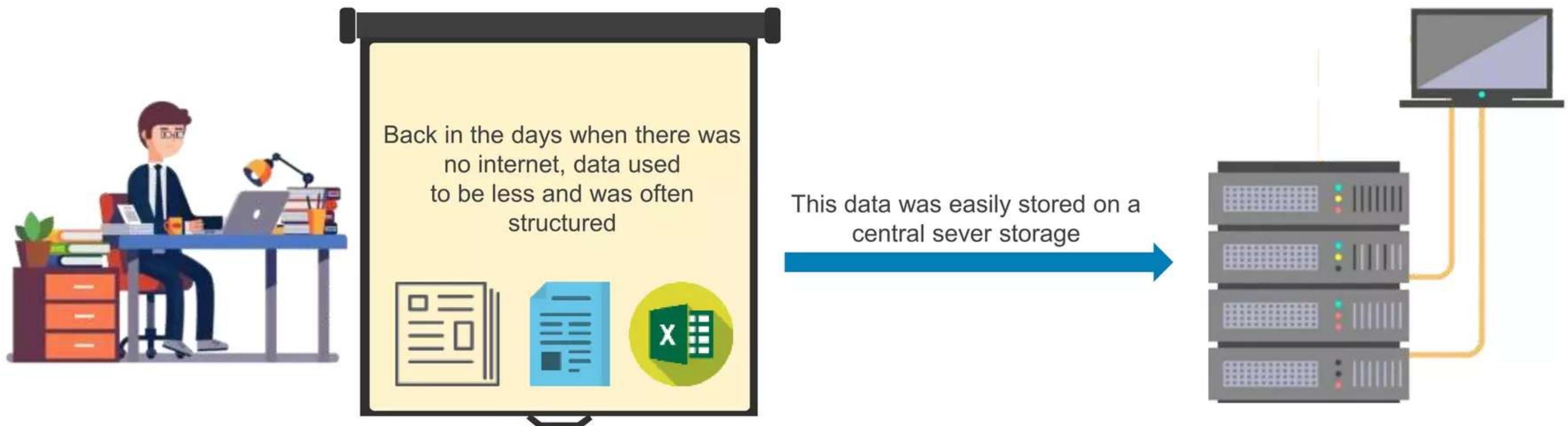


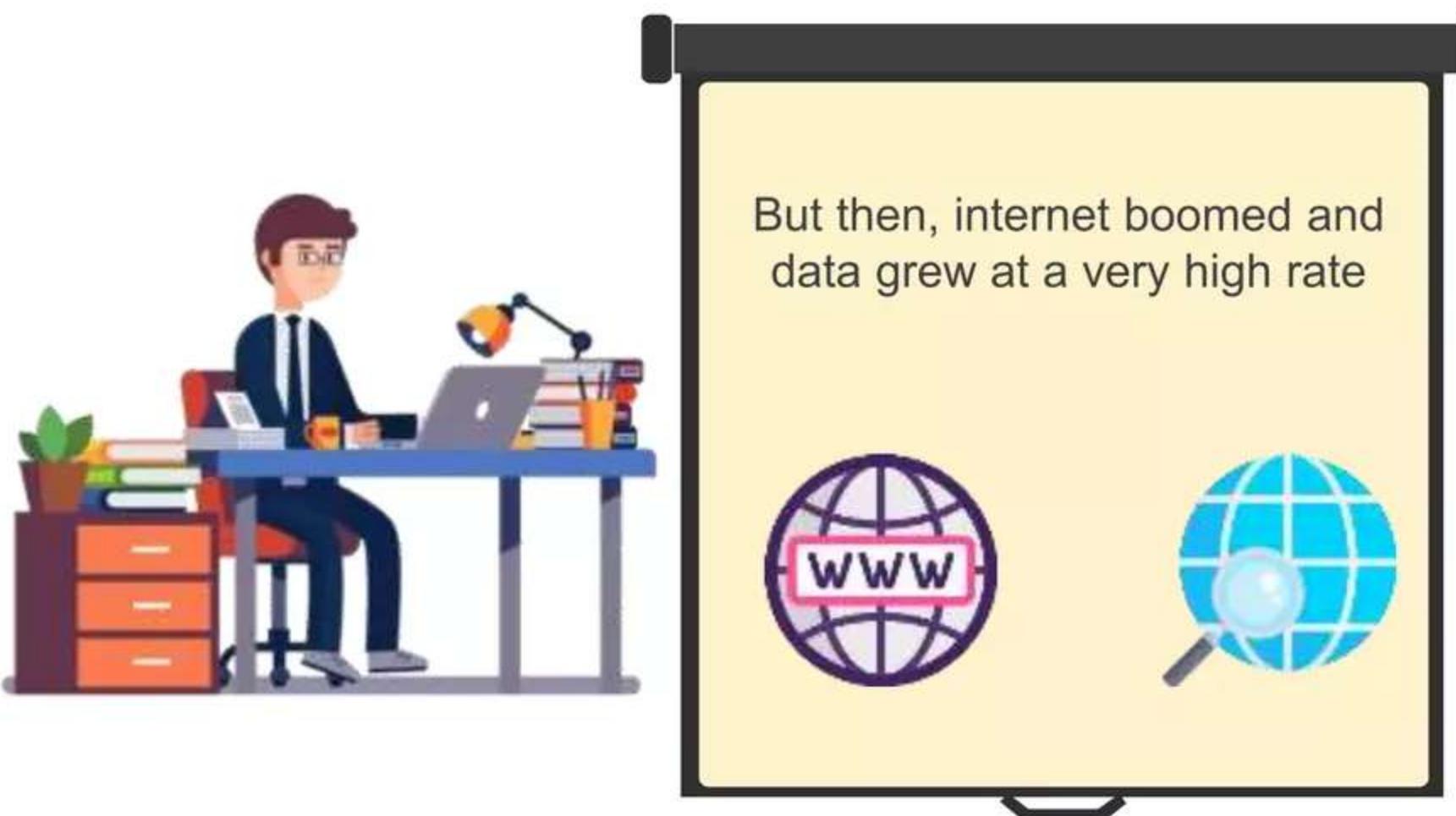
Hadoop Architecture

simplilearn

How Big Data evolved?

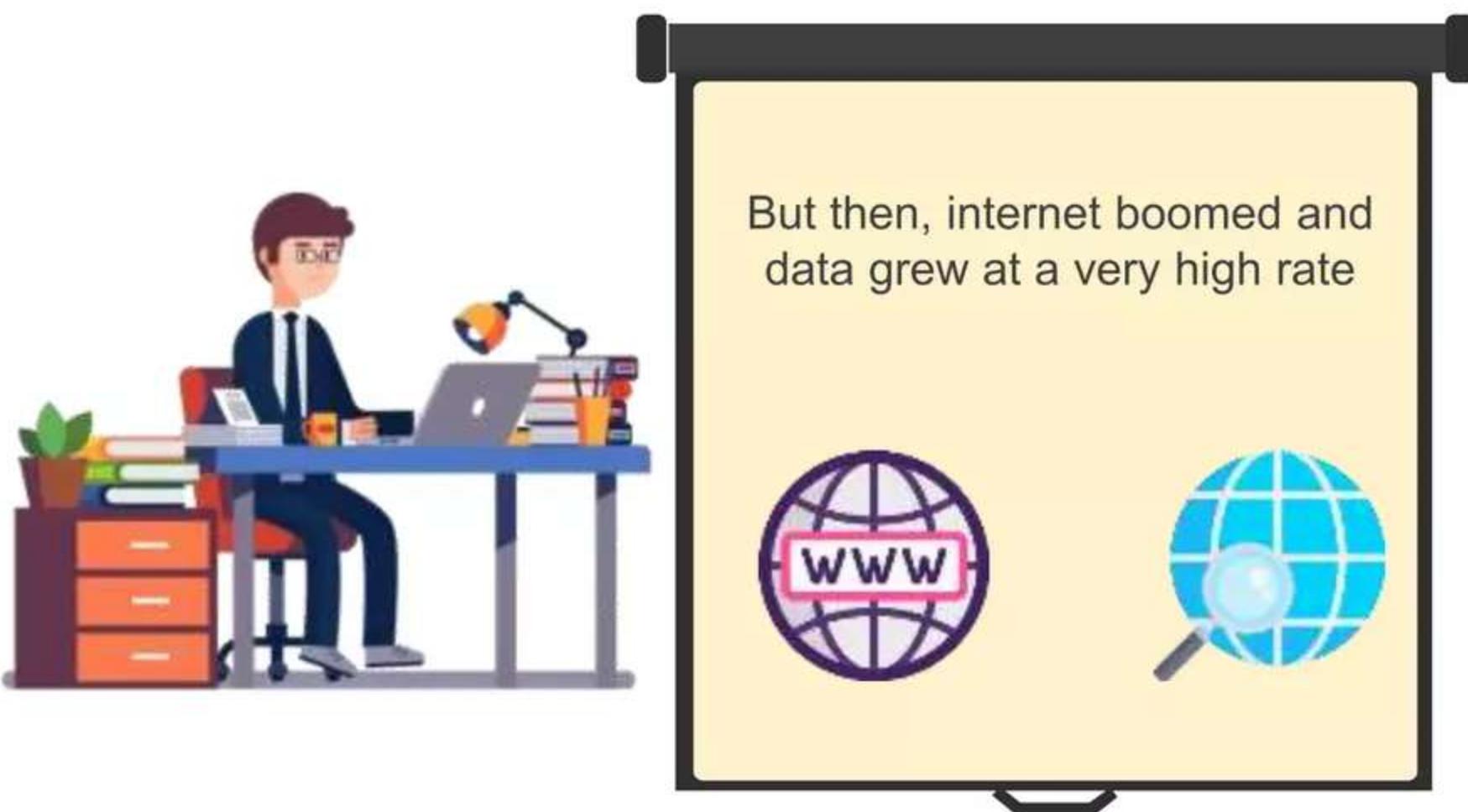


How Big Data evolved?



A lot of semi-structured and unstructured data was being generated

How Big Data evolved?



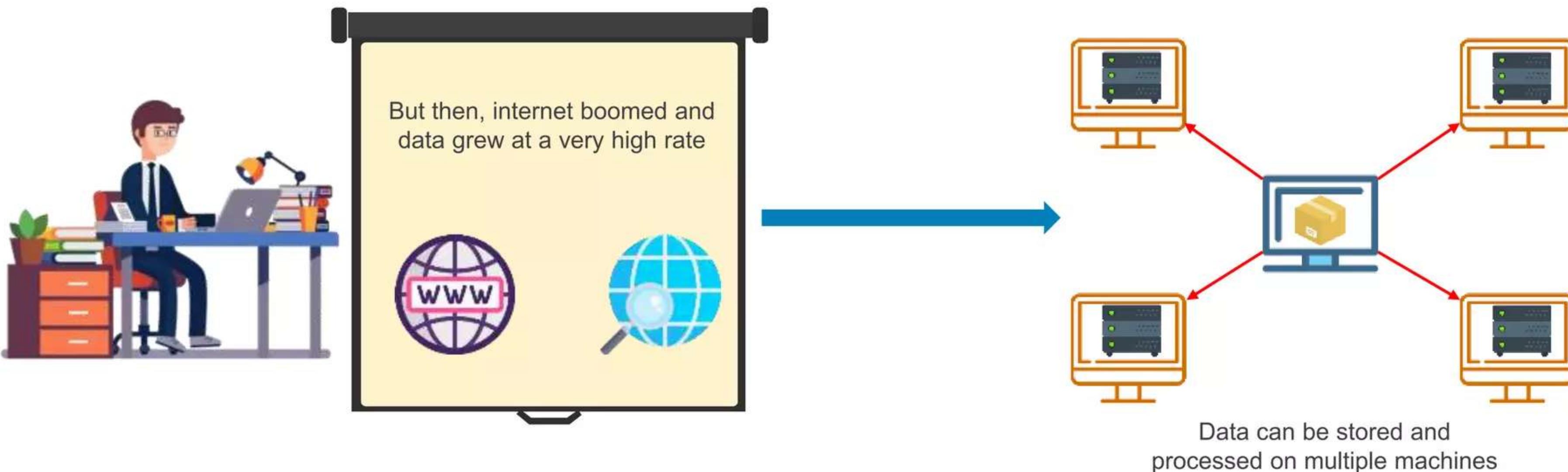
Storing such huge volumes of data on a single server was not an efficient way

How Big Data evolved?

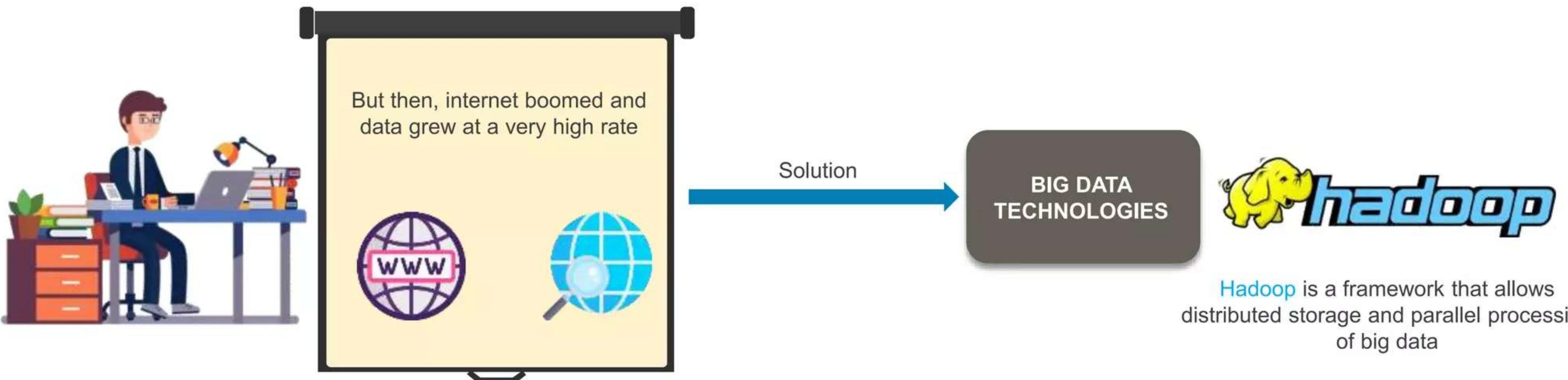


There was a need for **distributed storage machines** where data could be stored and processed parallelly

How Big Data evolved?



How Big Data evolved?



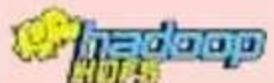
What's in it for you?



What is Hadoop?



Components of Hadoop



What is HDFS?



HDFS Architecture

Hadoop MapReduce



Hadoop MapReduce Example



Hadoop YARN



Demo on MapReduce



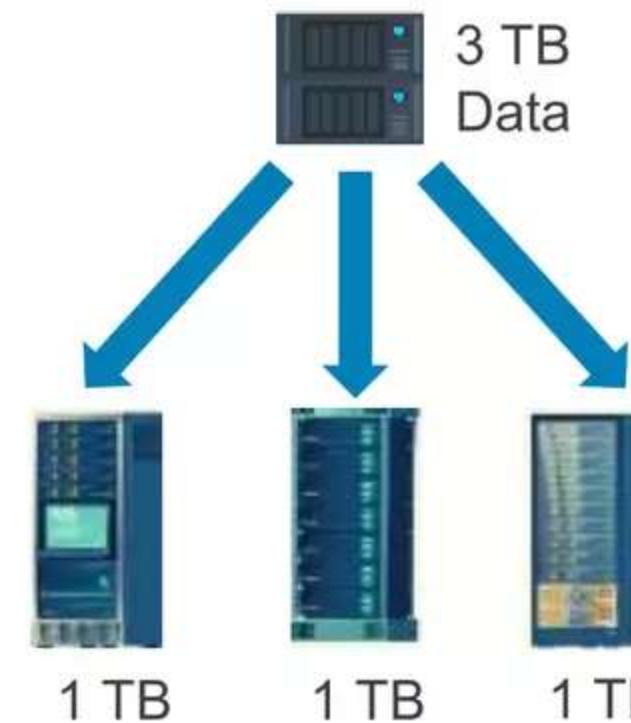
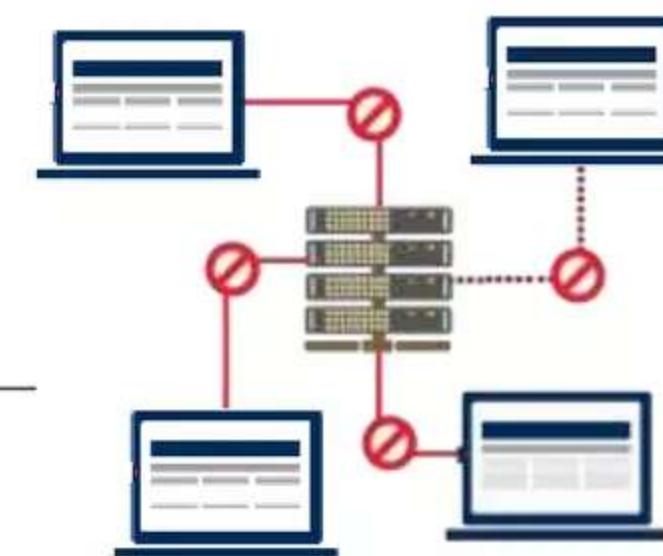
What is Hadoop?

What is Hadoop?

Hadoop is a framework that allows you to store large volumes of data on several node machines

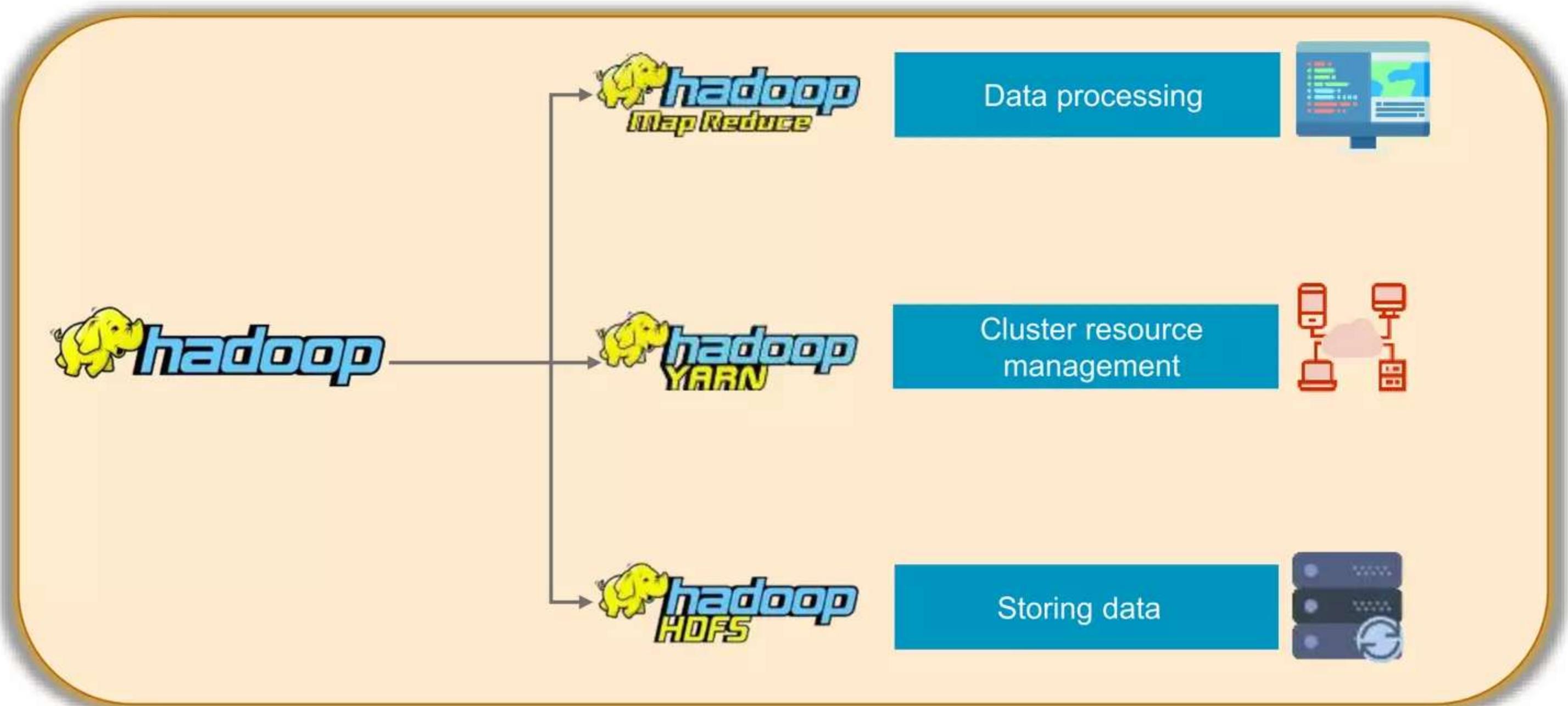


It also helps in processing the data in a parallel manner



Components of Hadoop

Components of Hadoop

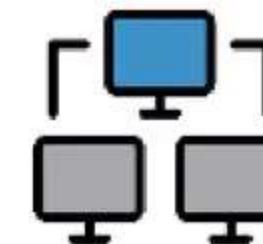
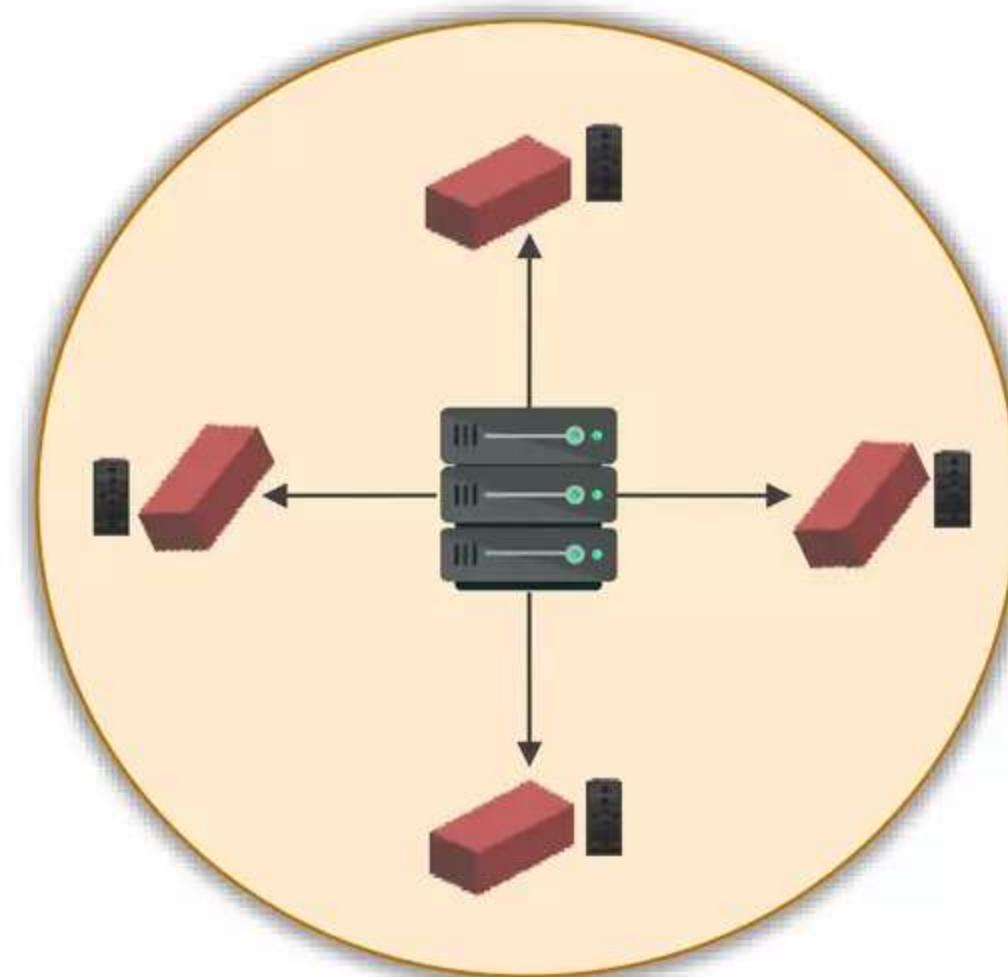


What is HDFS?

