkTM is a fast and general engine for large-scale data processing framework for massive parallel computing (cluster)
- Harnessing the power of cheap memory
- Written using Scala, Java, and Python languages
- High-level APIs support in Scala, Java and Python
- Has had 36,338 commits made by 1,319 contributors. Representing lakhs of lines of code with proper comments.
- Developers from 50+ companies including UC Berkeley, Cloudera, Yahoo, Databricks, Intel, Groupon, etc., Apache Committers from 16+ organizations

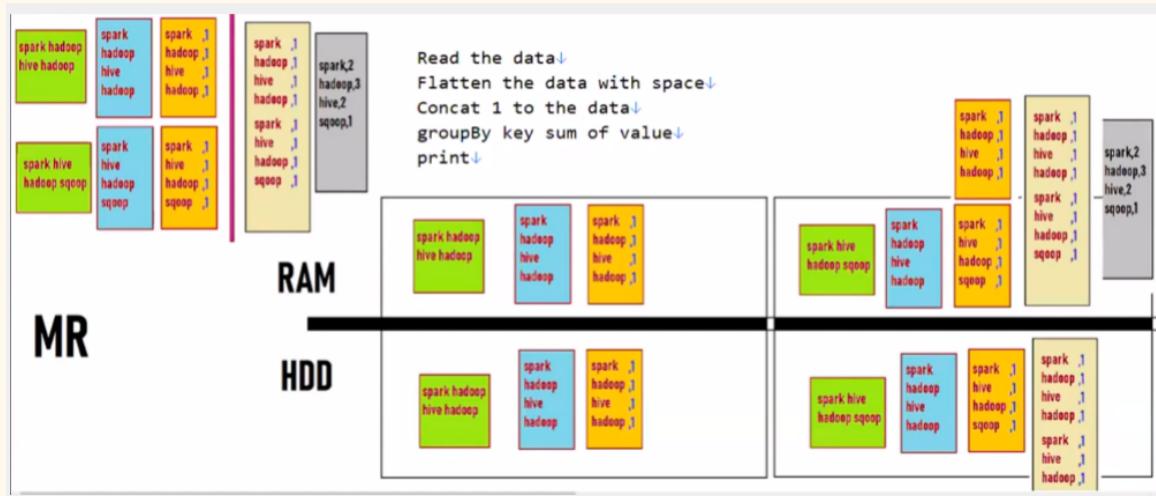
SPARK Vs MapReduce

Spark		MR	
In-memory processing		I/O processing	
Lazy evolution		Linear - Run steps one by one	
Till the Action is triggered none of the Transaction will work.		Operation run step by step	
Transactions	Actions		
Join	Count Select		

Spark - In-Memory Processing example -



MR - I/O Processing example -



Beyond MapReduce

MapReduce was great for batch processing, but users quickly needed to do more:

- More complex, multi-pass algorithms
- More interactive ad-hoc queries
- More real-time stream processing

Result: many *specialized* systems for these workloads

Spark Execution Model

Spark offers you

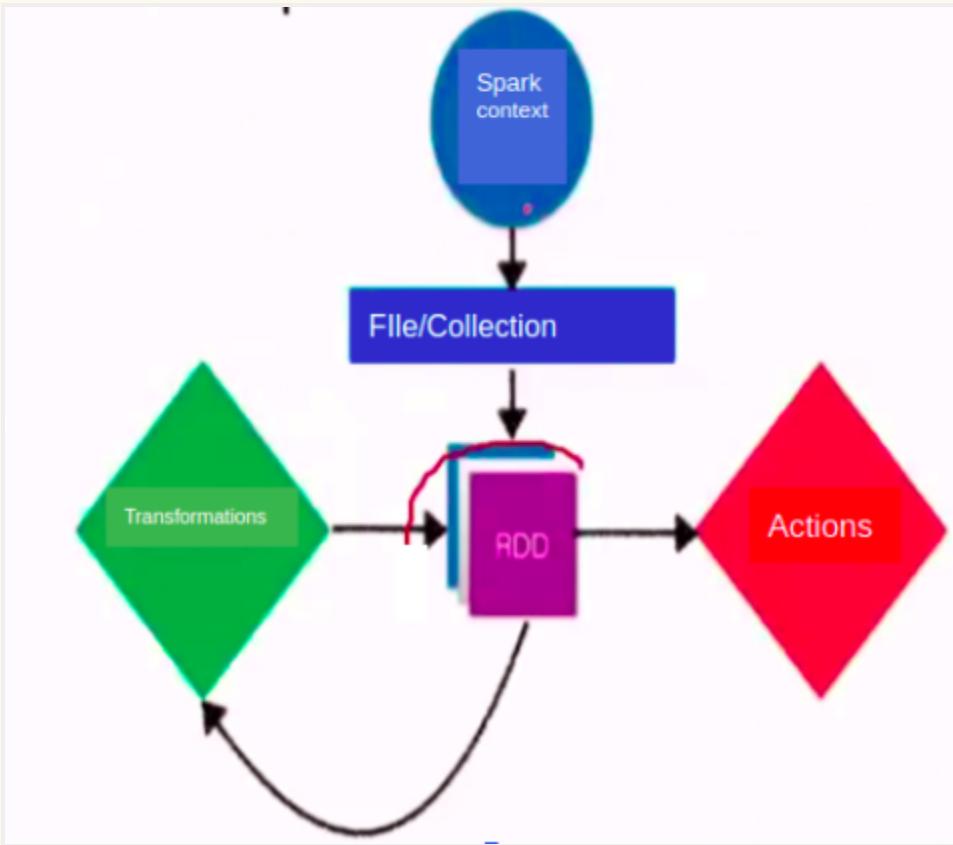
- Lazy Computations
 - Optimize the job before executing
- In-memory data caching
 - Scan HDD only once, then scan your RAM
- Efficient pipelining
 - Avoids the data hitting the HDD by all means

Pillars of Spark

Direct Acyclic Graph - sequence of computations performed on data

- Node - RDD partition
- Edge - transformation on top of data
- Acyclic - graph cannot return to the older partition
- Direct-transformation is an action that transitions data partition state (suppose from A node to B node)

Spark Internals



Spark Process the data in Two Ways

- ★ RDD and other one is Dataframe

RDD

1. RDD means -
 - ★ **Resilient** - if data in memory is lost, it can be recreated
 - ★ **Distributed** - stored in memory across the cluster
 - ★ **Dataset** - initial data can come from a file or be created programmatically
2. RDDs are the fundamental unit of data in Spark
3. Most Spark programming consists of performing operations on RDDs
4. It is fault-tolerant if you perform multiple transformations on the RDD and then due to any reason any node fails. The RDD, in that case, is capable of recovering automatically.
5. it does not have an optimizer like catalyst
6. it is recommended to use a Data frame because of its performance

7. it is available from 1.0 version
8. For the RDD process, we need a jar “**Spark-core.jar**” which contains the class name “**SparkContext**” which has many methods like **textFile()**, **filter**, **map()**, **flatMap()**, etc.
9. To import ‘**SparkContext**’ in scala file write →**import org.apache.spark.SparkContext**

Below is the example of the RDD query and output -

1. ----- open Spark shell -----
spark-shell
2. `sc.textFile("file:///home/cloudera/data.csv").filter(x => x.contains("US")).foreach(println)`

```
5,Jai,S,US
6,Swathi,S,US
```

===== How to Create SparkContext Object =====

```
import org.apache.spark.SparkContext

object SparkWordCount {
  def main(args: Array[String]) {

    /* "setAppName()" -> This is the name of the application that you want to run.
     "setMaster()" --> This parameter denotes the master URL to connect the spark application to.
    */
    val conf = new SparkConf().setAppName("first").setMaster("local[*]")

    /* Creating a [RDD] SparkContext Object */
    val sc = new SparkContext(conf)

    /* Read input file and store in variable */
    val input = sc.textFile("input.txt")
    val count = input.flatMap(line ⇒ line.split(" "))
      .map(word ⇒ (word, 1))
      .reduceByKey(_ + _)
    count.saveAsTextFile("outfile")
  }
}
```

- Use local[x] when running in Standalone mode. x should be an integer value greater than 0; this represents how many partitions it should create when using RDD,

DataFrame, and Dataset. Ideally, the x value should be the number of CPU cores you have.

