

# HDFS

Press **Esc** to exit full screen

Hadoop File System was developed using distributed file system design. It is run on commodity hardware. Unlike other distributed systems, HDFS is highly faulttolerant and designed using low-cost hardware.

HDFS holds very large amount of data and provides easier access. To store such huge data, the files are stored across multiple machines. These files are stored in redundant fashion to rescue the system from possible data losses in case of failure. HDFS also makes applications available to parallel processing.

*Dr. Neelesh Jain*

Dr. Neelesh Jain 8770193851 Follow me: Youtube/FB : DrNeeleshjain



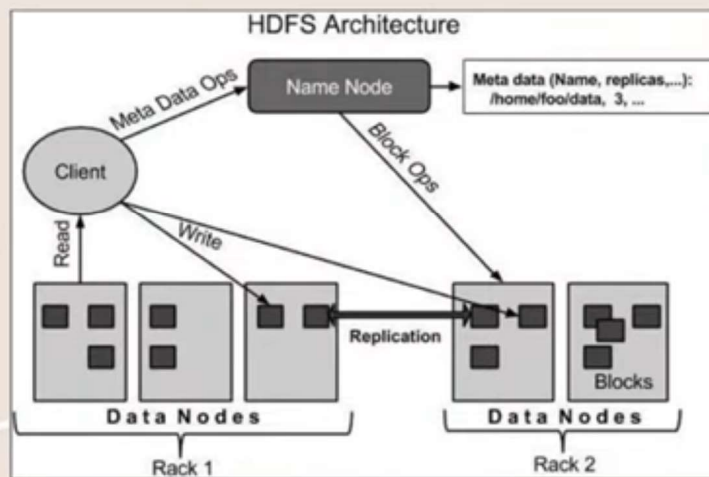
## Features HDFS

- It is suitable for the distributed storage and processing.
- Hadoop provides a command interface to interact with HDFS.
- The built-in servers of namenode and datanode help users to easily check the status of cluster.
- Streaming access to file system data.
- HDFS provides file permissions and authentication.

*Dr. Neelesh Jain*

Dr. Neelesh Jain 8770193851 Follow me: Youtube/FB : DrNeeleshjain

# HDFS Architecture



Dr. Neelesh Jain

Dr. Neelesh Jain 8770193851 Follow me: Youtube/FB : DrNeeleshjain

# HDFS Architecture

HDFS follows the master-slave architecture and it has the following elements.

## Namenode

The namenode is the commodity hardware that contains the GNU/Linux operating system and the namenode software. It is a software that can be run on commodity hardware. The system having the namenode acts as the master server and it does the following tasks –

- Manages the file system namespace.
- Regulates client's access to files.
- It also executes file system operations such as renaming, closing, and opening files and directories.



Dr. Neelesh Jain

Dr. Neelesh Jain 8770193851 Follow me: Youtube/FB : DrNeeleshjain

# HDFS Architecture

- **Block**

Generally the user data is stored in the files of HDFS. The file in a file system will be divided into one or more segments and/or stored in individual data nodes. These file segments are called as blocks. In other words, the minimum amount of data that HDFS can read or write is called a Block. The default block size is 64MB, but it can be increased as per the need to change in HDFS configuration.



*Dr. Neelesh Jain*

Dr. Neelesh Jain 8770193851 Follow me: Youtube/FB : DrNeeleshjain

## Goals of HDFS

**Fault detection and recovery** – Since HDFS includes a large number of commodity hardware, failure of components is frequent. Therefore HDFS should have mechanisms for quick and automatic fault detection and recovery.

**Huge datasets** – HDFS should have hundreds of nodes per cluster to manage the applications having huge datasets.

**Hardware at data** – A requested task can be done efficiently, when the computation takes place near the data. Especially where huge datasets are involved, it reduces the network traffic and increases the throughput.



*Dr. Neelesh Jain*

Dr. Neelesh Jain 8770193851 Follow me: Youtube/FB : DrNeeleshjain

Hadoop Distributed File System (HDFS) II Blocks II Name Node and Data Naode Explained in Hindi

5 Minutes Engineering

The diagram illustrates the HDFS architecture. At the top, a box labeled 'HDFS' contains a sequence of blocks numbered 1, 2, 3, and 4. Below this, a central box labeled 'Name node (meta data)' is connected to three 'Data node' boxes. The first Data node contains blocks 1, 3, and 4, with block 1 marked with a red 'X'. The second Data node contains blocks 2 and 4. The third Data node is empty. A '64MB' label is next to the Name node. The video player interface shows the video is paused at 5:48 / 6:33.