

---

2018

---

---

---

---



# Big Data ?

Data which can not be handled by traditional Databases.

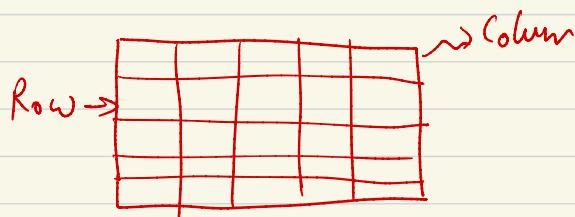
## Technical Definition

Any data related problem which satisfies 3 V's criteria.

① Volume: Huge amount of data  
(hundreds of GB's, TB, PB ..)

② Variety: Different format of data.

a.) Structured Data



b.) Semi-Structured Data

↳ Where schema is flexible

{ Ex: JSON, XML

{ "name": "Shashank",

  "age": 29

},

{ "name": "Rahul",

  "age": 29,

3 "Salary": 1000

### ④ Unstructured Data

↳ Doesn't have any fixed representation.

Ex: Image, audio, video files

### ③ Velocity: Speed of data generation

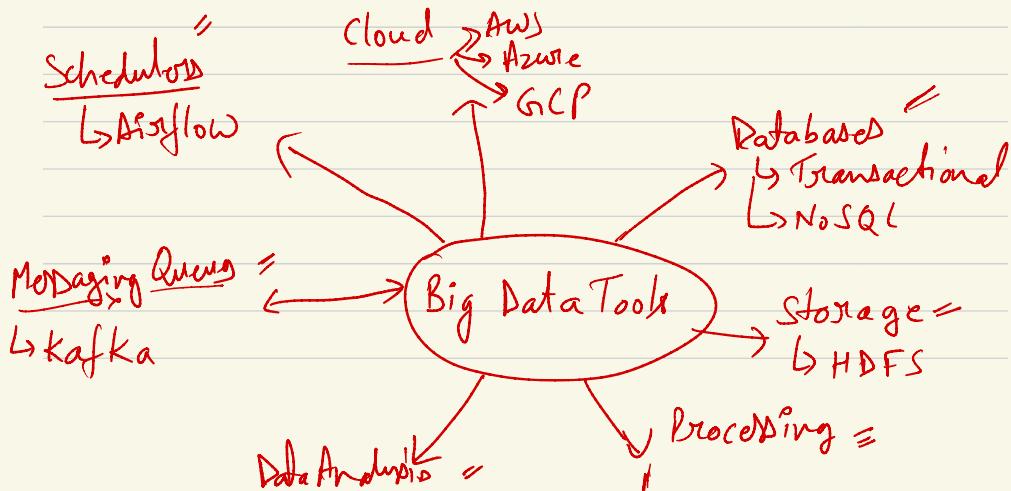
a) Batch Processing → Electricity Bill, Credit Card Bills

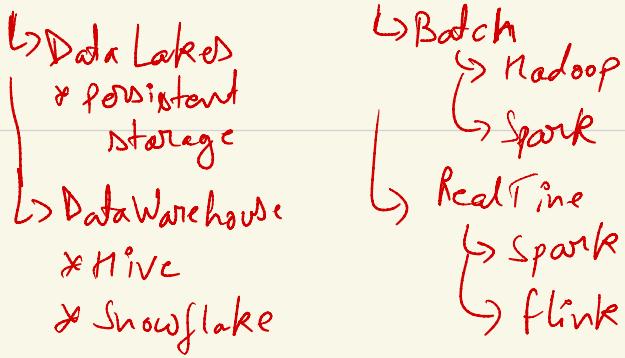
b) Real Time Processing

→ Twitter, Stock feed  
→ Streaming  
→ Live Gameplay

### ④ Value: Extracting meaningful information.

### ⑤ Voracity: Related to uncertainty in the Data.

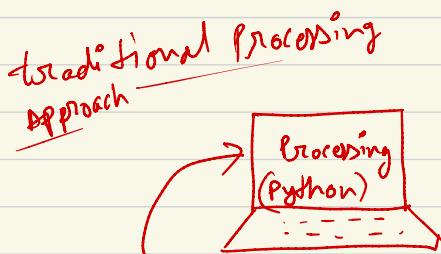




## ii) Hadoop

- ↳ Distributed Computation framework.
- ↳ Specially Designed for Batch Data Processing.