RDD Advantages

- In-Memory Processing
- Immutability
- Fault Tolerance
- Lazy Evolution
- Partitioning
- Parallelize

DataFrame → SPARK SQL

- ★ It is conceptually a Row Columnar format
- ★ It has a optimizer like catalyst
- ★ It is available from 1.3 version
- ★ Can read and write different file formats
- ★ For the DataFrame process, we need a jar “**Spark-sql.jar**” which contains the class name “**SparkSession**” which has many methods like **read()**, **filter**, etc.
- ★ To import **SparkSession** in scala file write →**import org.apache.spark.sql.SparkSession**

Below is the example of the RDD query and output -

1. **spark-shell**
2. **spark.read.csv("file:///home/cloudera/data.csv").filter("_C3='US'").show()**

_c0	_c1	_c2	_c3
5	Jai	S	US
6	Swathi	S	US

===== How to Create SparkContext Object =====

```
import org.apache.spark.sql.SparkSession

object SparkWordCount {
  def main(args: Array[String]) {
```

```
/* SparkSession.builder() – Return SparkSession.Builder class. This is a builder for SparkSession.  
master(), appName() and getOrCreate() are methods of SparkSession.Builder.
```

“appName()” -> This is the name of the application that you want to run.
“master()” --> This parameter denotes the master URL to connect the spark application to.

```
/  
Creating a [DataFrame] Spark Session Object */
```

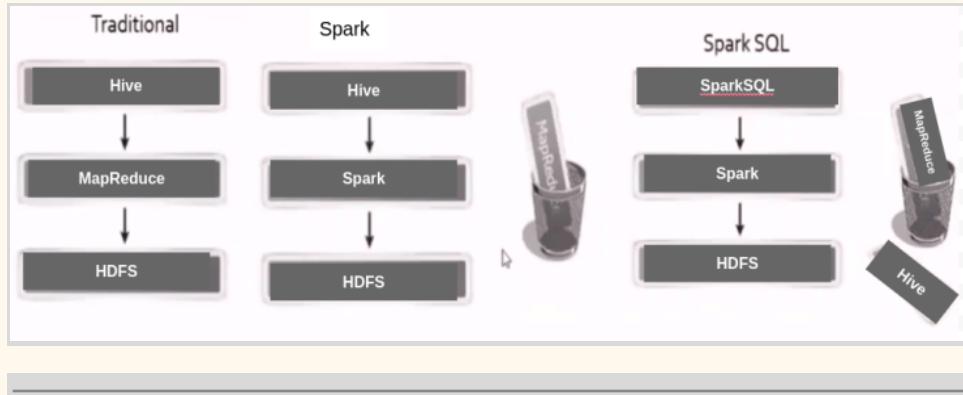
```
var sparkSession = SparkSession.builder()  
    .master("local[1]")  
    .appName("SparkByExamples.com")  
    .getOrCreate();
```

```
/* Read input file and store in variable */  
val input = sc.read.csv("input.csv")  
}  
}
```

RDD Vs DataFrame

RDD	DataFrame
Low Level -> If the requirement is small and easy, then RDD is a good choice.	High Level -> If the requirement is huge and big scenarios, then DF is a good choice. Like have semi-structured data
Support only Structured data	Support Structure and semi-structured
Through RDD, we can process structured as well as unstructured data. But, in RDD user need to specify the schema of ingested data, RDD cannot infer its own	In the data frame data is organized into named columns. Through dataframe, we can process structured and unstructured data efficiently. It also allows Spark to manage schemas.
Performance is slow as compared DF	Performance is faster as compared to RDD. DF has an inbuilt Tungston framework which makes DF to perform faster.
Less integration to read data from external source - It can only ready files, There are some extend to read data from Hbase, Hive, Cassandra, NoSQL but it is very very difficult	Very good integration to read data from external source due to Unified formula <code>(spark.read.format("csv").load("path"))</code>

Evolution of Spark SQL



Scala programs

Practice more programs here - <https://www.w3resource.com/scala-exercises/basic/index.php>

Word Count Program

```
import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark._

object SparkWordCount {
  def main(args: Array[String]) {

    /* "setAppName()" -> This is the name of the application that you want to run.
       "setMaster()" --> This parameter denotes the master URL to connect the spark application to.
    */
    val conf = new SparkConf().setAppName("first").setMaster("local[*]")

    /* Creating a [RDD] SparkContext Object */
    val sc = new SparkContext(conf)
```

```

/* Read input file and store in variable */
val input = sc.textFile("input.txt")
val count = input.flatMap(line => line.split(" "))
    .map(word => (word, 1))
    .reduceByKey(_ + _)
count.saveAsTextFile("outfile")
}
}

```

Program - How to use Filter, Map, FlatMap. Below is the List ⇒

```

"Amazon-Jeff-America",
'Microsoft-BillGates-America',
'TCS-TATA-india',
'Reliance-Ambani-INDIA"

```

Below are the operation to do on the above list ⇒

- Filter ‘india’ records (Ensure capital INDIA also filtered)
- Flatten(split) the result with Hyphen
- Then replace “india” with “local” value
- Then Concatenate “- Done” at the end of each item of the list then print it

```

object obj {

def main(args: Array[String]): Unit = {
var liststr = List("Amazon-Jeff-America",
  "Microsoft-BillGates-America",
  "TCS-TATA-india",
  "Reliance-Ambani-INDIA")

println("String List ==> " + liststr)

var list2 = liststr.filter(x => x.toLowerCase().contains("india"))
    .flatMap(x => x.split("-"))
    .map(x => x.replace("india", "local"))
    .map(x => x.concat(", Done"))

println("filter list ==> ")
list2.foreach(println)

===== OUTPUT =====
TCS- Done
TATA- Done
local- Done

```

Reliance- Done
Ambani- Done
INDIA- Done

Write a Scala program to compute the sum of the two given integer values. If the two values are the same, then return triples of their sum.

```
package pack

object obj2 {

    def main(args: Array[String]): Unit= {
        print(sum(1,2))
        print(sum(2,2))
        print(sum(4,4))
    }

    def sum(i: Int, i1: Int): Int= {
        if (i == i1) (i + i1) * 3 else (i + i1)
    }
}
```

RDD - Use Case #1

- Read the data file (data.csv) [just copy below data in a file]
- /*

```
first_name,last_name,company_name,address,city,county,state,zip,age,phone1,phone2,email,web
James,Butt,"Benton, John B Jr",6649 N Blue Gum St,New
Orleans,Orleans,LA,70116,9,504-621-8927,504-845-1427,jbutt@gmail.com,http://www.bentonjohnbjr.com
Josephine,Darakjy,"Chanay, Jeffrey A Esq",4 B Blue Ridge
Blvd,Brighton,Livingston,MI,48116,8,810-292-9388,810-374-9840,josephine_darakjy@darakjy.org,http://www.chanayjeffreyaesq.com
Art,Venere,"Chemel, James L Cpa",8 W Cerritos Ave
#54,Bridgeport,Gloucester,NJ,8014,7,856-636-8749,856-264-4130,art@venere.org,http://www.chemeljameslcpa.com
*/


- Filter rows greater than which length > 200
- Flatten(split) the Filtered data with comma
- Replace hyphen with Nothing (Remove the Hyphen)
- Concat string → ",zeyo" at the End of each element
- Store result in a file



```
=====
```


```

```

package pack
import org.apache.spark.{SparkConf, SparkContext}
object obj1 {

def main(args:Array[String]): Unit={

  val conf = new SparkConf().setAppName("Spark practice").setMaster("local[*]")
  val sc = new SparkContext(conf)

  /* Read the data file */
  var data = sc.textFile("file:///home/Practice/data/data.csv")

  data = data.filter(x => x.length > 200)      /* Filter rows greater than which length > 200 */
  .flatMap(x => x.split(","))           /* Flatten(split) the Filtered data with comma */
  .map(x => x.replace("-", ""))
  .map(x => x.concat(" -> zeyo done")) /* Replace hyphen with Nothing (Remove the Hyphen) */
  .map(x => x.concat(" -> zeyo done")) /* Concat string →",zeyo" at the End of each element */

  data.foreach(println)

  /* Data write to a file (provide non-exist folder name in below)*/
  data.coalesce(1).saveAsTextFile("file:///home/Practice/resultData")
}

}

```

RDD with DF- Use Case #2

1. Read data from text file
2. Split each column with comma like below

00000000,06-26-2011,Exercise,GymnasticsPro
00000002,06-01-2011,Exercise,GymnasticsPro
00000003,06-05-2011,Gymnastics,Rings
/* After split using map() */

index 0 1 2 3

00000000 06-26-2011 Exercise GymnasticsPro
00000002 06-01-2011 Exercise GymnasticsPro
00000003 06-05-2011 Gymnastics Rings

