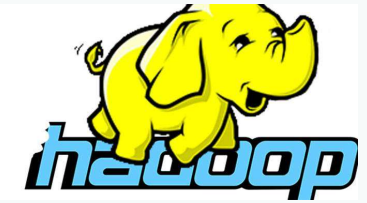HDFS

# What is Hadoop?

**An open-source framework**

**Created to make it easier to work with big data.**

**It provides a method to**

- Access data
- Process data
- manage resources across the computing and network resources

# Hadoop Core Modules

## Hadoop Distributed File System (HDFS)

- Provides access to application data.
- Can work with other file systems - FTP, Amazon S3 etc

## Hadoop YARN

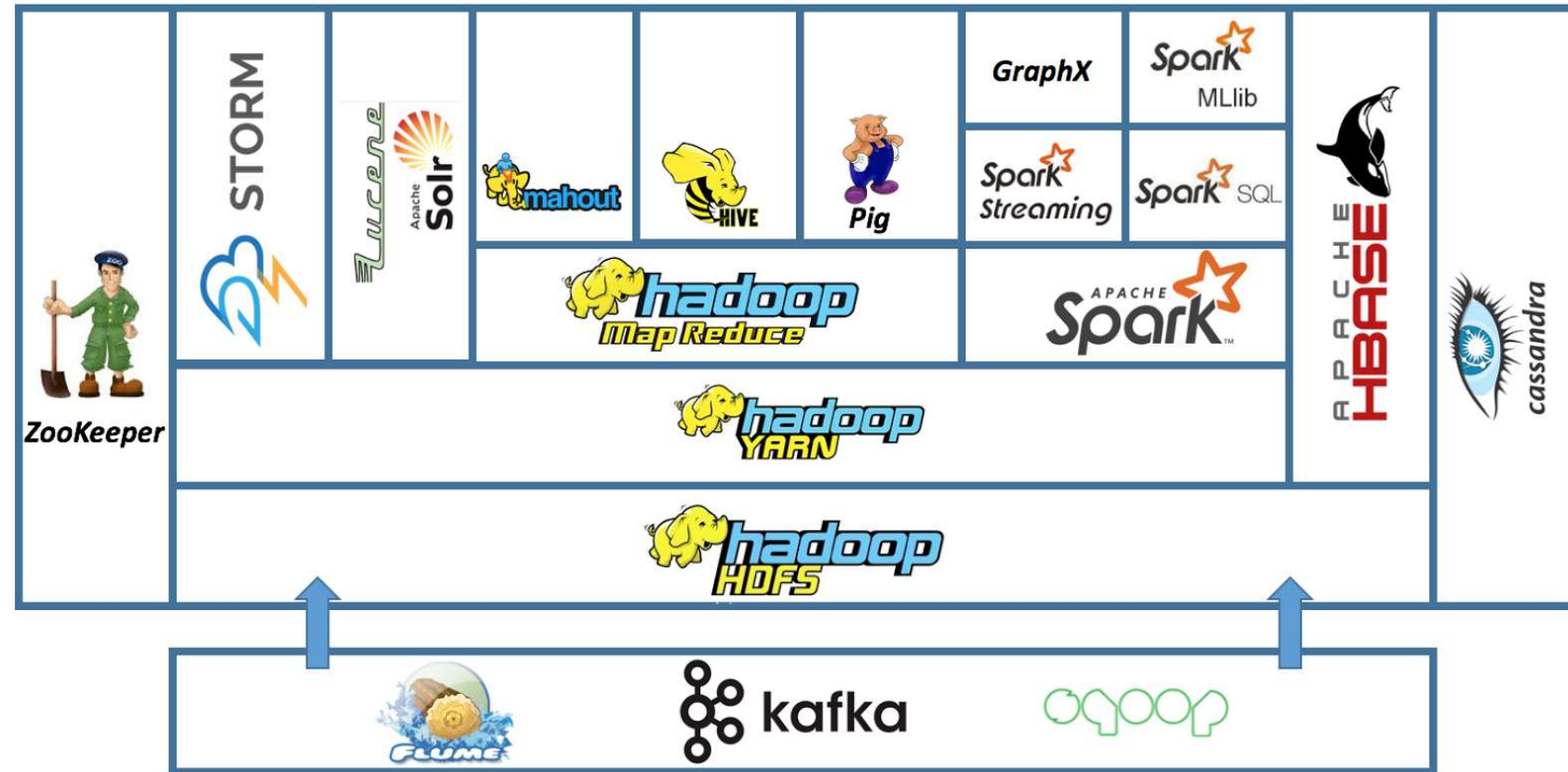- Provides the framework to schedule jobs and manage resources

## Hadoop MapReduce

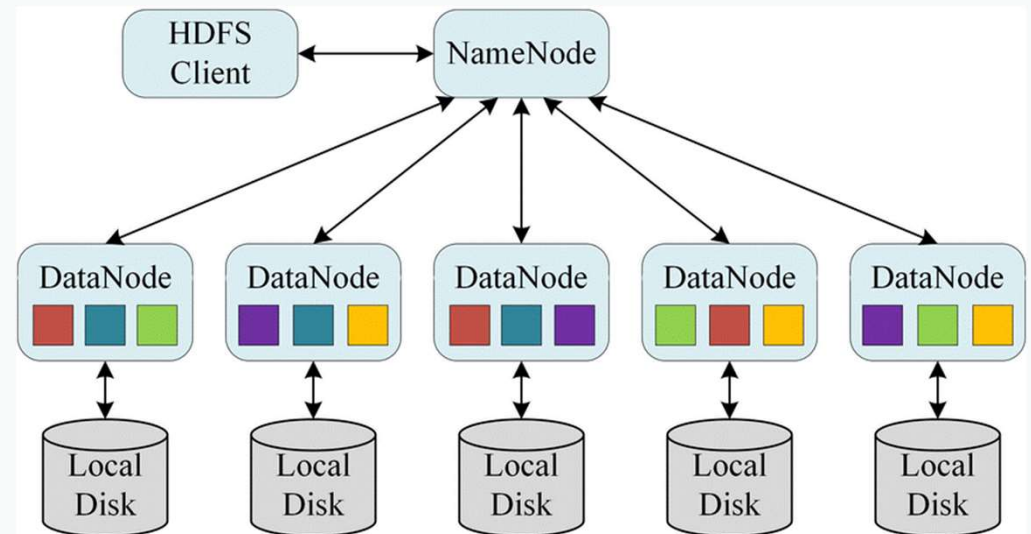- Parallel processing system for large data sets.

## Hadoop Common

- A set of utilities that supports the three other core modules.
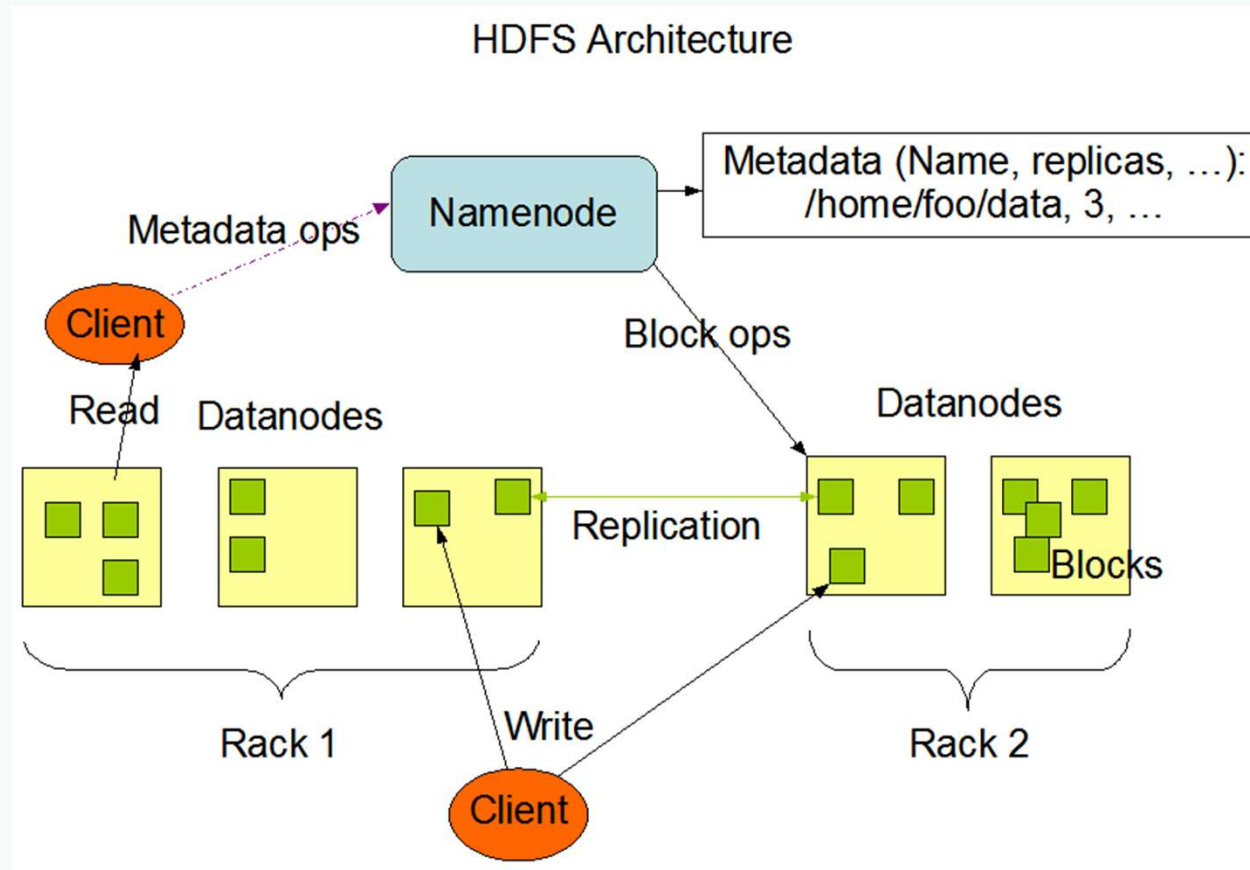
# Hadoop Ecosystem

# HDFS

- Manages how data files are divided and stored across the cluster.

- Data is divided into blocks

- Each server in the cluster contains data from different blocks

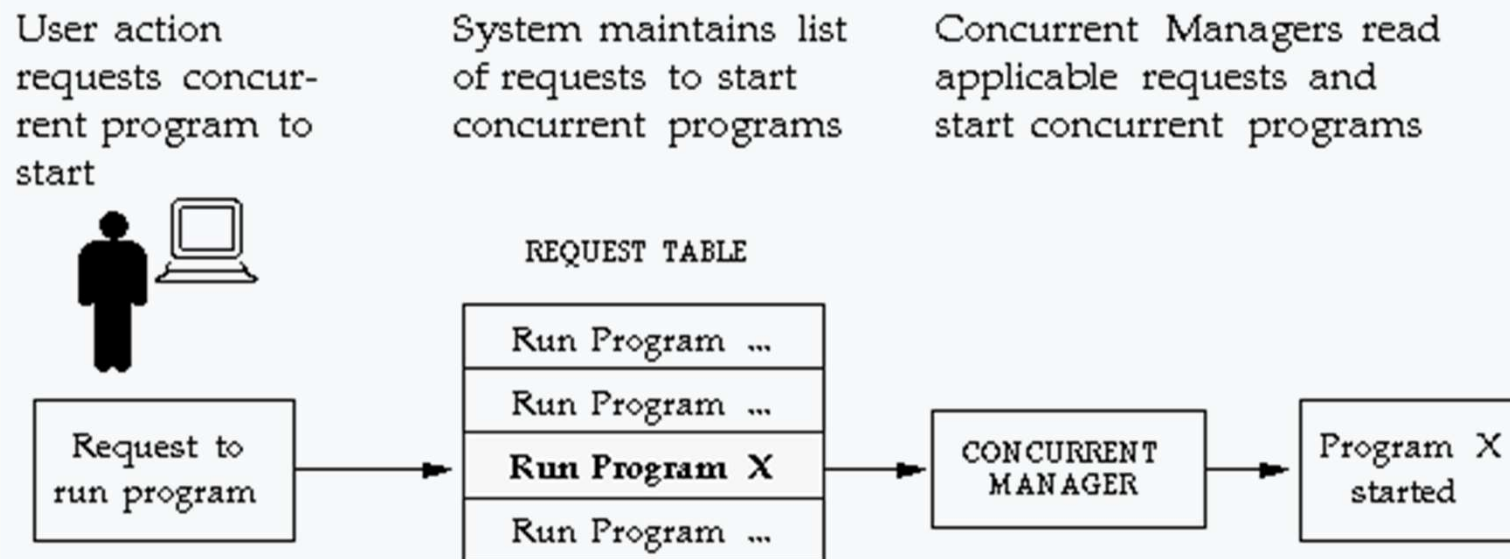- There is also some built-in redundancy.