What is HDFS?



Hadoop Distributed File System (HDFS) is the storage layer of Hadoop that stores data in multiple data servers



Data is divided into multiple blocks



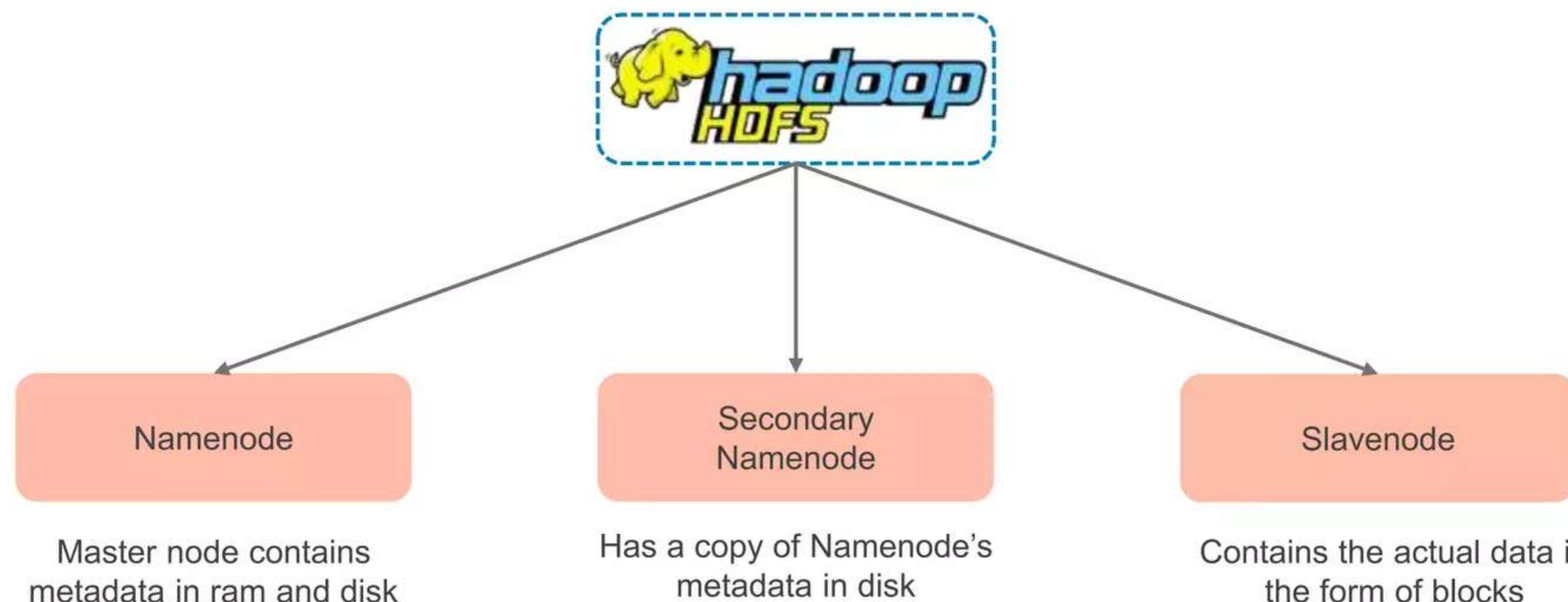
Stores them over multiple nodes of the cluster

What is HDFS?



Hadoop Distributed File System (HDFS) is the storage layer of Hadoop that stores data in multiple data servers

3 core components

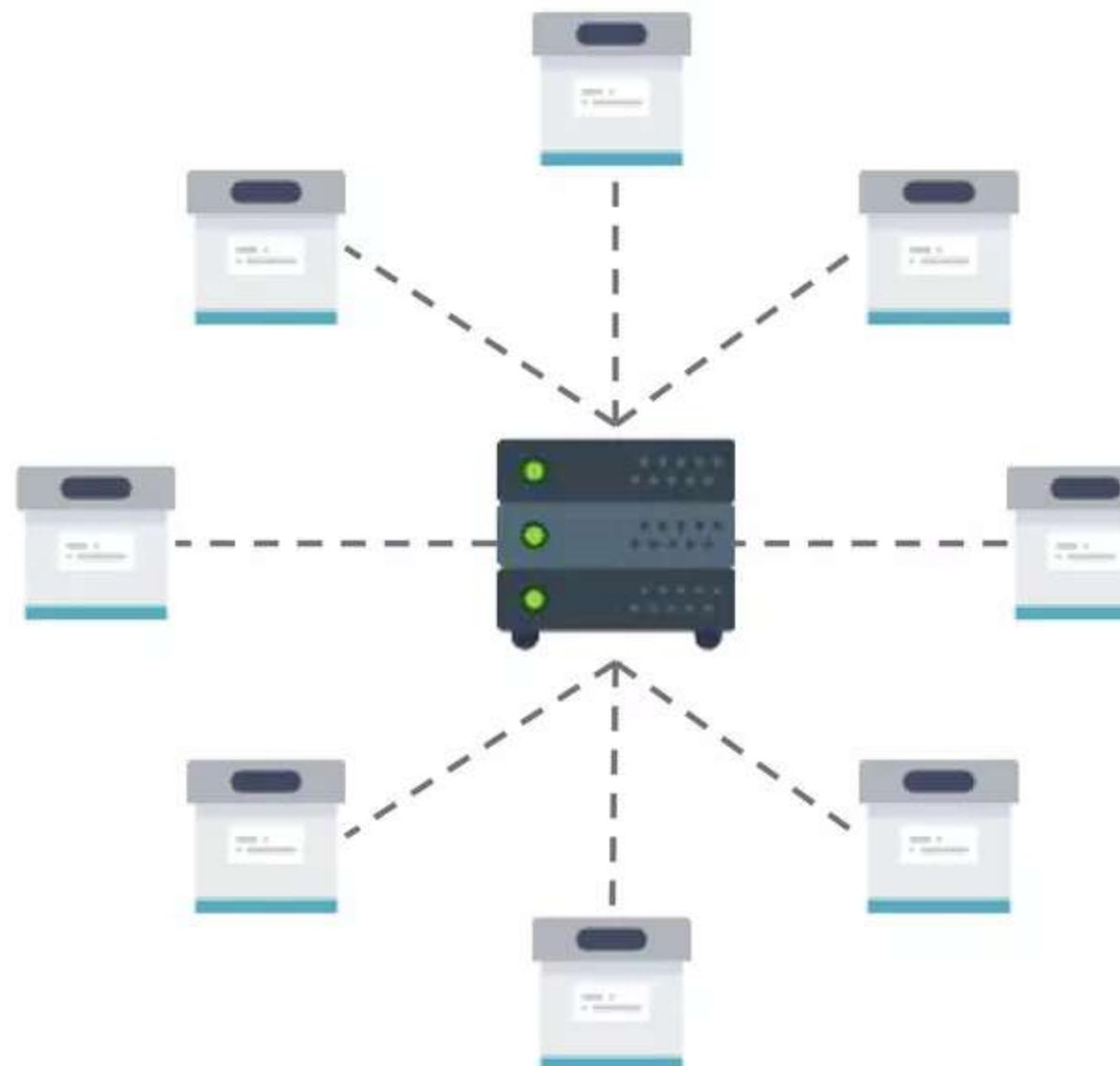


HDFS Blocks

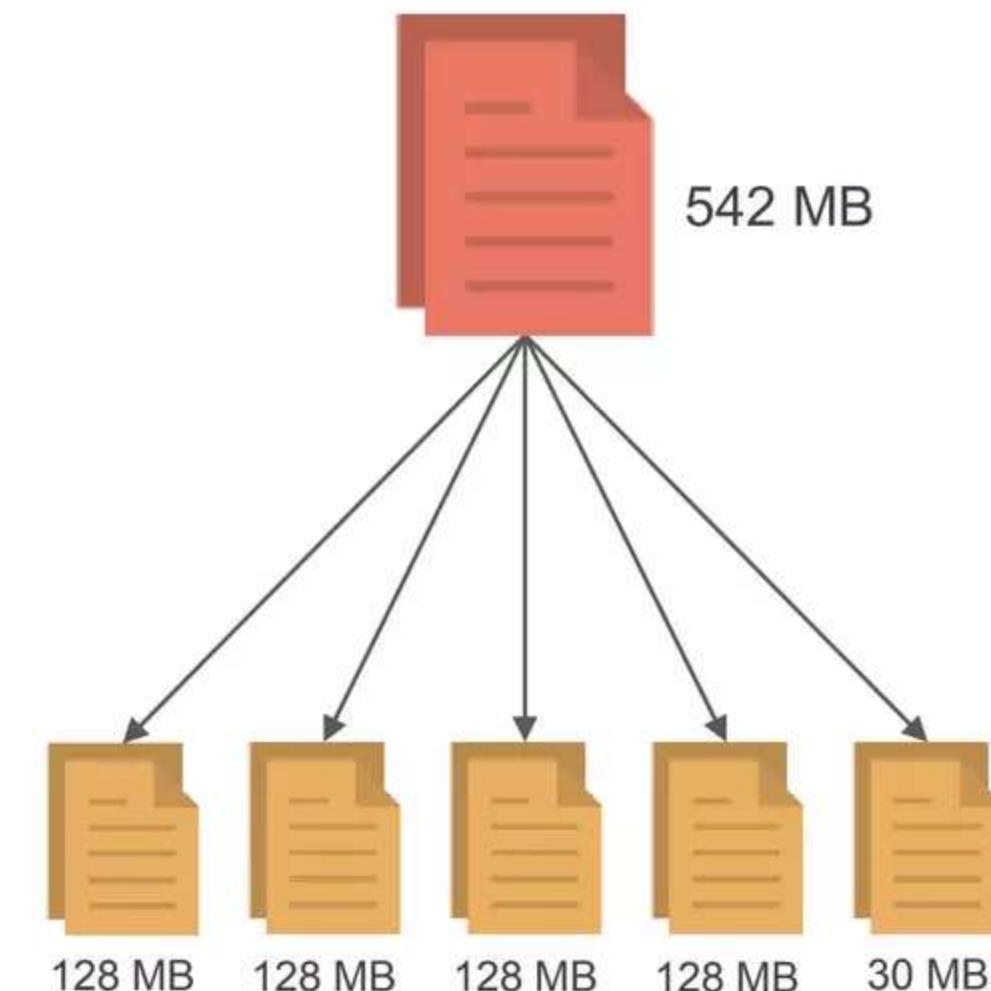
simplilearn

90

HDFS Blocks



Suppose, we have a 542 MB file



HDFS divides large data into different blocks

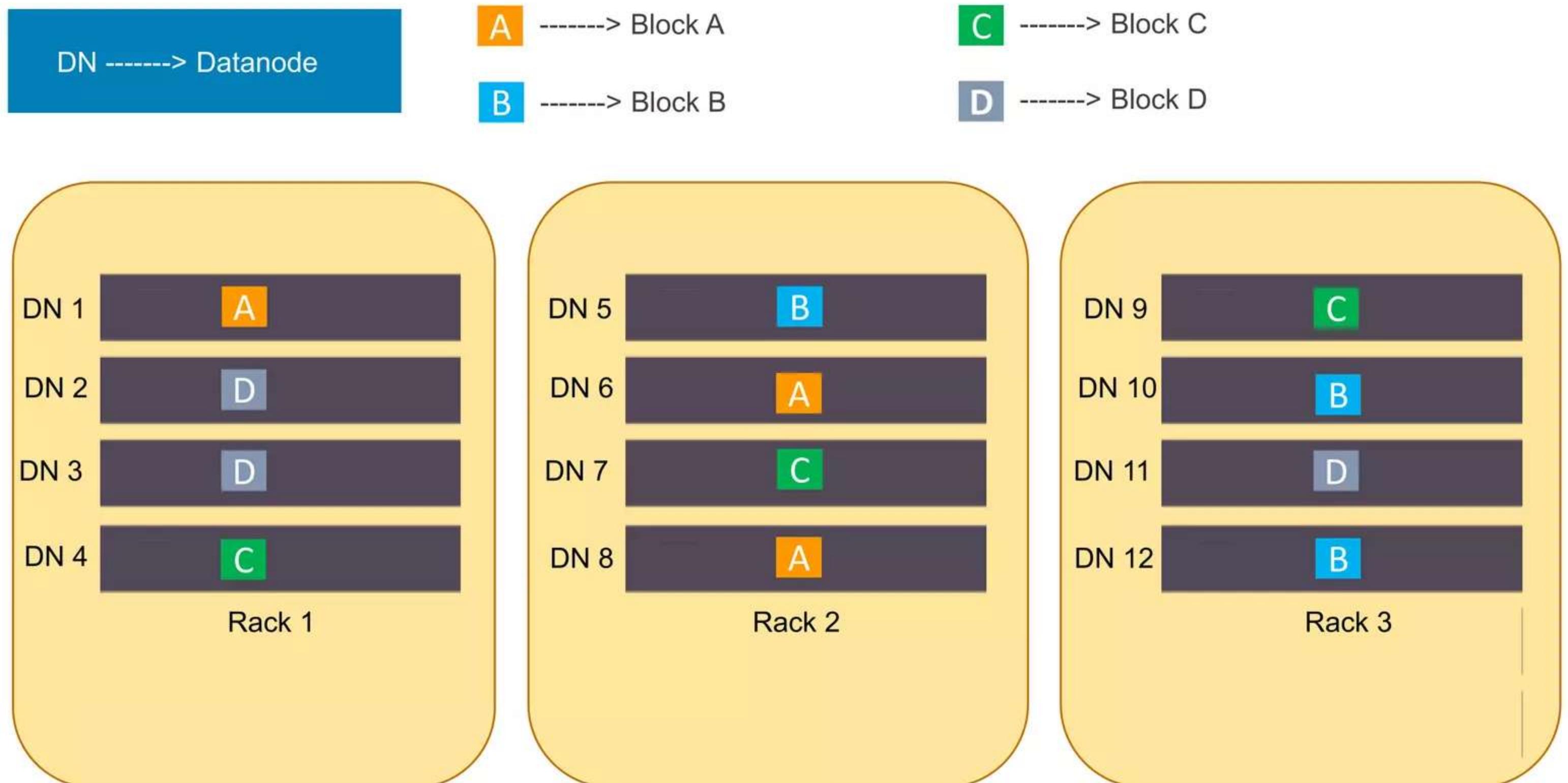
Each block by default has 128 MB's of data

Data Replication

simplilearn

90

Data Replication in HDFS

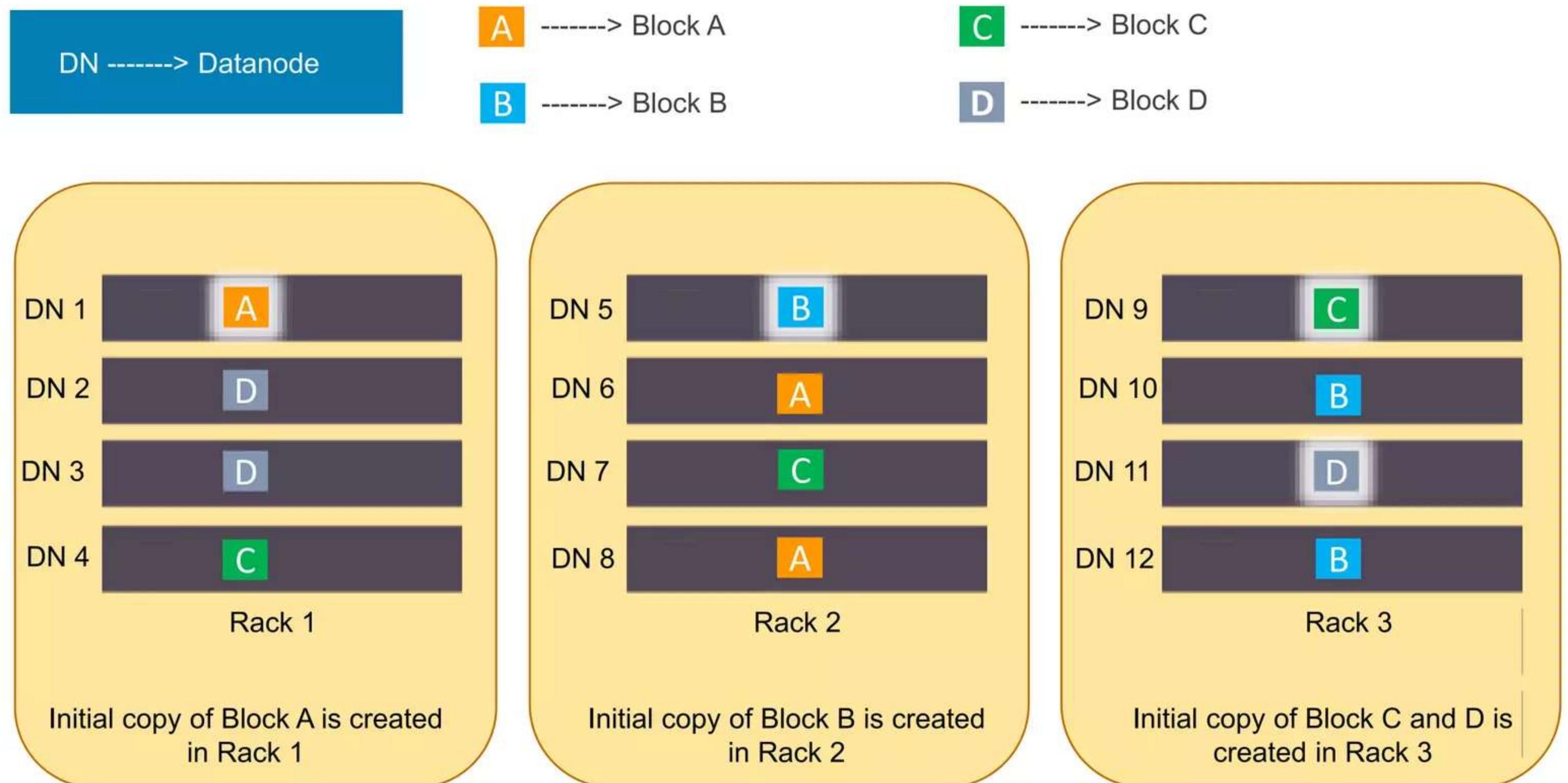


Each block of data is being replicated thrice on different datanodes present in different racks



