

# Hadoop File System

1

# Basic Features: HDFS

2

- Highly fault-tolerant
- High throughput
- Suitable for applications with large data sets
- Can be built out of commodity hardware

# Fault tolerance

3

- Failure is the norm rather than exception
- A HDFS instance may consist of thousands of server machines, each storing part of the file system's data.
- Since we have huge number of components and that each component has non-trivial probability of failure means that there is always some component that is non-functional.
- Detection of faults and quick, automatic recovery from them is a core architectural goal of HDFS.

# Data Characteristics

4

- Batch processing rather than interactive user access.
- Large data sets and files: gigabytes to terabytes size
- High aggregate data bandwidth
- Scale to hundreds of nodes in a cluster
- Tens of millions of files in a single instance
- Write-once-read-many: a file once created, written and closed need not be changed – this assumption simplifies coherency
- A map-reduce application or web-crawler application fits perfectly with this model.

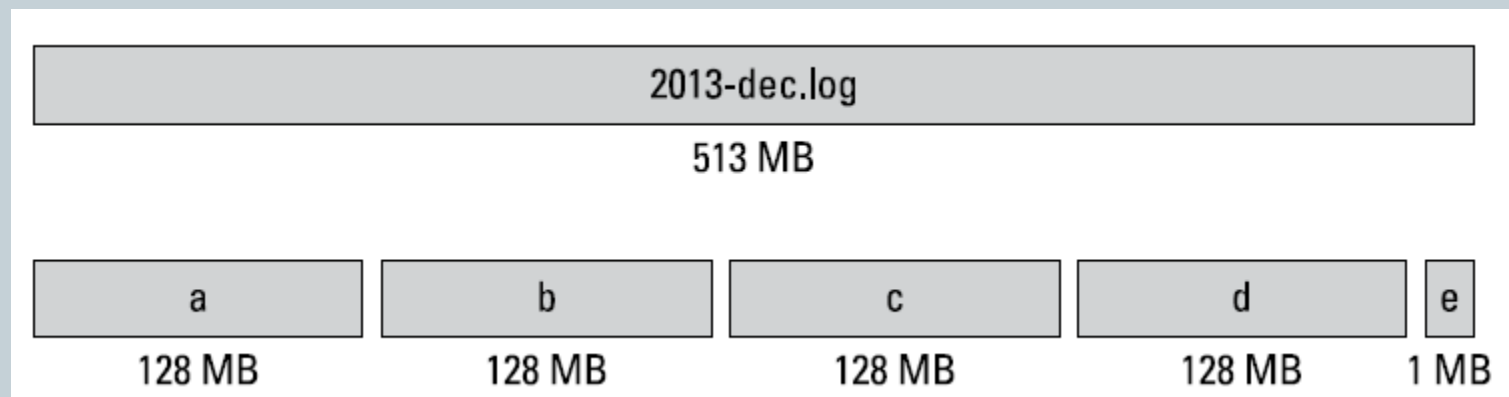
# Architecture

5

# Data Blocks

6

- HDFS support write-once-read-many with reads at high speeds.
- A data file is split into blocks.
- A typical block size is 128MB.
- A file is chopped into 128MB chunks and stored.

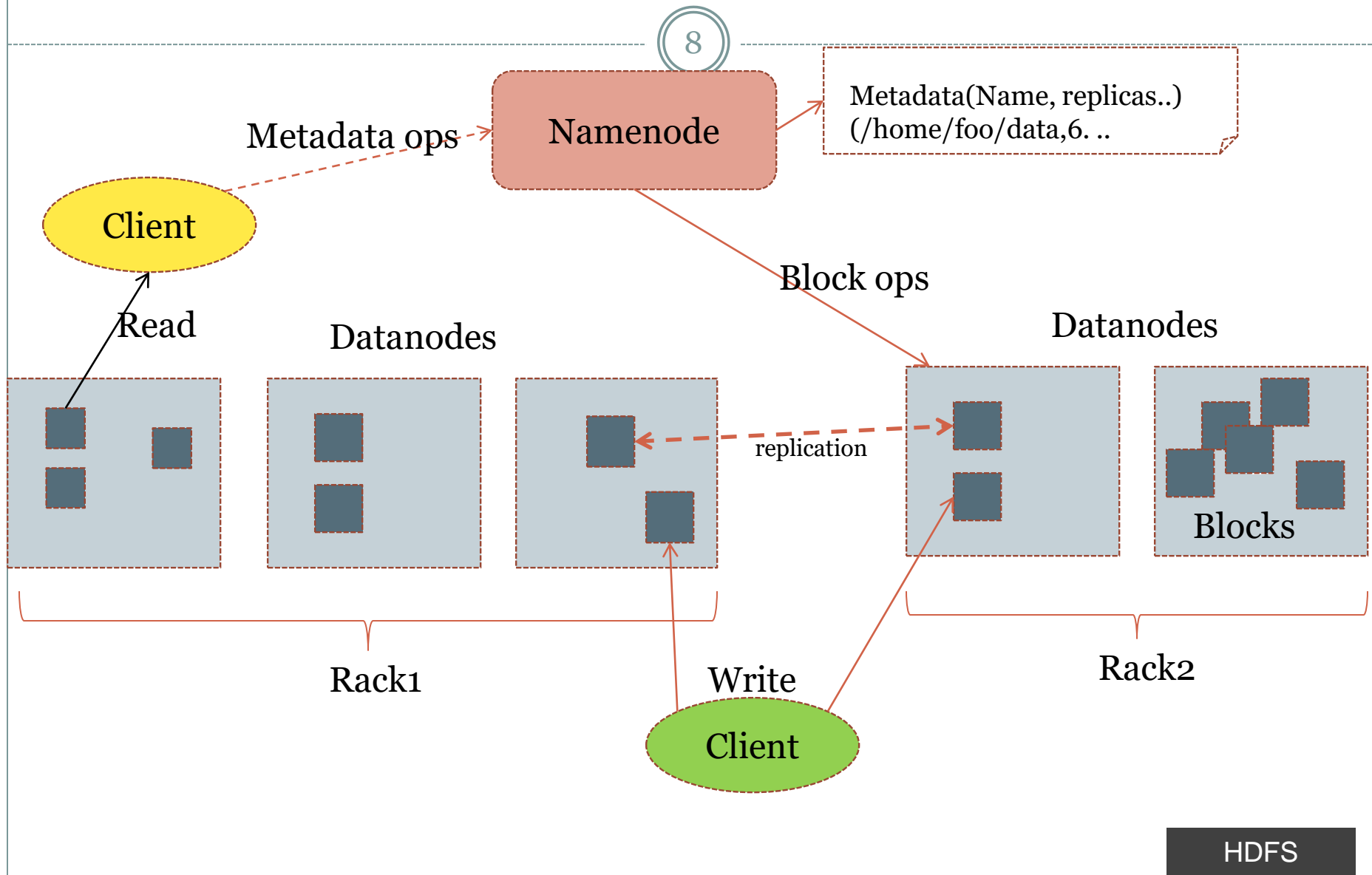


# Namenode and Datanodes

7

- Master/slave architecture
- HDFS cluster consists of a single **Namenode**, a master server that manages the file system namespace and regulates access to files by clients.
- There are a number of **DataNodes** usually one per node in a cluster.
- The DataNodes manage storage attached to the nodes that they run on.
- HDFS exposes a file system namespace and allows user data to be stored in files.
- A file is split into one or more blocks and set of blocks are stored in DataNodes.
- DataNodes: serves read, write requests, performs block creation, deletion, and replication upon instruction from Namenode.

