



SNS COLLEGE OF ENGINEERING



Kurumbapalayam(Po), Coimbatore – 641 107

Accredited by NAAC-UGC with 'A' Grade

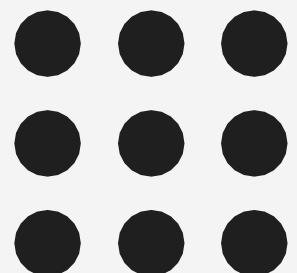
Approved by AICTE, Recognized by UGC & Affiliated to Anna University, Chennai

Department of Artificial Intelligence and Data Science

**Course Name – Big Data Analytics
III Year / V Semester**

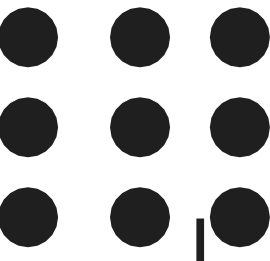
Unit 3 – DATA ANALYTICAL FRAMEWORKS

Topic - HDFS





HDFS



HDFS (Hadoop Distributed File System)

- The Hadoop Distributed File System (HDFS) is based on the Google File System (GFS) and provides a distributed file system that is designed to run on large clusters (thousands of computers) of small computer machines in a reliable, fault-tolerant manner.
- Unlike other distributed systems, HDFS is highly fault tolerant and designed using low-cost hardware.
- HDFS holds very large amount of data and provides easier access. To store such huge data, the files split into blocks 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



HDFS Architecture

- HDFS uses a master/slave architecture where master consists of a single **Name Node** that manages the file system metadata and one or more slave.
- **Data Nodes** that store the actual data.

Name Node

- NameNode is the master node in the Apache Hadoop HDFS Architecture that maintains and manages the blocks present on the DataNodes (slave nodes).
- NameNode is a very highly available server that manages the File System Namespace and controls access to files by clients.
- It Stores metadata for the files, like the directory structure of a typical File System.