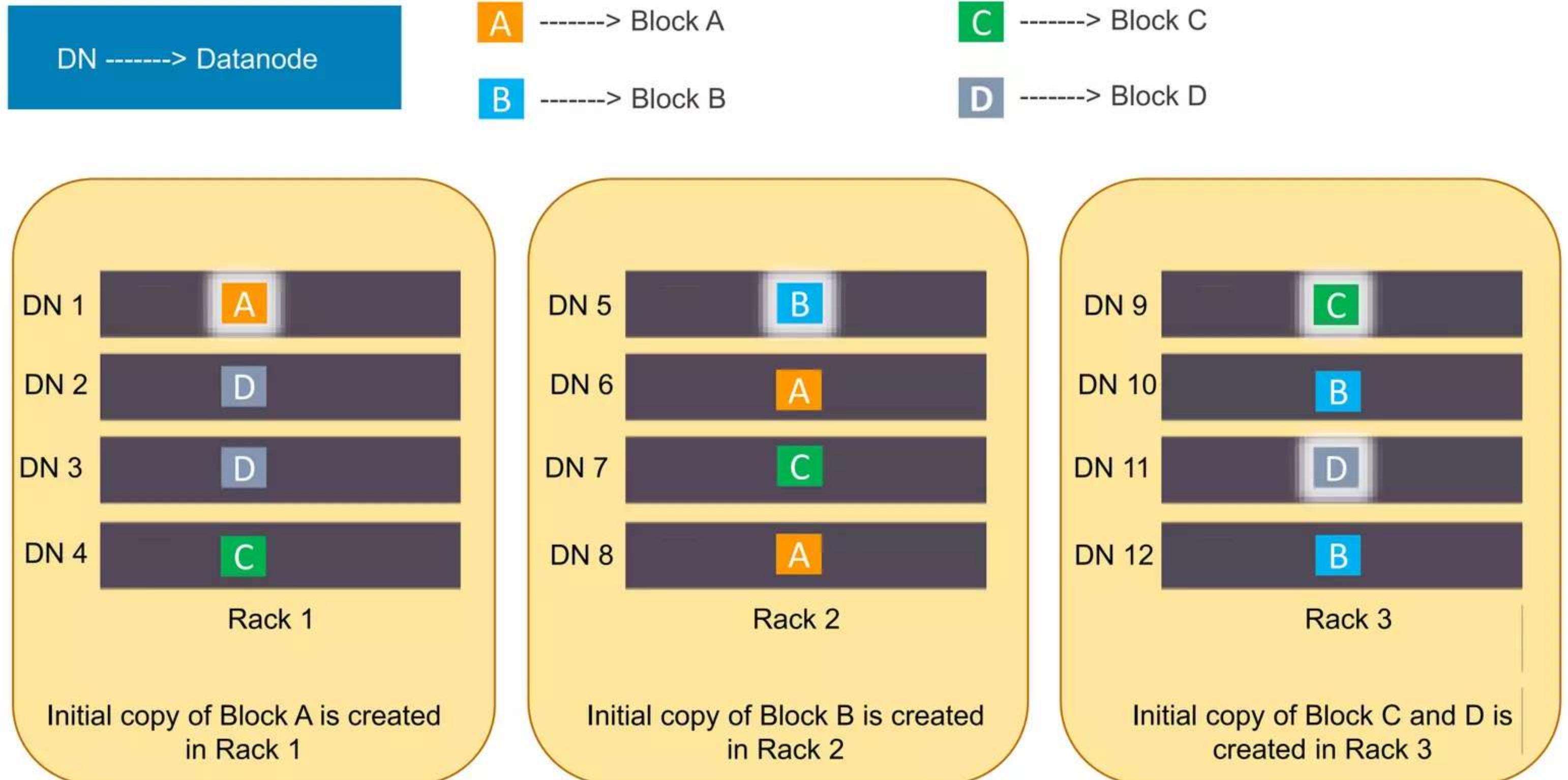
Do you understand what's happening here?

Data Replication in HDFS



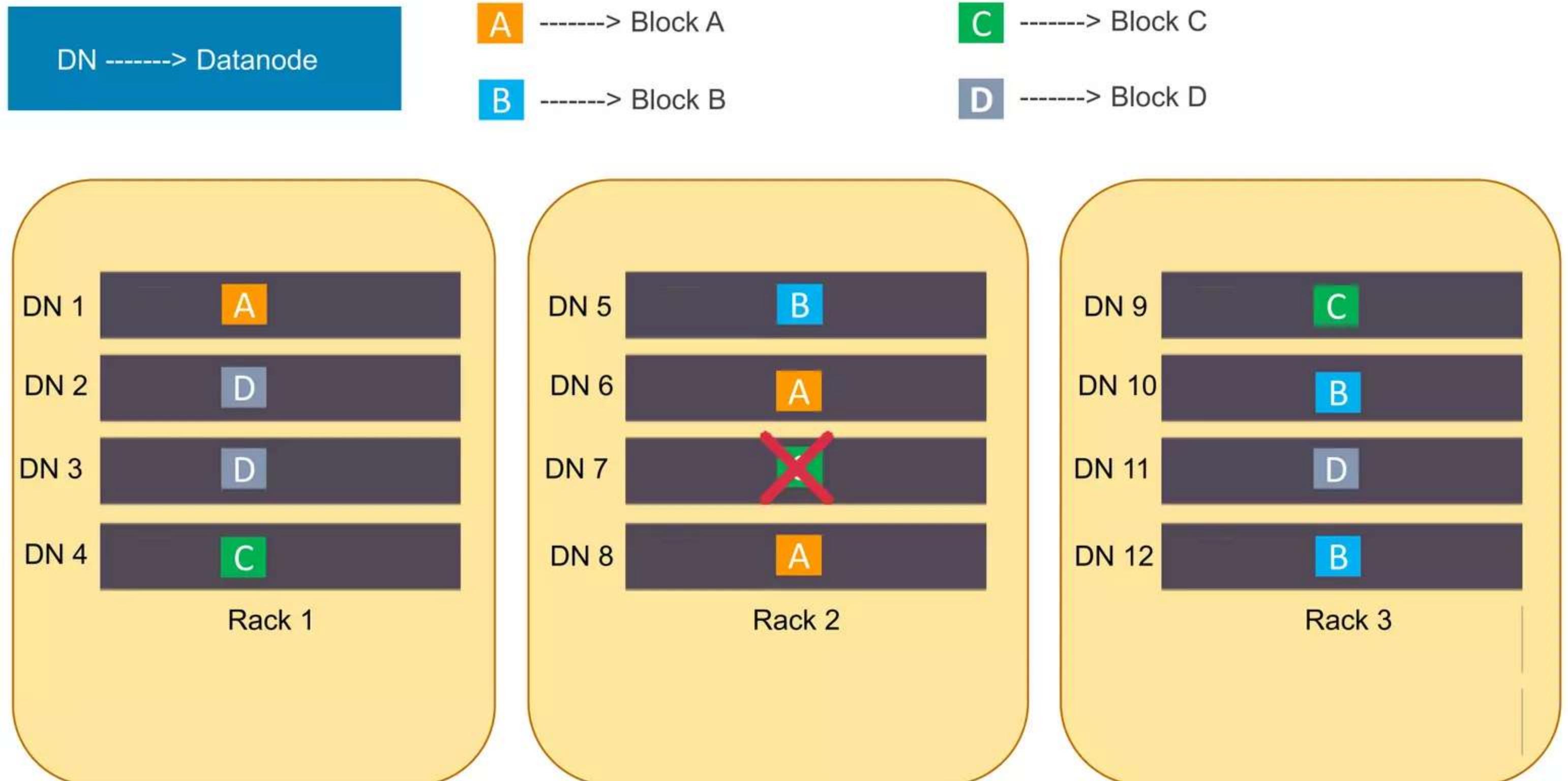
Two identical blocks cannot be placed on the same datanode

Data Replication in HDFS

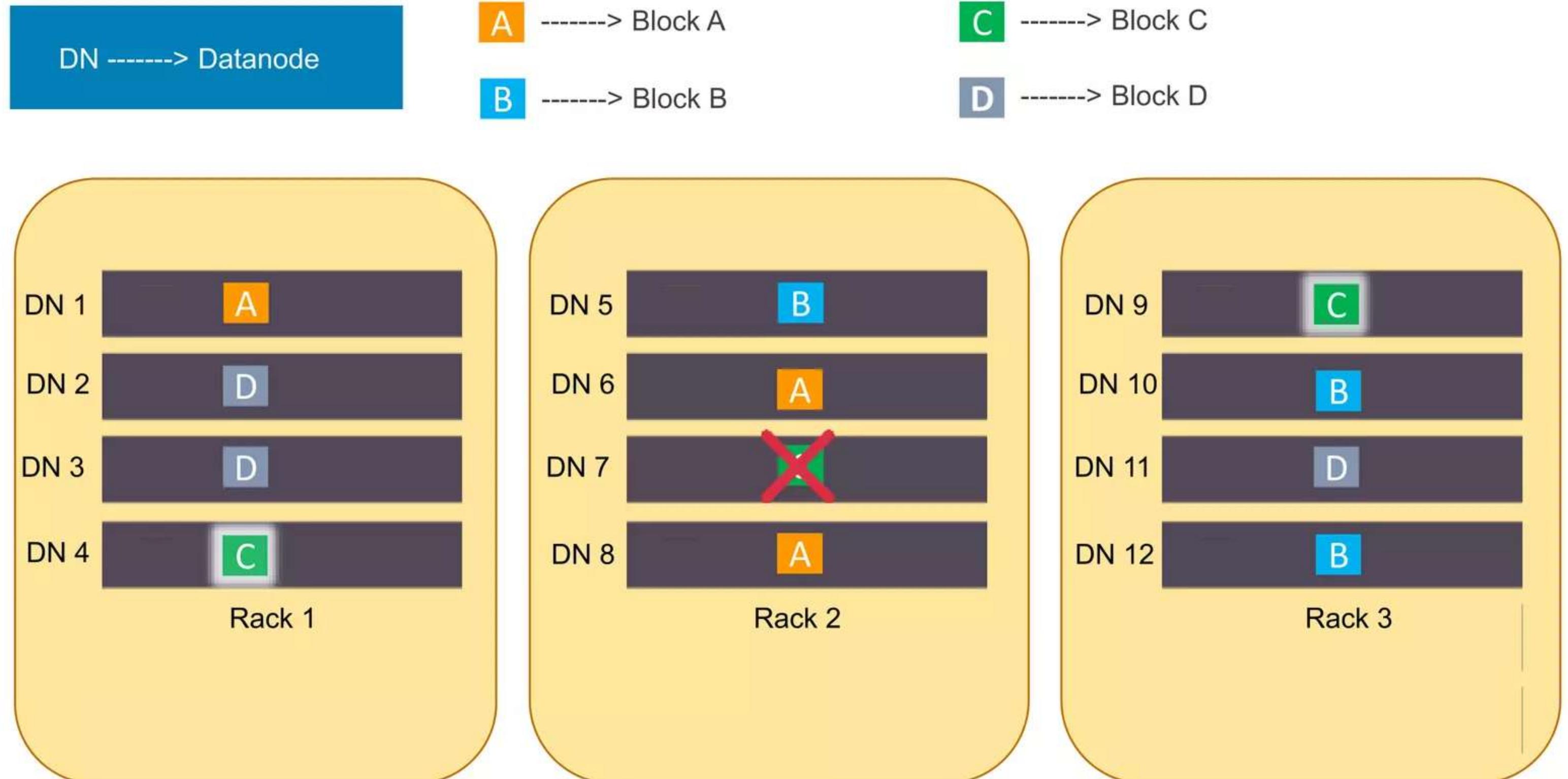


When cluster is rack aware, all the replicas of a block will not be placed on the same rack

Data Replication in HDFS



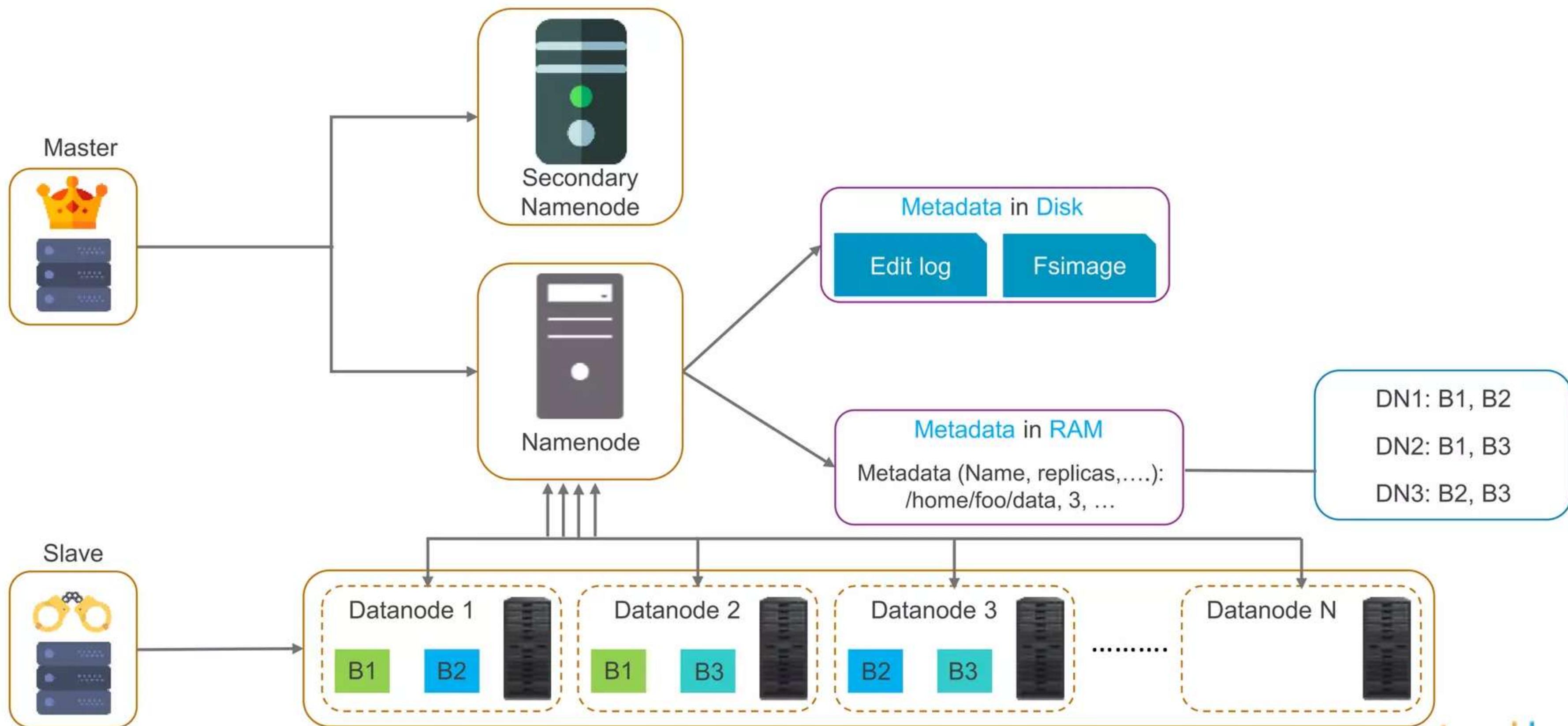
Data Replication in HDFS



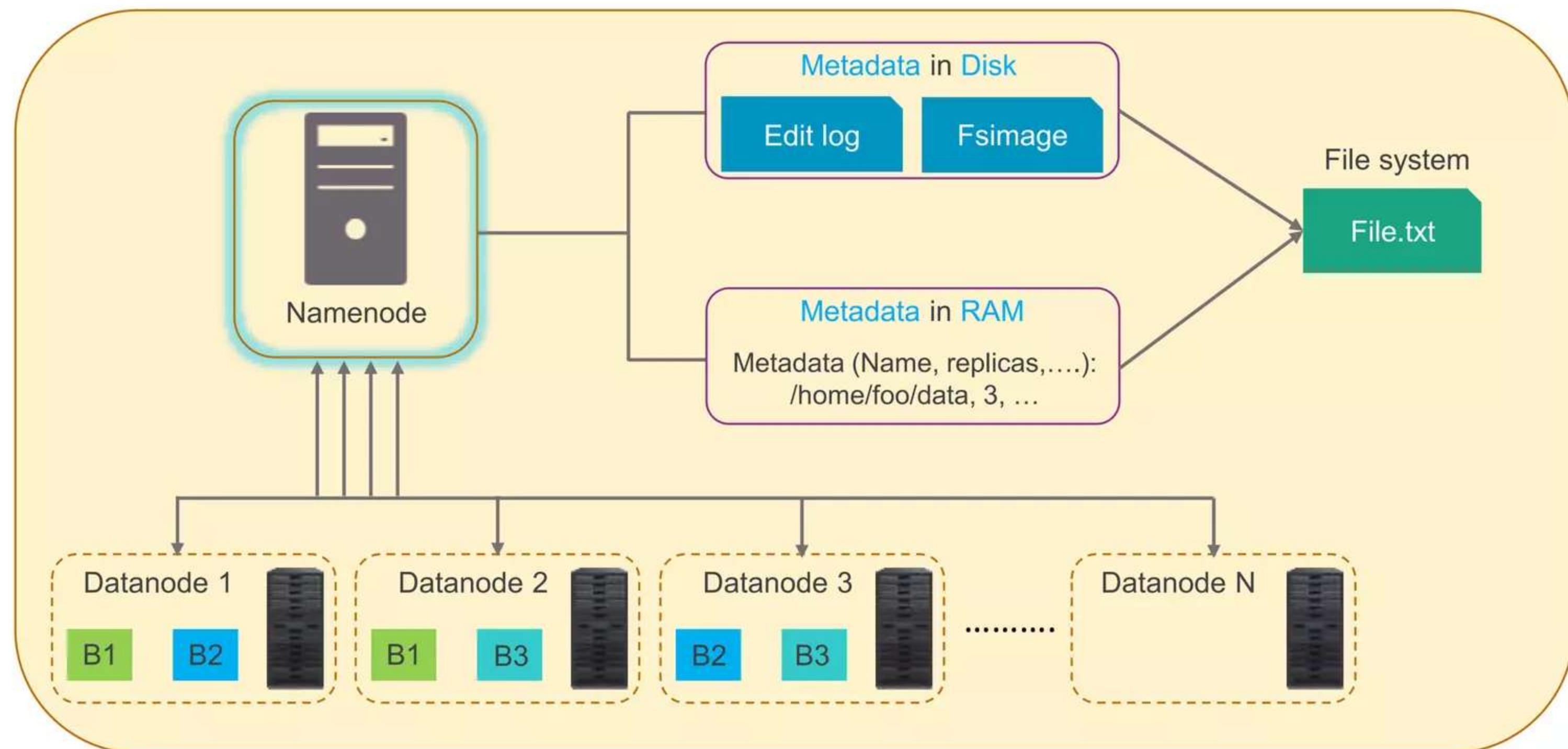
We will still have 2 copies of Block C data on DN 4 of Rack 1 and DN 9 of Rack 3

HDFS Architecture

HDFS Architecture

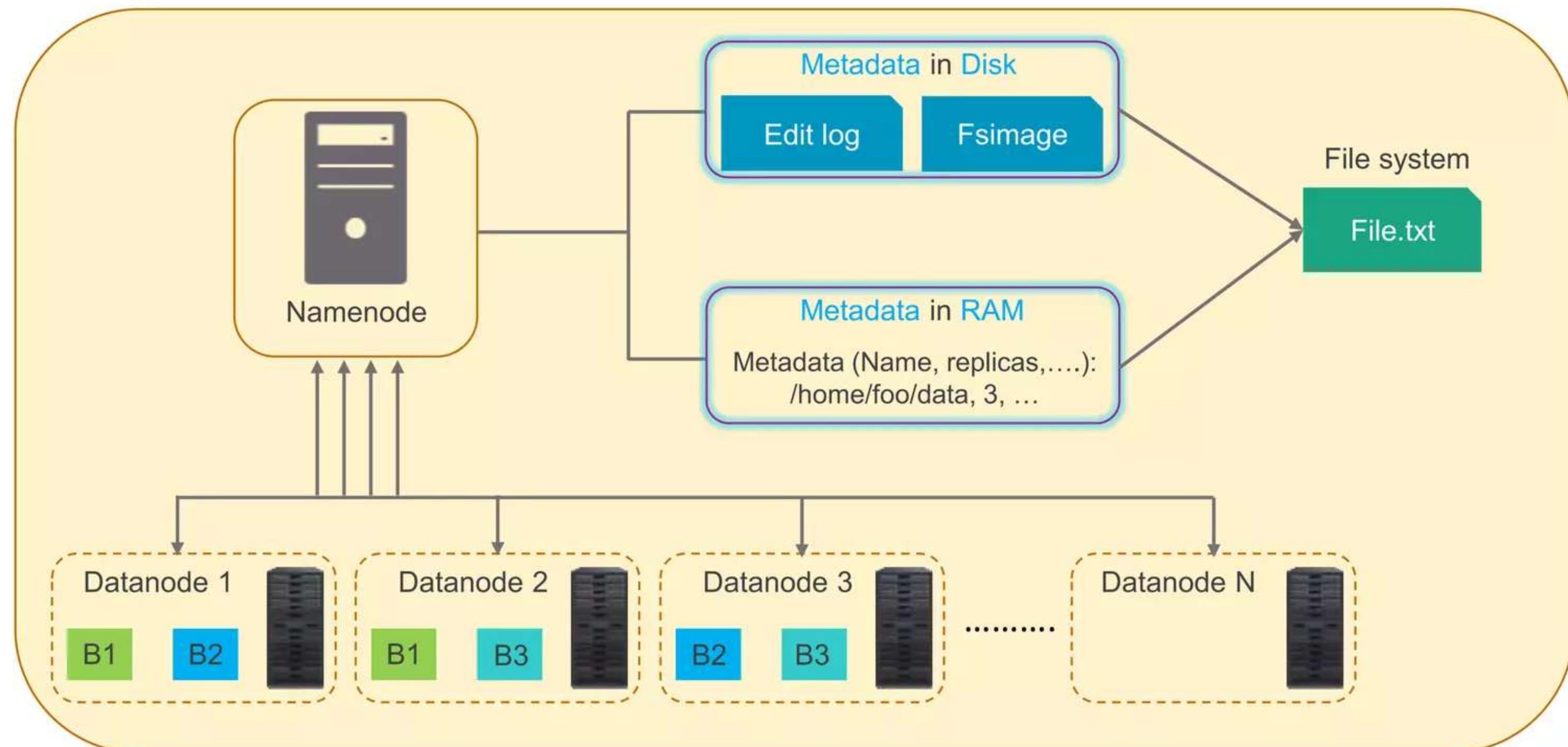


HDFS - Namenode



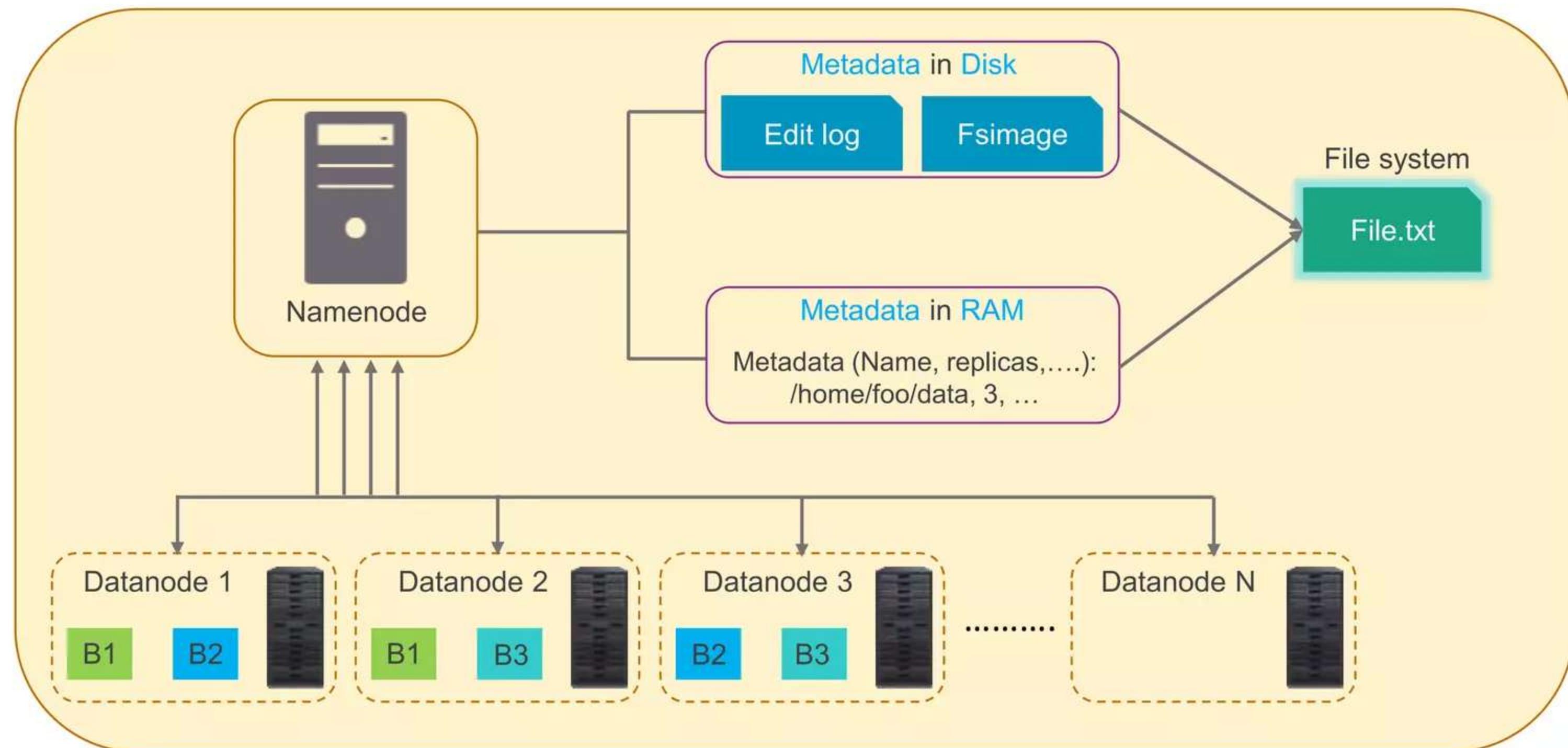
Namenode is the master server. In a non high availability cluster, there can be only one Namenode. In a Hadoop cluster, 2 Namenodes are possible

HDFS - Namenode



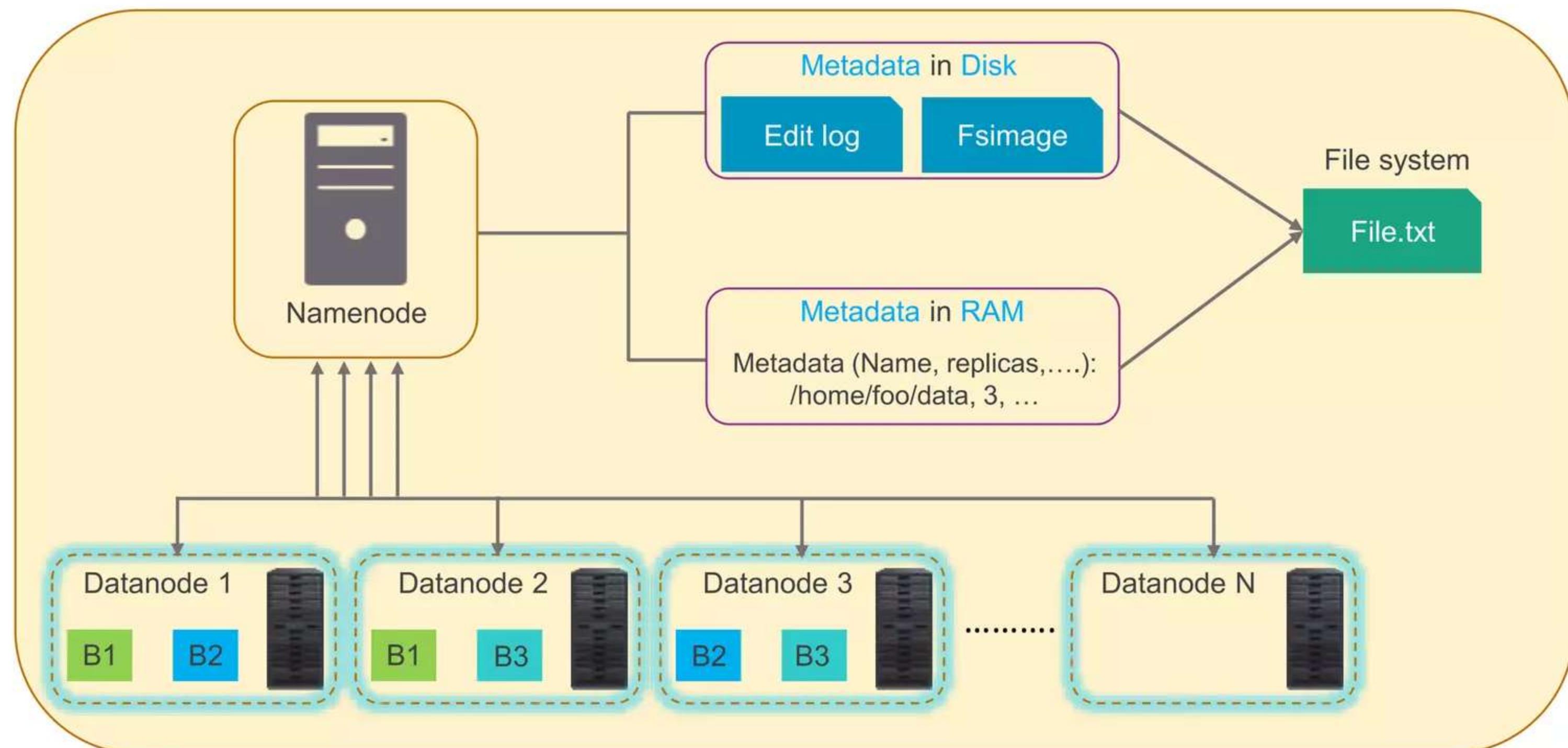
Namenode holds metadata information about the various Datanodes, their location, the size of each block, etc.

HDFS - Namenode



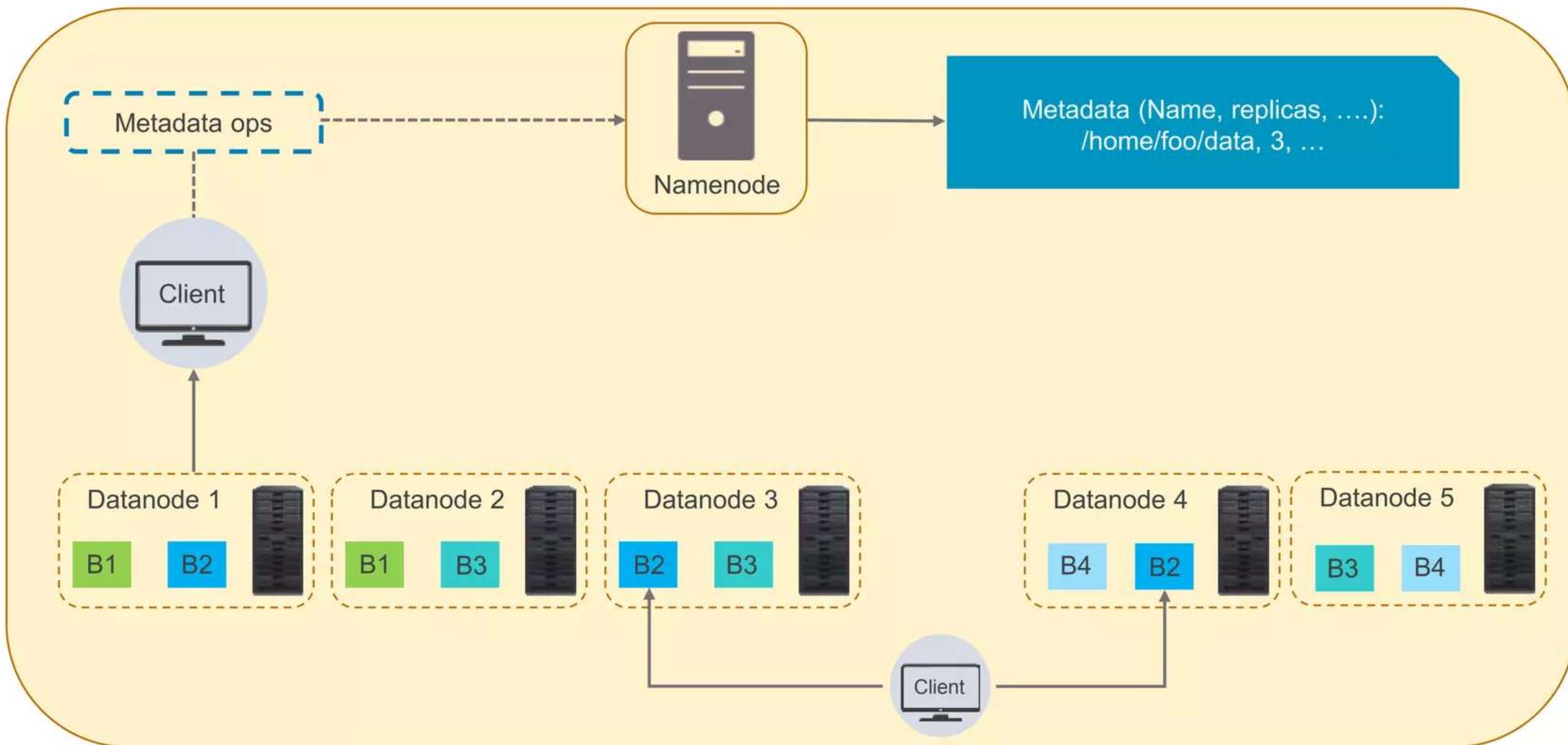
Helps to execute file system namespace operations –
opening, closing, renaming files and directories

HDFS - Namenode



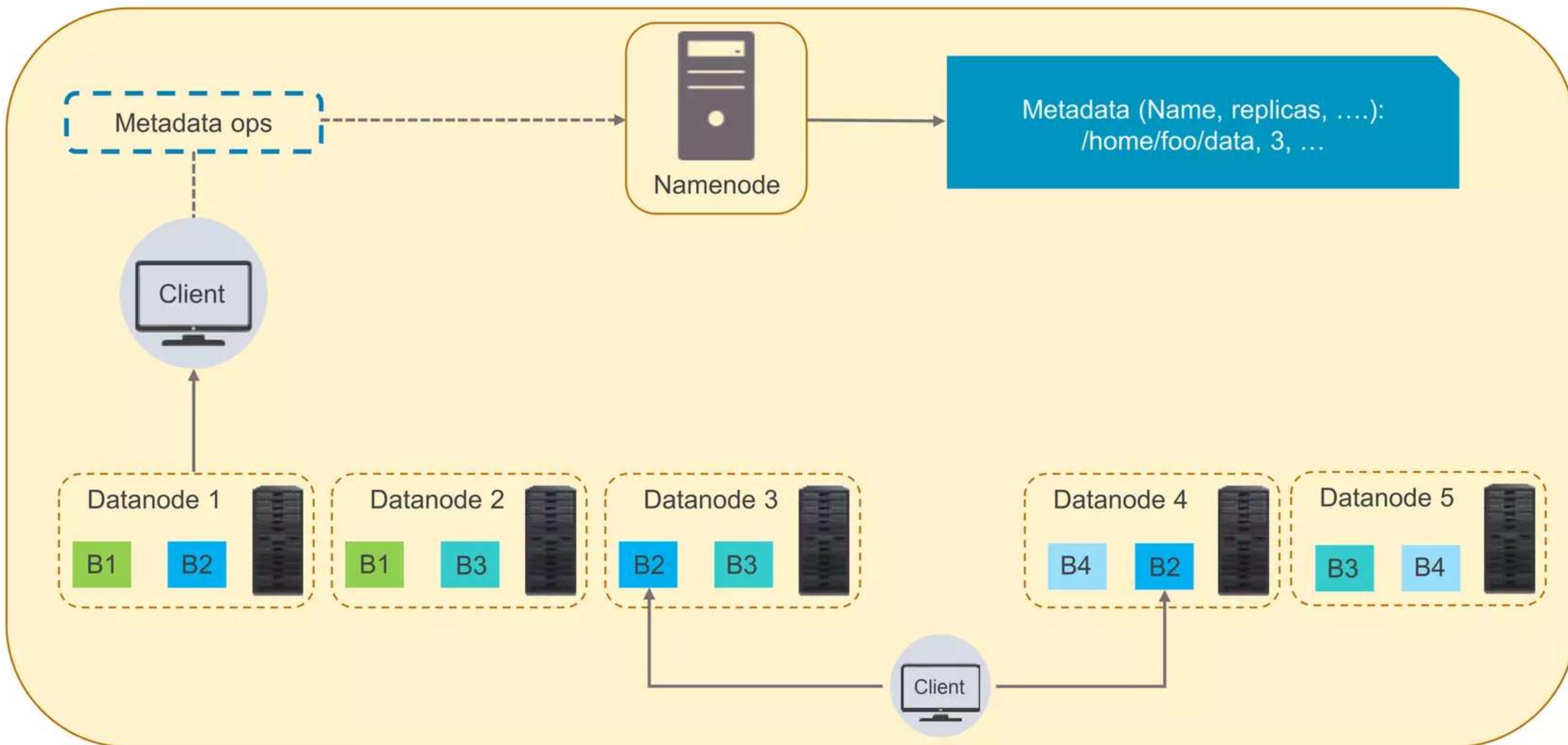
Datanodes send block reports to Namenode every 10 seconds

HDFS - Datanode

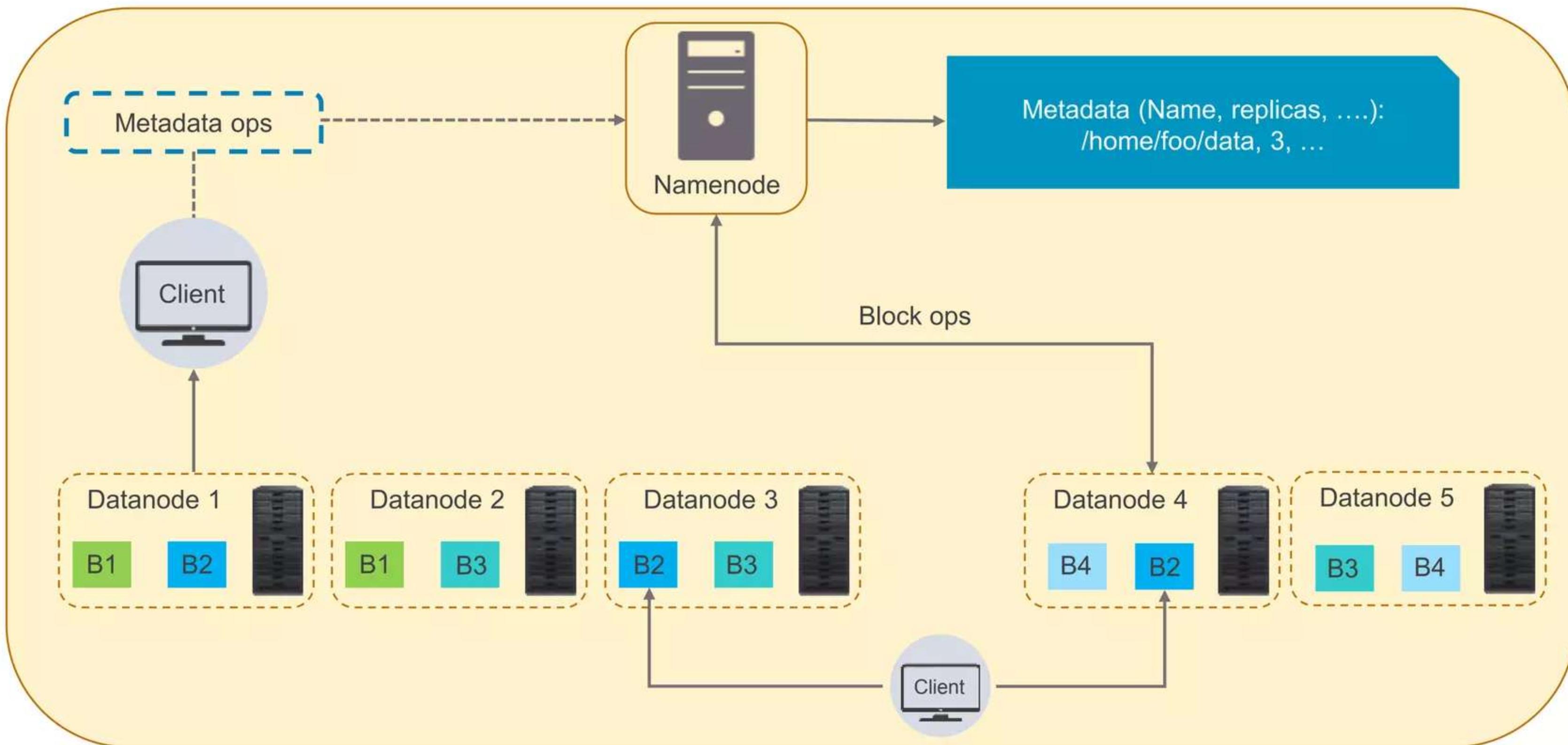


Datanode is a multiple instance server. There can be N number of Datanode servers

