

15/3/23

Hadoop

- It is an open source software framework used for storing data and running applications on a group of commodity hardware.
- open source.
- Distributed data processing
- map - Reduce programming.

Why hadoop

- handle any data type
 - structured / unstructured data.
- schema, schema less
 - ↳ transaction → job.

- high volume / low volume of data.
- All kind of analytic application
- * proven with petabyte scale.
- * Grows with business
- * capacity and performance grows.

features

- No vendor lock in
- Rich eco system and community development
- More affordable
- Scalable

- cost effective
- flexible
- fast
- Resilience to failure.

Hadoop Infrastructure

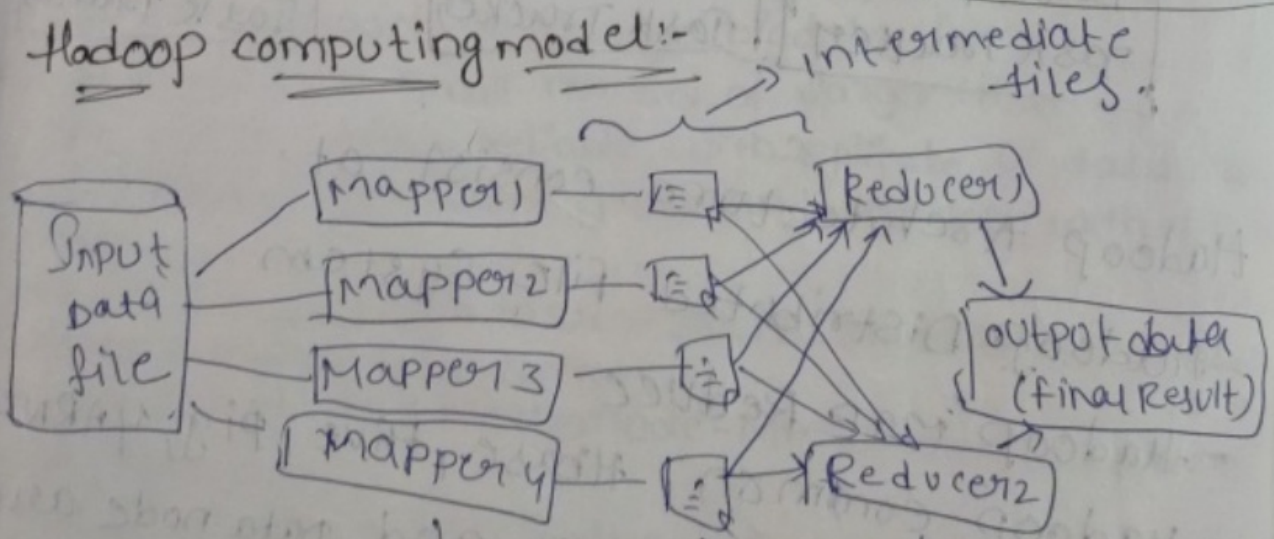
Two Infrastructure model of Hadoop

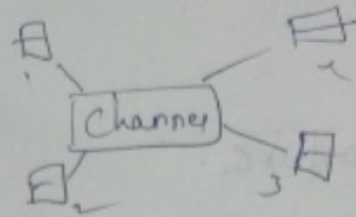
1. Data model
2. computing model

Distributed data model vs Hadoop data model

- | | |
|---------------------------------------|--|
| - deals with tables and relation | - deals with flat files in any format |
| - must have schema for data | - No schema |
| - data fragmentation and Partitioning | - files are automatically divided into blocks. |