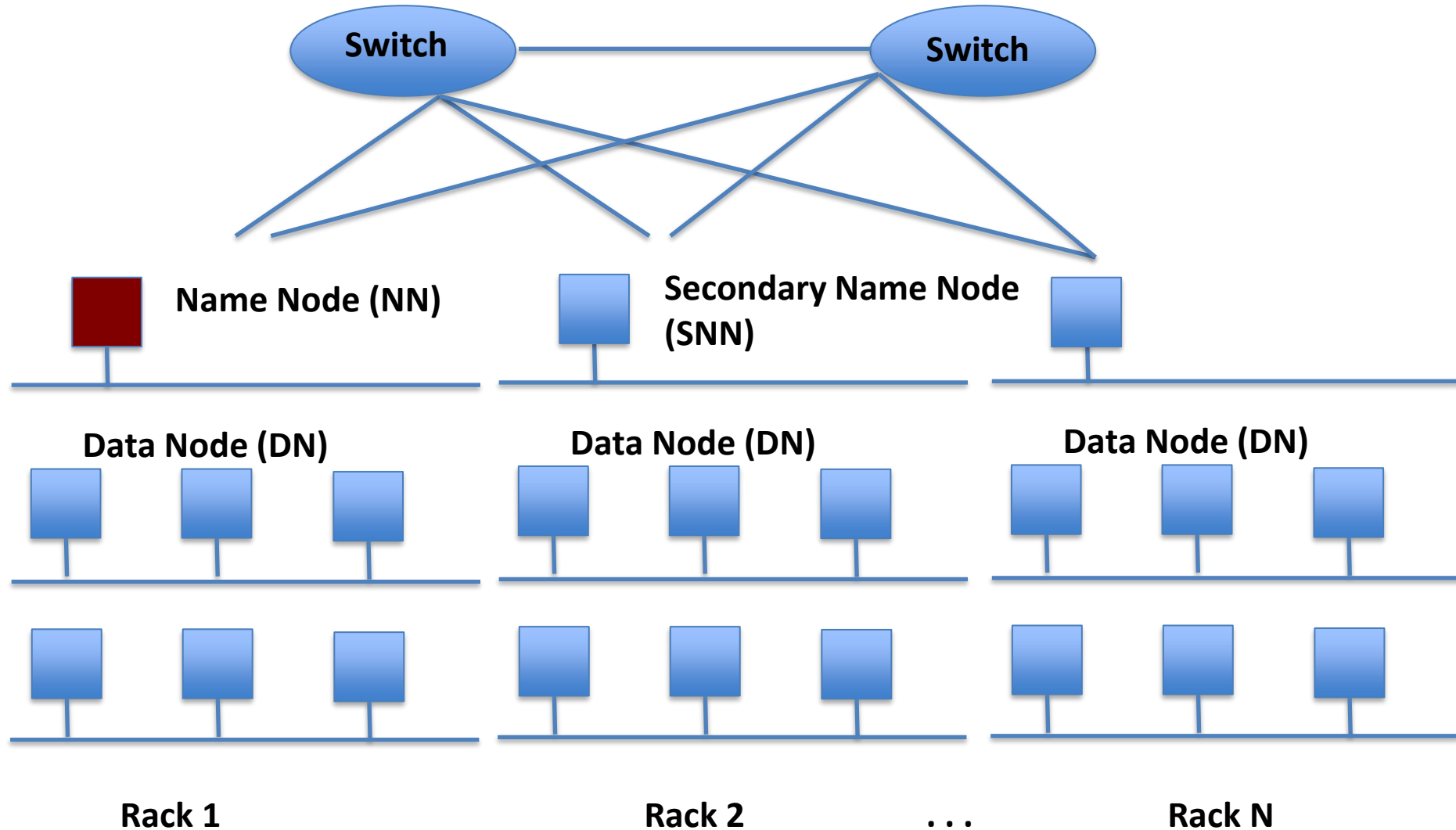
# HDFS Architecture





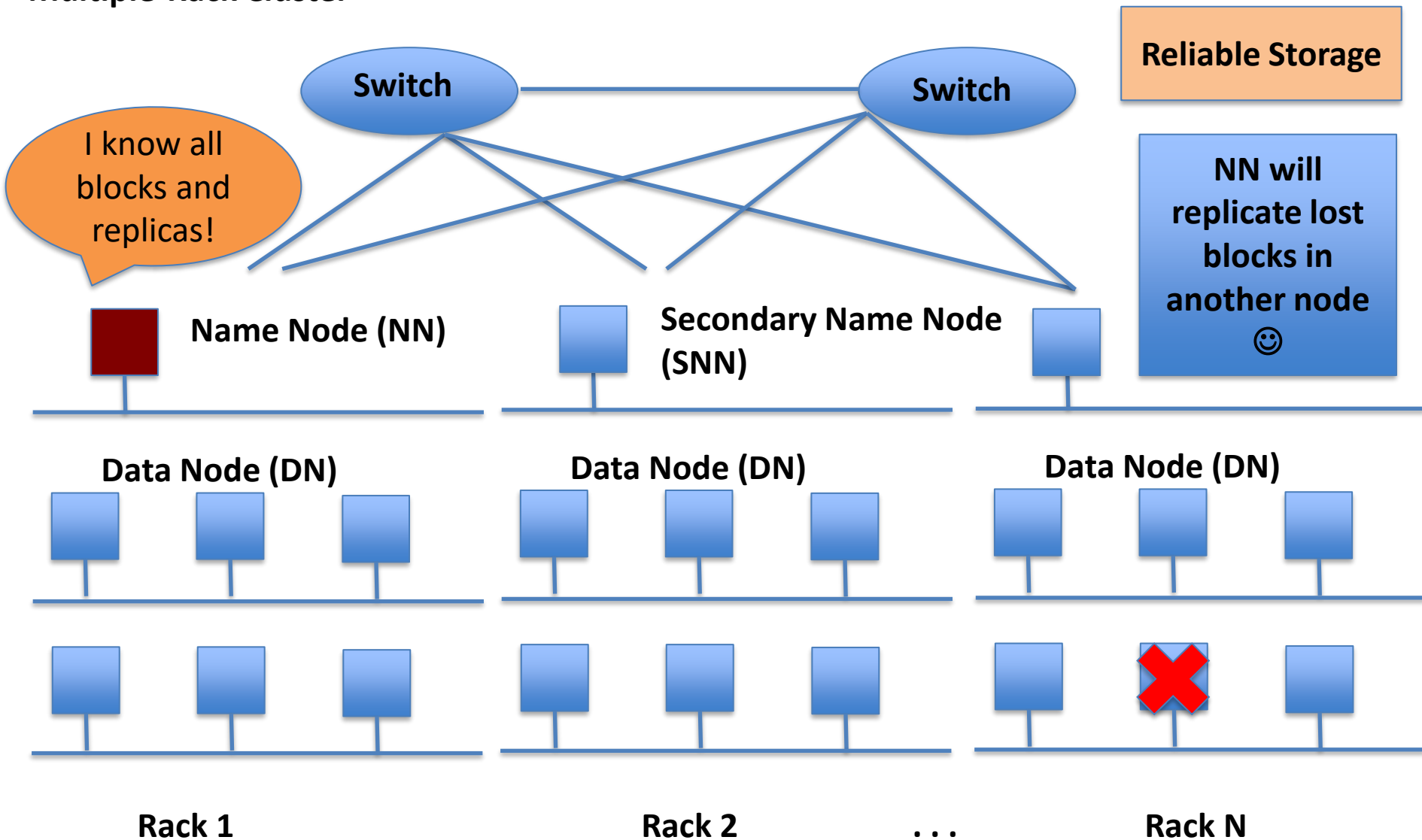
# HDFS Architecture: Master-Slave

## Multiple-Rack Cluster



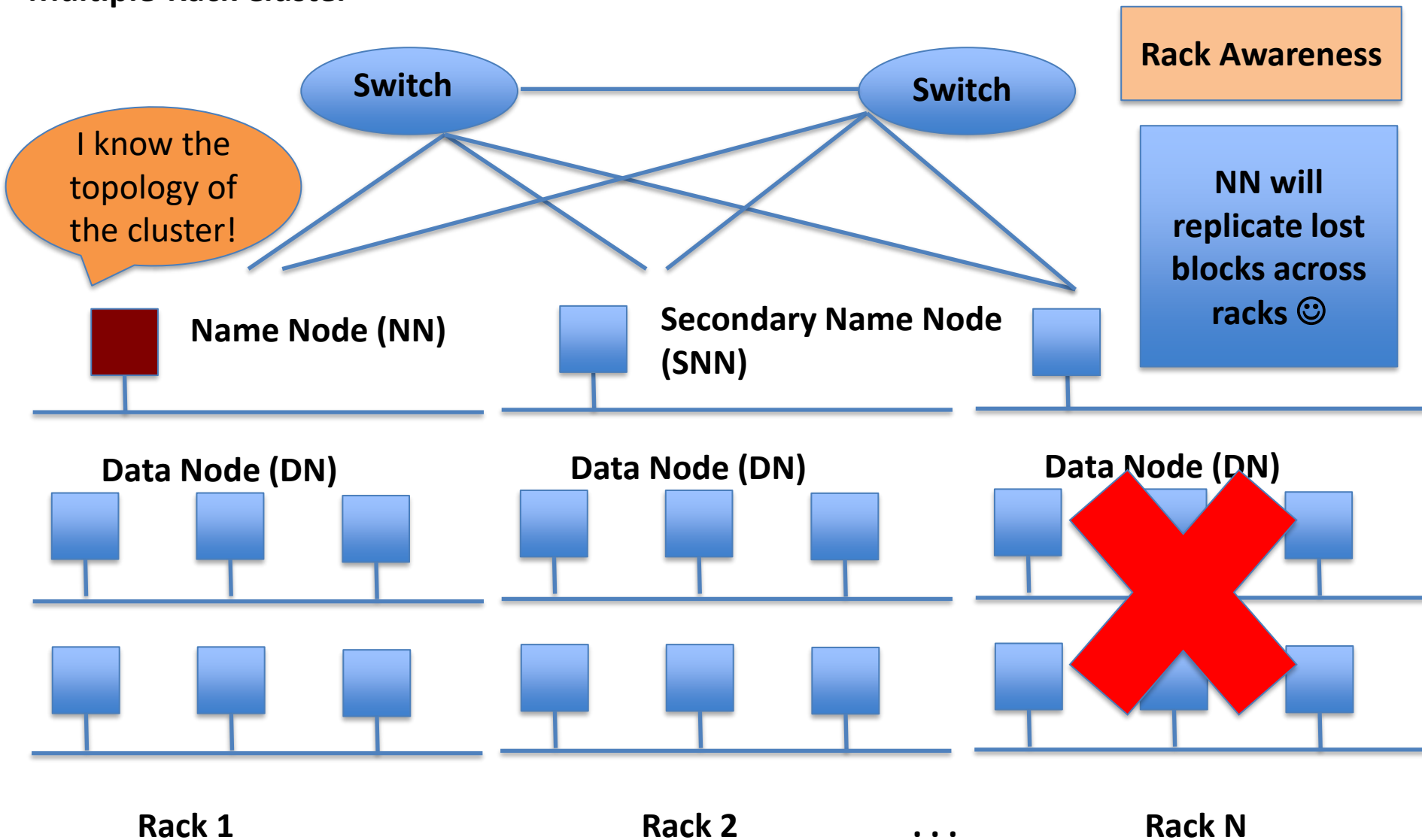
# HDFS Architecture: Master-Slave

## Multiple-Rack Cluster



# HDFS Architecture: Master-Slave

## Multiple-Rack Cluster



# File system Namespace

12

- Hierarchical file system with directories and files
- Create, remove, move, rename etc.
- Namenode maintains the file system
- Any meta information changes to the file system recorded by the Namenode.
- An application can specify the number of replicas of the file needed: replication factor of the file. This information is stored in the Namenode.

# Data Replication

13

- HDFS is designed to store very large files across machines in a large cluster.
- Each file is a sequence of blocks.
- All blocks in the file except the last are of the same size.
- Blocks are replicated for fault tolerance.
- Block size and replicas are configurable per file.
