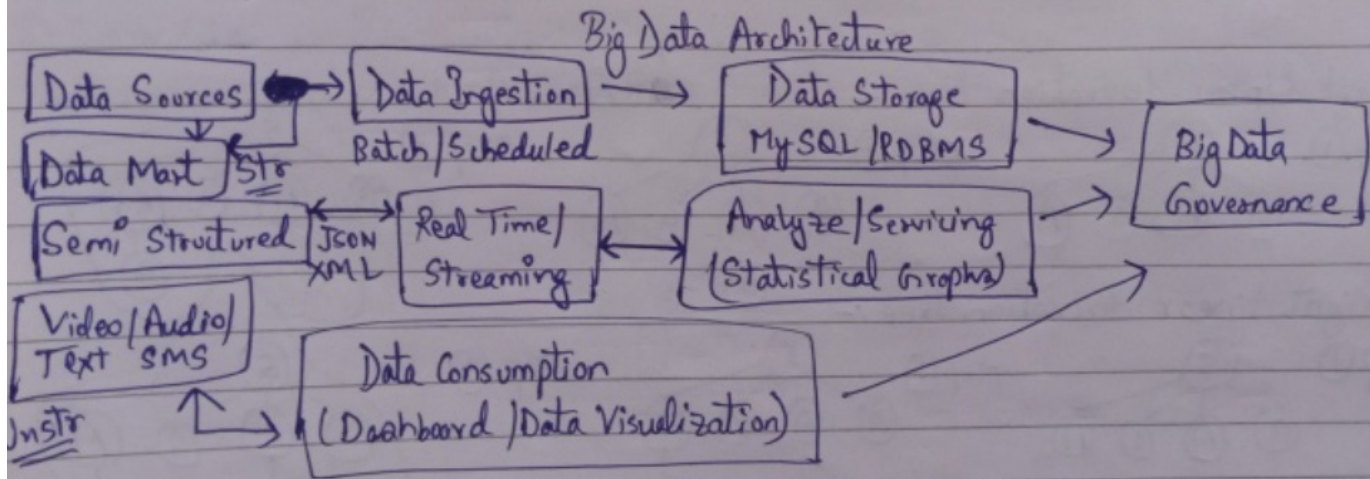


Left Linear

BD

Characteristics of Big Data [5 V's] :-

- ① ~~Velocity~~ :- Data is uncertain. Entries incomplete. NULL values possible.
Variety :-
- ③ Variety :- Data stored from different sources, eg: video, picture, voice.
- ② Volume :- Large no. of files.
- ④ Velocity :- Previous years, exponential increase in collection of data.
- ⑤ Value :- Useful data to make some analysis, valuable data.



BD

4 Components of Big Data :- / Flow of Data :-

- ① Capturing the data (Raw collection of facts or figure)
- ② Storing the data (Structured, Semi-structured or unstructured)
- ③ Data Ingestion / Processing / Analysis (Outliers, error data checking)
- ④ Data Visualization

⑤ Fields of Big Data :-

- ① Predictive Analysis
- ② Machine Learning / Artificial Intelligence { Pattern Recognition }
- ③ Natural Language Processing [NLP].
- ④ Computer Vision (AI)

⑥ Technologies related to Big Data :-

- ① Apache HADOOP :- It is open source framework. used in distributed computing of large data set.
- ② Apache SPARK :- used along with HADOOP, open source processing engine.
- ③ Apache FLINK :- open source stream processing framework, uses friendly API, livestreaming of real time data.
- ④ PRESTO :- Open source ^{supports interactive} SEQUEL engine, analysis of huge data set. Different systems are there for different analysis.
- ⑤ DRUID :- Open source Analytical data storage designed for queries on event based data. eg: log files.

(Stored in accumulator) (Accumulator - H)

BD

⑧ Features of Big Data :-

- ① Data Preparation :- used before ^{data} iterating ^{of} the model, capturing the data construction of model used at
- ② Data Exploration :- visualize the insights like pictorial ^{graphical} representation of data.
- ③ Scalability :- uses less network gear, use less energy (memory + processing should be fast) performance not degraded.
- ④ Supports various types of Analytics :- Numerical values,

IDS: Intrusion Detection System } use of certain firewalls.
for security purpose.

statistical values by pictorial representation. (pie chart, graphs, tabular charts) Categorical analysis.

⑤ Version Control :- particular version gets enhanced.

⑥ Data Management :- Capture, Storing in cost effective / space / Time.

⑦ Data Integration :- Collection of data and storing in database (integrating)

⑧ Data Governance :- Data should be accurate, available, reliable i.e. data should be easily governable.

⑨ Security :-

⑩ Visualization :- Data driven environment.

* Application of Big Data :-

→ Monitoring then Targeting :-

→ Recommendation of Application :-

→ Healthcare → Smartest Traffic System. →
(IOT) (Ola, Uber, Google Map)

PP

WAP to enter a number and check whether it is a perfect number.

```
print ("Enter a number")
num = int(input())
for i in range(1, num):
    if (num % i == 0):
        sum += i
if (sum == num):
    print ("Yes")
else:
    print ("No")
Q i = 1
while (i <= 5):
    print (i)
    i = i + 1
else:
    print ("Hello")
```


BD

Advantages of using Big Data in Auditing:-

- ① Reductions in operational cost :-
- ② Improved decision making :-
- ③ High Customer Retention :-
- ④ High Satisfaction Rate :-

Where are auditing used - banking, manufacturing, credit card, healthcare, weather prediction.

CD

Phases of a Compiler :-

→ No. of phases = 6, First 3 phases are called Analysis phase
Last 3 phases are called synthesis phase.

Phase

BD

UNIT :- 2 { HADOOP }

Hadoop provides access to the file system. Hadoop common package contains necessary JAR files and script.

