

Why Big Data now?

- More data are being collected and stored
- Open source code
- Commodity hardware / Cloud

7 V's

3 V's
gb/3B - (1/B')



- High-Volume
- High-Velocity
- High-Variety

Stream

Java RMI

RCS - iMessage / And

5 V's

→ Artificial Intelligence

Variability

Versatility →

Volume
Velocity
Variety } 3v's

So important

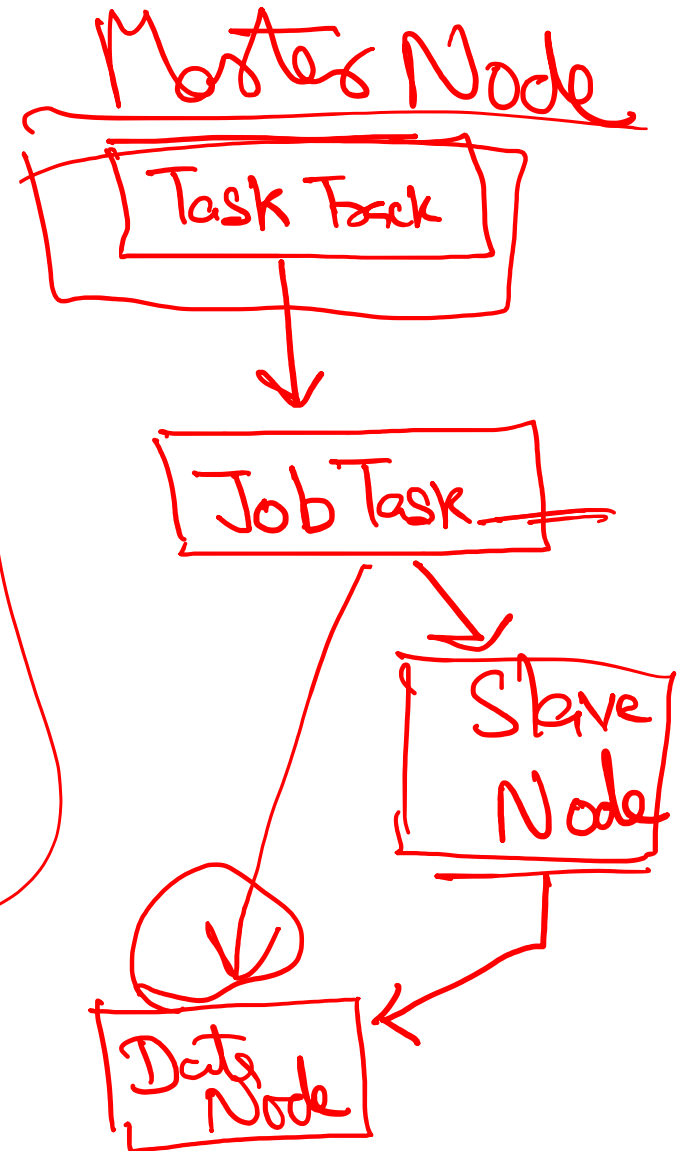
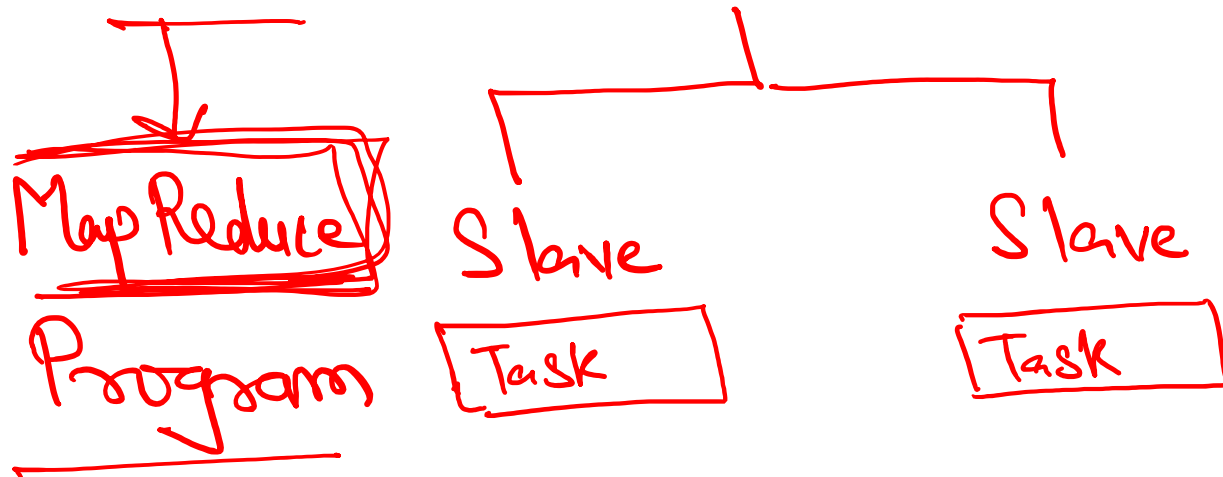
Variability
Versacity } 5v's

Visualization } 7v's

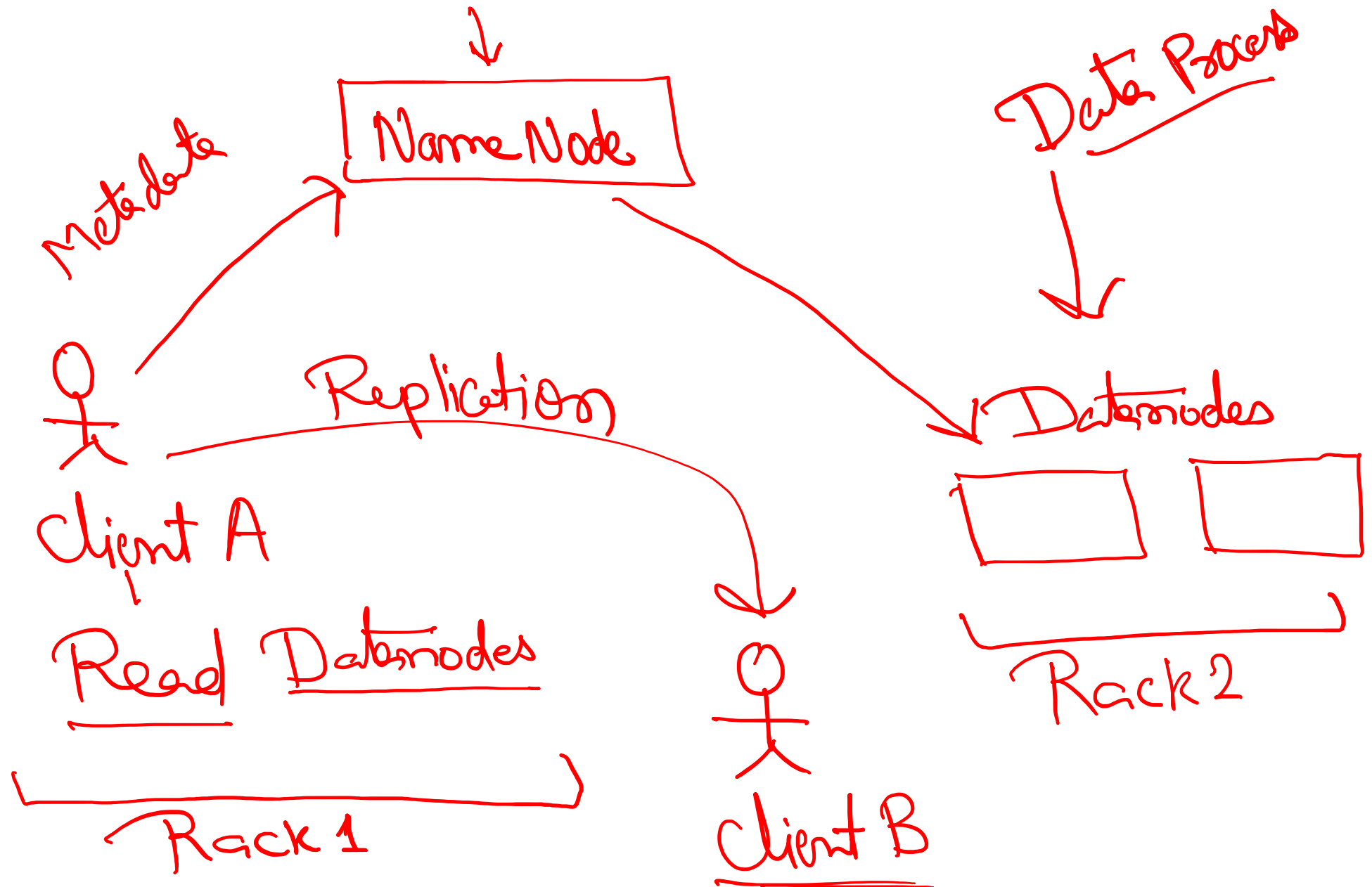
Value is the end game of everything.