HDFS - Datanode

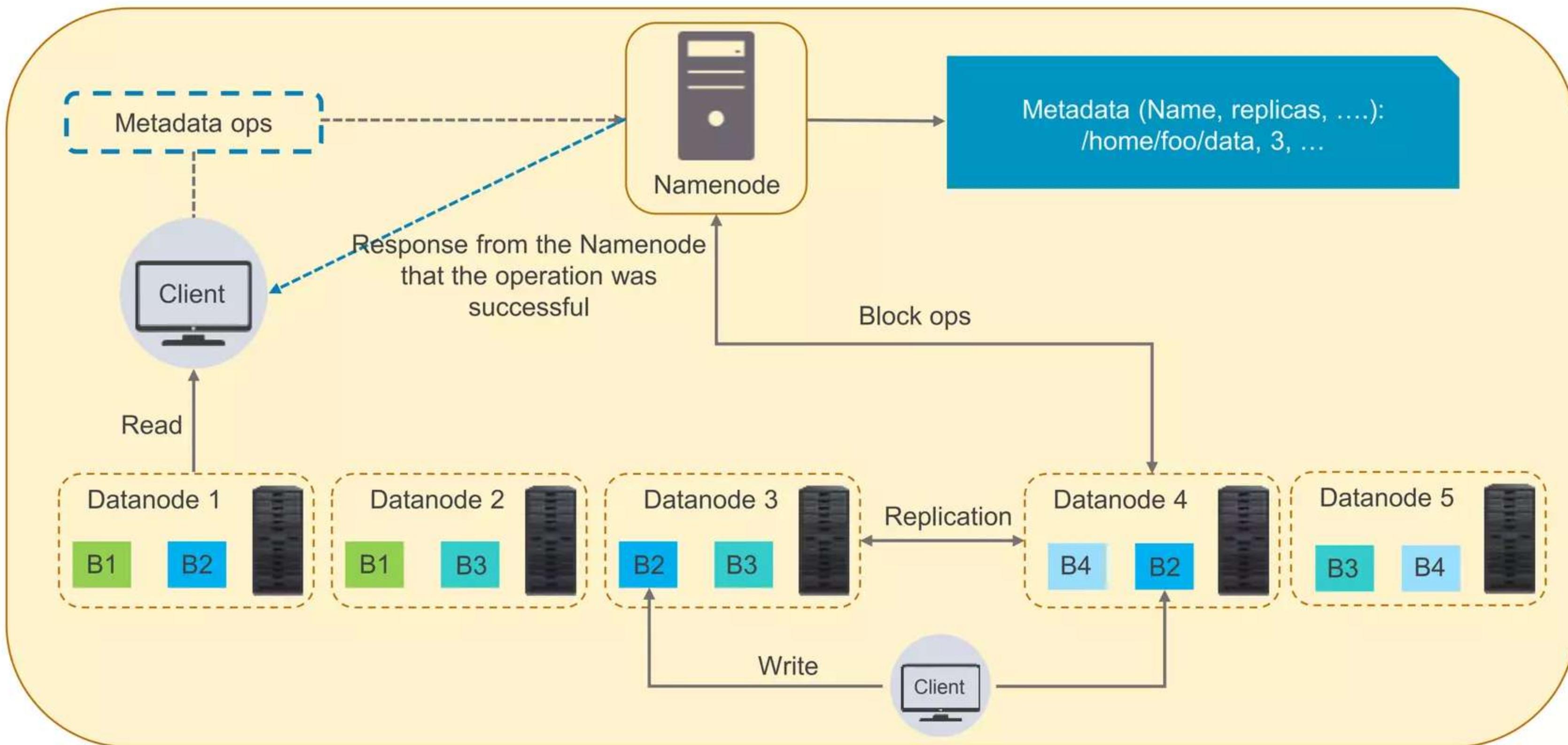


HDFS - Datanode



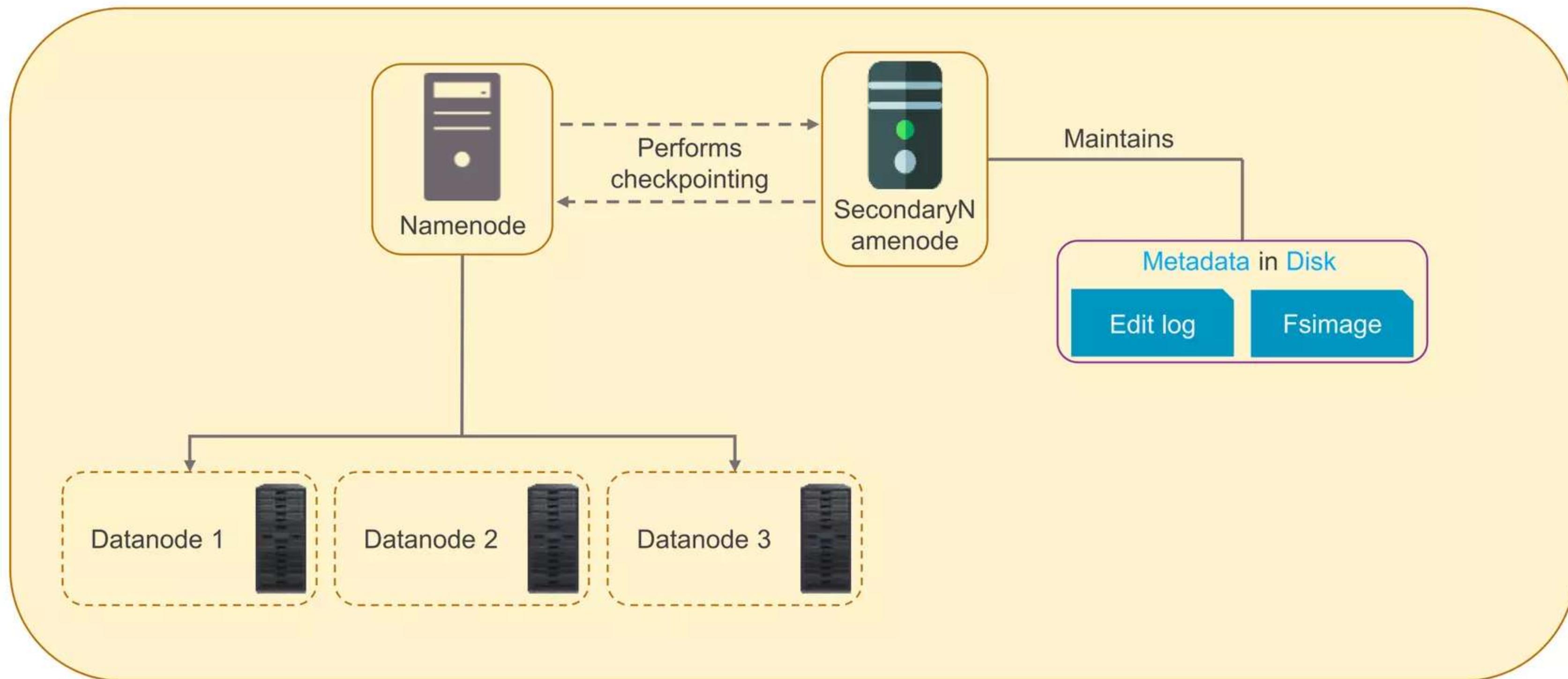
Datanode stores and retrieves the blocks when asked by the Namenode

HDFS - Datanode



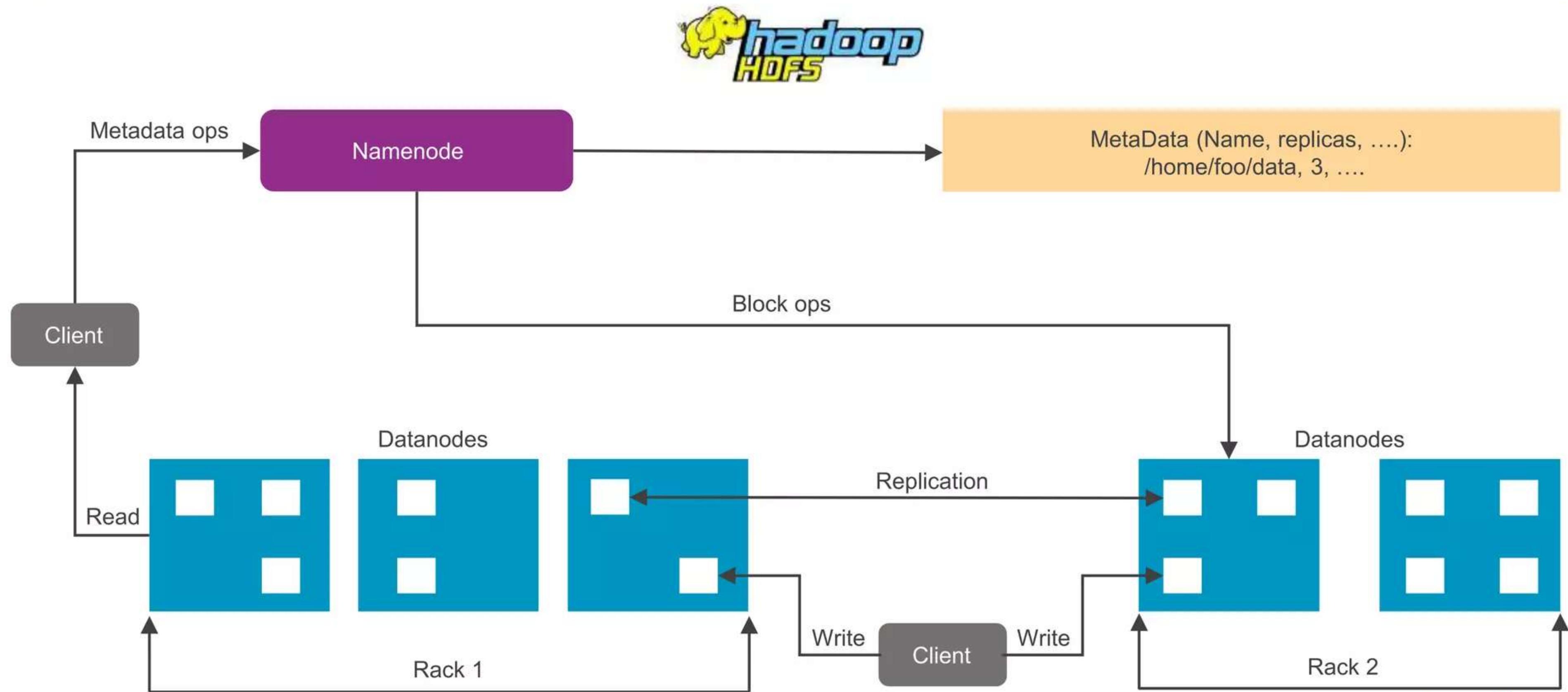
It reads and writes client's request and performs block creation, deletion and replication on instruction from the Namenode

HDFS – Secondary Namenode

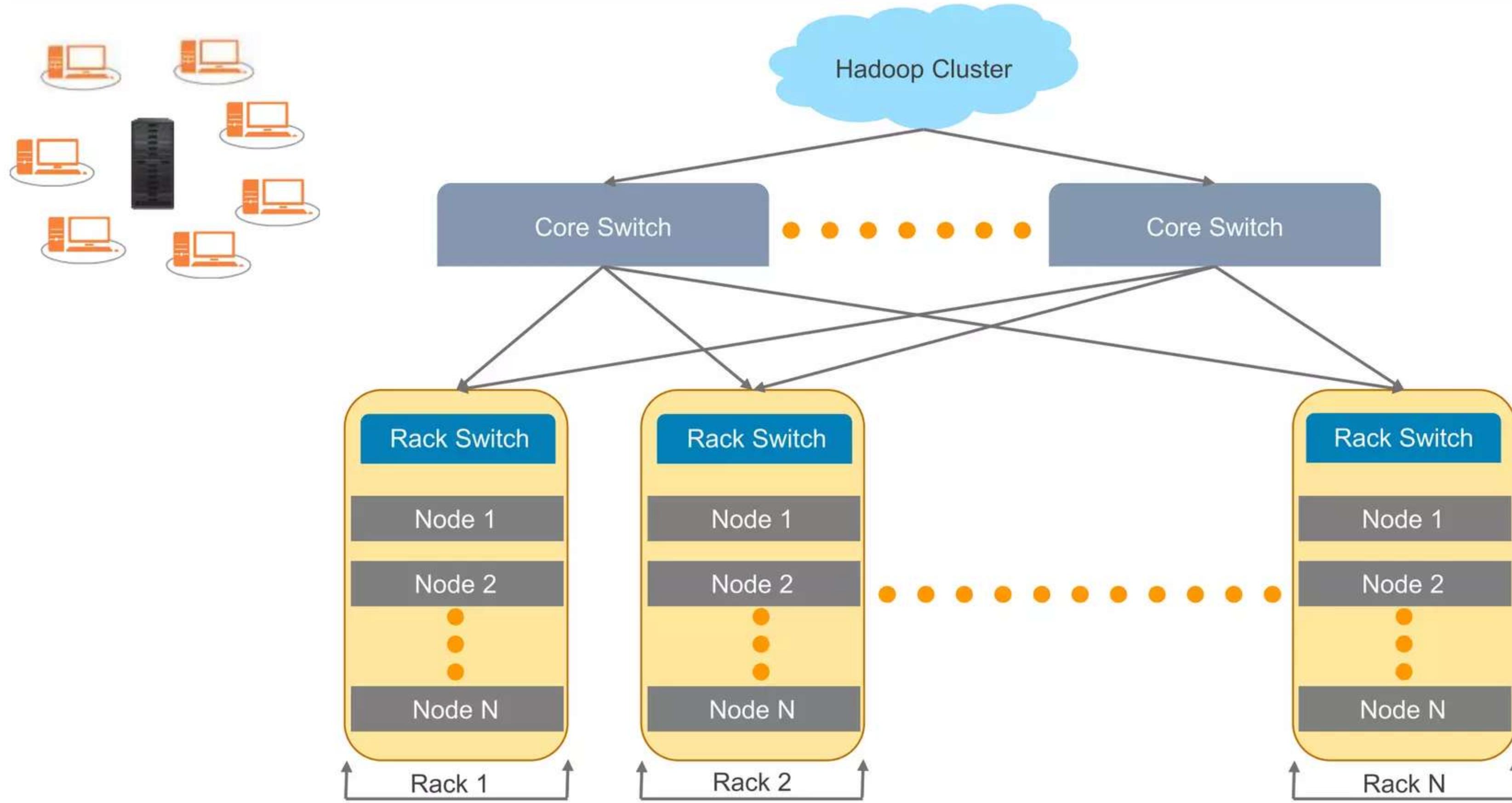


Secondary Namenode server is responsible for maintaining a copy of
Metadata in disk

HDFS Architecture

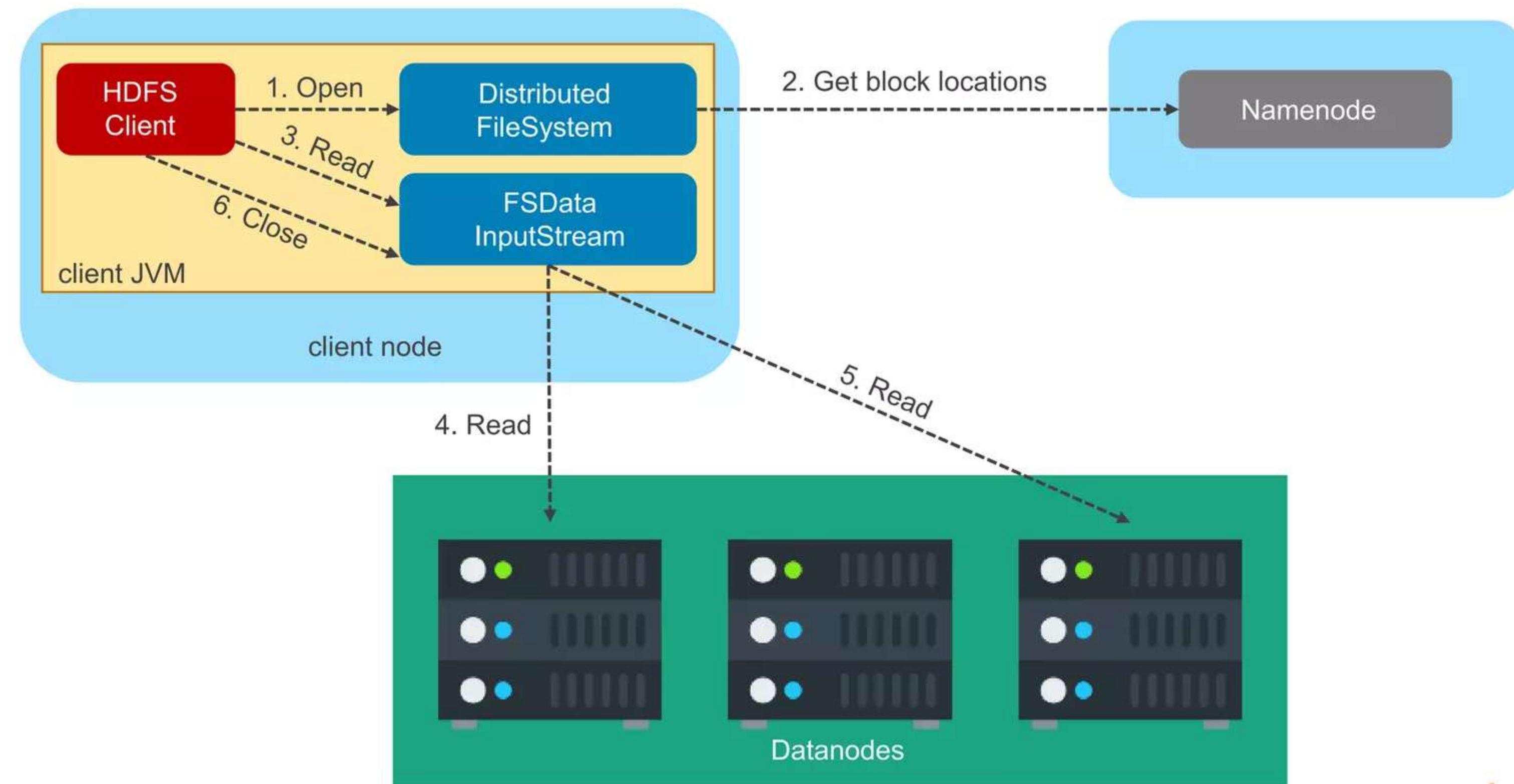


Hadoop Cluster – Rack Based Architecture

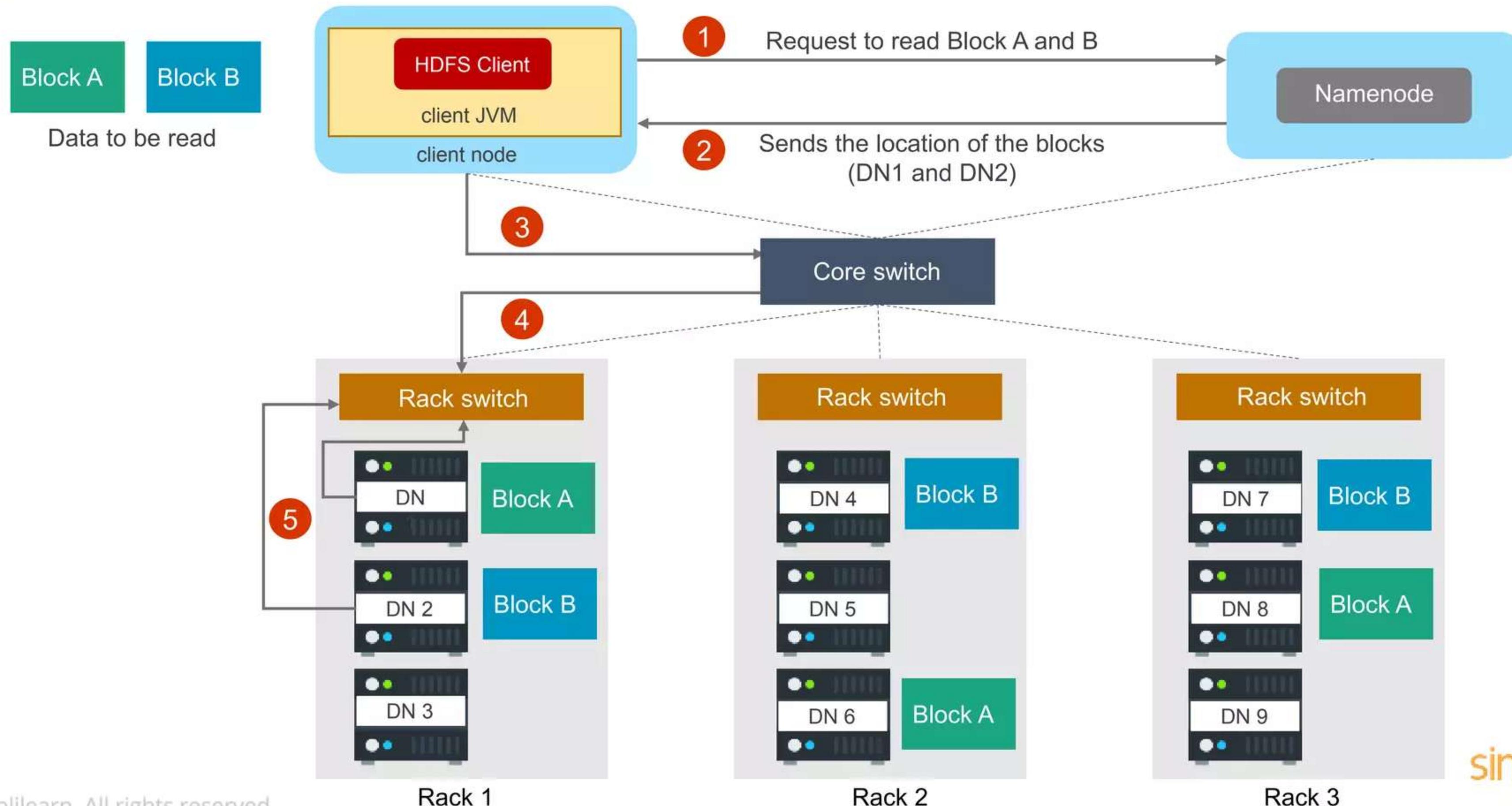


HDFS Read Mechanism

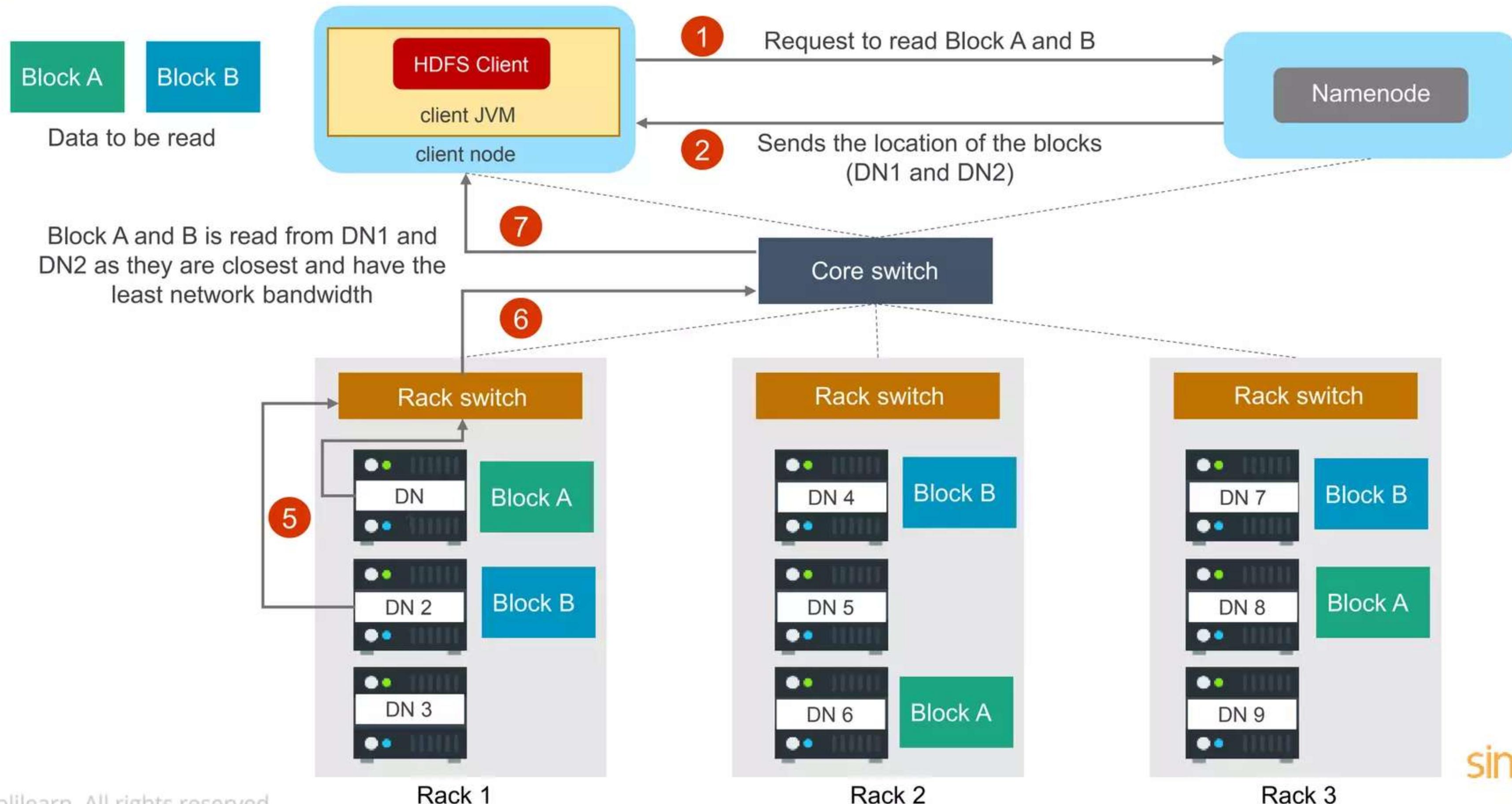
HDFS Read Mechanism



HDFS Read Mechanism

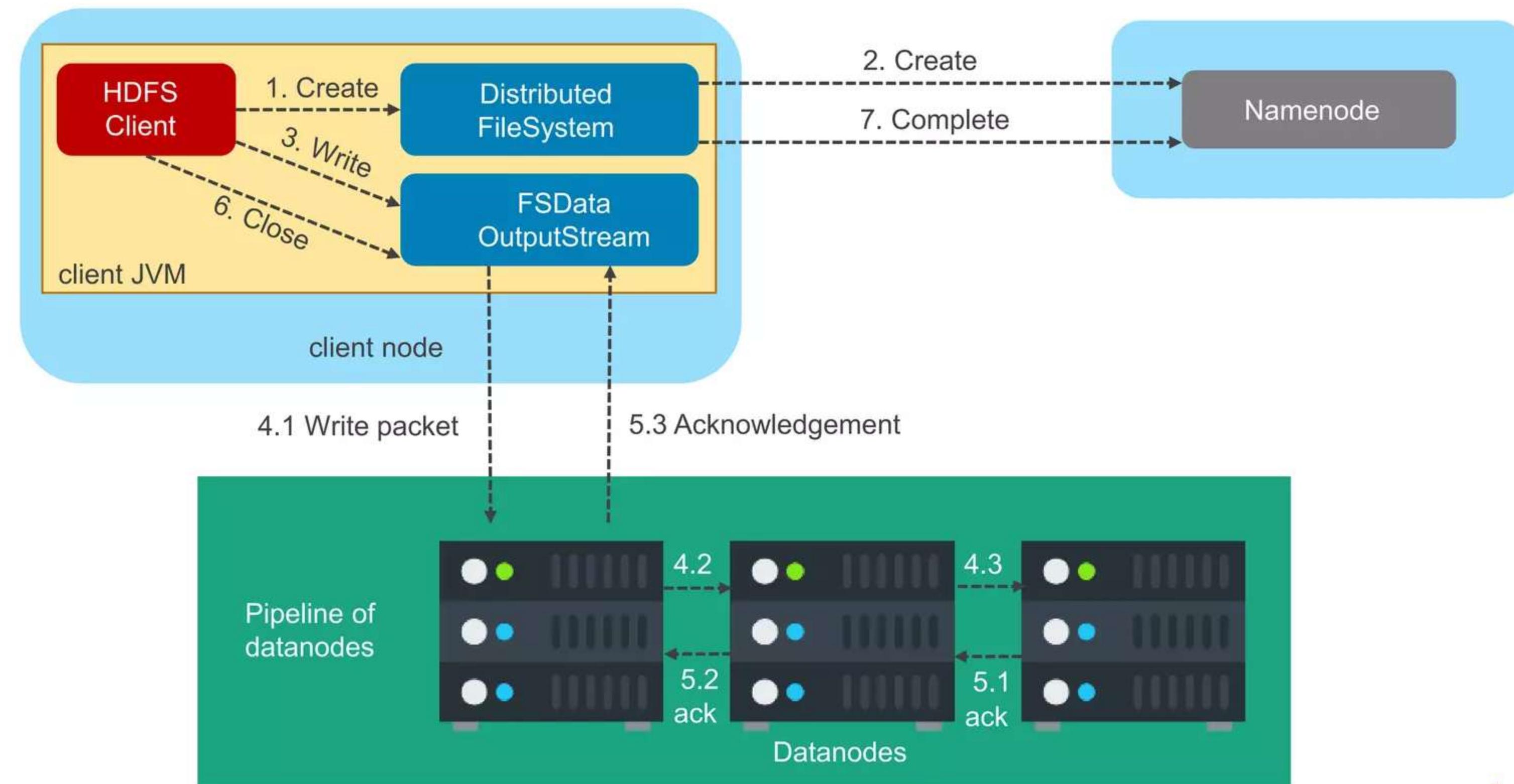


HDFS Read Mechanism

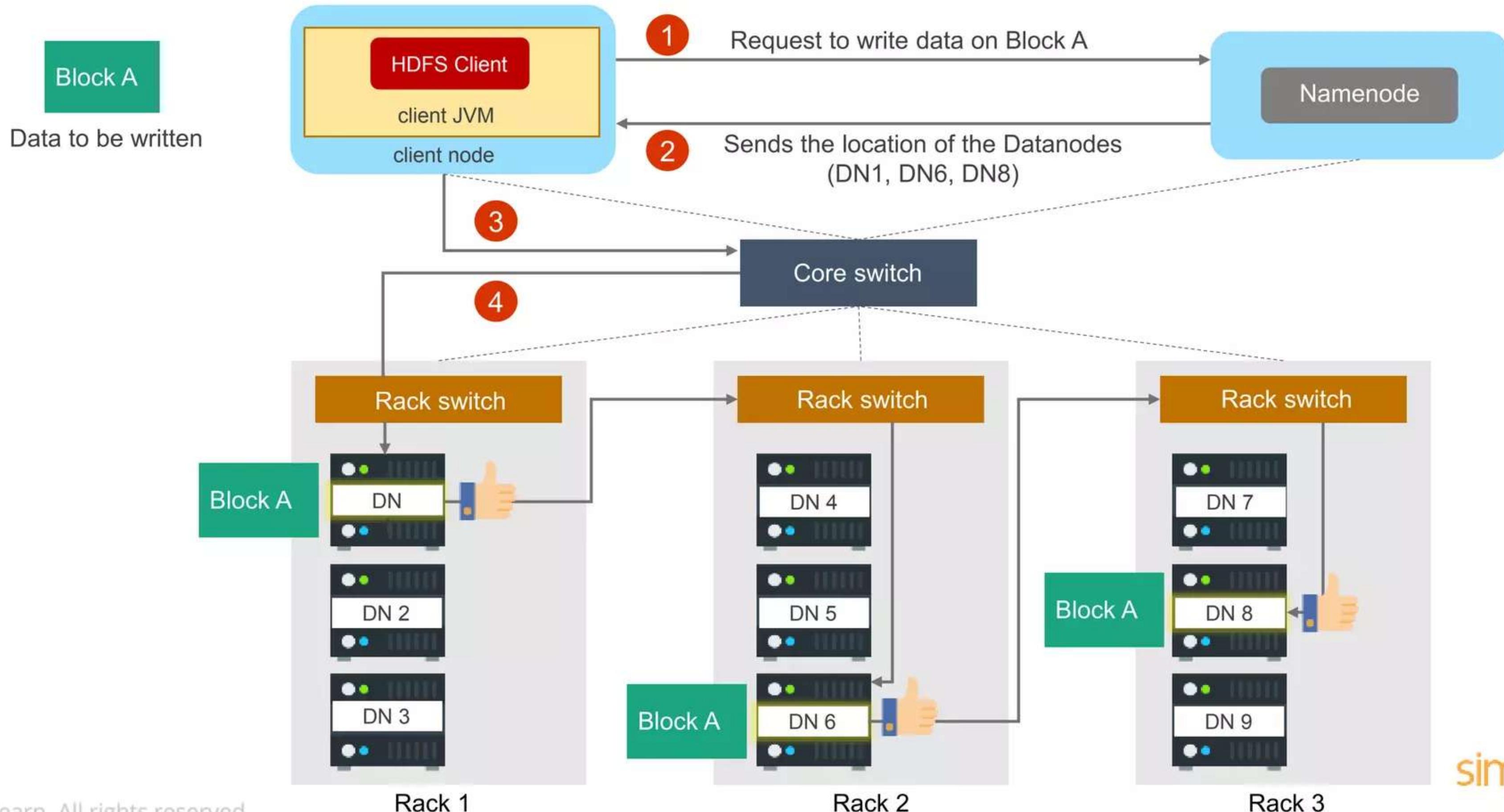


HDFS Write Mechanism

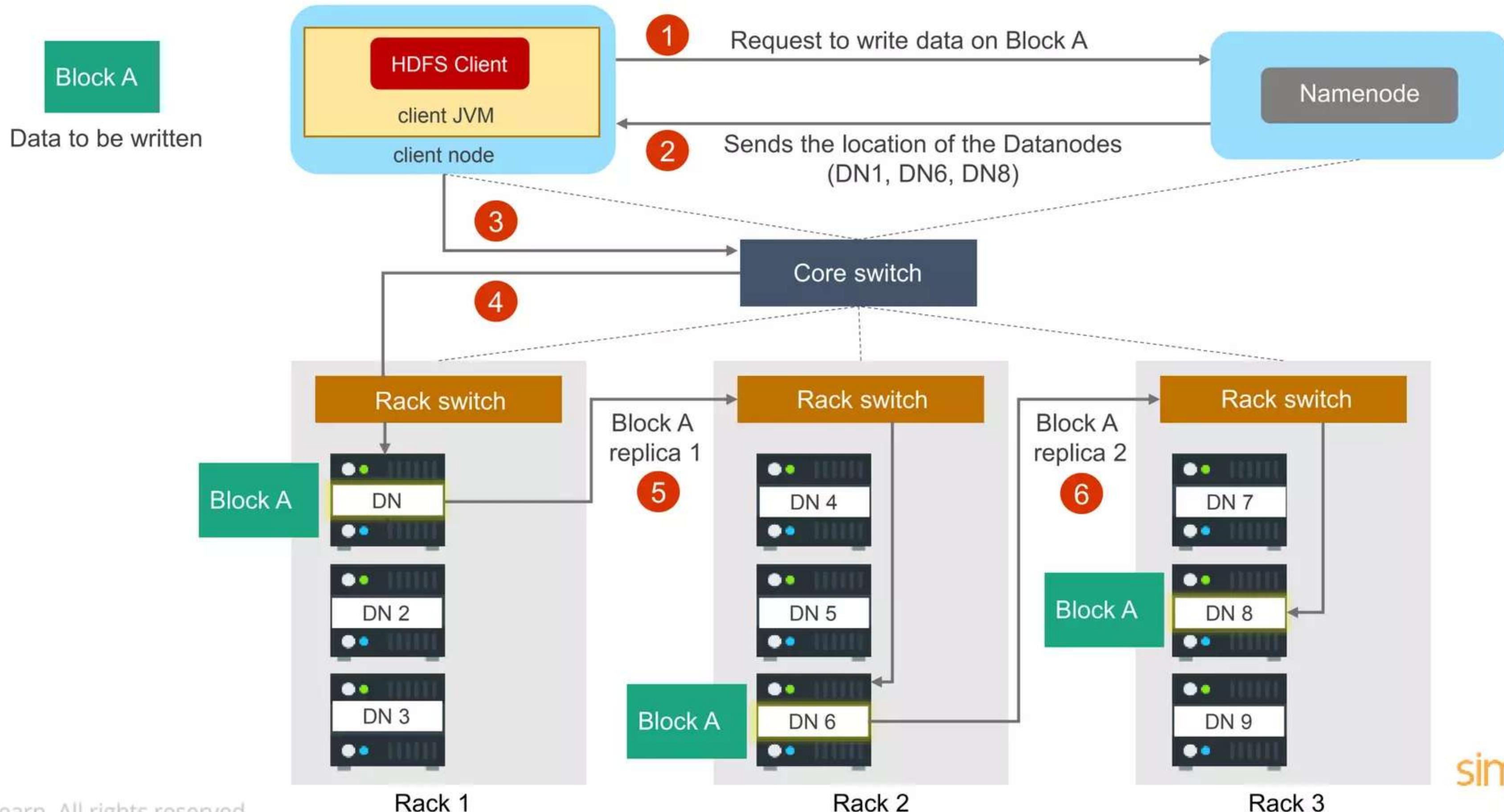
HDFS Write Mechanism



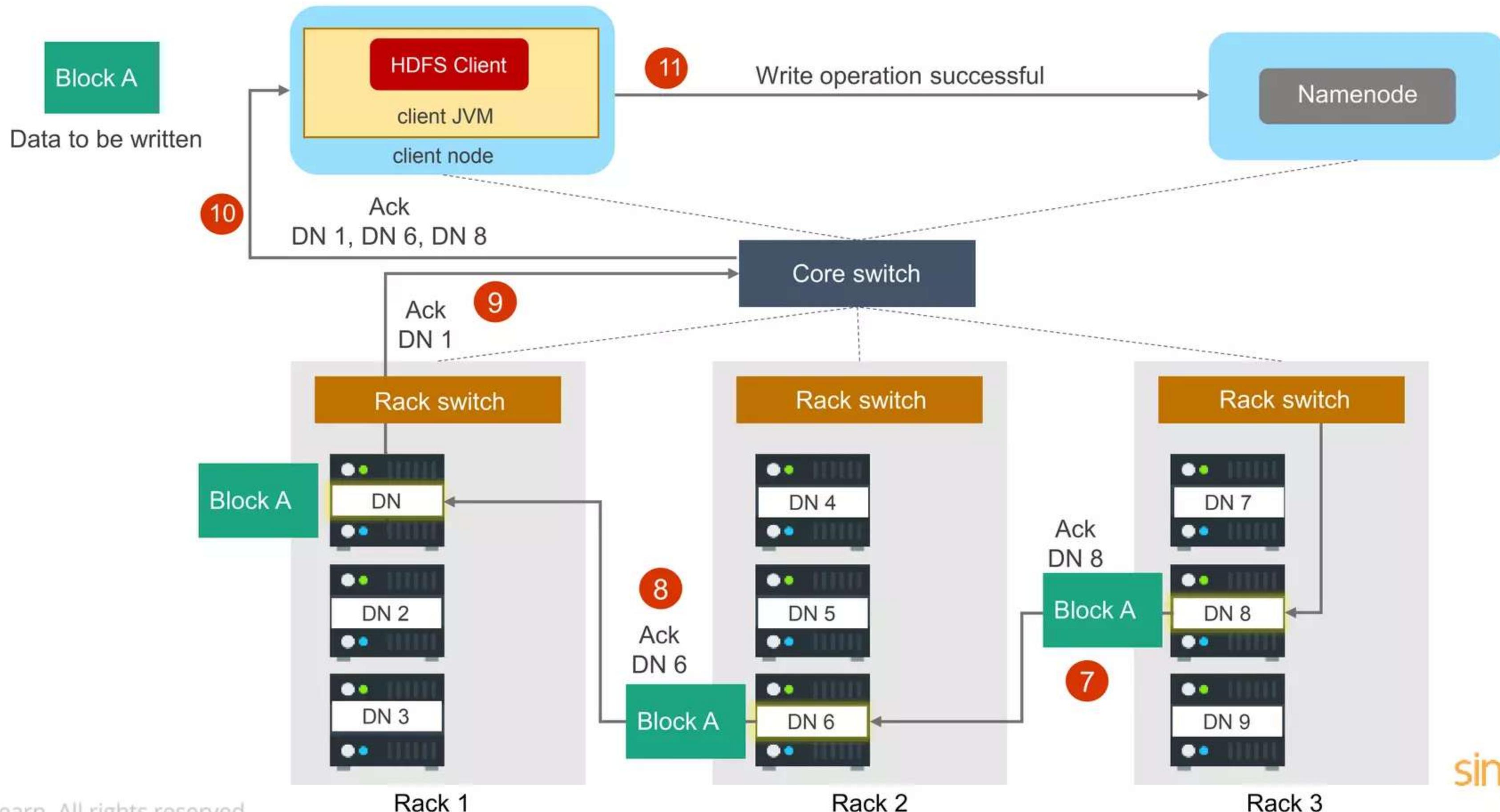
HDFS Write Mechanism



HDFS Write Mechanism



HDFS Write Mechanism



Hadoop MapReduce

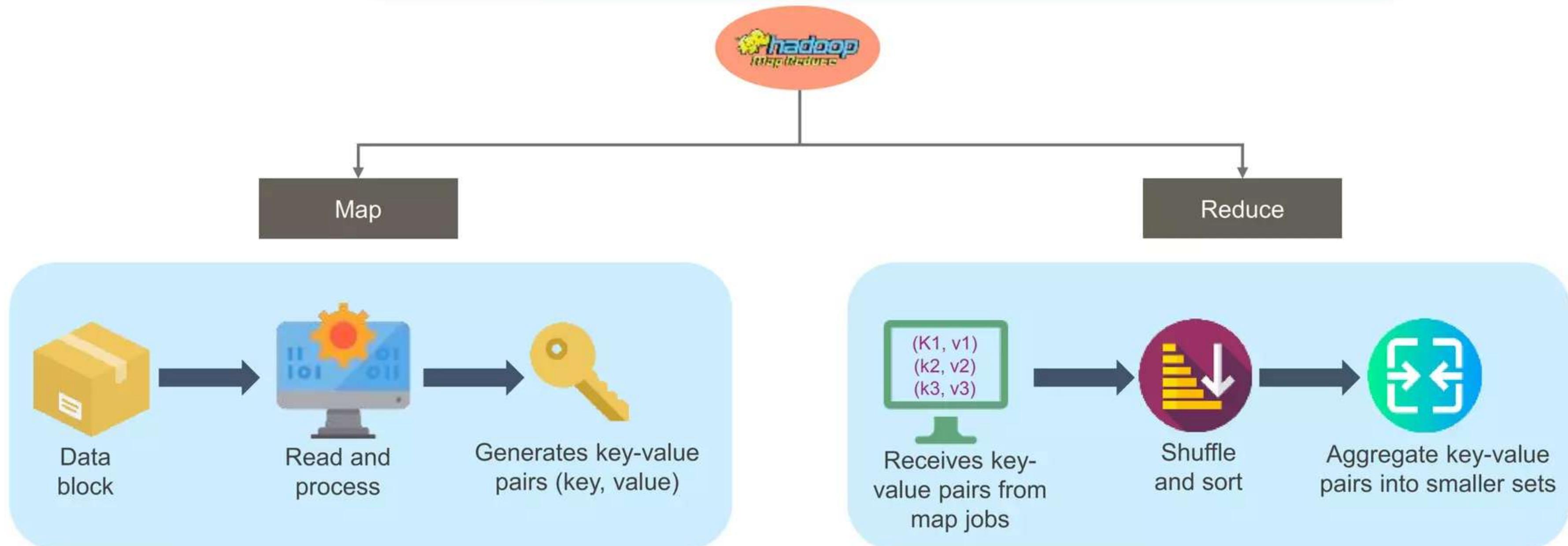
simplilearn

90

Hadoop MapReduce



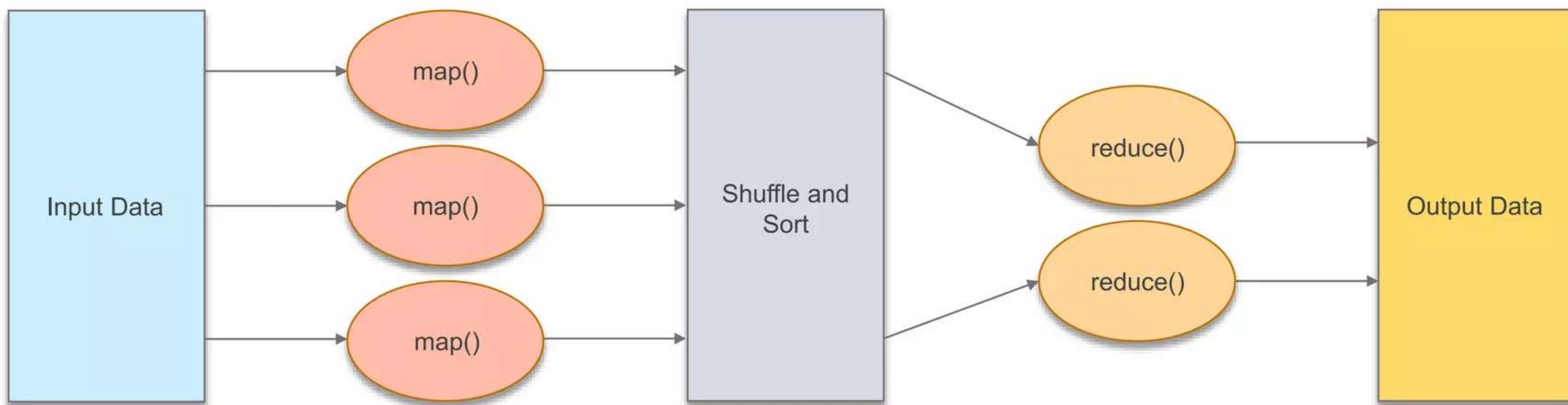
MapReduce is a framework that performs distributed and parallel processing of large volumes of data



