

HADOOP ASSIGNMENT

Ans 1. Characteristics of Big Data

(i) Volume:

Volume refers to the unimaginable amount of information generated every from social media, cell phones, cars, credit cards, M2M sensors, Images, video and ~~whatnot~~ etc.

(ii) Variety:

It refers to heterogeneous sources and the nature of data, both structured and unstructured. This variety of unstructured data poses certain issues for storage, mining and analyzing data.

(iii) Velocity:

The term 'velocity' refers to the speed of generation of data. How fast the data is generated and processed to meet the demands, determines real potential in the data.

(iv) Variability

This refers to the inconsistency which can be shown by the data at times, thus hampering the process of being able to handle and manage the data effectively.

Some of the techniques of data mining are used to analyze the big data such as clustering, prediction and classification and decision tree etc. Apache Hadoop, Apache Spark, Apache Storm,

MongoDB, NOSQL, HPC are the tools used to handle big data.

Ans 2. For many industries, big data isn't a choice but a naturally shaped reality, as the amount of structured and unstructured data is growing exponentially, along with a wide network of IoT devices capturing it.

The major business opportunities presented by big data for any industry include:

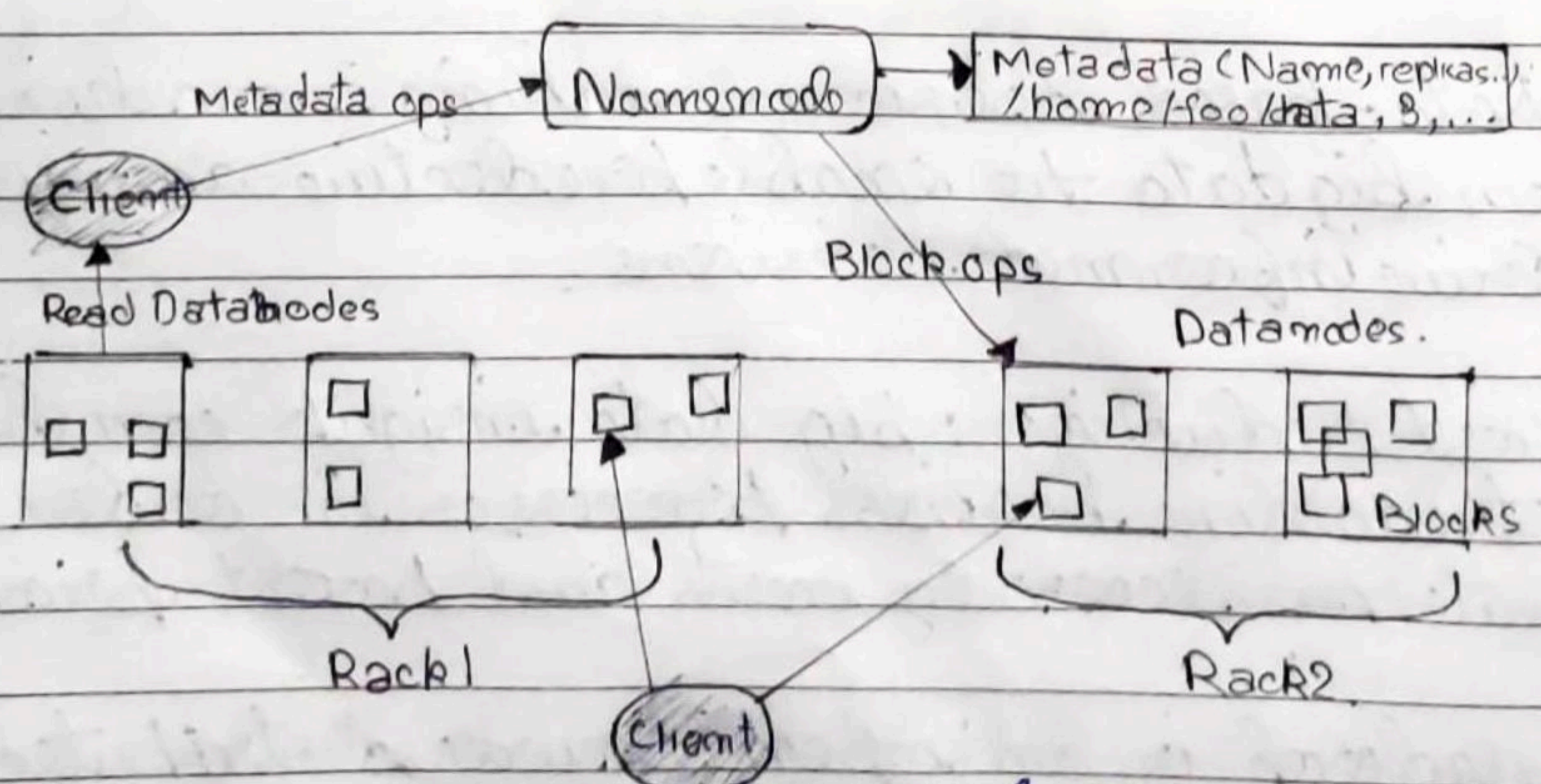
- **Automation:** data driven IT infrastructures allow businesses to automate time-consuming processes such as data collection and analysis.
- **Trends and insights:** big data reveals hidden opportunities and patterns that can be used to tailor products and services to end users' needs to increase operational efficiency.
- **Data-based decision-making:** machines learn on big data to enable predictive analysis and drive informed decisions.
- **Cost reduction:** big data insights can be used to streamline business processes in order to eliminate unnecessary costs and boost productivity.

Ans 3. Hadoop is an open source distributed processing framework that manages data processing and storage for big data applications in scalable clusters of computer servers.

Some myths about hadoop are:

- Hadoop is a database
- Hadoop is cheap.
- Hadoop requires MapReduce: Some users opt to deploy HDFS with Hive or HBase but not Map Reduce
- MapReduce only controls analytics: it also handles parallel programming, fault tolerance of wide variety of coded logics.
- Hadoop is too risky for enterprise use.
- Hadoop is a complete, single product: Hadoop is an ecosystem, a family of open source products.

Ans 4 HDFS is the storage system of Hadoop framework. It is a distributed file system that can conveniently run on commodity hardware for processing unstructured data.



Hadoop HDFS Architecture follows a Master/Slave Architecture where a cluster comprises of a single **NameNode (Master node)** and all the other nodes are

Data Nodes (Slave nodes). HDFS can be deployed on a broad spectrum of machines that support Java.

- Ans 5 Hadoop ecosystem includes multiple components that support each stage of Big Data processing.
- flume and Sqoop ingest data, #
 - spark and MapReduce process data.
 - Pig Hive and Impala analyze data.
 - Hue and Cloudera Search help to explore data.
 - Oozie manages the work flow of hadoop jobs.