

Big data.

Page No.

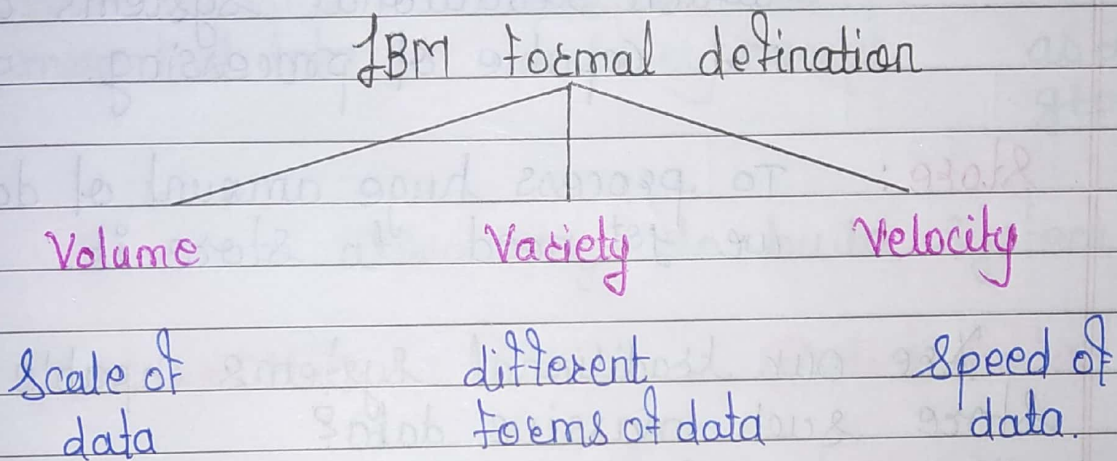
Date

What is big data?

"big data is a term that describes the large volume of data."

How to classify big data?

3V's of bigdata by IBM.



• Volume - 2.5 quintillion

• Variety:

① Structured data: RDBMS databases (Oracle & MySQL)

② Semi-structured data: CSV, XML, JSON.

③ Unstructured data: Audio, video, Image log files.

• velocity: 900 million photos on Facebook.
600 million tweets on twitter.

4th v of data:

Veracity : uncertainty of data
i.e poor quality data
unclean data.

Why big data?

process: To process huge amount of data which traditional systems are not capable of processing.

store: To process huge amount of data we 1st need to store it.

Are our traditional systems capable to store such massive data?

Big data System Requirements?

* Store huge amount of data :

• Traditional systems are NOT fit to store such huge amount of data.

• So store massive amount of data.

* process huge amount of data in a efficient and timely manner.

X traditional system are not capable to handle.

* scale easily to accomodate growing requirement

• traditional systems have serious limitations

store

Process

Scale

store massive amount of data

process it in a timely manner

scale easily as data grows.

Two ways to build system.

① Monolithic

② distributed

Monolithic: one powerful system with lot of resources.

distributed: many smaller systems come together.

Q. Monolithic or distributed?

① Monolithic is a single powerful server. Hard to do add resources after a certain limit.

② Resources:
① RAM - 8 GB (Memory)
② Hard disk - 1 TB (Storage)
③ CPU quad core (compute)

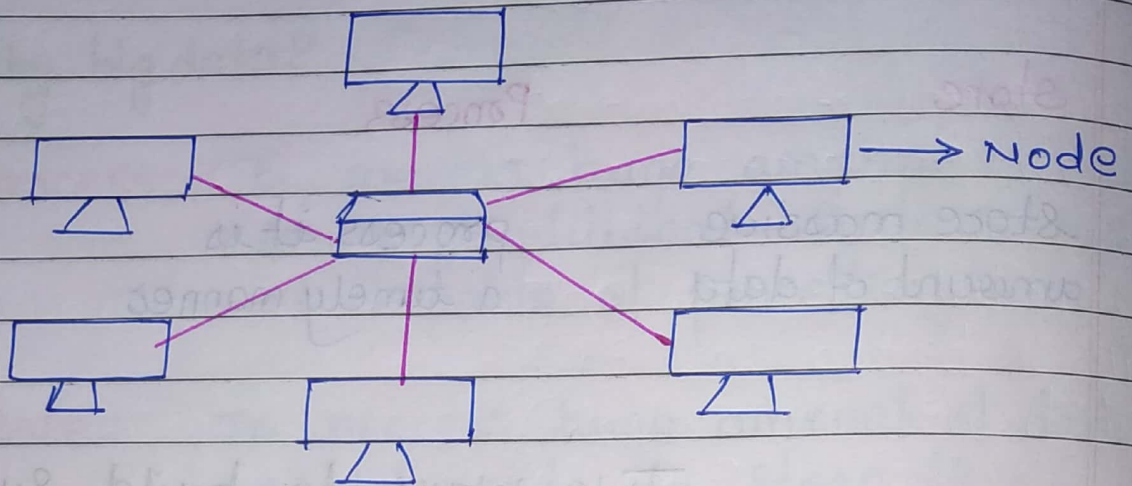
Monolithic

Is monolithic scalable?

No!

$2x$ resources $\neq 2x$ performance

Distributed:



6 Node cluster.

① The set of all nodes is known as cluster.

② Many small and cheap computers come together to act as a single entity.

Is distributed system scalable?

Yes, distributed systems are linearly scalable.

Distributed systems are scalable.

$2x$ resources = $2x$ performance speed.

Monolithic Architecture

↓
Vertically Scaling
(Not true Scaling)

Distributed Architecture

↓
Horizontal Scaling
(True Scaling)

So, Monolithic or distributed?

X Monolithic

distributed ✓

That's why all good big data systems are based on distributed architecture!

Hadoop

What is Hadoop?

Hadoop is a framework to solve big data problems.

Hadoop Evolution:

2003

Google released a paper to describe how to store large datasets.

2003

This paper was called as GFS (Google File System).

2004

Google released another paper to describe how to process large data.

2004

This paper was called as MapReduce.

2006

Yahoo took these papers & implemented it.

2006

The implementation of GFS was named as HDFS (Hadoop distributed file system).

The implementation of MapReduce was named as MapReduce (unchanged).

Hadoop Version ~~0.9~~ 1.0

HDFS - for distributed storage
MapReduce - for distributed processing.

2009

Hadoop comes under Apache Software Foundation & become open source.

2013

Apache released Hadoop 2.0 to provide major performance enhancement.

2003

google
GFS

2004

google
MR

2006

Yahoo
Implementation

2009

Hadoop
under
Apache

2013

Hadoop
2.0
released

Hadoop 1.0

MapReduce
HDFS

Hadoop 2.0

MapReduce YARN
HDFS