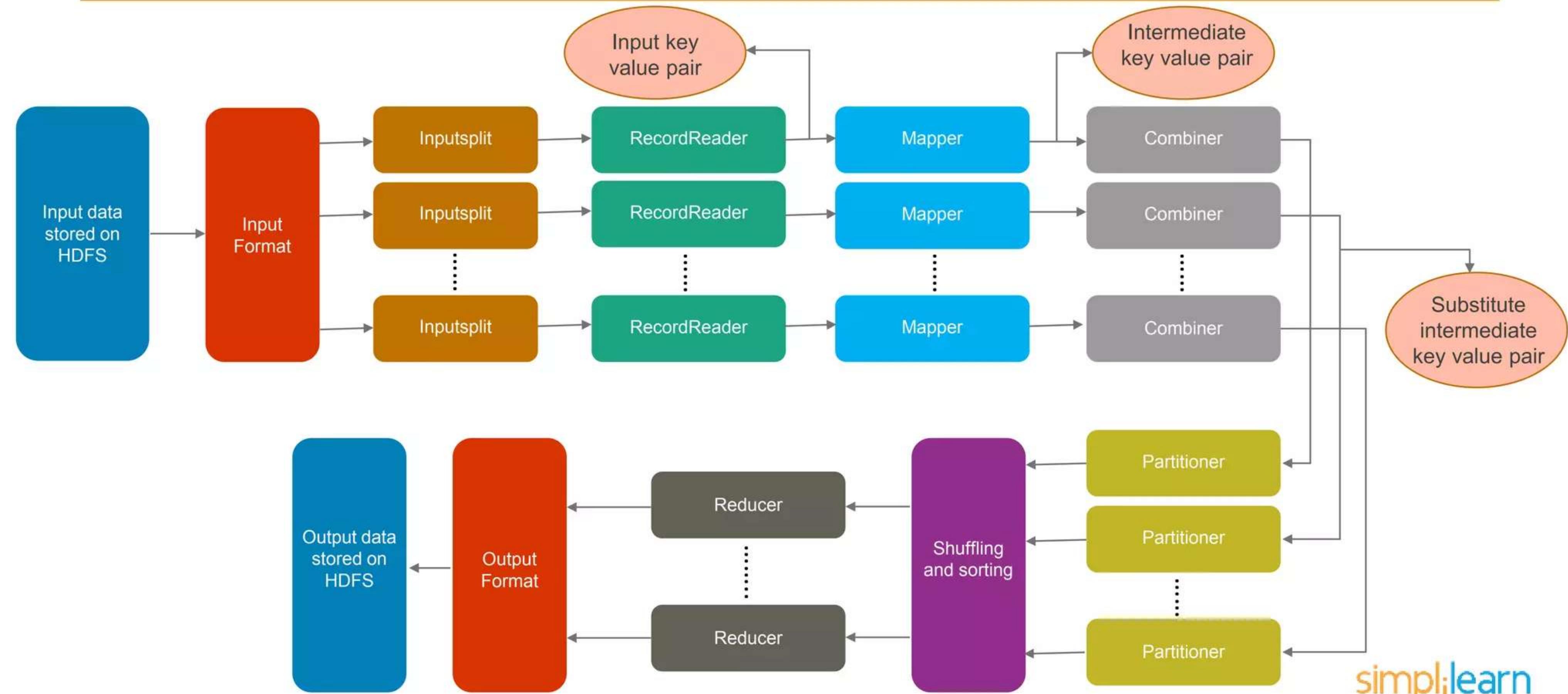
Hadoop MapReduce



MapReduce is a framework that performs distributed and parallel processing of large volumes of data



MapReduce Job Execution



MapReduce Example

Input

Big data comes in various formats. This data can be stored in multiple data servers

Split

Big data comes in various formats

Map

Big, 1
data, 1
comes, 1
in, 1
various, 1
formats, 1

This data can be stored in multiple data servers

This, 1
data, 1
can, 1
be, 1
stored, 1
in, 1
multiple, 1
data, 1
servers, 1

Shuffle

be, (1)
Big, (1)
be, (1)
can, (1)
comes, (1)
data, (1,1)
formats, (1)
in, (1,1)
multiple, (1)
servers, (1)
stored, (1)
This, (1)
various, (1)

Reduce

be, (1)
Big, (1)
be, (1)
can, (1)
comes, (1)
data, (2)
formats, (1)
in, (2)
multiple, (1)
servers, (1)
stored, (1)
This, (1)
various, (1)

Hadoop YARN

simplilearn

90

Hadoop YARN

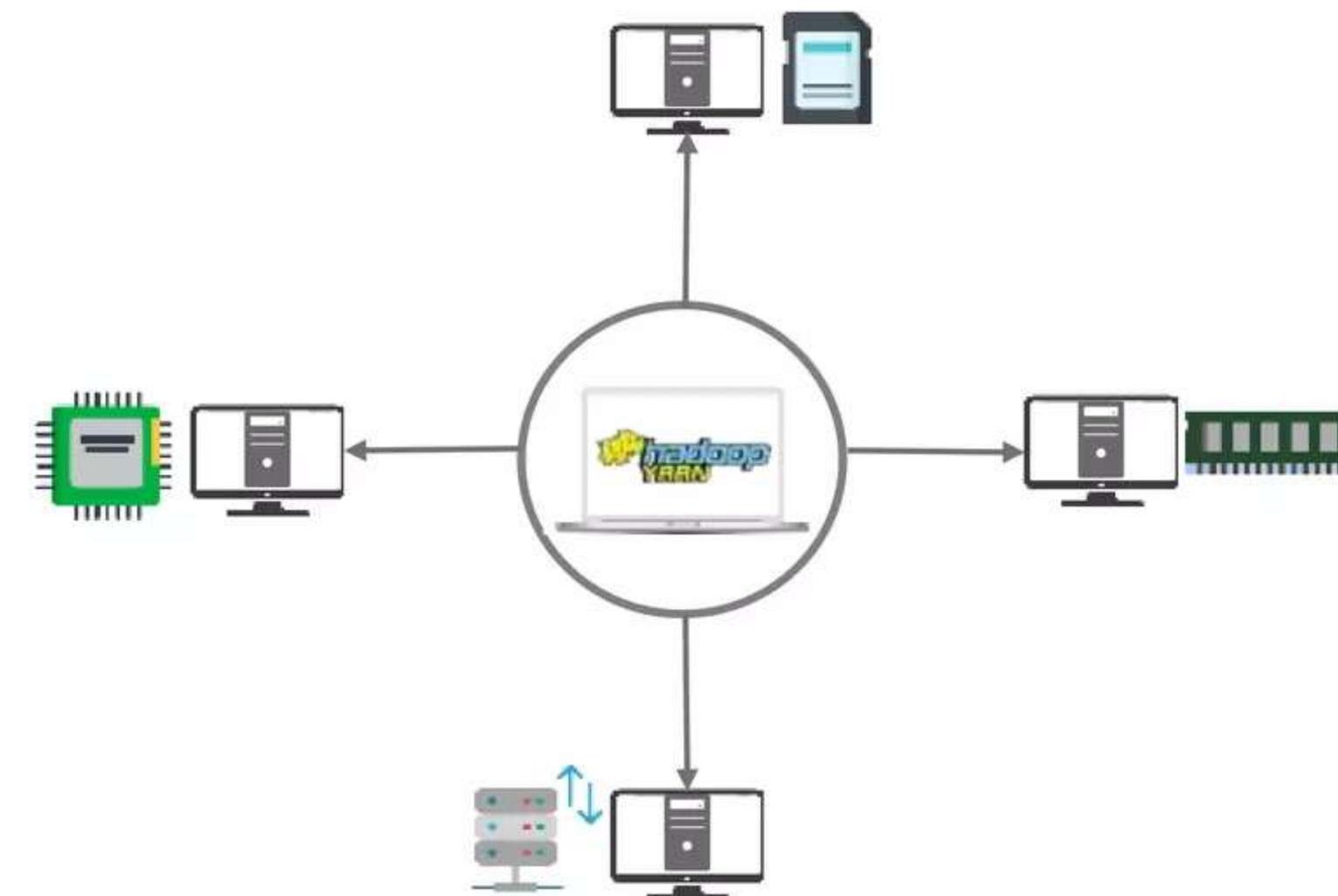


YARN -----> Yet Another Resource Negotiator

Introduced in Hadoop 2.0 version

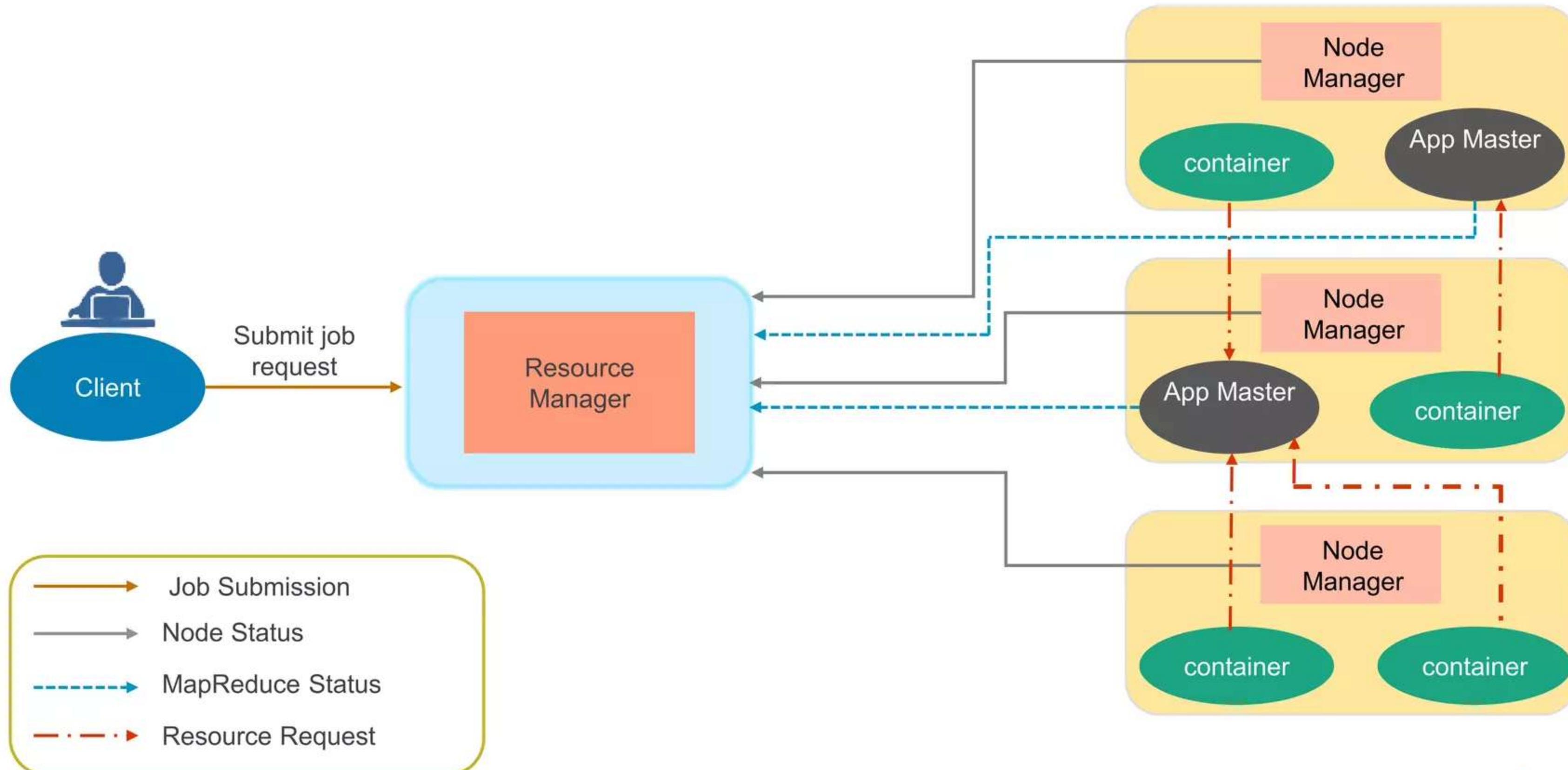


It is the middle layer between HDFS
and MapReduce

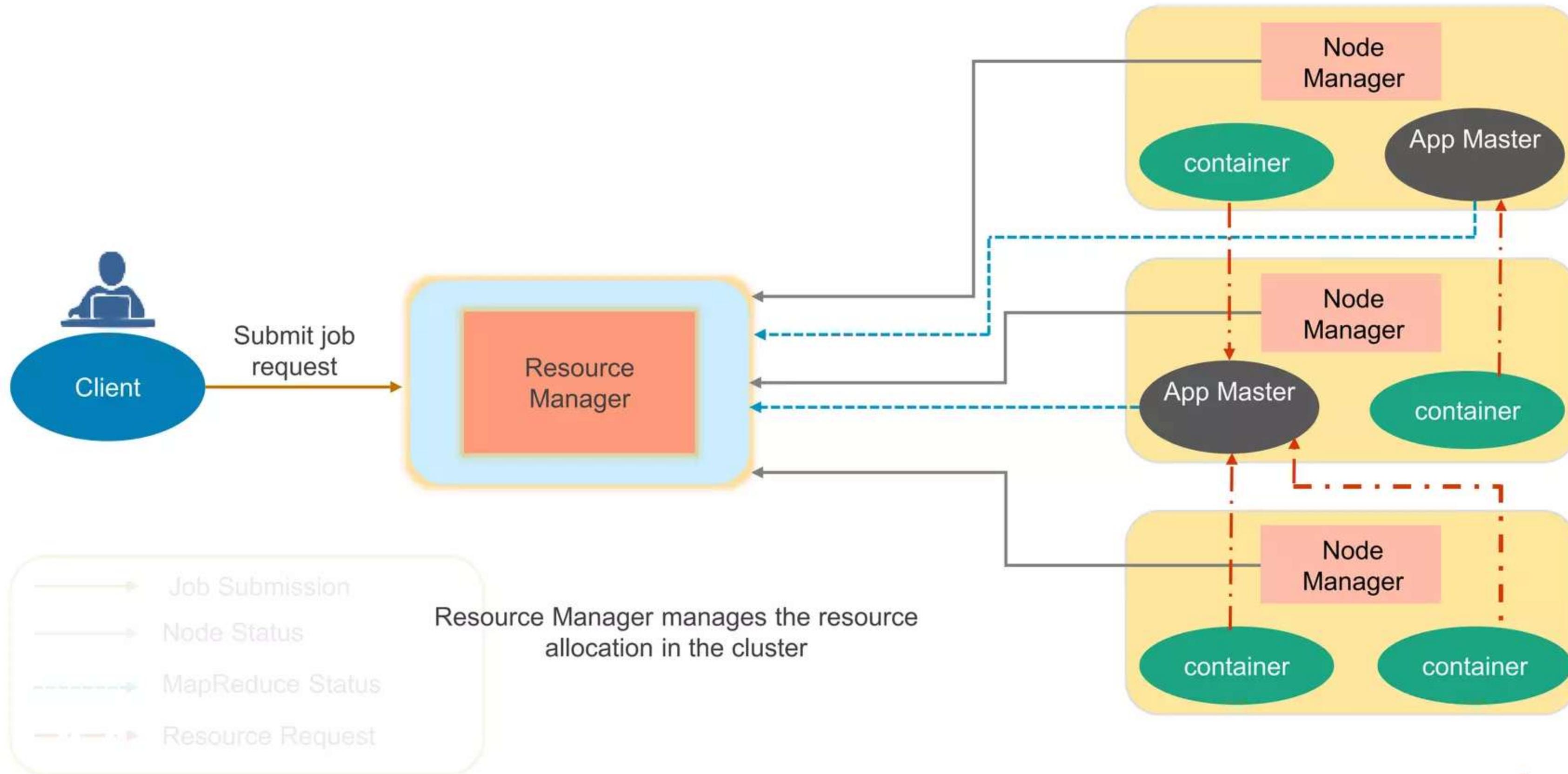


Manages cluster resources (memory,
network bandwidth, disk IO, CPU)