- The Namenode receives a Heartbeat and a BlockReport from each DataNode in the cluster.
- BlockReport contains all the blocks on a Datanode.

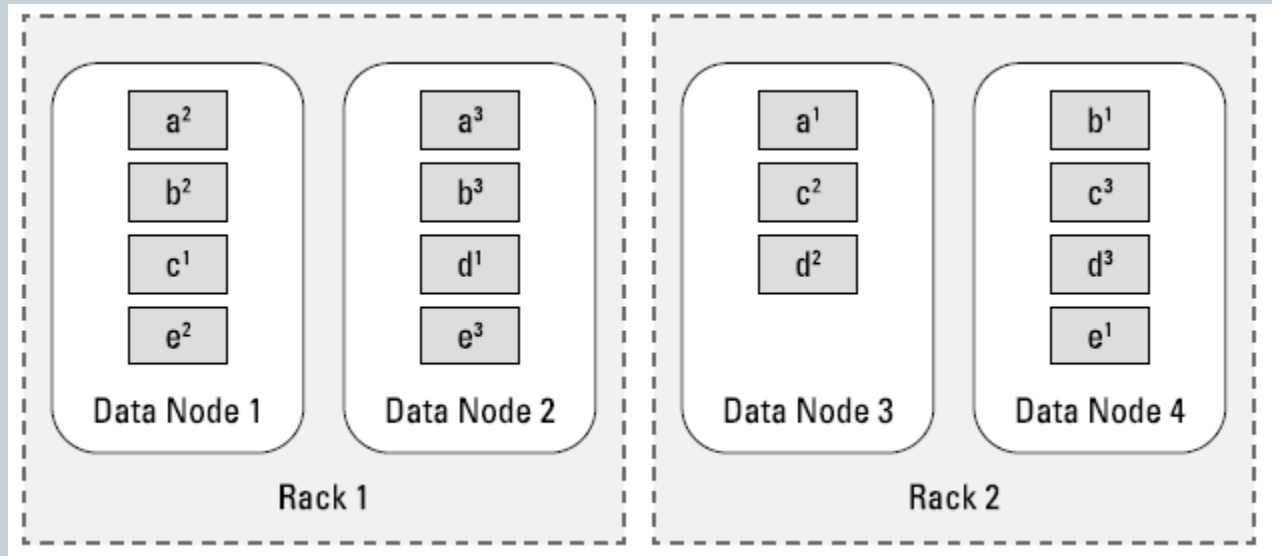
# Replica Placement

14

- The placement of the replicas is critical to HDFS reliability and performance.
- Optimizing replica placement distinguishes HDFS from other distributed file systems.
- Rack-aware replica placement:
  - Goal: improve reliability, availability and network bandwidth utilization
- Many racks, communication between racks are through switches.
- Network bandwidth between machines on the same rack is greater than those in different racks.
- Namenode determines the rack id for each DataNode.
- Replicas are typically placed on unique racks
  - Simple but non-optimal
  - Writes are expensive
  - Replication factor is 3
- Replicas are placed: one on a node in a local rack, one on a different node in the local rack and one on a node in a different rack.
- 1/3 of the replica on a node, 2/3 on a rack and 1/3 distributed evenly across remaining racks.