Hadoop Ecosystem

HDFS — Master Node



HDFS Architecture



Course Main Thrust 1: Apache Hadoop and Big Data



The Apache™ Hadoop® project develops open-source software for reliable, scalable, distributed computing.

The Apache Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using simple programming models. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage. Rather than rely on hardware to deliver high-availability, the library itself is designed to detect and handle failures at the application layer, so delivering a highly-available service on top of a cluster of computers, each of which may be prone to failures.

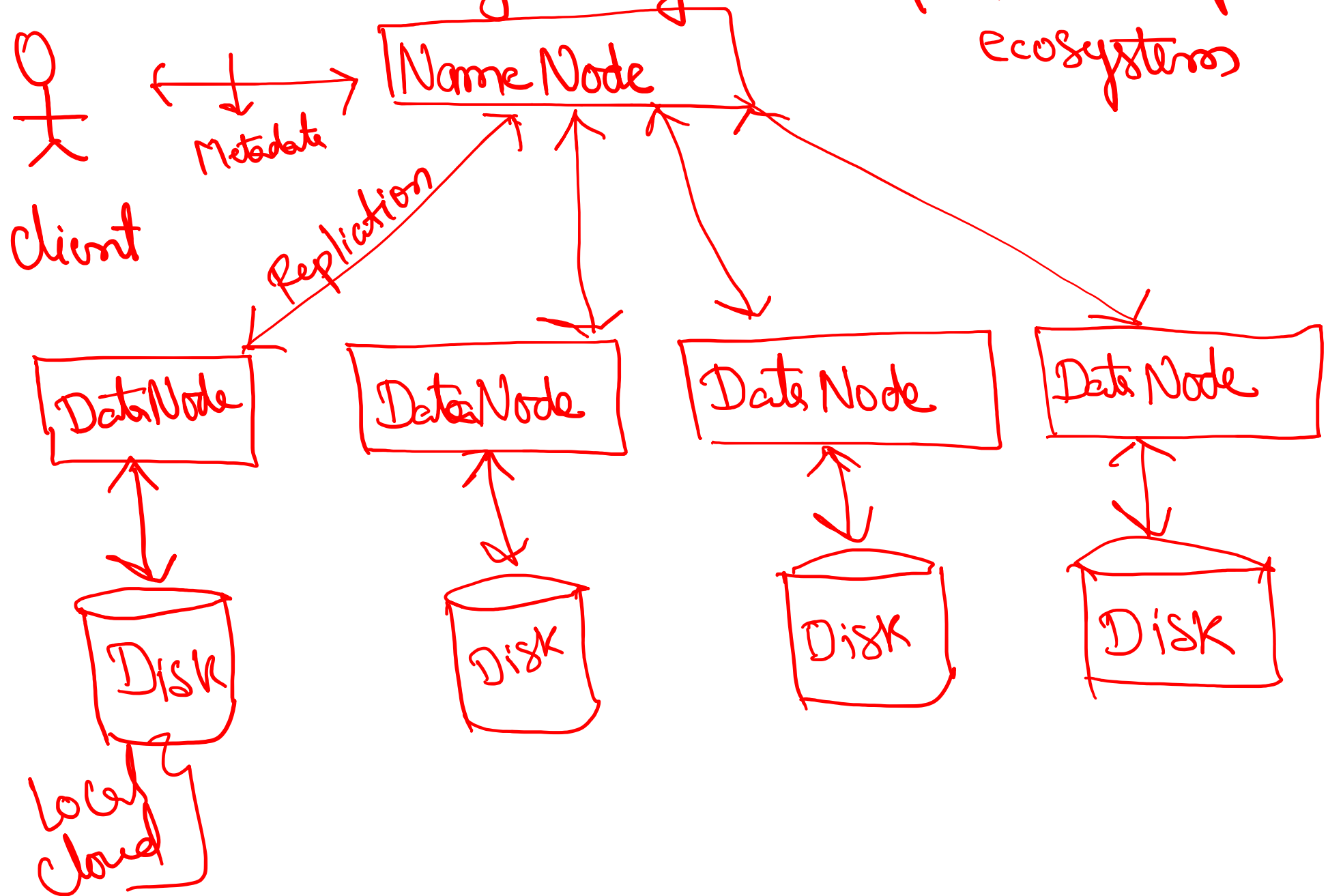
The project includes these modules:

- **Hadoop Common:** The common utilities that support the other Hadoop modules.
- **Hadoop Distributed File System (HDFS™):** A distributed file system that provides high-throughput access to application data.
- **Hadoop YARN:** A framework for job scheduling and cluster resource management.
- **Hadoop MapReduce:** A YARN-based system for parallel processing of large data sets.