3. **Impose(create) schema using case class**

4. Convert it into schemaRDD

5. Filter products(column/index 4) element contains "Gymnastics"

6. Convert it into DataFrame

7. Save result as Parquet

Using SchemaRDD

```
package pack
import org.apache.spark.sql.SparkSession
import org.apache.spark.{SparkConf, SparkContext}
object obj1 {

  case class schema(id: String, date: String, category: String, product: String) /* define schema using case class */

  def main(args: Array[String]): Unit = {
    val conf = new SparkConf().setAppName("spark practice").setMaster("local[*]")
    val sc = new SparkContext(conf)

    val data = sc.textFile("file:///home/Practice/data/datatxns.txt")

    val splitData = data.map(x => x.split(",")) /* split data */

    /* Using SchemaRDD*/
    val rddSchemaResult = splitData.map(x => schema(x(0), x(1), x(2), x(3))) /* convert data into schema(case class)*/
      .filter(x => x.product.contains("Gymnastics")) /* filter data */

    println("===== RDD result for 'Gymnastics' product =====")
    rddSchemaResult.foreach(println)

    println("===== Start DataFrame =====")
    val ss = new SparkSession.Builder().getOrCreate()
    import ss.implicits._

    val df = rddSchemaResult.toDF() /* convert rdd result into dataframe */
    df.createOrReplaceTempView("filterTable") /* create sql view to perform more operation on data*/
    val df_result = ss.sql("SELECT * from filterTable where date='06-26-2011'") /* run sql query on created view */

    println("===== Filter more data using DataFrame =====")
    df_result.show() /* print data */
    df_result.coalesce(1).write.parquet("file:///home/Practice/data/parquetData") /* save data in parquet file*/
  }
}
```

Solution 2 → Run same able scenario but instead of using schemaRDD(case class) use Row rdd

1. Step (1) & 2 is same as above solution
2. After split data now CONVERT it to ROW rdd

3. Filter products(column/index 4) element contains “Gymnastics”

4. **Impose(create) schema using StructType()**

5. Convert it into DataFrame

6. Save result as Parquet

Using RowRDD

```
package pack
import org.apache.spark.sql.types.{StringType, StructField, StructType}
import org.apache.spark.sql.{Row, SparkSession}
import org.apache.spark.{SparkConf, SparkContext}

object obj1 {
  def main(args: Array[String]): Unit = {
    val conf = new SparkConf().setAppName("spark practice").setMaster("local[*]")
    val sc = new SparkContext(conf)
    val data = sc.textFile("file:///home/Practice/data/datatxns.txt")

    val splitData = data.map(x => x.split(",")) /* split data */

    /* Using Row RDD*/
    val rowData = splitData.map(x => Row(x(0), x(1), x(2), x(3))) /* convert data into Row rdd */
    .filter(x => x(3).toString.contains("Gymnastics")) /* filter data */

    println("===== RDD result for 'Gymnastics' product =====")
    rowData.foreach(println)

    val structSchema = StructType(Array(
      StructField("id", StringType, true),
      StructField("date", StringType, true),
      StructField("category", StringType, true),
      StructField("product", StringType, true)
    ))
    println("===== Start DataFrame =====")
    val ss = new SparkSession.Builder().getOrCreate()

    val df = ss.createDataFrame(rowData, structSchema) /* create frame using rdd data and struct schema */
    df.createOrReplaceTempView("filterTable")
    /* create sql view to perform more operation on data*/
```

```

val df_result = ss.sql("SELECT * from filterTable where date='06-26-2011'") /* run sql query on created view */

println("===== Filter more data using DataFrame =====")
df_result.show() /* print data */
df.coalesce(1).write.parquet("file:///home/Practice/data/parquetData") /* save data in parquet file*/
}
}

```

Ways of Creating RDD

```

# 1. parallelizing data collection
my_list = [1, 2, 3, 4, 5]
my_list_rdd = sc.parallelize(my_list)

## 2. Referencing to external data file
file_rdd = sc.textFile("path_of_file")

```

Ways of Create Dataframe

```

=====
Create DataFrame from RDD
=====

val columns = Seq("language", "users_count")
val data = Seq(("Java", "2000"), ("Python", "100000"), ("Scala", "3000"))
val rdd = sc.parallelize(data)
val dfFromRD1 = rdd.toDF()
dfFromRDD.printSchema()

```

```

=====
Using createDataFrame() with the Row type
=====

val columns = Seq("language", "users_count")
val data = Seq(("Java", "2000"), ("Python", "100000"), ("Scala", "3000"))
val rdd = sc.parallelize(data)
val schema = StructType( Array(
    StructField("language", StringType, true),
    StructField("users", StringType, true)
))

```

```
        ))
val rowRDD = rdd.map(x => Row(x._1, x._2))
val dfFromRDD = spark.createDataFrame(rowRDD,schema)
```

```
=====
Using createDataFrame() from SparkSession
```

```
=====
val columns = Seq("language","users_count")
val data = Seq(("Java", "2000"), ("Python", "100000"), ("Scala", "3000"))
val rdd = sc.parallelize(data)
var dfFromData = spark.createDataFrame(data).toDF(columns:_*)
```

Read data from CSV/TEXT/JSON/Parquet/ORC (These file type are supported by Spark by default)

```
val df = spark.read.csv("/src/resources/file.csv")
val df = spark.read.text("/src/resources/file.txt")
val df = spark.read.json("/src/resources/file.json")
val df = spark.read.format("csv").options("header", true).load("file:///C:/data/datatxns.txt")
val df = spark.read.format("json").load("file:///C:/data/devices.json")
val df = spark.read.format("orc").load("file:///C:/data/data.orc")
val df = spark.read.format("parquet").load("file:///C:/data/data.parquet")
```

Read data from Avro

(need to add jara or dependency -> `spark-avro_2.12`)

```
=====
val df = spark.read.format("avro").load("file:///home/Practice/data.avro")
```

Read data from XML

(need to add jara or dependency -> `spark-xml_2.12`)

```
=====
val df = spark.read.format("xml")
    .option("rootTag", "POSLog")
```

```
.option("rowTag", "Transaction")
.load("file:///home/Practice/data/transactions.xml")

----- Below is xml -----
<?xml version="1.0" encoding="UTF-8"?>
<POSLog>
  <Transaction>
    <RetailStoreID>48</RetailStoreID>
    <WorkstationID>6</WorkstationID>
    <SequenceNumber>73</SequenceNumber>
    <BusinessDayDate>2014-09-30</BusinessDayDate>
    <EndDateTime>2014-09-30T06:20:14</EndDateTime>
    <OperatorID OperatorName="KERRY P">48237</OperatorID>
    <CurrencyCode>USD</CurrencyCode>
    <RetailTransaction Version="2.1">
      ...
    </RetailTransaction>
  </Transaction>
  <Transaction>
    ...
  </Transaction>
</POSLog>
```

Read data from RDBMS Database

(need to add jara or dependency -> mysql-connector-java)

```
=====
val df_mysql = spark.read.format("jdbc")
  .option("url", "jdbc:mysql://localhost:port/db")
  .option("driver", "com.mysql.jdbc.Driver")
  .option("dbtable", "tablename")
  .option("user", "user")
  .option("password", "password")
  .load()
```

Read data from AWS s3 location file

(need to add jara or dependency -> aws-java-sdf)

```
=====
/* Read file from AWS S3 location */
```

```
val df = spark.read.format("csv")
  .option("header", true)
  .option("fs.s3a.access.key", "access-key")
  .option("fs.s3a.secret.key", "secret-key")
  .load("s3a://zeyodevbb/datatxns.txt")

df.show()
```

Read data from Hive Table

```
val df = sparkSession.read
  .options(Map(HBaseTableCatalog.tableCatalog -> catalog))
  .format("org.apache.spark.sql.execution.datasources.hbase")
  .load()
```

Spark File Mode/ Write file in Spark

There are 4 modes → error(default), append, overwrite, ignore

```
=====
Write new data in existing folder but remove all old data files
=====

object obj {
  def main(args: Array[String]): Unit = {
    val conf = new SparkConf().setAppName("practice").setMaster("local[*]")
    val sc = new SparkContext(conf)
    val spark = SparkSession.builder().getOrCreate()

    val data = spark.read.orc("file:///home/Practice/data/devices.csv")
    data.createOrReplaceTempView("dataView")
    val filterData = spark.sql("SELECT * from dataView where lat>40")

    /* 1st way to create write orc file */
    filterData.write.mode("overwrite").orc("file:///home/Practice/data/d1")

    /* 2nd way to create write orc file */
    filterData.write.format("orc")
      .mode("overwrite")
      .save("file:///home/Practice/data/d2")
```

```
}

=====
Write/Add new data in existing folder but also keep all old data
=====
filterData.write.format("orc")
    .mode("append")
    .save("file:///home/Practice/data/d2")

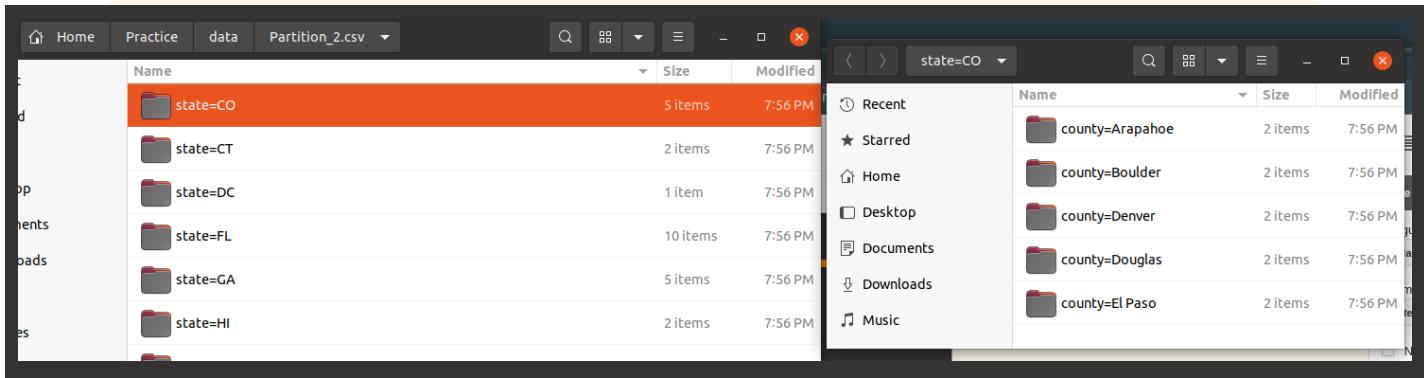
=====
Write new data in folder but if folder already exists don't write data
=====
filterData.write.format("orc")
    .mode("ignore")
    .save("file:///home/Practice/data/d2")
```