# HDFS Architecture



HDFS Architecture

Metadata ops → Namenode → Metadata (Name, replicas, …): /home/foo/data, 3, …

Block ops

Client

Read    Datanodes

Datanodes

Replication

Blocks
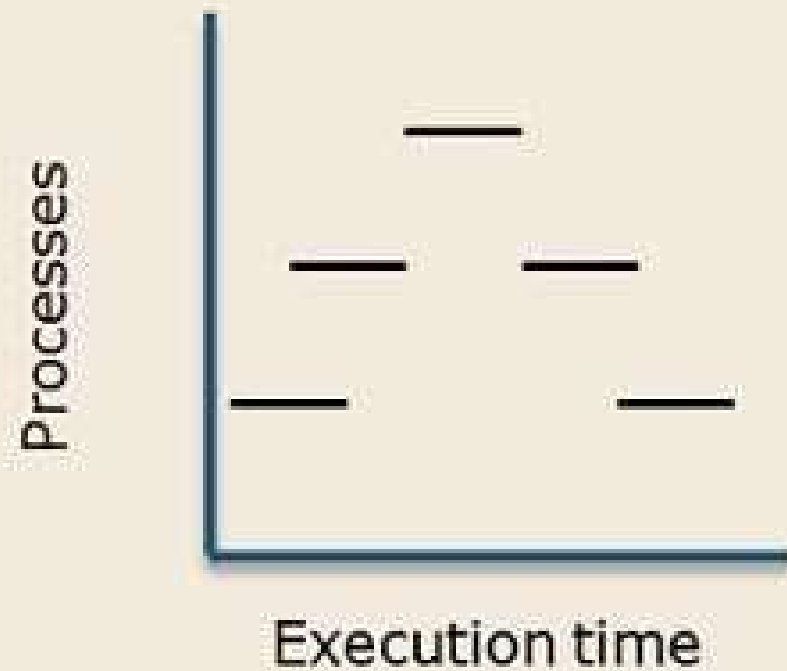
Rack 1    Write    Rack 2

Client

# Concurrent Processing

- The simultaneous execution of several interrelated computer programs