# Data Replication Patterns

15



# Replica Selection

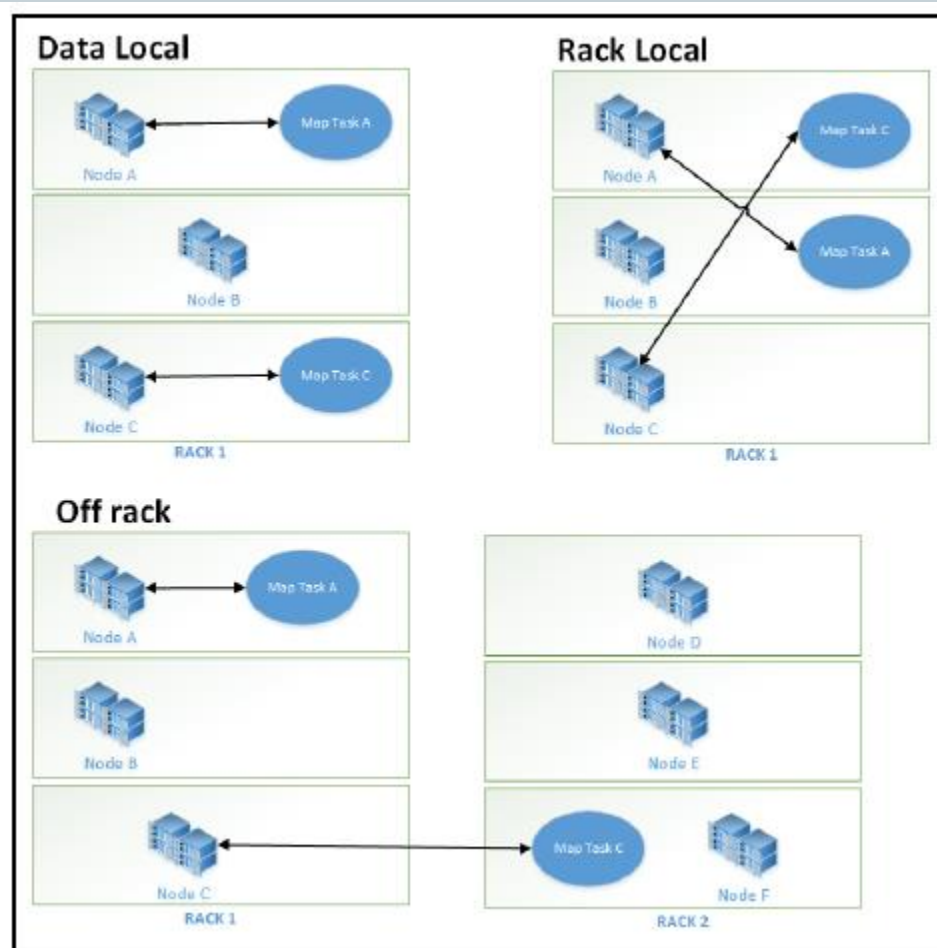
16

- Replica selection for READ operation: HDFS tries to minimize the bandwidth consumption and latency.
- Replica selection criteria
  - Data Local: Replica available on the same node
  - Rack Local: Replica available on the same rack
  - Off Rack: Replica available on different rack
- Replica selection preference:  
Data Local > Rack Local > Off Rack (Worst Case)
- HDFS cluster may span multiple data centers:  
replica in the local data center is preferred over the remote one.



# Replica Selection

17



# Safemode Startup

18

- On startup Namenode enters Safemode.
- Replication of data blocks do not occur in Safemode.
- Each DataNode checks in with Heartbeat and BlockReport.
- Namenode verifies that each block has acceptable number of replicas
- After a configurable percentage of safely replicated blocks check in with the Namenode, Namenode exits Safemode.
- It then makes the list of blocks that need to be replicated.
- Namenode then proceeds to replicate these blocks to other Datanodes.

# Filesystem Metadata

19

- The HDFS namespace is stored by Namenode.
- Namenode uses a transaction log called the EditLog to record every change that occurs to the filesystem meta data.
  - For example, creating a new file.
  - Change replication factor of a file
  - EditLog is stored in the Namenode's local filesystem
- Entire filesystem namespace including mapping of blocks to files and file system properties is stored in a file FsImage. Stored in Namenode's local filesystem.

# Name Node

20

- Keeps image of entire file system namespace and file Blockmap in memory.
- 4GB of local RAM is sufficient to support the above data structures that represent the huge number of files and directories.
- When the Namenode starts up it gets the FsImage and Editlog from its local file system, update FsImage with EditLog information and then stores a copy of the FsImage on the filesystem as a checkpoint.
- Periodic checkpointing is done. So that the system can recover back to the last checkpointed state in case of a crash.