## What is distributed Computation?

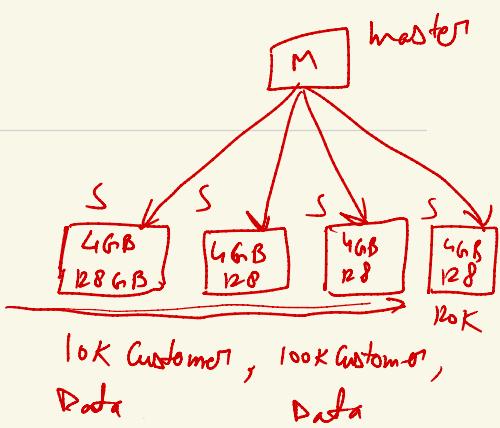
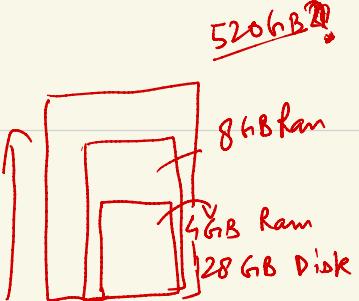


8GB RAM  
520GB External Disk

5GB Data to process \* I am sending Data to the Code ?

> Python procod.py

\* Difference Between Vertical & Horizontal Scaling!



## Vertical Scalability

- ↳ Adding Extra capacity in existing machine.

Horizontal Scaling

- ↳ Adding more machines in the system.

## \* Commodity Hardware

- ↳ Simple machine which has storage and processing capacity.

- ✓) Distributed Storage →

Break data into small pieces and store it on different machine

- ✓) Distributed Computation

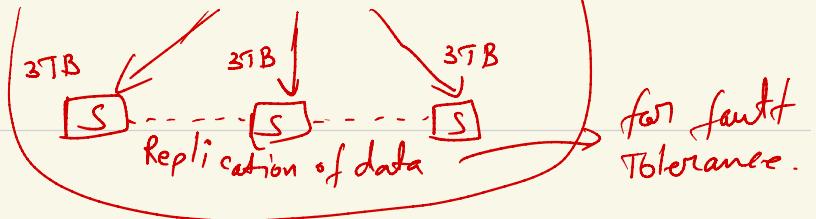
→ Process multiple parts of data on different machines at the same time.

## \* Distributed Storage

Master-Slave

ATB Data

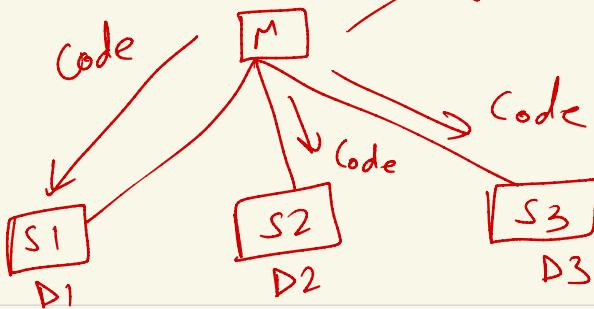




cluster: Combination of multiple machines.

## Distributed Computation

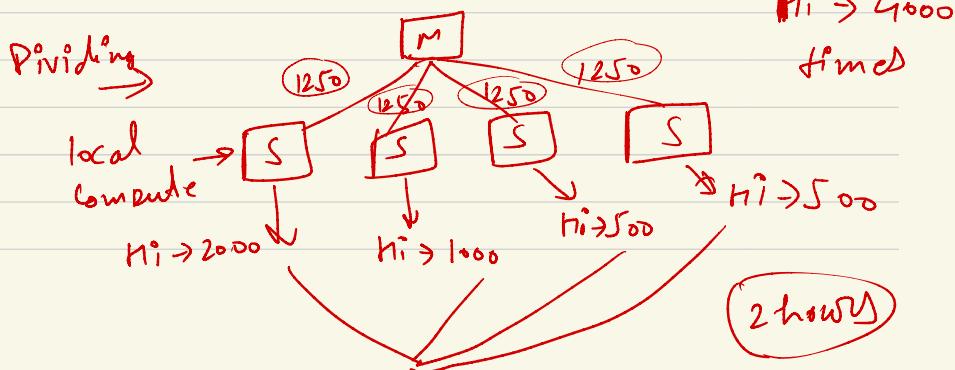
We are sending code towards the data.

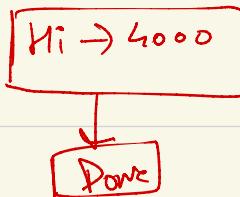


Q.) 5000 page of Book. How much time you need to count frequency of each word?

1 day

↳ You have 4 friends





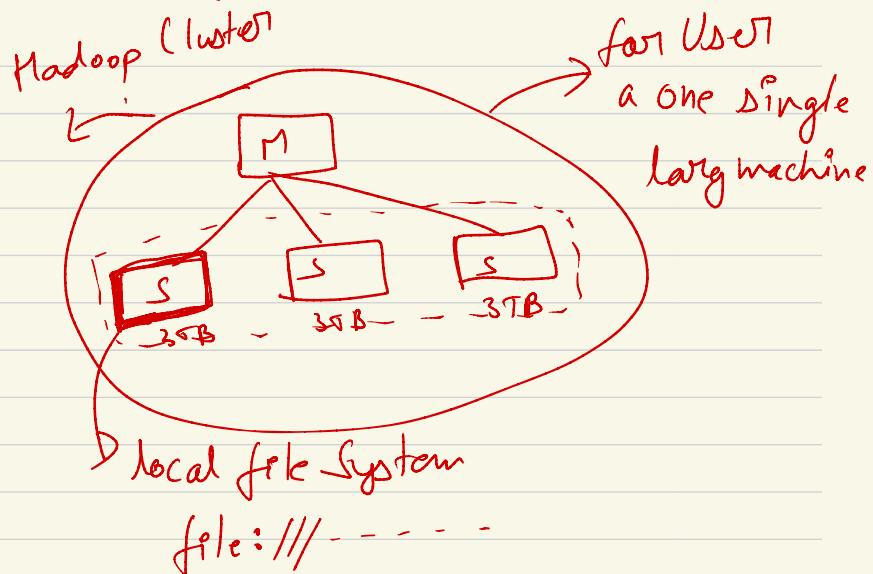
times

## Hadoop -

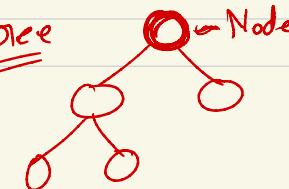
- ↳ HDFS (Storage)
- ↳ MAP-REDUCE (Processing)
- ↳ YARN (Resource Management)

## HDFS (Hadoop file System)

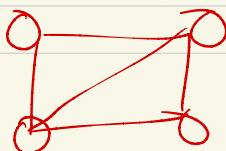
↳ Distributed file storage system.