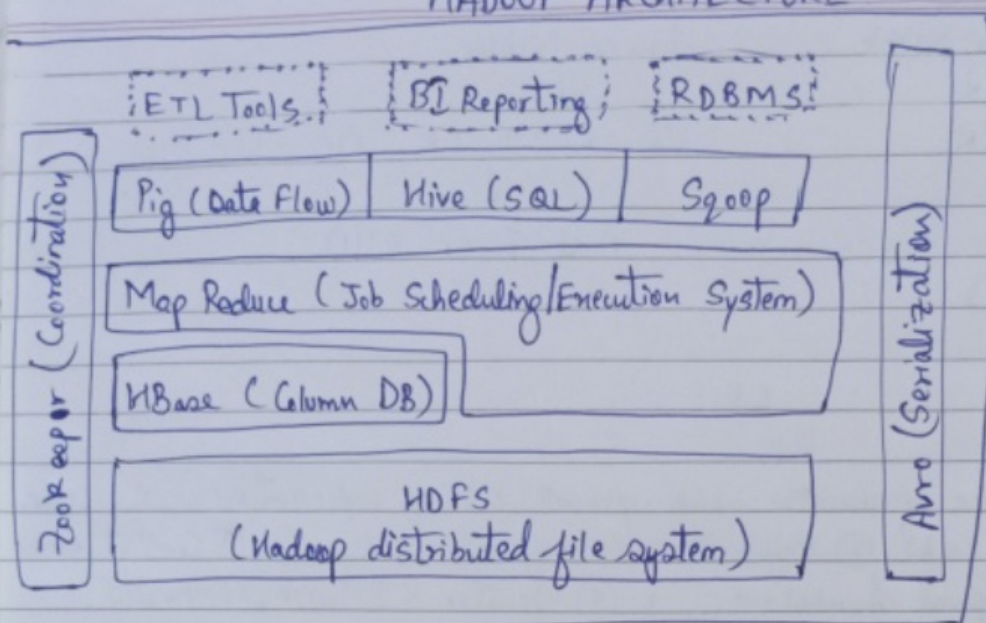
The package also provides source code, documentation and a contribution section that includes projects from the Hadoop Community.

Extract Transport Load (ETL)

Business Intelligence (BI)

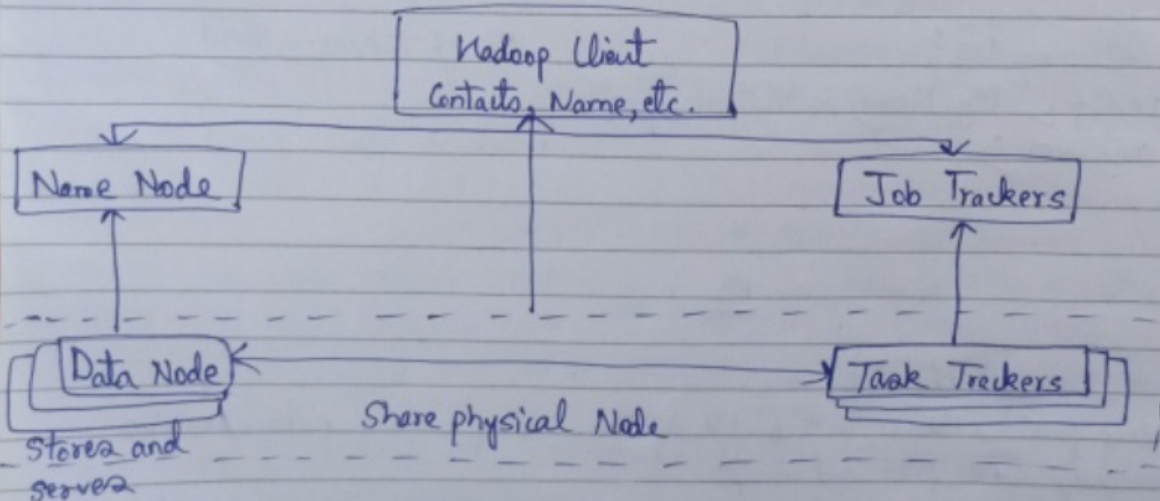
Command Line Interface
= Sqoop.

HADOOP ARCHITECTURE



Q What is HDFS? → Distributed file system
→ Traditional hierarchical file organization. → Single namespace for entire cluster. → Write once - read many access model → Aware of network topology

High Level Architecture

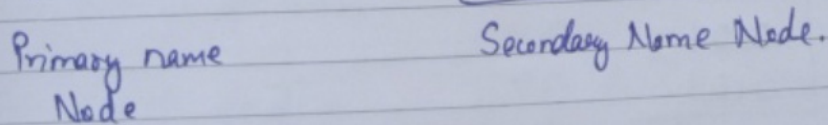


The diagram illustrates the data flow in a log-structured file system. It consists of two main components: a left box representing the main storage and a right box representing the log storage.

- Left Box:** Contains two sub-boxes labeled "RAM" and "HDD".
- Right Box:** Contains two sub-boxes labeled "RAM" and "HDD". Below the "HDD" sub-box, the text "Stores to HDD" is written.

Arrows indicate the direction of data flow:

- A top arrow labeled "pull transaction" points from the left "RAM" to the right "RAM".
- A middle arrow labeled "log" points from the left "RAM" to the right "RAM".
- A bottom arrow labeled "push" points from the right "HDD" back to the left "HDD".



MPJ

Logical Instructions

1) ANA R $A \leftarrow A \wedge R$ (Register)
AND operation

2) ANA M $A \leftarrow A \wedge [M]$ (Indirect)

AND Accumulator

AND 25H $A \leftarrow A \wedge 25H$ (Immediate)