HDFS

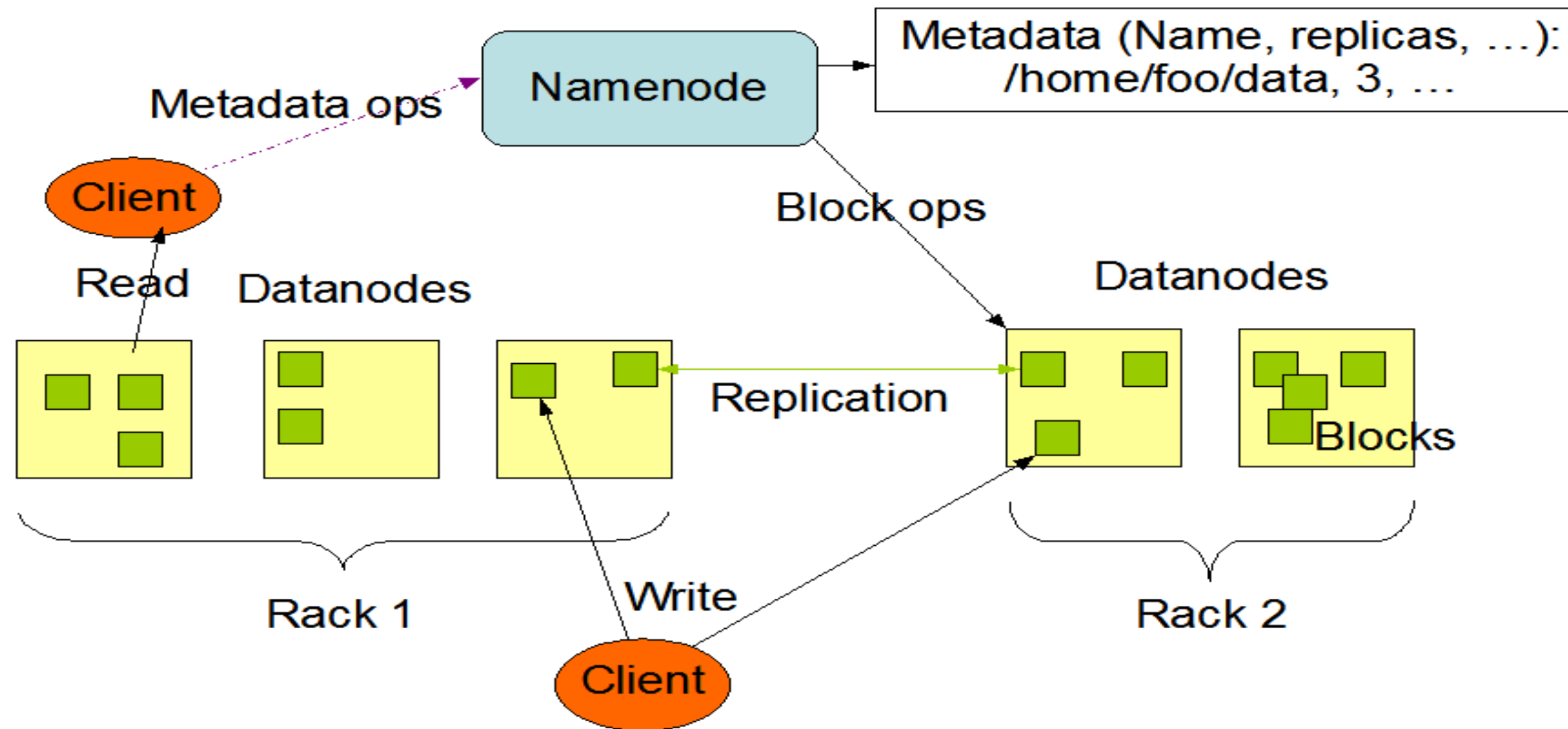


Functions of Name Node

- It Regulates client's access to files.
- It also executes file system operations such as renaming, closing, and opening files and directories. It also determines the mapping of blocks to DataNodes.
- It is the master daemon that maintains and manages the DataNodes (slave nodes)
- It records the metadata of all the files stored in the cluster, e.g. The location of blocks stored, the size of the files, permissions, hierarchy, etc.
- There are two files associated with the metadata:
 - FsImage: It contains the complete state of the file system namespace since the start of the NameNode.
 - EditLogs: It contains all the recent modifications made to the file system with respect to the most recent FsImage
- It regularly receives a Heartbeat and a block report from all the DataNodes in the cluster to ensure that the DataNodes are live.
- The NameNode is also responsible to take care of the replication factor of all the blocks which we will discuss in detail later in this HDFS tutorial blog.

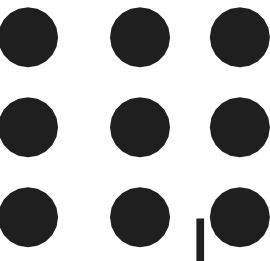
HDFS

HDFS Architecture





HDFS

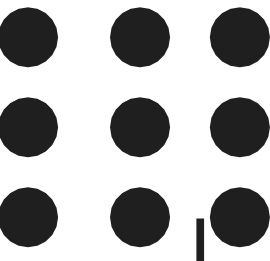


Data Node:

- The datanode is a commodity hardware having the GNU/Linux operating system and datanode software. These nodes manage the data storage of their system.
- A file in an HDFS namespace is split into several blocks and those blocks are stored in a set of Data Nodes.
- Data nodes store and retrieve blocks when they are requested by client or name node.
- They report back to name node periodically, with list of blocks that they are storing.
- The data node also perform operations such as block creation, deletion and replication as stated by the name node.
- They send heartbeats to the NameNode periodically to report the overall health of HDFS, by default, this frequency is set to 3 seconds.



HDFS



Secondary NameNode:

- Apart from these two daemons, there is a third daemon or a process called Secondary NameNode.
- The Secondary NameNode works concurrently with the primary NameNode as a helper daemon. And don't be confused about the Secondary NameNode being a backup NameNode because it is not.

Functions of Secondary NameNode:

- The Secondary NameNode is one which constantly reads all the file systems and metadata from the RAM of the NameNode and writes it into the hard disk or the file system.
- It is responsible for combining the EditLogs with FsImage from the NameNode.
- It downloads the EditLogs from the NameNode at regular intervals and applies to FsImage. The new FsImage is copied back to the NameNode, which is used whenever the NameNode is started the next time.



HDFS

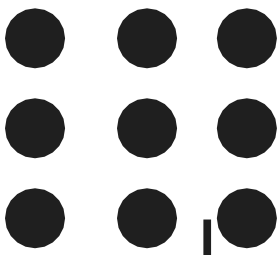


Secondary NameNode:

- Apart from these two daemons, there is a third daemon or a process called Secondary NameNode.
- The Secondary NameNode works concurrently with the primary NameNode as a helper daemon. And don't be confused about the Secondary NameNode being a backup NameNode because it is not.

Functions of Secondary NameNode:

- The Secondary NameNode is one which constantly reads all the file systems and metadata from the RAM of the NameNode and writes it into the hard disk or the file system.
- It is responsible for combining the EditLogs with FsImage from the NameNode.
- It downloads the EditLogs from the NameNode at regular intervals and applies to FsImage. The new FsImage is copied back to the NameNode, which is used whenever the NameNode is started the next time.
- Hence, Secondary NameNode performs regular checkpoints in HDFS. Therefore, it is also called CheckpointNode



THANK YOU