# HDFS Inside: Name Node

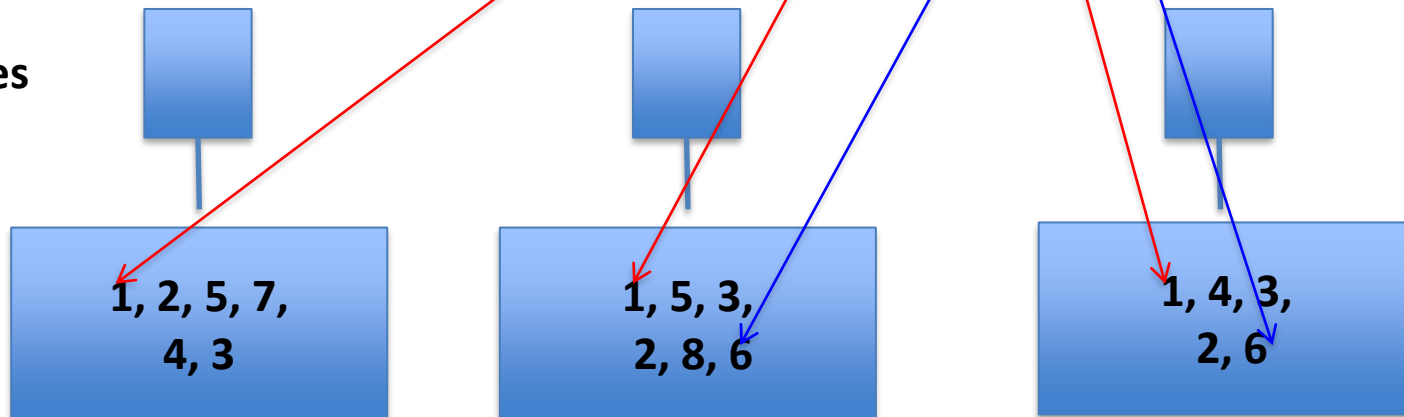
Name Node

Snapshot of FS

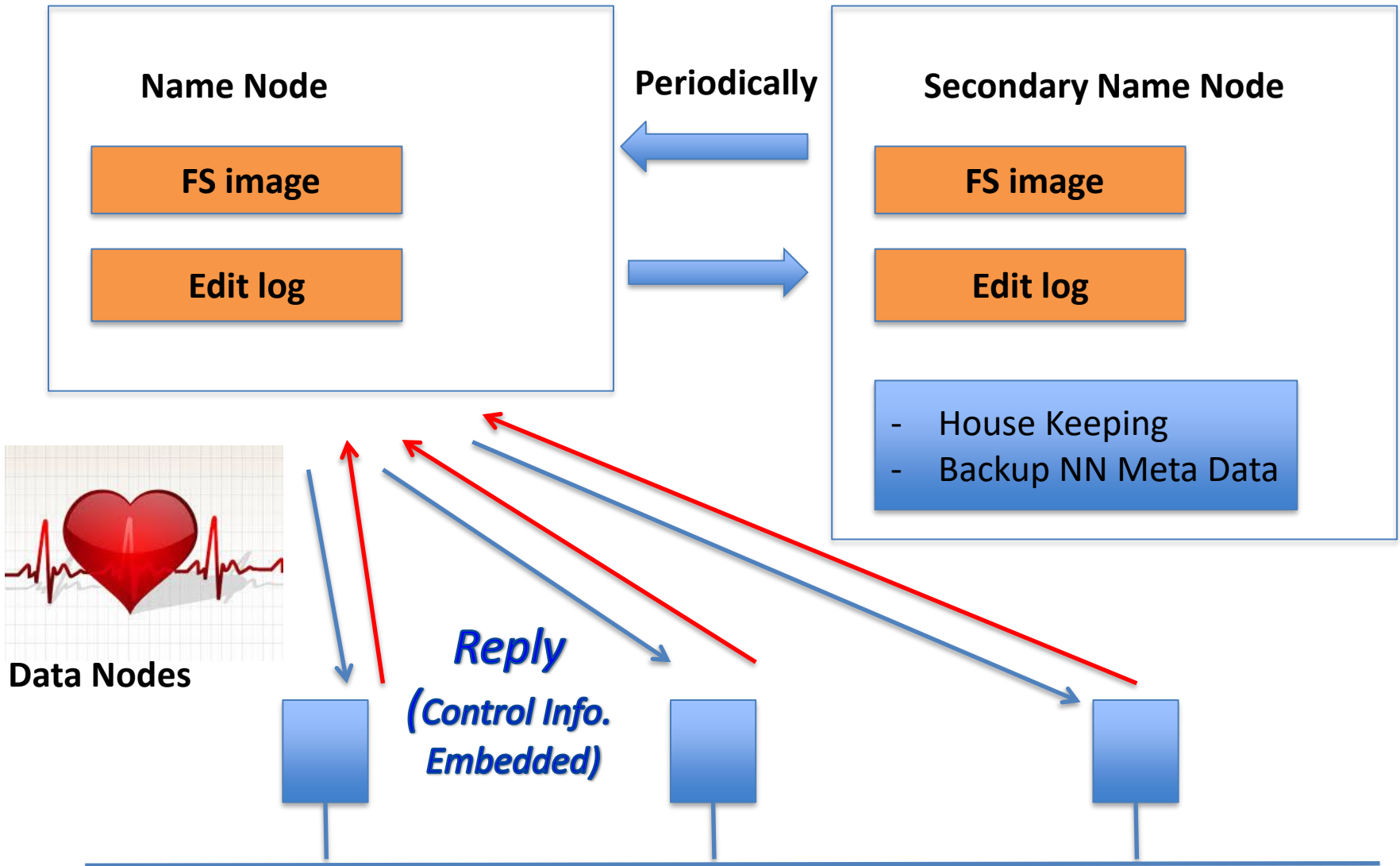
Edit log: record  
changes to FS

Filename	Replication factor	Block ID
File 1	3	[1, 2, 3]
File 2	2	[4, 5, 6]
File 3	1	[7, 8]

Data Nodes



# HDFS Inside: Name Node



# Secondary Name Node

23

- NOT a backup name node
- Merging FsImage and EditLog is memory intensive
- Secondary Name Node periodically downloads current Name Node FsImage and EditLog, merges them and uploads the new image back

# Datanode

24

- A Datanode stores data in files in its local file system.
- Datanode has no knowledge about HDFS filesystem
- It stores each block of HDFS data in a separate file.
- Datanode does not create all files in the same directory.
- When the filesystem starts up it generates a list of all HDFS blocks and send this report to Namenode: Blockreport.



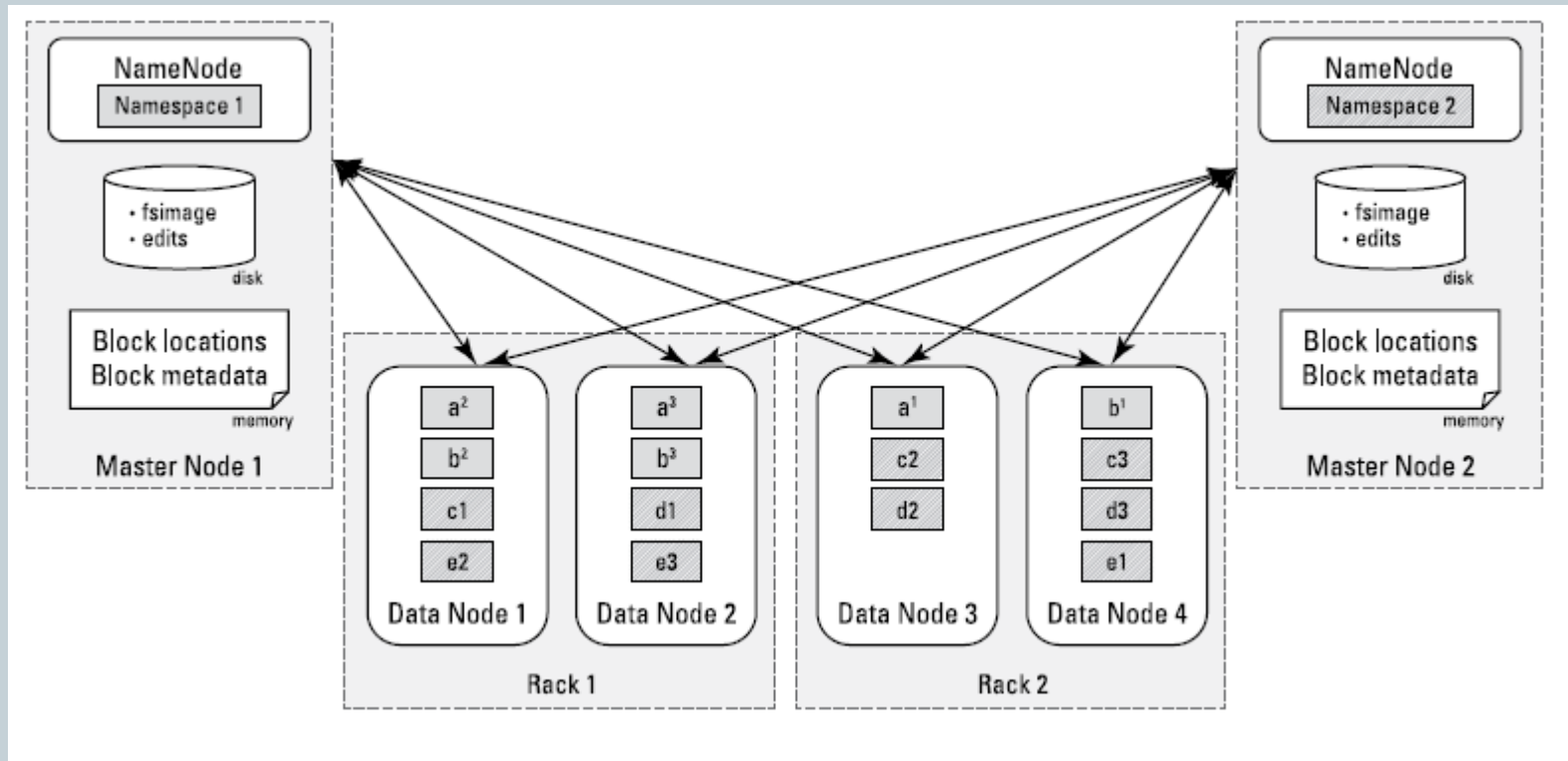
# HDFS Federation

25

- Before Hadoop 2, a single Name Node had to store metadata for every bloc of data
- Difficult to manage increasing number of blocks
- HDFS Federation –
  - Multiple Name Nodes
  - Each Name Node responsible for a particular Namespace (set of blocks)
  - Only individual blocks are partitioned among Name Nodes. Data Nodes can be shared.

# HDFS Federation

26



# Robustness

27

# Objectives

28

- Primary objective of HDFS is to store data reliably in the presence of failures.
- Three common failures are: Namenode failure, Datanode failure and network partition.

# DataNode failure and heartbeat

29

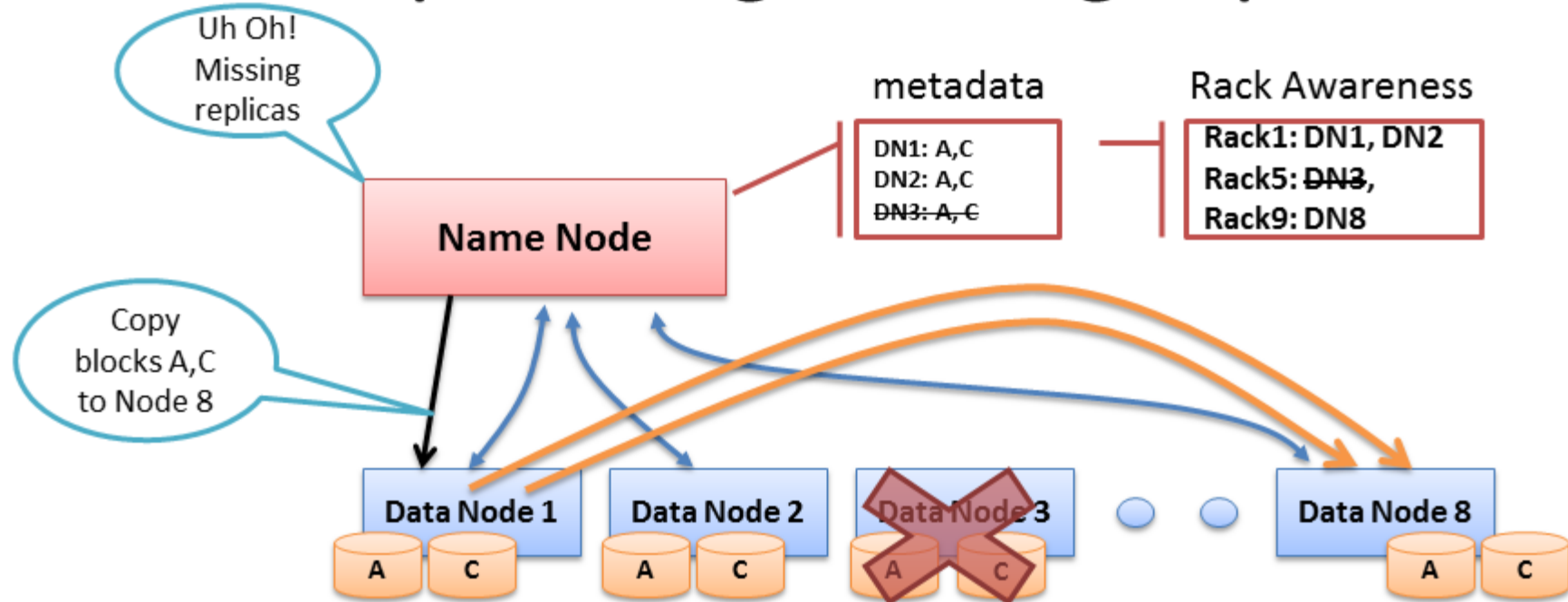
- A network partition can cause a subset of Datanodes to lose connectivity with the Namenode.
- Namenode detects this condition by the absence of a Heartbeat message.
- Namenode marks Datanodes without Heartbeat and does not send any IO requests to them.
- Any data registered to the failed Datanode is not available to the HDFS.
- Also the death of a Datanode may cause replication factor of some of the blocks to fall below their specified value.

# Re-replication

30

- The necessity for re-replication may arise due to:
  - A Datanode may become unavailable,
  - A replica may become corrupted,
  - A hard disk on a Datanode may fail, or
  - The replication factor on the block may be increased.

# Re-replicating missing replicas



- Missing Heartbeats signify lost Nodes
- Name Node consults metadata, finds affected data
- Name Node consults Rack Awareness script
- Name Node tells a Data Node to re-replicate

# Cluster Rebalancing

32

- HDFS architecture is compatible with data rebalancing schemes.
- A scheme might move data from one Datanode to another if the free space on a Datanode falls below a certain threshold.
- In the event of a sudden high demand for a particular file, a scheme might dynamically create additional replicas and rebalance other data in the cluster.



# Data Integrity

33

- Consider a situation: a block of data fetched from Datanode arrives corrupted.
- This corruption may occur because of faults in a storage device, network faults, or buggy software.
- A HDFS client creates the checksum of every block of its file and stores it in hidden files in the HDFS namespace.
- When a client retrieves the contents of file, it verifies that the corresponding checksums match.
- If it does not match, the client can retrieve the block from a replica.

# Metadata Disk Failure

34

- FsImage and EditLog are central data structures of HDFS.
- A corruption of these files can cause a HDFS instance to be non-functional.
- For this reason, a Namenode can be configured to maintain multiple copies of the FsImage and EditLog.
- Multiple copies of the FsImage and EditLog files are updated synchronously.
- Meta-data is not data-intensive.