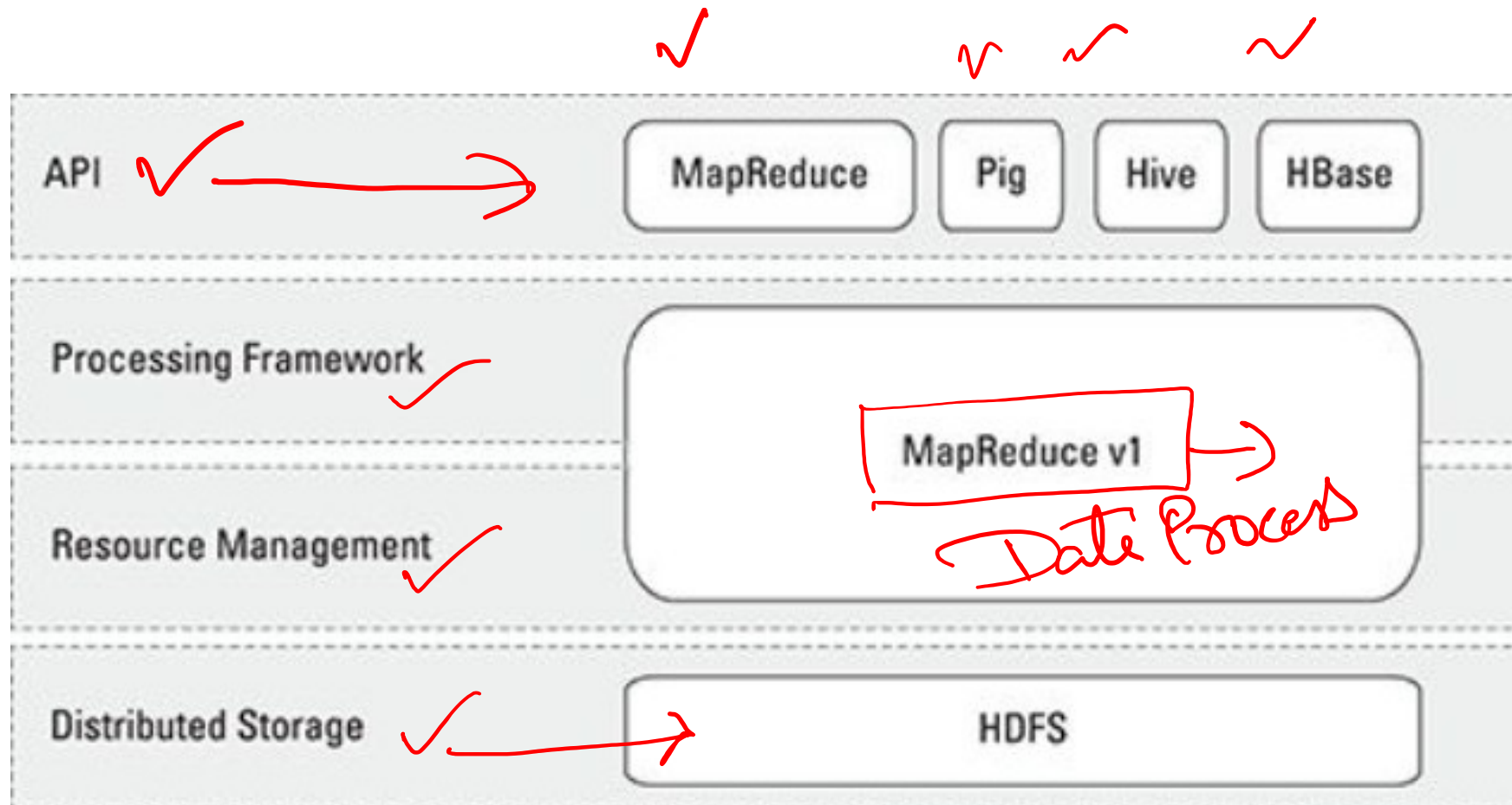


<http://hadoop.apache.org>

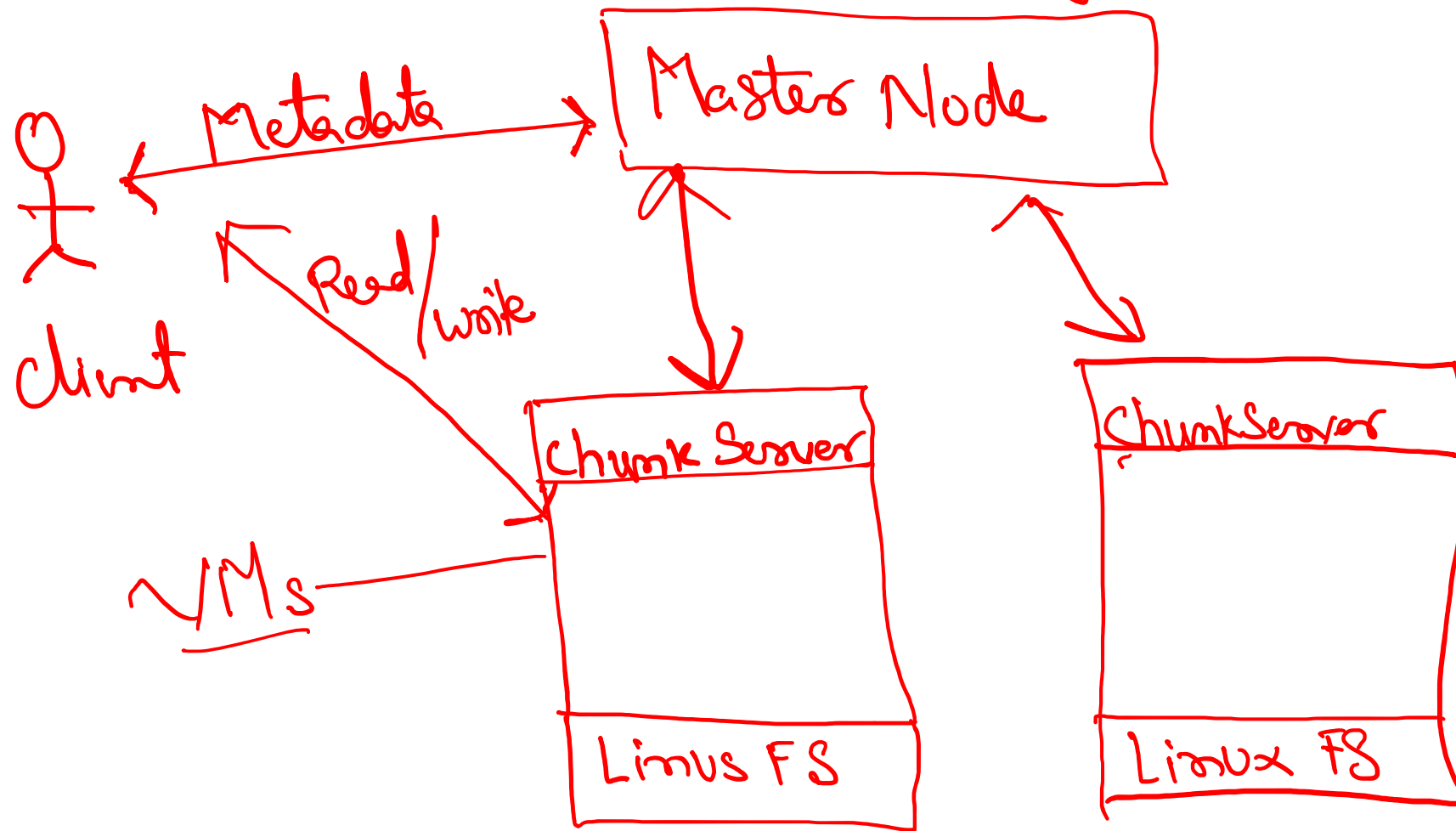
HDFS - Storage layers -> Apache Hadoop ecosystem



Four distinctive layers of Hadoop



Google file System → GCP → google cloud Platform



Apache Hadoop Architecture Eco System Diagram.