Create Partitioned data using scala?

```
/* Read file form Local */
val df = spark.read.format("csv").option("header", true)
    .load("file:///home/Practice/data/usdata.csv")

/* data Partition by 'state' */
df.write.format("csv").partitionBy("state")
    .mode(SaveMode.Overwrite).save("file:///home/Practice/data/Partition_1.csv")

/* data Partition by 'state' and "country" */
df.write.format("csv").partitionBy("state", "county")
    .mode(SaveMode.Overwrite).save("file:///home/Practice/data/Partition_2.csv")
```



Options use cases

Delimiter - it will split data from that ‘~’ into columns

```
val csvdf = spark.read.format("csv")
    .option("header", "true")
    .option("delimiter", "~")
    .load("file:///C:/data/txns_head_pipe")
```

Spark SQL

1. Read data from the file and save it in variable
2. Create SQL view (the view is like a table for temporary to execute options that should not affect the original table)
3. Run SQL command using spark.sql(command) on that view just like a table.

```

val jsondf = spark.read.format("json")
.load("file:///C:/data/devices.json")

jsondf.show()

jsondf.createOrReplaceTempView("jdf")
val finaldf = spark.sql("select * from jdf where lat>40")

finaldf.show()

```

Spark DSL(Domain specific language)

Filters with column conditions and operations

To run any operation over any column value, use col(columnName) to take columns from data.

```

val df = spark.read.format("csv").option("header", value = true)
.load("file:///D:/BigData/data/dt.txt")

println("===== Or '||' filter =====")
df.filter(col("category") === "Exercise"
    || col("spendby") === "cash").show()

println("===== And '&&' filter =====")
df.filter(col("category") === "Exercise"
    && col("spendby") === "cash").show()

println("===== Multi value filter =====")
df.filter(col("category") isin ("Exercise", "Gymnastics")).show()
println("===== not in Multi value filter =====")
df.filter(!(col("category") isin ("Exercise", "Gymnastics"))).show()

println("===== 'like' operator using filter =====")
df.filter(col("product") like "%Gymnastics%").show()

println("===== 'null' value filter =====")
df.filter(col("id")isNull).show()

println("===== not 'null' value filter =====")
df.filter(col("id") isNotNull).show()

//SQL Expression
df.filter("gender == 'M'").show(false)

```

```

//Array condition - array_contains() is inbuilt SQL pack function
df.filter(array_contains(df("languages"), "Java")).show(false)

//Struct condition
df.filter(df("name.lastname") === "Williams").show(false)

-----
Few below alternatives for filter the column where "state" is column name
-----

df.filter('state' === "OH").show()
df.filter($state === "OH").show()
df.filter(col("state") === "OH").show()
df.where(df("state") === "OH").show()
df.where('state' === "OH").show()
df.where($state === "OH").show()
df.where(col("state") === "OH").show()

```

withColumn use case

withColumn (colName , Expr)



Give me a column as well as expression, I will perform expression and assign the result to that Column

If you give the column which does not exist in colName, I will still perform the expression but I will create a new column at the End

- To run the condition on only one condition and without bothering other columns use this method, as the result expression will run on only that column and all other columns will also be in the result with any change.
- For example, using selectExpr() we have to give all column names that exist in the table with the expressional column, using withColumn we can reduce that work.

===== Using withColumn =====

```
println("==== using withColumn select only year from date and renamed column with  
'withColumnRenamed' ===")  
  
df.withColumn("tdate", expr("split(tdate, '-')[2]"))  
.withColumnRenamed("tdate", "YEAR")  
.withColumn("status", expr("case when spendby='cash' then 0 else 1 end"))  
.show()
```

	id	YEAR	amount	category	product	spendby
1	00000000	2011	200	Exercise GymnasticsPro	cash	
2	00000001	2011	300	Exercise Weightlifting	credit	
3	00000002	2011	100	Exercise GymnasticsPro	cash	
4	00000003	2011	100	Gymnastics	Rings	credit
5	00000004	2011	300	Team Sports	Field	cash
6	00000005	2011	200	Gymnastics	Ring	cash

===== Using selectExpr =====

```
println("===== use of 'selectExpr' example =====")  
df.selectExpr("id",  
    "split(tdate, '-')[2] as year", /* divide date with '-' and select  
    "amount",  
    "concat(category, ', done')",  
    "product").show()
```

- If in **colName** value ‘column’ exist in table than expression value will be override otherwise if column not exist but expression is correct new column will be created in the end of table with expressions value, like below -

```
println("==== using withColumn when column not exist ===")  
  
df.withColumn("nonExistColumn", expr("split(tdate, '-')[2]")).show()
```

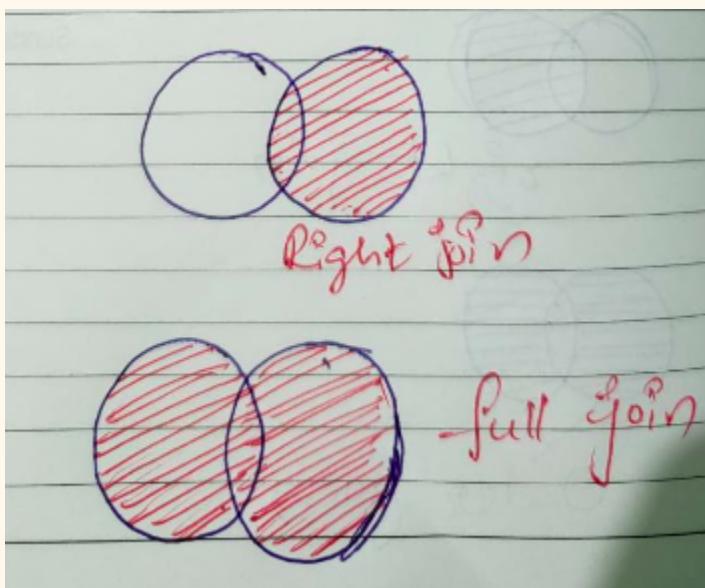
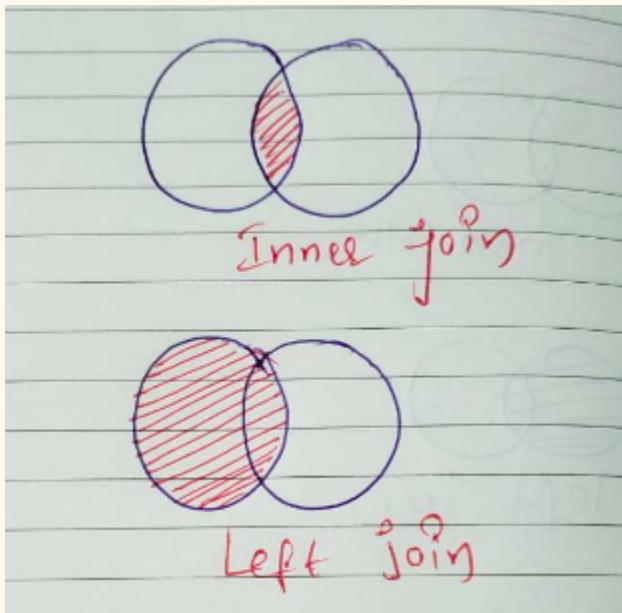
==== using withColumn when column not exist ===							
	id	tdate	amount	category	product	spendby	nonExistColumn
1	00000000	06-26-2011	200	Exercise GymnasticsPro	cash		2011
2	00000001	05-26-2011	300	Exercise Weightlifting	credit		2011
3	00000002	06-01-2011	100	Exercise GymnasticsPro	cash		2011
4	00000003	06-05-2011	100	Gymnastics	Rings	credit	2011
5	00000004	12-17-2011	300	Team Sports	Field	cash	2011
6	00000005	02-14-2011	200	Gymnastics	Ring	cash	2011

DSL Functions Queries

```
println("==== GroupBy and Aggregate method use =====")  
  
/* find total amount of category */  
df.groupBy("category")  
  .agg(sum("amount").as("Total Amount"))  
  .orderBy(col("category").desc)  
  .show()  
  
/* find count of category */  
df.groupBy("category")  
  .agg(count("category").as("count"))  
  .orderBy(col("category").desc)  
  .show()
```

Joins

Inner, Left, Right, Full joins, left_anti joins



```

val df1 = spark.read.format("csv").option("header", "true")
    .load("file:///home/Practice/data/join1.csv")

val df2 = spark.read.format("csv").option("header", "true")
    .load("file:///home/Practice/data/join2.csv")

println("===== Inner Join =====")
df1.join(df2, df1("txnno") === df2("tno"), "inner")
    .drop("tno")
    .show()

println("===== Left Join =====")
df1.join(df2, df1("txnno") === df2("tno"), "left")
    .drop("tno")

```

```

.show()

println("===== Right Join =====")
df1.join(df2, df1("txnno") === df2("tno"), "right")
.drop("txnno")
.select("tno", "txndate", "amount", "custno", "spendby")
.show()

println("===== Full Join =====")
df1.join(df2, df1("txnno") === df2("tno"), "full")
.withColumn("txnno", expr("case when txnno is null then tno else txnno end"))
.drop("tno")
.show()

println("===== Left-Anti Join =====")
df1.join(df2, df1("txnno") === df2("tno"), "left_anti")
.drop("tno")
.show()

```

Tasks -

```

=====
Task #1. Read data.csv file and only print specific columns in dataframe
=====

/* read file */
val df = spark.read.format("csv").option("header", true)
.load("file:///home/Practice/data/usdata.csv")

/* print only specific columns */
df.select("first_name", "last_name").show()

=====

Task #2. Select only year from date column value
  Using selectExpr, we can run any sql function over any column
=====
println("===== use of 'selectExpr' example =====")
df.selectExpr("id",
  "split(tdate, '-')[2] as year", /* divide date with '-' and select 2nd index */
  "amount",
  "concat(category, ', done')",
  "product").show()

```

```
=====
Task #3. Add new column and provide '0' or '1' as per spendBy column values
Using selectExpr,we can run any sql function over any column
=====
println("===== use of 'selectExpr' example =====")
df.selectExpr("id",
  "split(tdate, '-')[2] as year", /* divide date with '-' and select 2nd index */
  "amount",
  "concat(category, ', done')",
  "Product",
  "case when spendby = 'case' then 0 else 1 end as STATUS").show()
```

```
=====
Task #4 - Get Day of the week from Date and Timestamp columns
=====
df.withColumn("date", to_date(col("date"), "d/M/yyyy")) /* Converts the column into a `DateType` with a
specified format */
.withColumn("Day", date_format(col("date"), "EE")) /* Converts a date/timestamp/ to a value of string in the
format specified by the date format given by the second argument. */
.groupBy("dispatching_base_number", "Day")
.agg(sum("trips").as("Total Trips"))
.orderBy(col("Total Trips"))
.show()
```

Scala Technical Questions

1. Difference between sc.textfile() and sc.wholeTextFiles()?

sc.textfile() → Each row will be a record

sc.wholeTextFiles() → read the entire content of a file as a single record

```
val readDataUsingTextFile = sc.textFile("file:///home/Practice/data/data2.txt")
val readDataUsingWholeTextFiles = sc.wholeTextFiles("file:///home/Practice/data/data2.txt")
val result = readDataUsingTextFile.map(x => (x,1))
val result2 = readDataUsingWholeTextFiles.map(x => (x,1))

println("==== result using text file =====")
result.take(2).foreach(println) // '1' is added after every row record
println("==== result using whole text file =====")
result2.take(2).foreach(println) // "1" is added after all content of file
-----
sc.textFile() result      ||      sc.wholeTextFiles() result
```

```
((file:/home/qualitest/Practice/data/data2.txt,ApacheJMeter  
KatalonStudio  
(ApacheJMeter,1)  
(KatalonStudio,1)  
AI  
,1)
```

2. Function return type is Unit, what does it mean?

Unit type is the same as void in JAVA, which means the function doesn't return anything

3. What is closure?

Closure is like a Function whose return type is defined on the expression. Closure can be defined in any scope (class, Object, function)

```
var a = (i:Int, j:Int) => i * j      /* return type is Int */  
var a2 = (i:Int, j:Int) => "Hello"    /* return type is String */  
var a3 = (i:Int, j:Int) => List("Hello", "world", "Ayushi")  /* return type is List */
```

4. What is String Interpolation?

- i. **The 's' String Interpolator**

The literal 's' allows the usage of variable directly in processing a string, when you prepend 's' to it

```
val name = "James"  
println(s "Hello, $name") //output: Hello, James
```

- ii. **The 'f' Interpolator**

The literal 'f' interpolator allows to create a formatted String, similar to printf in C language

```
val height = 1.9d  
val name = "James"  
println(f"$name%s is $height%2.2f meters tall") //James is 1.90 meters tall
```

- iii. **'raw' Interpolator**

The ‘raw’ interpolator is similar to ‘s’ interpolator except that it performs no escaping of literals within a string.

```
object Demo {  
    def main(args: Array[String]) {  
        println(raw"Result = \n a \n b")  
    }  
}
```

Output –

```
Result = \n a \n b
```

5. What is Pattern Matching?

It is the same as the Switch case in Java.

```
def matchTest(x: Int): String = x match {  
    case 1 => "one"  
    case 2 => "two"  
    case _ => "many" /* mend of pattern match*/  
  
    print( "this is done")  
    return String  
}
```

6. Exceptional Handling in Scala?

```
object Demo {  
    def main(args: Array[String]) {  
        try {  
            val f = new FileReader("input.txt")  
        } catch {  
            case ex: FileNotFoundException => {  
                println("Missing file exception")  
            }  
  
            case ex: IOException => {  
                println("IO Exception")  
            }  
        } finally {  
            println("Exiting finally...")  
        }  
    }  
}
```

7. What is a Match expression?

```
def main(args: Array[String]) {  
    val age = 100;  
  
    age match {  
        case 20 => println(age);  
        case 18 => println(age);  
        case 30 => println(age);  
        case 40 => println(age);  
        case 50 => println(age);  
  
        case _ => println("default");  
    }  
  
    val i = 7;  
    i match {  
        case 1 | 3 | 5 | 7 | 9 => println("odd");  
        case 2 | 4 | 6 | 8 | 10 => println("even");  
    }  
}
```

8. Basic ways to write functions in Scala?

```
def add(x: Int, y: Int): Int = {
    return x + y;
}

def subtract(x: Int, y: Int): Int = {
    x - y;
}

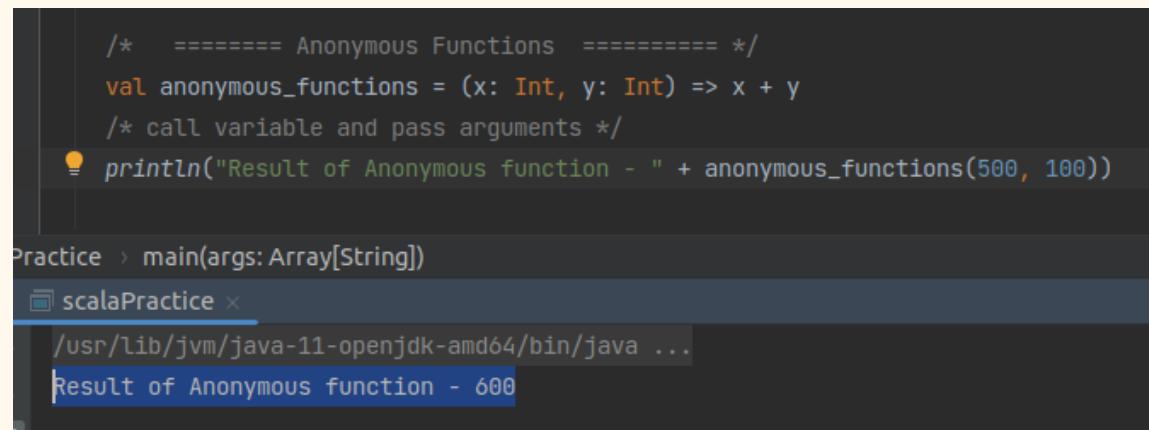
def multiply(x: Int, y: Int): Int = x * y;

def divide(x: Int, y: Int) = x / y;

def main(args: Array[String]) {
    println(add(45, 15));
}
```

9. What are Anonymous functions in Scala?

Without giving any name to a function we can assign it to a **Variable**, it's called Anonymous function.



The screenshot shows a Scala code editor with the following content:

```
/* ===== Anonymous Functions ===== */
val anonymous_functions = (x: Int, y: Int) => x + y
/* call variable and pass arguments */
println("Result of Anonymous function - " + anonymous_functions(500, 100))
```

Below the code, the terminal window shows the output:

```
Practice > main(args: Array[String])
scalaPractice x
/usr/lib/jvm/java-11-openjdk-amd64/bin/java ...
Result of Anonymous function - 600
```

10. What are Higher order functions?

Higher order functions are those functions which are able to **take functions as arguments** and are able to **return functions as result**.

```
object scalaPractice {  
    /* ===== higher-order Functions ===== */  
    def higherOrder_functions(x: Int, y: Int, f: (Double, Double) => Double): Double = {  
        f(x, y)  
    }  
  
    def main(args: Array[String]): Unit = {  
        /* ===== call higher-order Functions, result: 5000.0 ===== */  
        println("Result of Anonymous function - " + higherOrder_functions(50, 100, (x,y) => x * y))  
    }  
}
```

11. Define Array in scala?

```
/* 1st way using new Keyword */  
val arr1: Array[Int] = new Array[Int](5)  
arr1(0) = 10  
arr1(1) = 10  
arr1(2) = 10  
  
/* 2nd way to create array by direct initialize */  
val arr2: Array[Int] = Array(1, 2, 3, 4, 5)
```

12. How to concat two Arrays in Scala?

```
/* Concat two Arrays */  
import Array.concat  
val result = concat(arr1, arr2)  
}
```

13. How to define **For Loop** in Scala?

```

/* 1st way loop within range */
/* for(initialize <- range) */
for(i <- 1 ≤ to ≤ 5) {
    println("number:" + i)
}

/* another way to define loop using dot(.) */
for(i <- 1 ≤ .to( ≤ 5)) {
    println("number:" + i)
}

/* 2nd way loop until range, */
for(i <- 1 ≤ until < 5) {
    println("number:" + i)
}

/* 3rd way loop with multiple ranges, below variable j range will work as nested loop */
for(i <- 1 ≤ until < 5 ; j <- 1 ≤ to ≤ i) {
    println("number i:" + i + ", number j:" + j)
}

/* 4th way with List */
var list = List(1,2,3,4,5)
for(i <- list) {
    println("number:" + i)
}

/* 5th way with Array */
val arr = Array(10,20,30,40,50)
for (i <- 0 ≤ to ≤ (arr.length-1)) {
    println(arr(i))
}

/* 6th way only go inside loop if condition satisfy */
for (i <- range(1,5); if i > 2) {
    println("hello -> " + i)
}

```

14. Define List in Scala?

```

val mylist: List[Int] = List(1,2,5,8,9,6,4);
val names: List[String] = List("Max", "Tom", "John");

def main(args: Array[String]) {
    println(0 :: mylist);
    println(mylist);
    println(names);
    println(1 :: 5 :: 9 :: Nil);
    println(mylist.head);
    println(names.tail);
    println(mylist.tail);
    println(names.isEmpty);
    println(mylist.reverse);
    println(List.fill(5)(2));
    println(mylist.max);

    mylist.foreach( println );
    var sum : Int = 0;
    mylist.foreach(sum +=_);
    println(sum);

    for (name <- names) {
        println(name);
    }

    println(names(0));
}

```

<terminated> Demo\$ [Scala Application] C:\Program File

```

List(0, 1, 2, 5, 8, 9, 6, 4)
List(1, 2, 5, 8, 9, 6, 4)
List(Max, Tom, John)
List(1, 5, 9)
1
List(Tom, John)
List(2, 5, 8, 9, 6, 4)
false
List(4, 6, 9, 8, 5, 2, 1)
List(2, 2, 2, 2, 2)
9
1
2
5
8
9
6
4
35
Max
Tom
John

```

15. How to define sets in Scala?

```

// Scala - Sets
object Demo {
    val myset: Set[Int] = Set(1,2,5,8,9,6,4);
    val myset2: Set[Int] = Set(4,2,9,18,19,16,14);
    val names: Set[String] = Set("Max", "Tom", "John");
    def main(args: Array[String]) {
        println(myset + 10);
        println(names("Maxxxx"));
        println(myset.head);
        println(myset.tail);
        println(myset.isEmpty);

        println(myset ++ myset2);
        println(myset.++(myset2));

        println(myset.&(myset2));
        println(myset.intersect(myset2));
        println(myset.max);

        println(myset.min);
    }
}

```

<terminated> Demo\$ [Scala Application] C:\Program Files\Java\jdk8\bin\javav

```

Set(5, 10, 1, 6, 9, 2, 8, 4)
false
5
Set(1, 6, 9, 2, 8, 4)
false
Set(5, 14, 1, 6, 9, 2, 18, 16, 8, 19, 4)
Set(5, 14, 1, 6, 9, 2, 18, 16, 8, 19, 4)
Set(9, 2, 4)
Set(9, 2, 4)
9
1

```

16. Map in scala?

```

/* ===== Map ===== */
val map1: Map[Int, String] = Map(1 -> "Ayushi", 2 -> "Rubal", 3 -> "xyz 123")

println(map1) /* print Map */
println(map1(1)) /* get value of '1' key */
println(map1.keys) /* get Map Keys */
println(map1.values) /* get Map values */
println(map1.isEmpty) /* check if Map is empty */
println(map1(50)) /* get value of invalid key will throw exception */

object scalaPractice {
    def main(args: Array[String]): Unit = {
        val map1: Map[Int, String] = Map(1 -> "Ayushi", 2 -> "Rubal", 3 -> "xyz 123")
        println(map1)
        println(map1(1))
        println(map1.keys)
        println(map1.values)
        println(map1.isEmpty)
        println(map1(50))
    }
}

```

17. Tuples in scala?

Tuples can contain elements of different/heterogeneous data-types. Tuples are immutable and they are fixed size.

```

/* ===== Tuples ===== */
val tuple = (1, "hello", 12.5, 'a', false, (10, "world"))

/* Access tuple specific element */
println("Tuple 1st element- " + tuple._1)
println("Tuple 2nd element- " + tuple._2)
/* Get nested tuple element */
println("Tuple 6th element- " + tuple._6._1)

/* Iterate Tuple */
tuple.productIterator.foreach(x => println("element value- " + x))

```

There is another way to define Tuples using new keyword, depend on size of tuple elements us that Tuple, for example if Tuple have 2 elements use Tuple2 -

```

/* Below tuple have 2 elements, so use "new Tuple2" */
val tuple1 = new Tuple2(1, "hello")

/* Below tuple have 5 elements, so use "new Tuple2" */
val tuple2 = new Tuple5(1, "hello", 12.50, 'a', false)
println(tuple1)

```

18. Reduce, fold and scan methods?

The screenshot shows a Scala IDE interface. On the left, the code editor displays `Demo.scala` with the following content:

```

// Scala - Reduce, fold or scan (Left/Right)
//reduceLeft, reduceRight, foldLeft, foldRight, ...
object Demo {
  val lst = List(1, 2, 3, 5, 7, 10, 13);
  val lst2 = List("A", "B", "C");
  def main(args: Array[String]) {
    println(lst.reduceLeft(_ +_));
    println(lst2.reduceLeft(_ +_));
    println(lst.reduceLeft((x,y) => {println(x + y); y}));

    println(lst.reduceLeft(_ -_));
    println(lst.reduceRight(_ -_));
    println(lst.reduceRight((x,y) => {println(x - y); y}));

    println(lst.foldLeft(100)(_ +_));
    println(lst2.foldLeft("z")(_ +_));
    println(lst.scanLeft(100)(_ +_));
    println(lst2.scanLeft("z")(_ +_));
  }
}

```

On the right, the console window shows the output of the code execution:

```

41
ABC
1 , 2
3 , 3
6 , 5
11 , 7
18 , 10
28 , 13
41
-39
7
10 , 13
7 , -3
5 , 10
3 , -5
2 , 8
1 , -6
7
141
zABC
List(100, 101, 103, 106, 111, 118, 128, 141)
List(z, zA, zAB, zABC)

```

19. Scala classes and constructors ?

Singleton class - If define a class with “object” keyword, it means the class is a singleton and can’t create an instance of this class.

```

package pack

/* This is singleton class, declare with "object" keyword */
object scalaPractice {

  def main(args: Array[String]): Unit = {
  }
}

```

Class -

```
/* Normal class without constructor */
class test1 {
}
```

20. Primary constructor and Auxiliary Constructors?

Primary constructor ->

```
/* class with primary constructor with private members */
class test2(private var id: Int, private val name: String) {
}

/* class with constructor when there is no keyword like val or var in-front of it
 * then those members can't be access outside/inside of the class.
 * Just initialize once creating object of this class, like --> new test3(12, "Ayushi") */
class test3(id: Int, name: String) {
}
```

Auxiliary constructor ->

- **Auxiliary** constructors are the alternative constructors for a class and are defined as methods in the class with the name "this".
- The auxiliary constructor in Scala is used for constructor overloading.
- A class can have many auxiliary constructors but they should have different signatures and must call any previously defined constructors.

```
/* class with Auxiliary constructor --> declare with this() method */
class test4(private var id: Int, private val name: String) {
    /* Auxiliary constructor with No Arguments */
    def this() {
        /* Call previous or primary constructor*/
        this(123, "Ayushi")
    }

    /* Auxiliary constructor with Arguments */
    def this(age: Int) {
        /* Call previous constructor*/
        this()
    }
}

object testObj {
    def main(args:Array[String]): Unit = {
        /* Initialize class with constructors */
        val obj1 = new test4()
        val obj2 = new test4( id = 123, name = "Ayushi")
        val obj3 = new test4( age = 20)
    }
}
```

21. Inheritance in Scala?

Scala doesn't support multiple inheritances just like java

Super class

```

/* ===== Super class ===== */
class A {
    def printMethod(): Unit = {
        println("This is super class method")
    }
}

object testInheritance {
    def main(args: Array[String]): Unit = {
        val obj_B = new B()
        val obj_C = new C( text= "Hello World!!")

        obj_B.printMethod()
        obj_C.printMethod()
    }
}

```

Subclass - Use “override” keyword to override the parent method.

```

/* Subclass inherit parent class "A" */
class B extends A {
}

```

```

/* Subclass inherit parent class "A" and overriding parent class methods */
class C (val text:String) extends A {
    /* Override parent method */
    override def printMethod(): Unit = println("This is override in subclass-" + text)
}

```

22. Abstract class in Scala?

An abstract can be initialized. Abstract methods should be implemented in the subclass.

```

val obj_1 = new A()
val obj_B = new B
val obj_C = new C

```

Class 'A' is abstract; cannot be instantiated

Super class

```
/* ===== Super Abstract class ===== */
abstract class A {
    def printMethod(): Unit = {
        println("This is super class method")
    }

    /* 1st abstract method */
    def abstractMethod(): Unit

    /* 2nd abstract method */
    def abstractMethod2(): String
}
```

Subclass -

```
/* Subclass inherit parent class "A" */
class B extends A {
    override def abstractMethod(): Unit = ???

    override def abstractMethod2(): String = ???
}
```

23. Scala lazy Evolution?

Scala supports strict evolution and lazy evolution

```
//Scala: Lazy Evaluation

class strict {
  val e = {
    println("strict");
    9
  }
}
class LazyEval {
  lazy val l = {
    println("lazy");
    9
  }
}

object Demo {
  def main(args: Array[String]) {
    val x = new strict;
    val y = new LazyEval;

    println(x.e);
    println(y.l);
  }
}
```

//Scala: Lazy Evaluation

```
object Demo {  
  def method1(n: Int) {  
    println("Method 1");  
    println(n);  
  }  
  def method2(n: => Int) {  
    println("Method 2");  
    println(n);  
  }  
  def main(args: Array[String]) {  
    val add = (a : Int, b: Int) => {  
      println("Add");  
      a + b  
    }  
  
    method1(add(5,6));  
  
    method2(add(5,6));  
  }  
}
```

```
<terminated> Der  
Add  
Method 1  
11  
Method 2  
Add  
11
```

24. Scala Trait?

Traits are like interfaces, they are partially implemented interfaces. We can inherit Trait using the “with” keyword just like “implements” keyword in java.

```
/* Trait example, it can have both type of methods - with definition or without definition */  
trait trait1 {  
  def traitPrintMethod(): Unit = {  
    println("This is super class method")  
  }  
  
  /* without definition method */  
  def traitMethod(): Unit  
}  
  
trait trait2 {  
  def trait2PrintMethod(): Unit = {  
    println("This is super class method")  
  }  
}
```

Just add “with” to inherit as many Traits you want, just like below we have added “with” two times to inheriting each Trait.

```
/* Subclass inherit parent class "A" */
class B extends A with trait1 with trait2 {
    def abstractMethod(): Unit = ???

    def abstractMethod2(): String = ???

    def traitMethod(): Unit = ???
}
```

- 25.
 - 26.
 - 27.
 - 28.
 - 29.
 - 30.
-

INTERVIEW QUESTIONS -

1. What is Hadoop?

Ans - [Hadoop](#)

Hadoop is a framework that uses distributed storage and parallel processing to store and manage big data. It is the software most used by data analysts to handle big data, and its market size continues to grow.

2. What is the mechanism for distributed Storage and distributed processing?

Ans -[distributed Storage](#) → DFS and [distributed processing](#) → MapReduce

3. What is the block size in 1.0 and 2.0?

Ans → Block size in Hadoop 1.0 ⇒ 64 MB

Block size in Hadoop 2.0 ⇒ 120 MB

4. What are name nodes, edge Node, secondary name, and data node?

Ans → name node ⇒ Master Node

edge Node ⇒ Client Node

secondary name ⇒ Where we do checkpointing

data node ⇒ Where store our data

5. What is the heartbeat interval? Ans - 30 sec

6. What is stale/grace time? Ans - 10 min

7. What is under replication and over replication?

Ans - When data nodes have more than 3 replicas of the same data it is called

Over replication and when data nodes have less than 3 replicas of the same data
it is called under replica.

8. Do name nodes do checking points? Ans - No, but once it goes to safemode and becomes a secondary node it will do checkpointing.

9. What is safe mode? Ans - Read [Hadoop 2.0 Architecture](#)

10. Difference between edit log and FS Image?

Ans - Edit log have metadata of current Transaction while FS Image has metadata of all Transaction.

11. Drawbacks of Hadoop 1.0? Ans - [Defects of Hadoop 1.0](#)

12. Who monitors nodes in 2.0? Ans - ZooKeeper

13. What are journal nodes? Ans - We write EL in journal node

14. What is the split-brain scenario? Ans - Read [Hadoop 2.0 Architecture](#)

15. What is the in_use.lock file? Ans - Read [Hadoop 2.0 Architecture](#)

16. What is fencing? Ans - Read [Hadoop 2.0 Architecture](#)

17. Insert 10 records in the RDBMS table and then Sqoop imports only Chennai records with columns only firstname, lastname?

----- 1st Way using Where -----

```
sqoop import --connect jdbc:mysql://localhost/sqlDB --username root --password cloudera --table customer --m 1  
--where "city='chennai'" --columns firstname,lastname --delete-target-dir --target-dir /user/cloudera/readdir
```

----- 2nd Way using Query -----

```
sqoop import --connect jdbc:mysql://localhost/sqlDB --username root --password cloudera
```

```
--query "SELECT firstname, lastname from customer where city='chennai' and \$CONDITIONS" --m 1 --target-dir /user/cloudera/revdir
```

18. Some questions in [Tasks/Interview questions](#)

19. What do you mean by the term or concept of Big Data?

Ans - [What is Big Data?](#)

20. What are the characteristics of Big Data?

Ans - [5 V's - \(now its 6 V's\)](#)

21. Deploying a Big Data solution includes the following steps?

- **Data Ingestion:** As a first step, the data is drawn out or extracted from various sources so as to feed it to the system.
- **Data Storage:** Once data ingestion is completed, the data is stored in either HDFS or NoSQL database.
- **Data Processing:** In the final step, the data is processed through frameworks and tools such as Spark, MapReduce, Pig, etc.

22. What is the reason behind using Hadoop in Big Data analytics?

Ans - Businesses generate a lot of data in a single day and the data generated is unstructured in nature. Data analysis with unstructured data is difficult as it renders traditional big data solutions ineffective. Hadoop comes into the picture when the data is complex, large, and especially unstructured. Hadoop is important in Big Data analytics because of its characteristics:

- Data storage
- Data processing
- Collection plus the extraction of data

23. What do you understand by fsck in Hadoop?

fsck stands for **file system check** in Hadoop, and is a command that is used in HDFS. fsck checks any and all data inconsistencies. If the command detects any inconsistency, HDFS is notified regarding the same.

24. Can you explain some of the important features of Hadoop?

Some of the important features of Hadoop are:

- **Fault Tolerance:** Hadoop has a high-level of fault tolerance. To tackle faults, Hadoop, by default, creates three replicas for each block at different nodes. This number can be modified as per the requirements. This helps to recover the data from another node if one node has failed. Hadoop also facilitates the automatic recovery of data and node detection.
- **Open Source**

- **Distributed Processing:** Hadoop stores the data in a distributed manner in HDFS. Distributed processing implies fast data processing. Hadoop also uses MapReduce for the parallel processing of the data.
- **Reliability:** One of the benefits of Hadoop is that the data stored in Hadoop is not affected by any kind of machine failure, which makes Hadoop a reliable tool.
- **Scalability:** Scalability is another important feature of Hadoop. Hadoop's compatibility with other hardware makes it a preferred tool. You can also easily add new hardware to the nodes in Hadoop.
- **High Availability:** Easy access to the data stored in Hadoop makes it a highly preferred Big Data management solution. Not only this, the data stored in Hadoop can be accessed even if there is a hardware failure as it can be accessed from a different path.