- A sequential computer program consists of a series of instructions to be executed one after another
- A concurrent program consists of several sequential programs to be executed in parallel
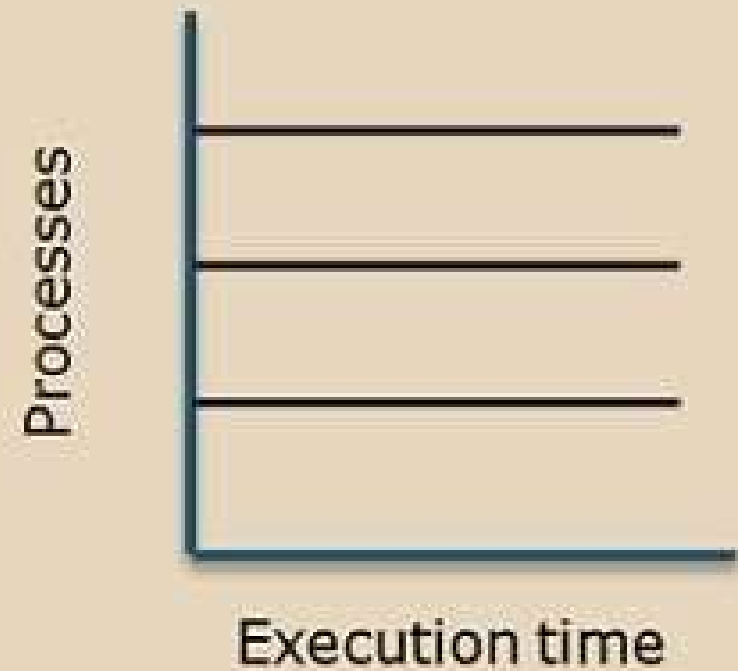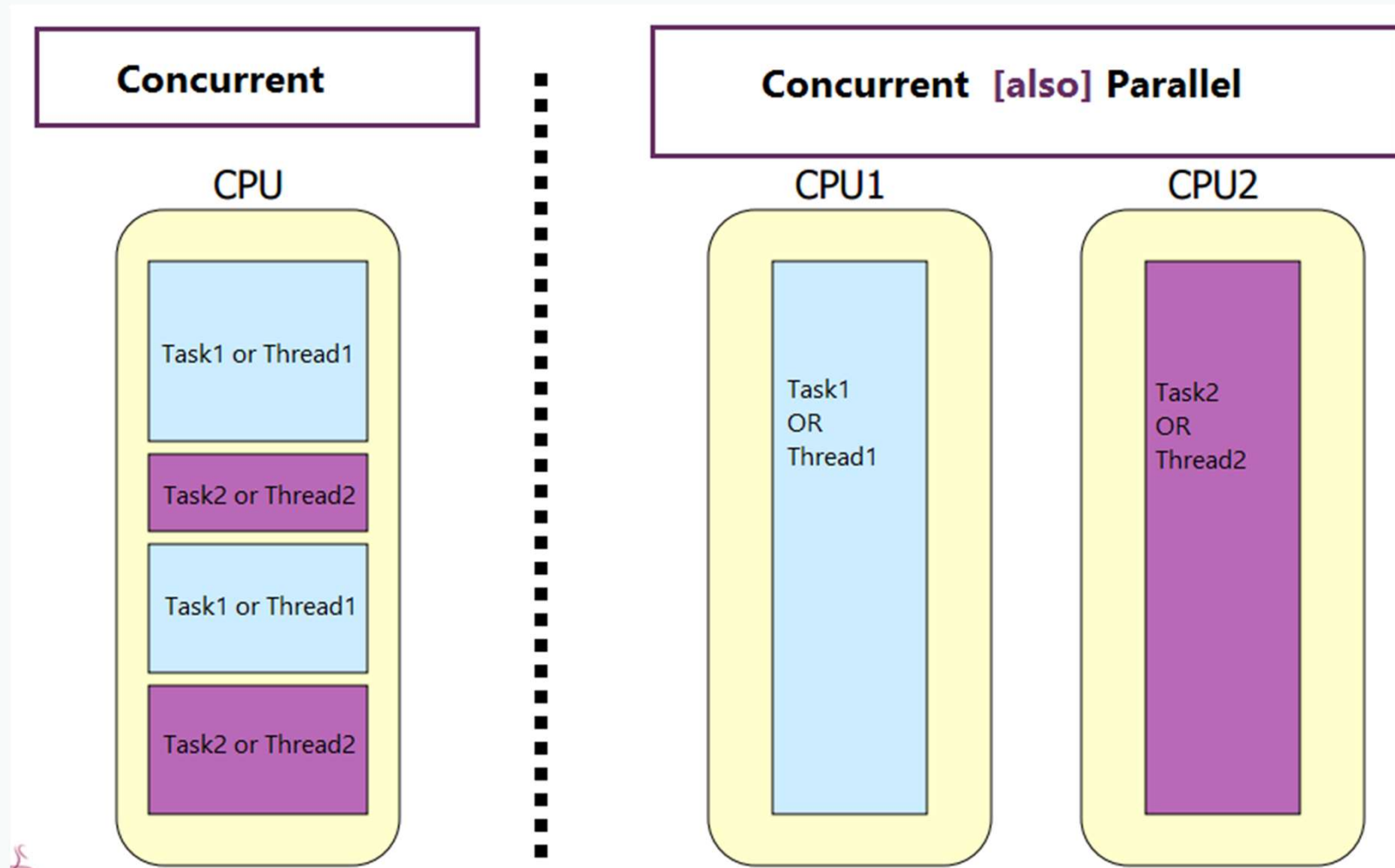- Each of the concurrently executing sequential programs is called a process