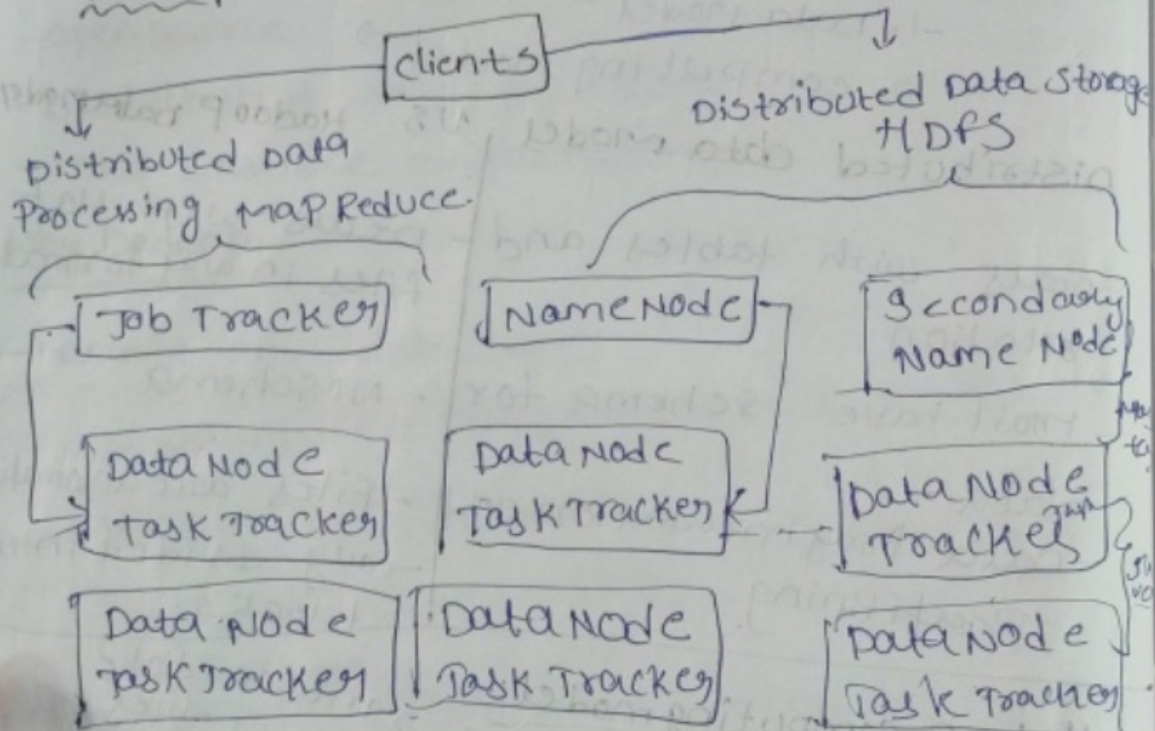
Hadoop computing model:-





16/3/23

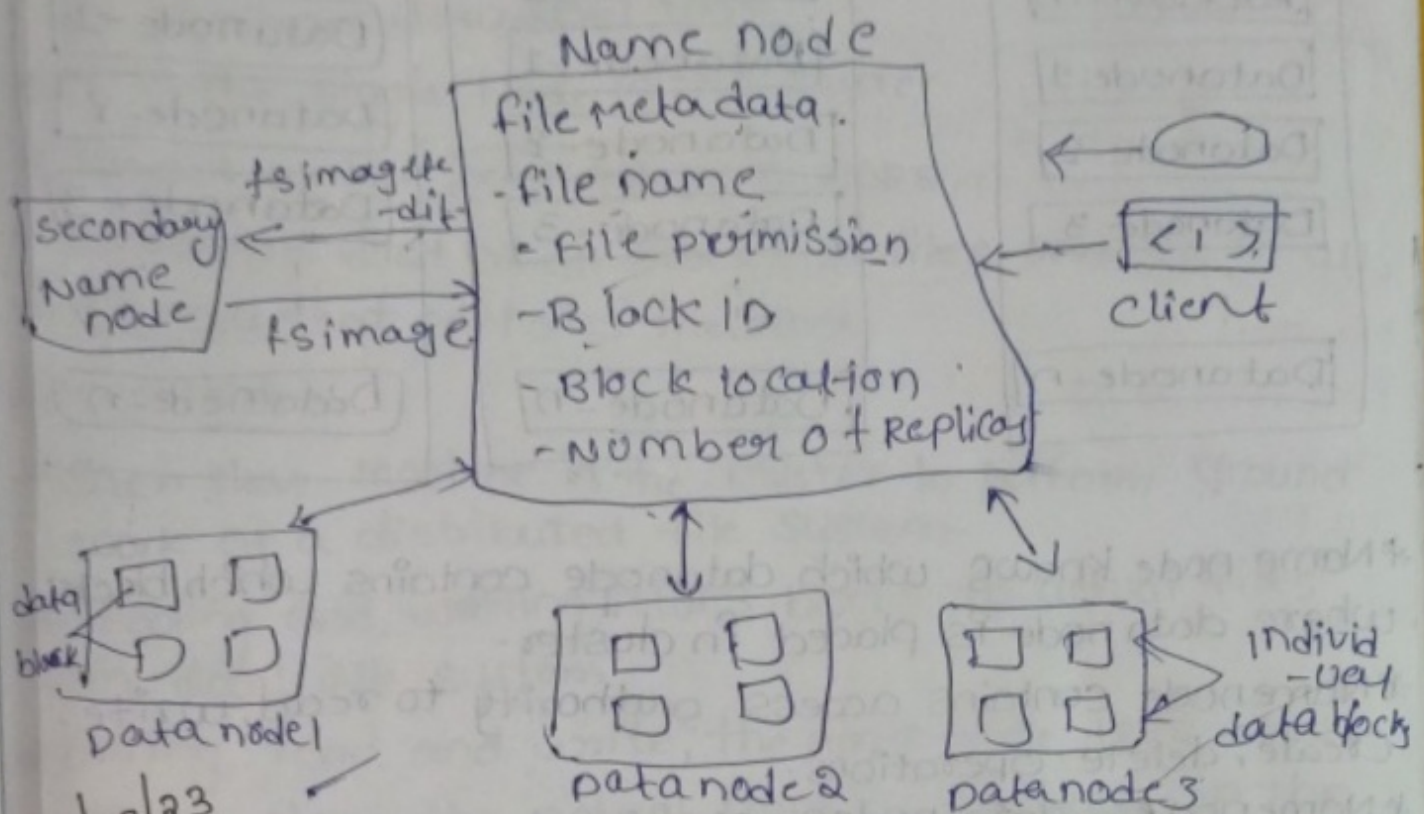
Hadoop Architecture:-



Hadoop Architecture consist of

- Hadoop Distributed file System
- Hadoop Map Reduce
- Hadoop common, Hbase, Hive, Pig, YARN etc,
- * Name Node is master and data node are Slaves
- * Job Tracker is the master and task tracker is the slave

* every slave node come with task track daemon and data node synchronize with the job tracker and namenode respectively.



17/03/23

* HDFS Working model: