# SESSION 3

# ASSIGNMENT 1

Different components of Hadoop 2.x are as follows:
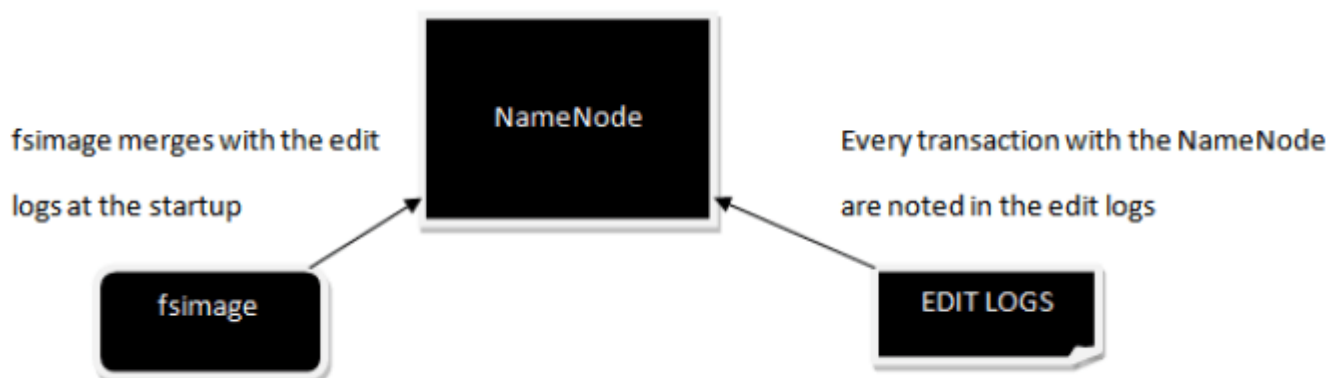
- **NameNode**
- **DataNode**
- **Secondary NameNode**
- **Resource Manager**
- **Node Manager**



# HDFS Daemons

**NameNode**

NameNode holds the meta data for the HDFS like Namespace information, block information etc. When in use, all this information is stored in main memory. But these information also stored in disk for persistence storage.



There are two different files associated with the NameNode

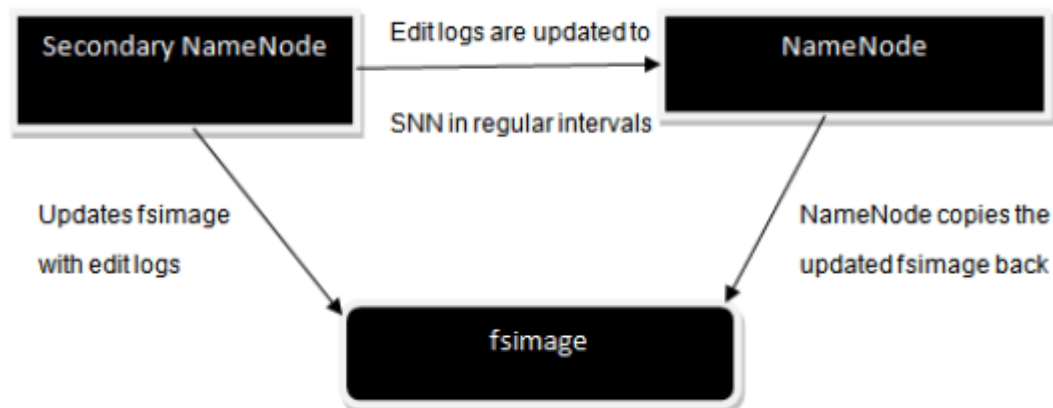1. fsimage– It is the snapshot of the filesystem

2. Edit logs – After the start of NameNode it maintains the every transaction that happened with NameNode.

The NameNode maintains the namespace tree and the mapping of blocks to DataNodes.

The Client communicates with the NameNode and provides data to the HDFS through it. The HDFS then stores data as blocks inside DataNodes.

By default, the block size in a hadoop cluster is 64MB. NameNode maintains the meta data information of all the blocks present inside the hadoop cluster like permissions, modification and access times, namespace and disk space quotas . This is the reason NameNode is also called as the Master node.

**SecondaryNameNode**



The above figure gives us a clear idea of how Secondary NameNode works

1. It asks the NameNode for its edit logs in regular intervals and copies them into the fsimage.

2. After updating the fsimage, NameNode copy back that fsimage

3. NameNode uses this fsimage when it starts, this eventually will reduce the startup time.

The main theme of secondary NameNode is to maintain a checkpoint in HDFS. When a failure occurs SNN won't become NameNode it just helps NameNode in bringing back its data. SNN is also called as checkpoint node in hadoop's architecture.

**DataNode**

This is the place where actual data is stored in the hadoop cluster in a distributed manner.

Every DataNode will have a block scanner and it directly reports to the NameNode about the blocks which it is handling.

The DataNodes communicate with the NameNode by sending Heartbeats for every 3seconds and if the NameNode does not receive any Heartbeat for 10 minutes, then it treats the DataNode as a dead node and re-replicates the blocks.

During startup each DataNode connects to the NameNode and performs a handshake. The purpose of the handshake is to verify the namespace ID and the software version of the DataNode. If either does not match that of the NameNode, the DataNode automatically shuts down.

The namespace ID is assigned to the file system instance when it is formatted. The namespace ID is persistently stored on all nodes of the cluster.

# YARN Daemons (Master/ResourceManager, Worker/NodeManager)

In a YARN cluster, there are two types of hosts:

- The ResourceManager is the master daemon that communicates with the client, tracks resources on the cluster, and orchestrates work by assigning tasks to NodeManagers.
- A NodeManager is a worker daemon that launches and tracks processes spawned on worker hosts.