25. What is Hadoop and what are its components?

Apache Hadoop is the solution for dealing with Big Data. Hadoop is an open-source framework that offers several tools and services to store, manage, process, and analyze Big Data. This allows organizations to make significant business decisions in an effective and efficient manner, which was not possible with traditional methods and systems.

Components of Hadoop

26. Explain Hadoop Architecture. [Hadoop 2.0 Architecture](#)

27. In what all modes can Hadoop be run?

ANS - Hadoop can be run in three modes:

- **Standalone Mode:** The default mode of Hadoop, standalone mode uses a local file system for input and output operations. This mode is mainly used for debugging purposes, and it does not support the use of HDFS. Further, in this mode, there is no custom configuration required for mapred-site.xml, core-site.xml, and hdfs-site.xml files. This mode works much faster when compared to other modes.
- **Pseudo-distributed Mode (Single-node Cluster):** In the case of pseudo-distributed mode, you need the configuration for all the three files mentioned above. All daemons are running on one node; thus, both master and slave nodes are the same.
- **Fully distributed mode (Multi-node Cluster):** This is the production phase of Hadoop, what it is known for, where data is used and distributed across several nodes on a Hadoop cluster. Separate nodes are allotted as master and slave nodes.

28. What is Hadoop federation?

29. What is the hive file format? Explain the Use case. (for example And why are you using Parquet, why not ORC?)

Ans - [Serialization File Types in SQOOP](#)

Hive File Formats

30. Difference between Parquet and AVRO file format?

Avro Data file -

Features -

- It can handle the Schema evolution, like delete a column, add new column--- It can map data to right columns
- 40-50% Compression
- Row format storage

PARQUET -

- i. Reduce compression from 60-80%
- ii. It saves the file data in **columnar file format**
- iii. **Faster query file format**
- iv. Parquet uses Predicate pushdown mechanism

31. how to put data from hive to Hadoop?

Ans - Data in Hive tables reside on HDFS, Hive is only a metadata layer that helps you access the HDFS data in form of tables and rows.

So you don't transfer data from Hive to HDFS, as the data is already on HDFS.

32. Partition vs Bucketing? **Ans** - Partition Vs Bucketing

33. Default Partition Algorithm used in MapReduce and Spark?

Or

Hive Bucket Internal formula?

Ans - How Bucket decides which part file records will go

34. How can I execute external files inside the hive?

35. Window function in Hive.

36. How can you transfer data from Hive to HDFS?

```
insert overwrite directory '/' select * from emp;
```

37. How does MapReduce work?

38. What is the difference between Tez and map-reduce?

39. What are the various steps involved in deploying a Big Data solution?

40. Spark Structure (How the job is executed behind the scene.

41. Explain Spark Architecture?

42. Repartition and Coalesce in SPARK in respective to Performance?

43. Transformation in Spark. data like -

- 1, Ram, 20000, "Ayodhya, Srilanka, India, Varanasi", Yes
- 2, Shyan 30000, "Bangalore, China, Canada, Hyderabad", No

file name test.csv

Write into hive table as parquet file

44. Broadcast variable

45. Lazy evaluation and shuffling

46. Write Scala code for the following scenario.

- a. JSON file HDFS location - /user/warehouse/input/employee.json
- b. Columns in Json file -id, name, salary, deptname
- c. Dataframe API to get the second highest salary.
- d. Save output dataframe in Parquet format in HDFS - /user/warehouse/output

47. Difference between Rank & Dense_rank?

48. What is scala file format?

[JSON/TEXT/CSV/XML/Avro/parquet](#)

Scala by default support JSON/text/parquet/CSV but for XML/Avro need to add jars or dependency for support these files

49. Write a word count program in spark.

```

import org.apache.spark.SparkContext
import org.apache.spark.SparkContext._
import org.apache.spark._

object SparkWordCount {
  def main(args: Array[String]) {

    /* "setAppName()" -> This is the name of the application that you want to run.
       "setMaster()" --> This parameter denotes the master URL to connect the spark application to.
    */
    val conf = new SparkConf().setAppName("first").setMaster("local[*]")

    /* Creating a [RDD] Spark Context Object */
    val sc = new SparkContext(conf)

    /* Read input file and store in variable */
    val input = sc.textFile("input.txt")
    val count = input.flatMap(line => line.split(" "))
      .map(word => (word, 1))
      .reduceByKey(_ + _)
    count.saveAsTextFile("outfile")
  }
}

```

50. Which technique did you use to display your code?

51. Write a Query to bind duplicate numbers?

52. What are nil and null?

53. What is transformation, action?

54. What is the access modifier in scala?

There is no public modifier, if no modifier is defined it is public by default.

Protected can we access either in subclass or sub-package.

Modifier	Class	Companion	Subclass	Package	World
No Modifier/Public	Yes	Yes	Yes	Yes	Yes
Protected	Yes	Yes	Yes	No *	No
Private	Yes	Yes	No	No *	No

55. how to debug in Hadoop?

56. scala version? spark version?

- Scala - 2.11.12
- Spark - 2.4.7

57. What package is used in scala?

- Org.apache.spark.Sparkcontext
- org.apache.spark.SparkSession
- org.apache.spark.sql

58. Once the file is loaded from HDFS to Hive, Why is the file moved completely to the Hive location(/user/hive/warehouse/)?

Ans -Although there is no documentation of why it happens. But I think HDFS shouldn't keep copies of the same data in two places(HDFS and HIVE)

59. What will happen if we don't set “**var conf = new**

SparkConf().setAppName("").setMaster("") to “**var sc = SparkContext()**”?

```
Exception in thread "main" org.apache.spark.SparkException Create breakpoint : A master URL must be set in your configuration
```

60. Difference Between var and val in Scala?

var	val
var varName="hello World"	val sName = "tutorialspoint.com"
Mutable	Immutable
Can be assigned multiple times.	Cannot be assigned multiple times.
Can be reassigned	Cannot be reassigned

61. In which Spark version was SparkSession introduced?

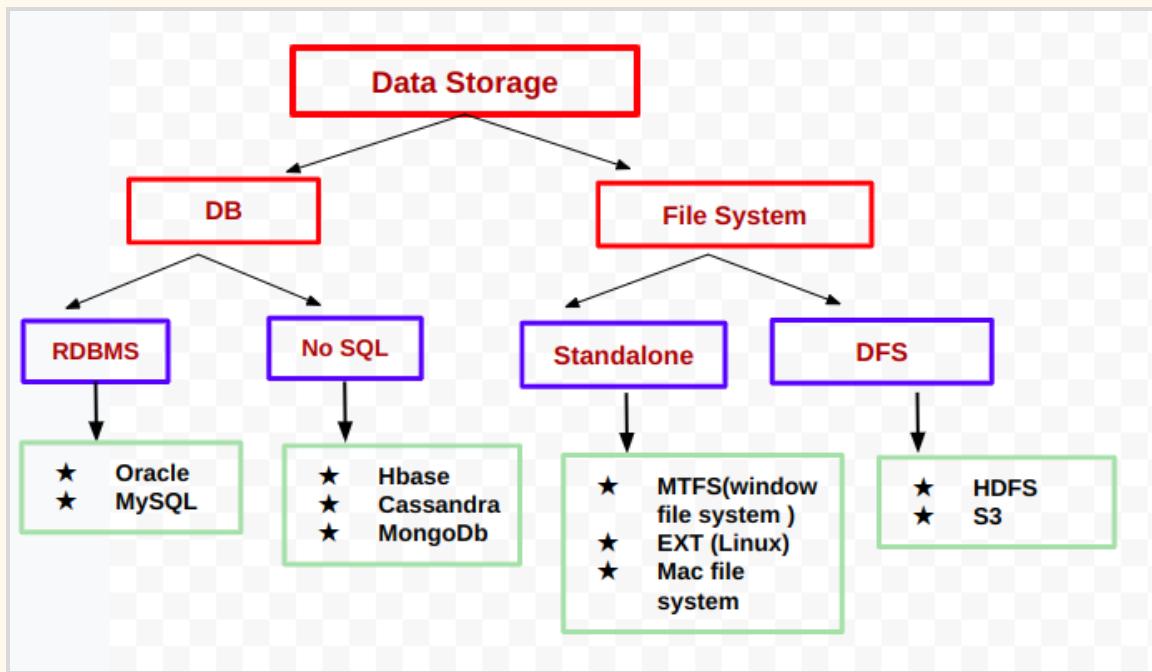
Ans - SparkSession was introduced in version Spark 2.0

62. What is Batch processing and stream processing?

Batch Processing - Your input data is stored in a file or database on which you are doing some processing and getting output. Using Hadoop, we can do batch processing

Stream Processing - Data coming live means not stored anywhere in the system and you are processing live data and getting the output, which is called Stream Processing. Using Spark, we can do Stream processing

63. How to migrate or take data from one place to another place in Spark?



64.

65.

66.

67.

68.

69.

70.