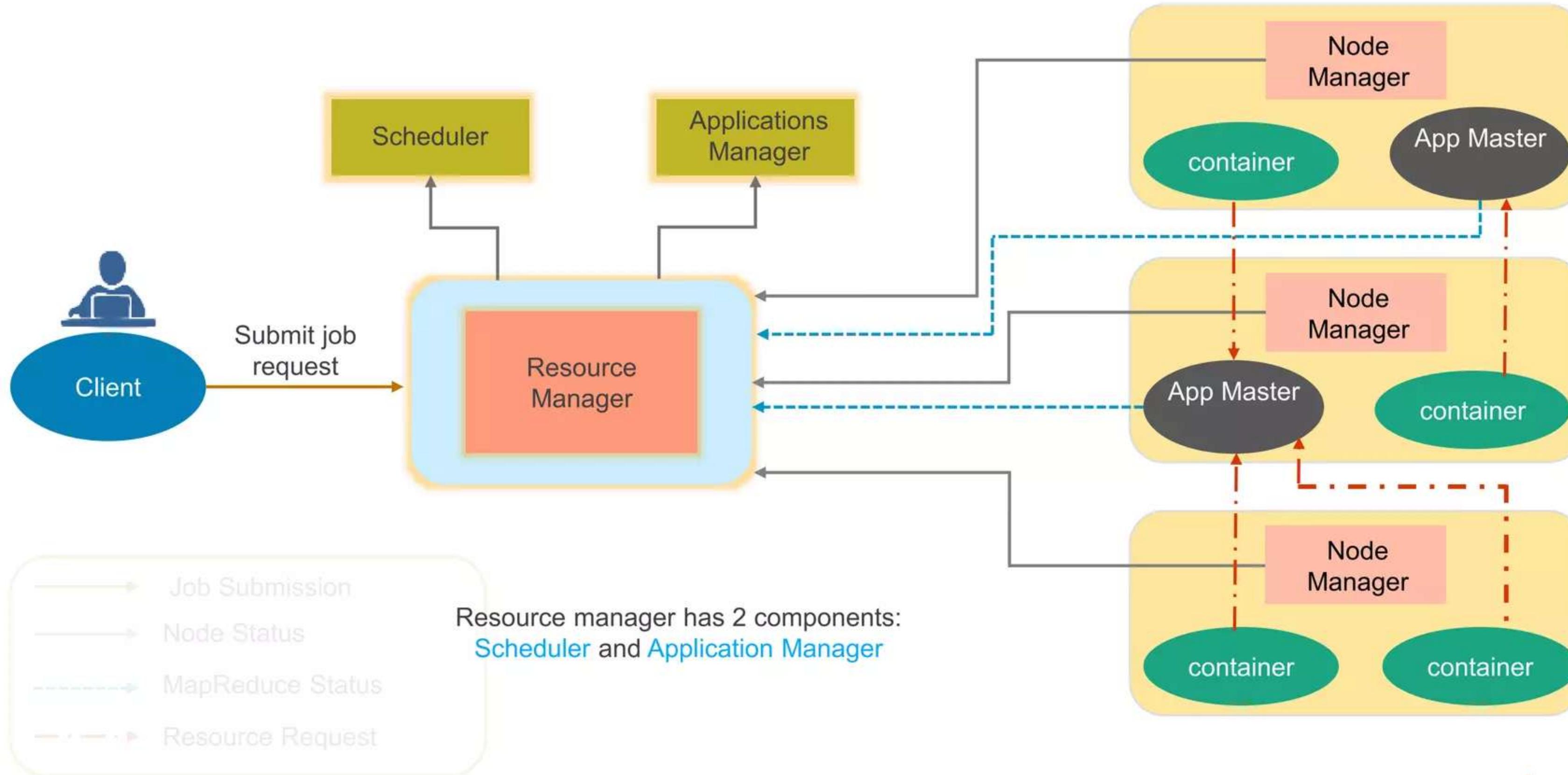
YARN Architecture



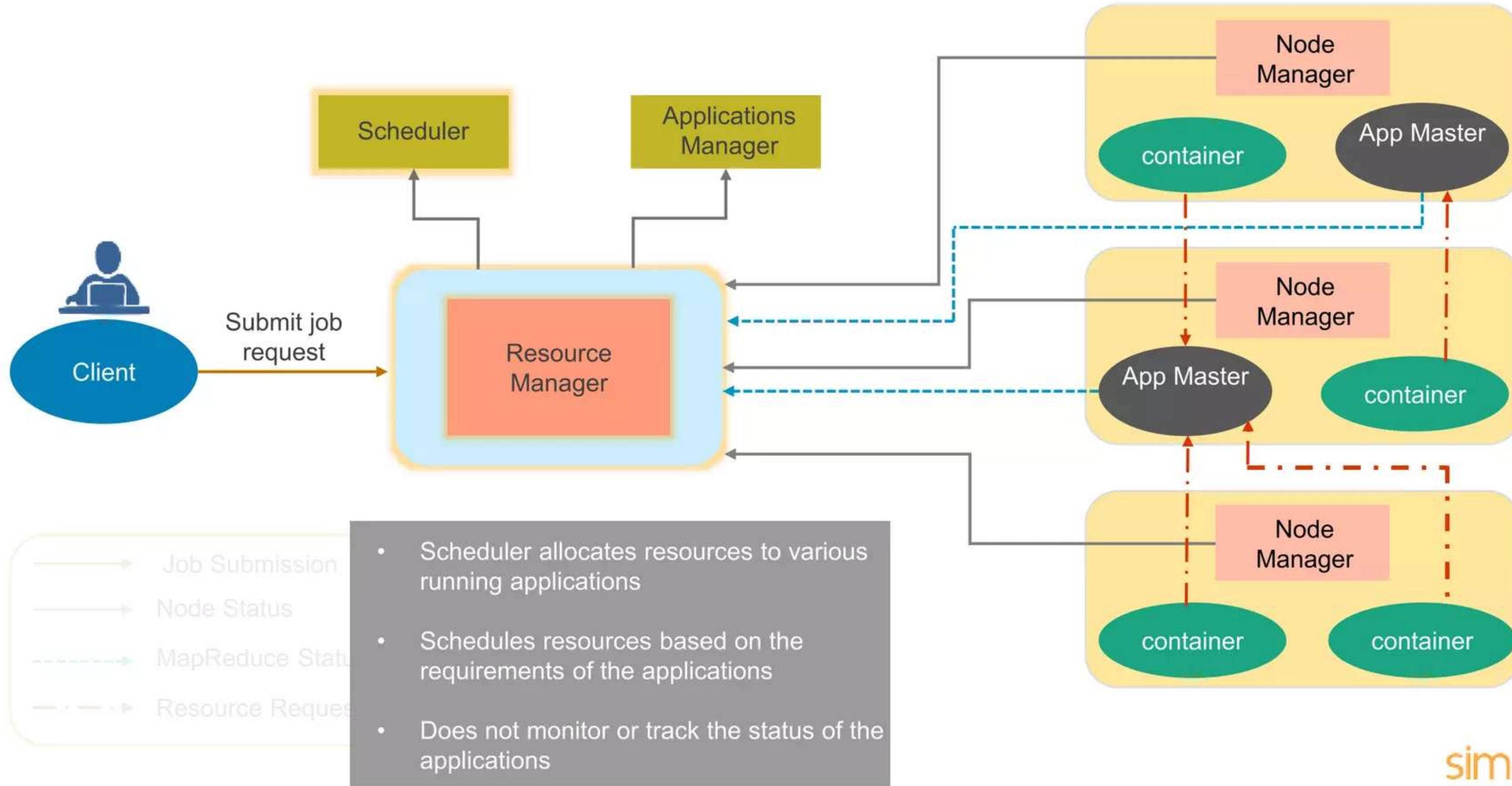
YARN Architecture – Resource Manager



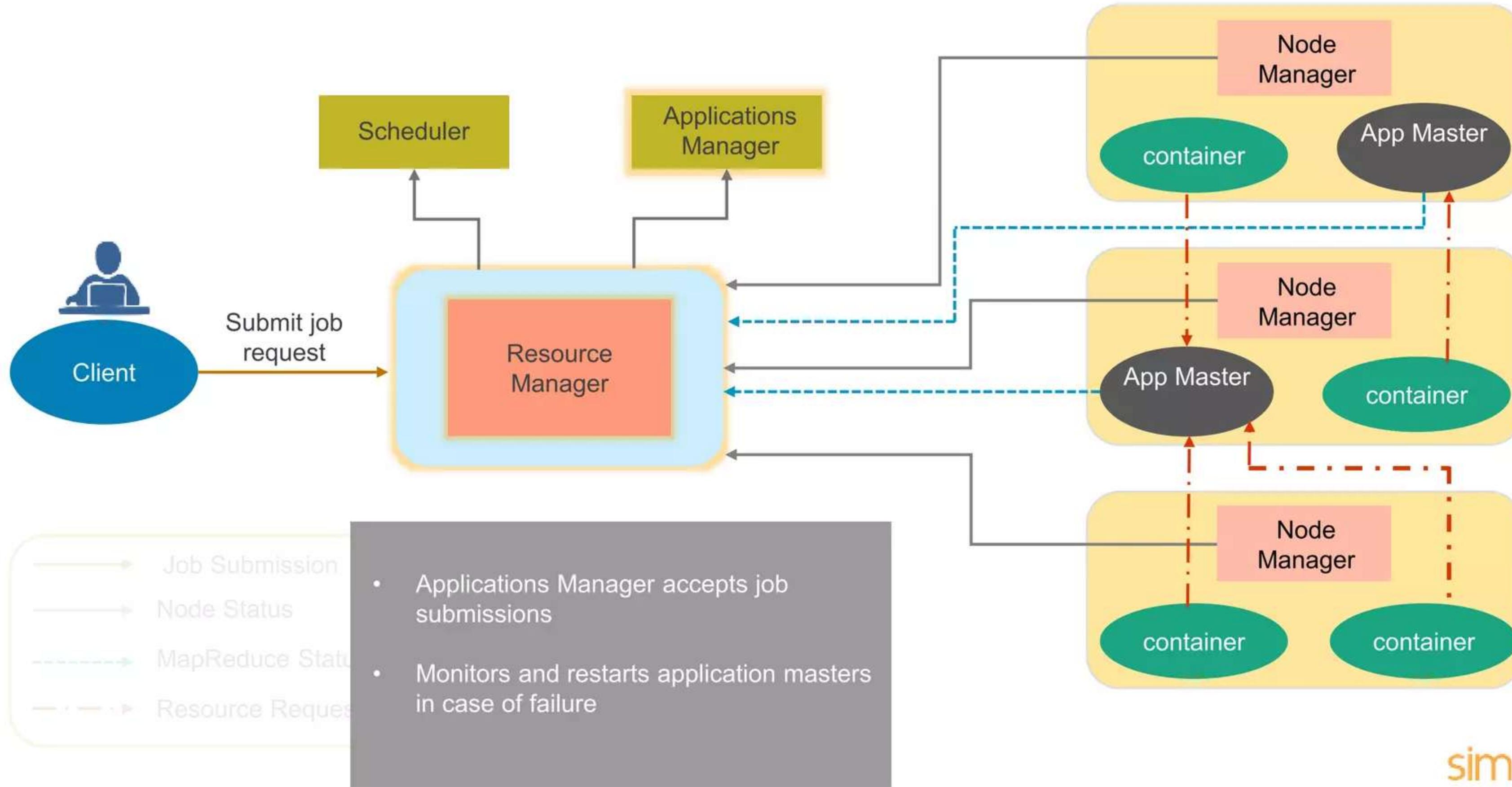
YARN Architecture – Resource Manager



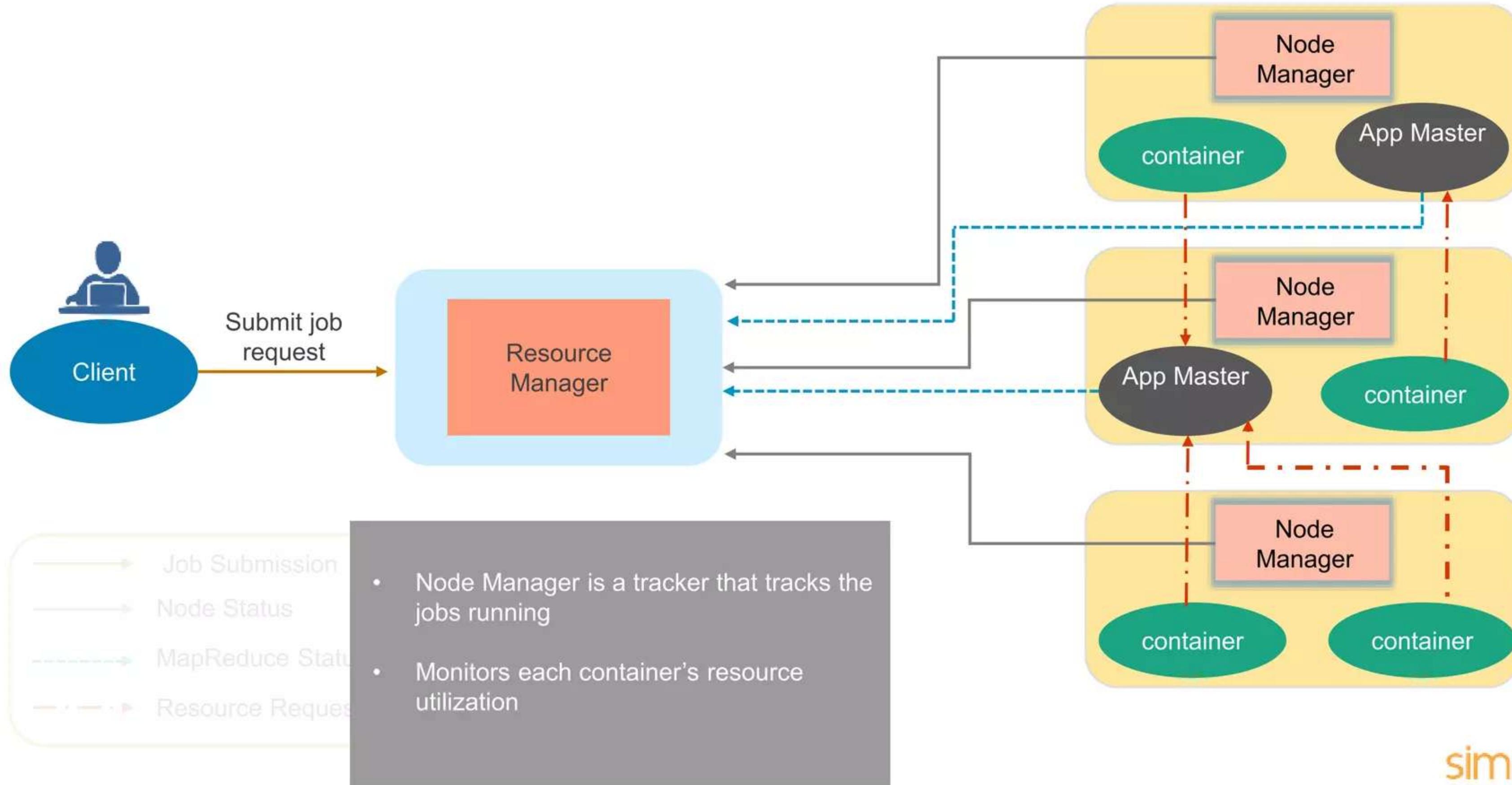
YARN Architecture – Scheduler



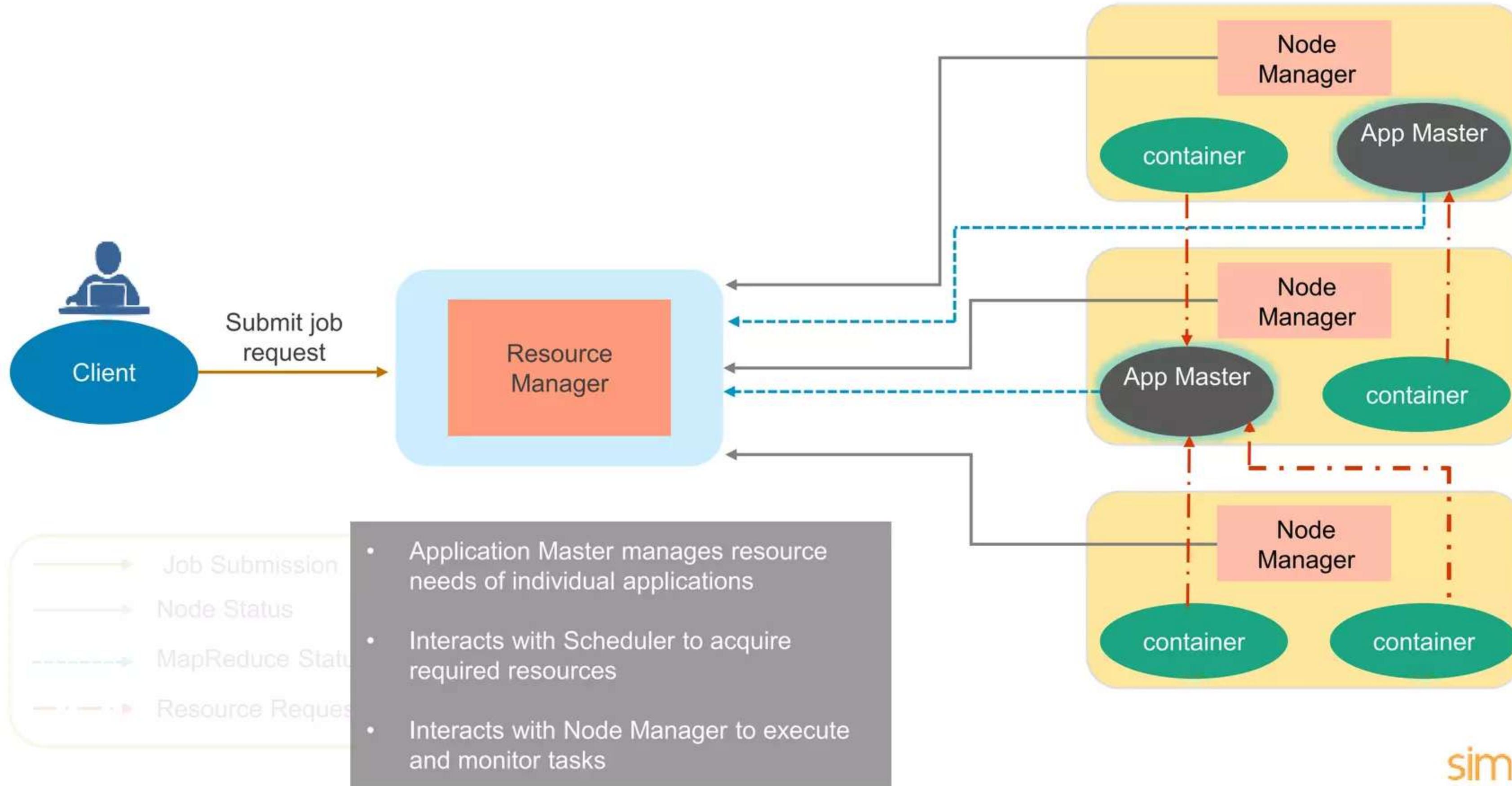
YARN Architecture – Applications Manager



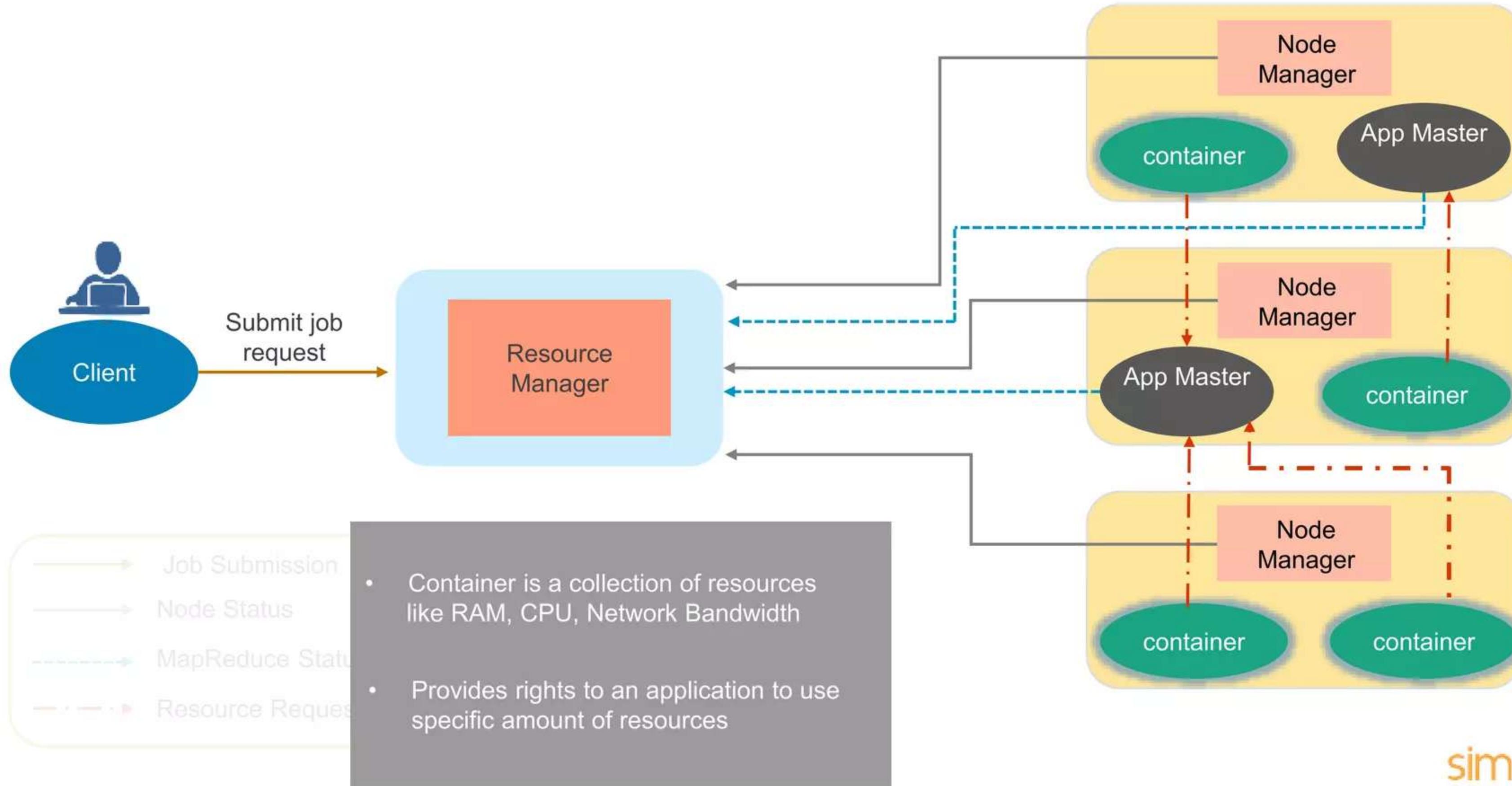
YARN Architecture – Node Manager



YARN Architecture – App Master



YARN Architecture - Container



Use case – Word Count using MapReduce



A blurred background image shows a person sitting at a desk, working on a laptop. The person's hands are visible, one resting on the keyboard and the other near the trackpad. A stack of papers or books is on the desk to the right of the laptop.

THANK YOU

For more information, visit

www.simplilearn.com

simplilearn