# HDFS High Availability

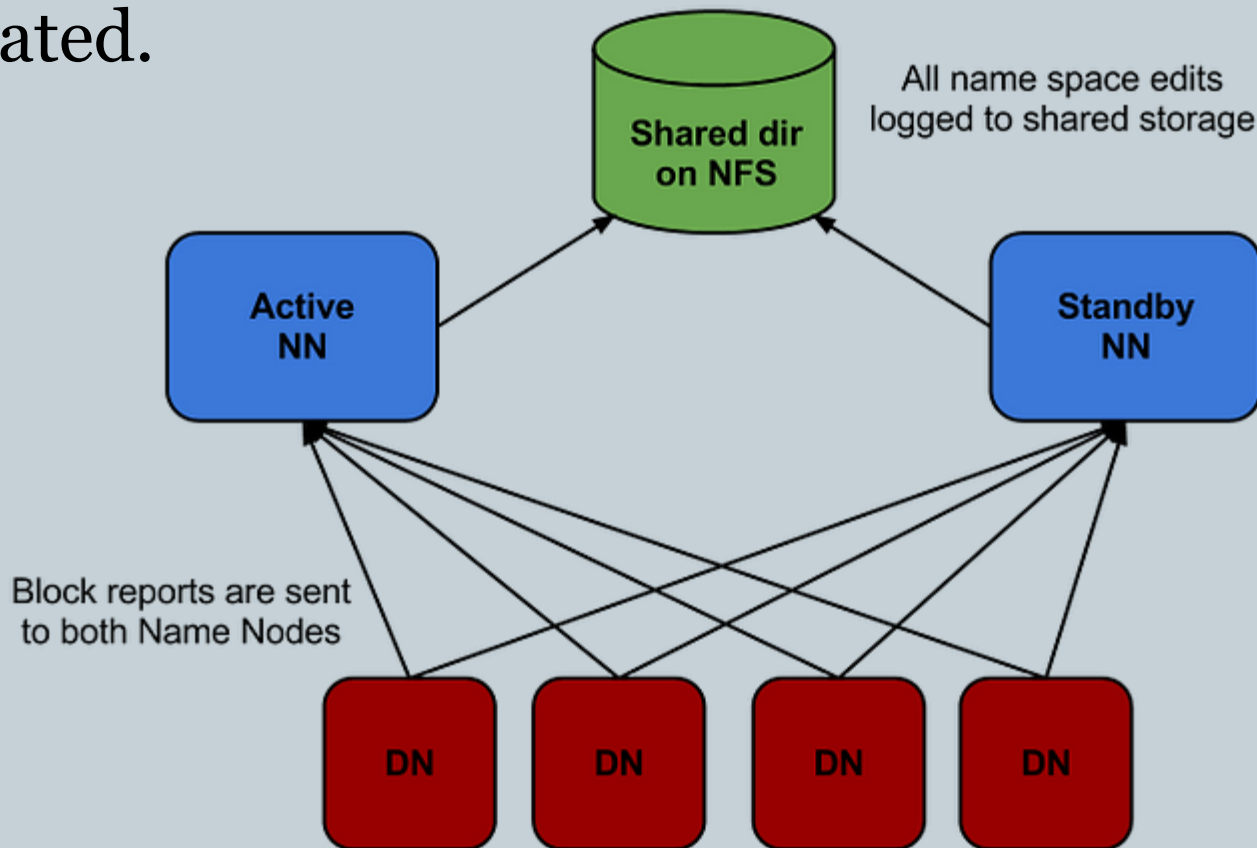
35

- Originally, Name Node was single point of failure.
- HDFS High Availability introduced in Hadoop 2.
- A dedicated Standby Name Node is created which has access to the up-to-date version of the current state (FsImage and EditLogs) of the file system.
- In case of Active Name Node failure, the standby Name Node takes over.
- Apache Zookeeper daemon detects Active Name Node failure and coordinates the shift.

# HDFS High Availability

36

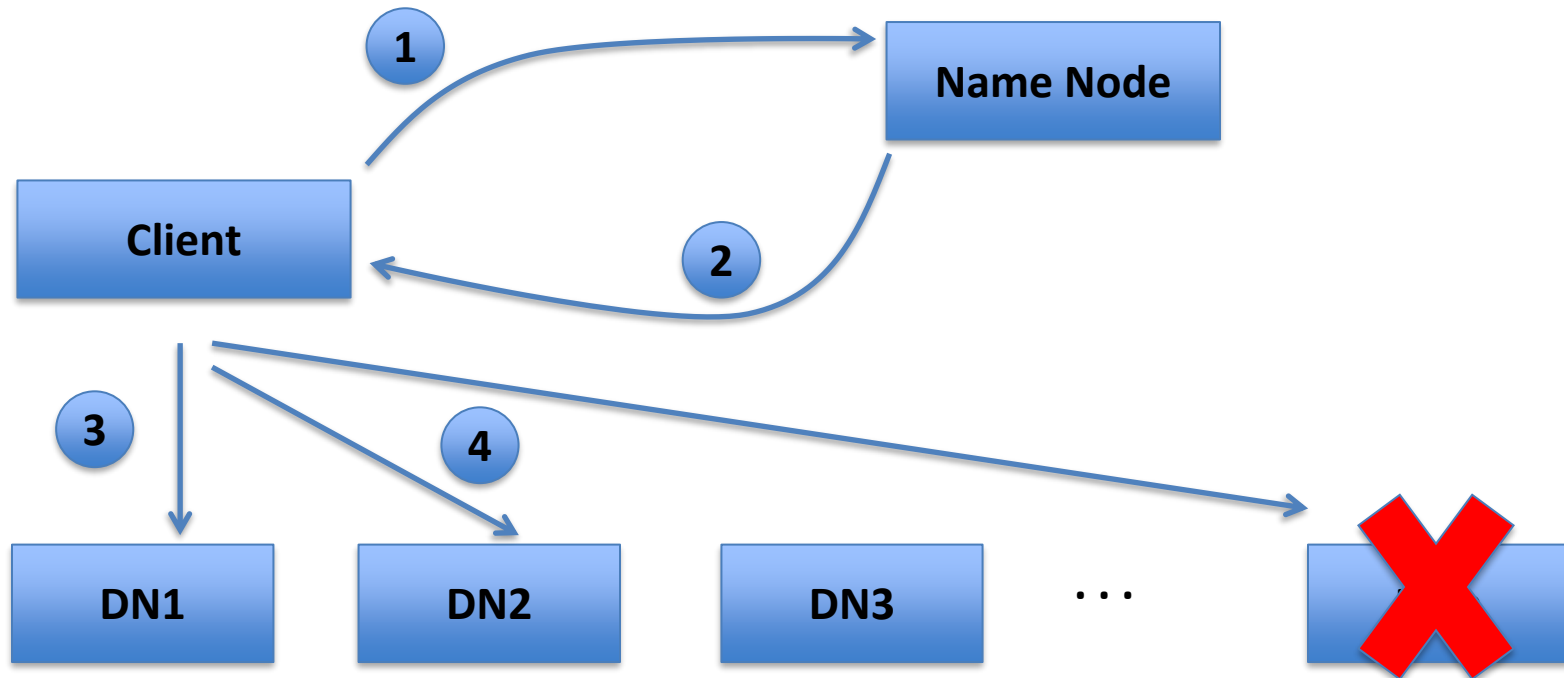
- Hadoop 2 only allowed one standby Name Node.
- In Hadoop 3, multiple standby Name Nodes can be created.