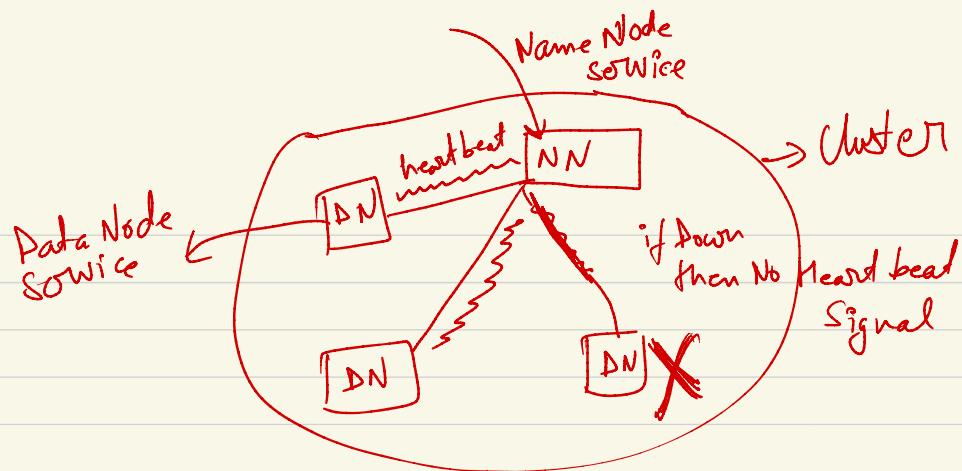
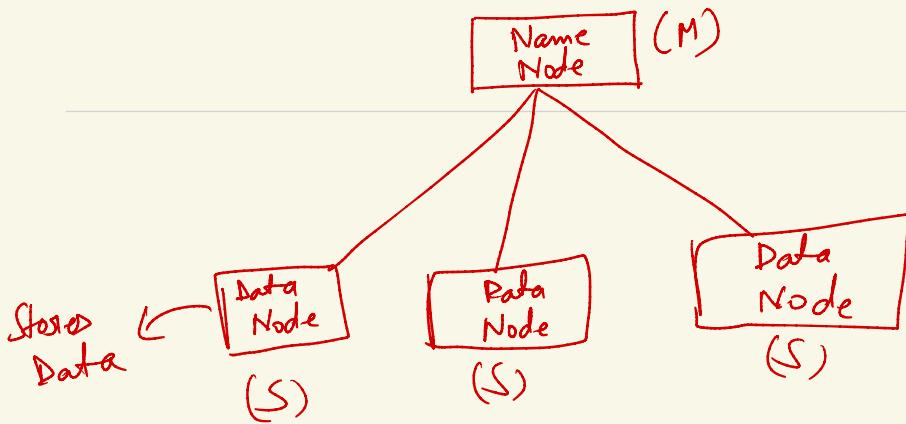


Tree



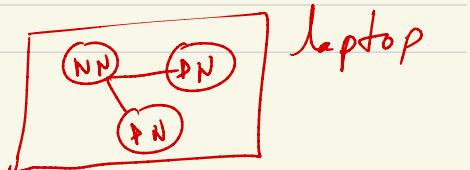
Graphs





Service → Set of Instructions / ~~Program~~  
which runs as a backend process.

→ Stand alone mode



## Data Node

- It is a daemon (background) process which runs on each slave node.
- The actual data is stored on Data Nodes.
- Data nodes will also help in computing.
- Data Nodes sends heartbeat to Name Node periodically to report the overall health of the system.

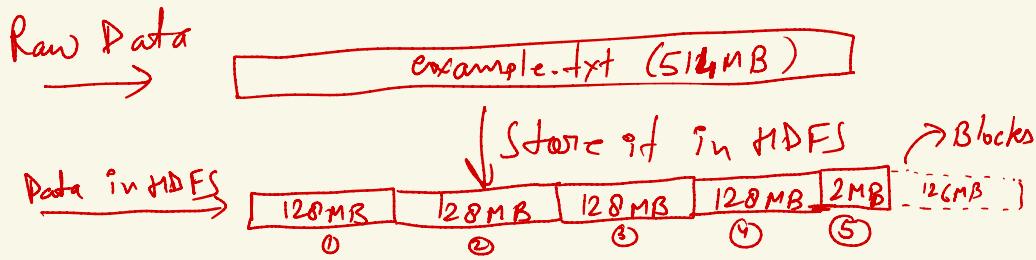
## Name Node functionalities

- It acts like a Master and manages Data Nodes.
- It records the metadata of all the files in the cluster ex: The location of the Blocks, the size of the files, hierarchy, where replicated data is stored.

\* What is Block?

Block is the smallest unit of physical memory where data is stored.

Ex If we want to store 514 MB data in HDFS, file example.txt



Hadoop 2.X → 128 MB

Hadoop 1.X → 64 MB } → Block Size (Predefined size)

Q.) file size is 514 MB

then how many blocks will be created inside HDFS?

$$514 / 128 = 4 \text{, remainder } 2 \text{ MB}$$

↓  
1

$$= 4 + 1 = 5$$

Q.) Store example.txt which is of 514 MB in HDFS

Calculation: Size of 1 Block = 128 MB

Total file size = 514 MB

Total block needed = 5

} 4 blocks can store = 128 MB  
1 block can = 2 MB