# Data Organization

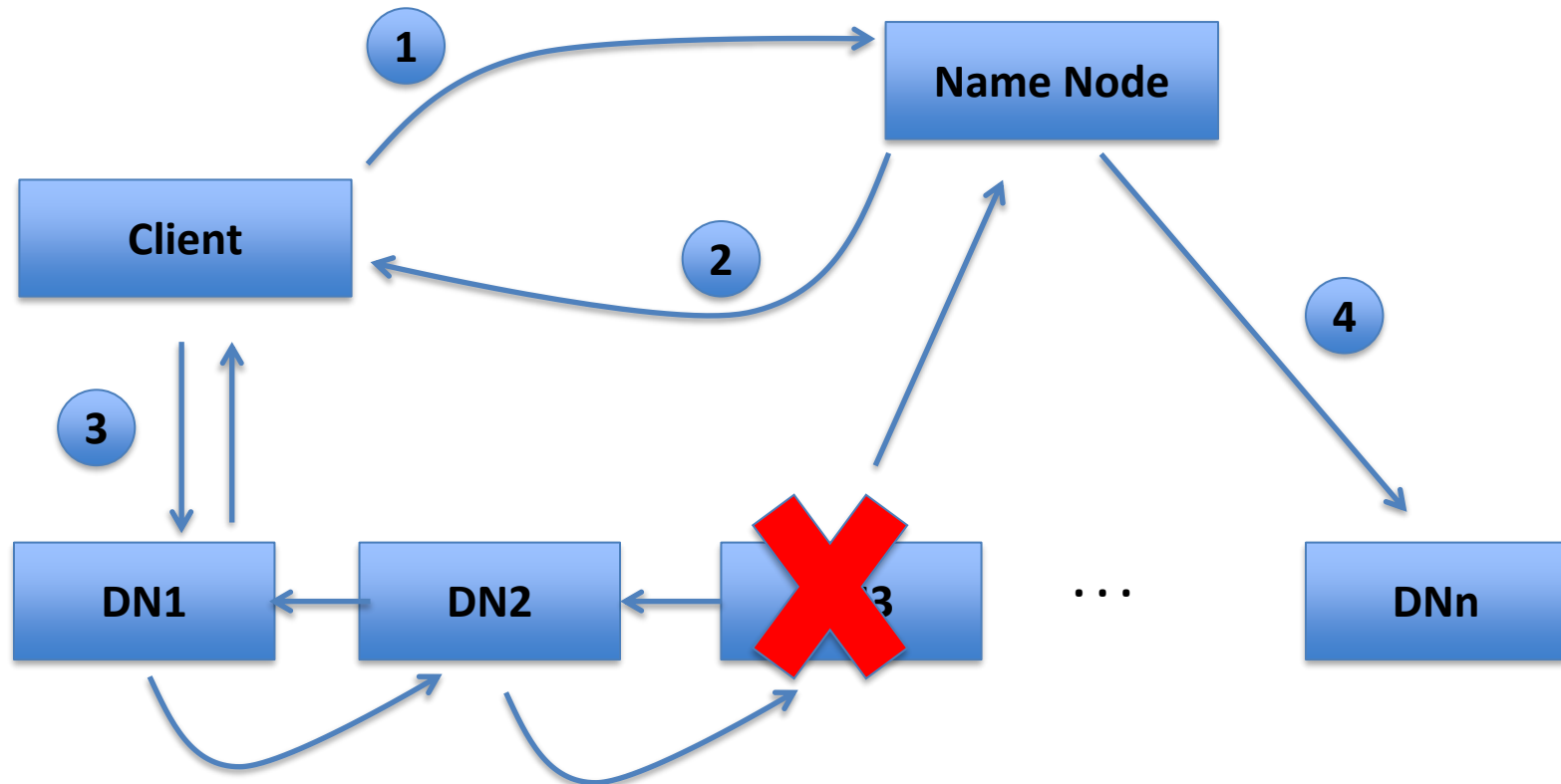
37

# HDFS Inside: Read Workflow



1. Client connects to NN to read data
2. NN tells client where to find the data blocks
3. Client reads blocks directly from data nodes (without going through NN)
4. In case of node failures, client connects to another node that serves the missing block

# HDFS Inside: Write Workflow



1. Client connects to NN to write data
2. NN tells client write these data nodes
3. Client writes blocks directly to data nodes with desired replication factor
4. In case of node failures, NN will figure it out and replicate the missing blocks

# Staging

40

- A client request to create a file does not reach Namenode immediately.
- HDFS client caches the data into a temporary file. When the data reached a HDFS block size the client contacts the Namenode.
- Namenode inserts the filename into its hierarchy and allocates a data block for it.
- The Namenode responds to the client with the identity of the Datanode and the destination of the replicas (Datanodes) for the block.
- Then the client flushes it from its local memory.



# Staging (contd.)

41

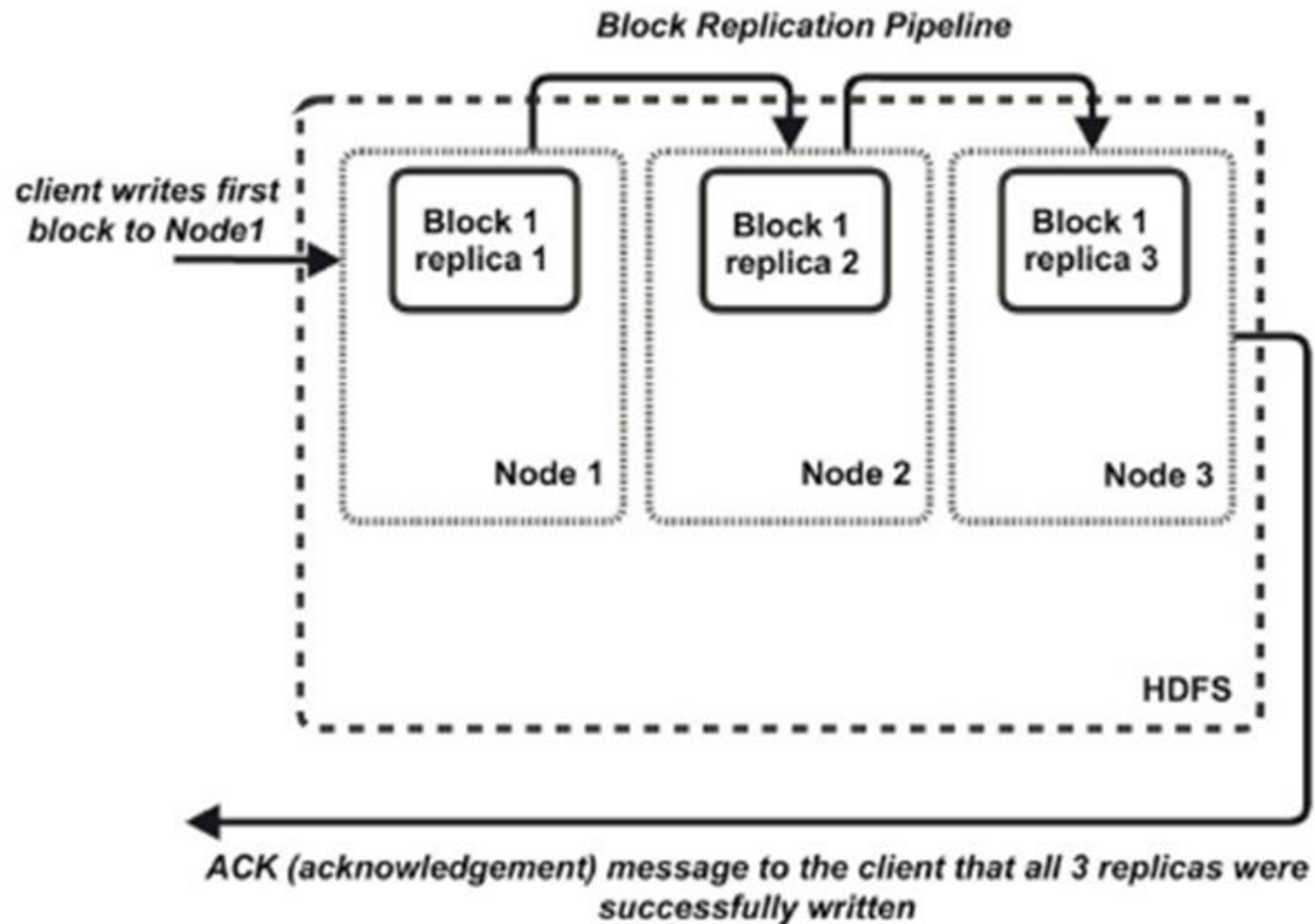
- The client sends a message that the file is closed.
- Namenode proceeds to commit the file for creation operation into the persistent store.
- If the Namenode dies before file is closed, the file is lost.
- This client side caching is required to avoid network congestion

# Replication Pipelining

42

- When the client receives response from Namenode, it flushes its block in small pieces (4K) to the first replica, that in turn copies it to the next replica and so on.
- Thus data is pipelined from Datanode to the next.

# Replication Pipelining



# Reference

44

- [The Hadoop Distributed File System: Architecture and Design by Apache Foundation Inc.](#)