**CONCURRENCY** — Processes / Execution time
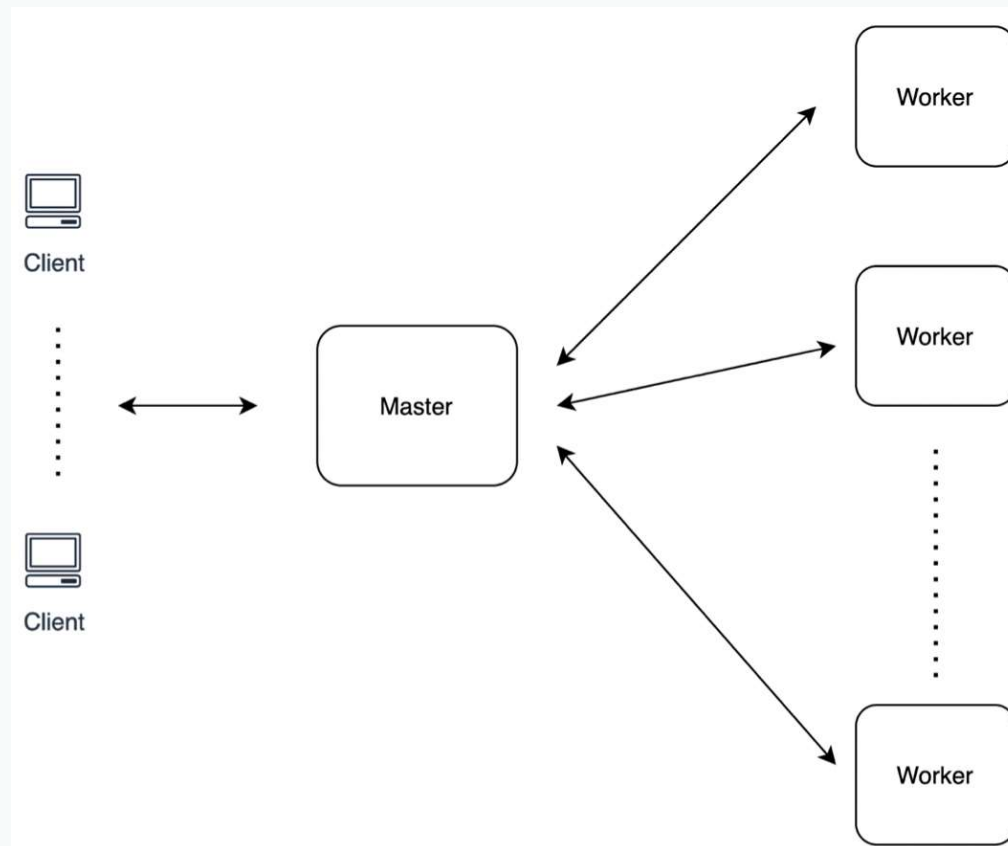
**VS**

**PARALLELISM** — Processes / Execution time
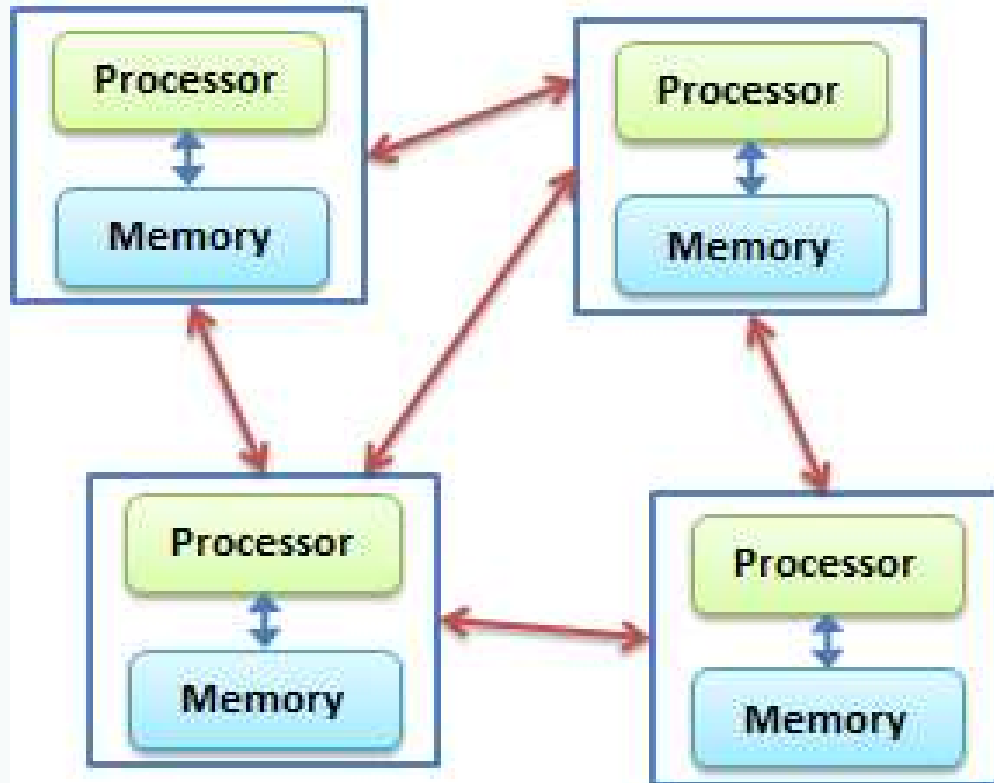
# Concurrency and Parallalism

# Distributed Processing

- Is the technique of linking together multiple computer servers over a network into a cluster

# Distributed Computing vs Parallel Computing

# HDFS Architecture

# HDFS Read Flow

# HDFS Write Flow

# Secondary Name Node

File system metadata

Name Node

File.txt = A,C

Its been an hour, give me your metadata

Secondary Name Node

- Not a hot standby for the Name Node
- Connects to Name Node every hour*
- Housekeeping, backup of Name Node metadata
- Saved metadata can rebuild a failed Name Node

# HDFS Redundancy

# HDFS Integrity



```
                              ┌──────────┐
                              │   File   │
                              └──────────┘
                                    │
        ┌───────────────────────────┼───────────────────────────┐
        │                           │                           │
┌───────────────────┐    ┌───────────────────┐    ┌───────────────────┐
│ HDFS Configuration A│    │ HDFS Configuration B│    │ ▢ Cloud Storage   │
│ (Block size n)      │    │ (Block size q)      │    │                   │
│   ┌───────────┐     │    │   ┌───────────┐     │    │   ┌───────────┐   │
│   │   File    │     │    │   │   File    │     │    │   │  Object   │   │
│   └───────────┘     │    │   └───────────┘     │    │   └───────────┘   │
└───────────────────┘    └───────────────────┘    └───────────────────┘
        │                           │                           │
┌───────────────────┐    ┌───────────────────┐    ┌───────────────────┐
│ Checksum:         │    │ Checksum:         │    │ Checksum:         │
│ [...]cec2f8b7     │    │ [...]911959f9     │    │ c517d290          │
└───────────────────┘    └───────────────────┘    └───────────────────┘
```
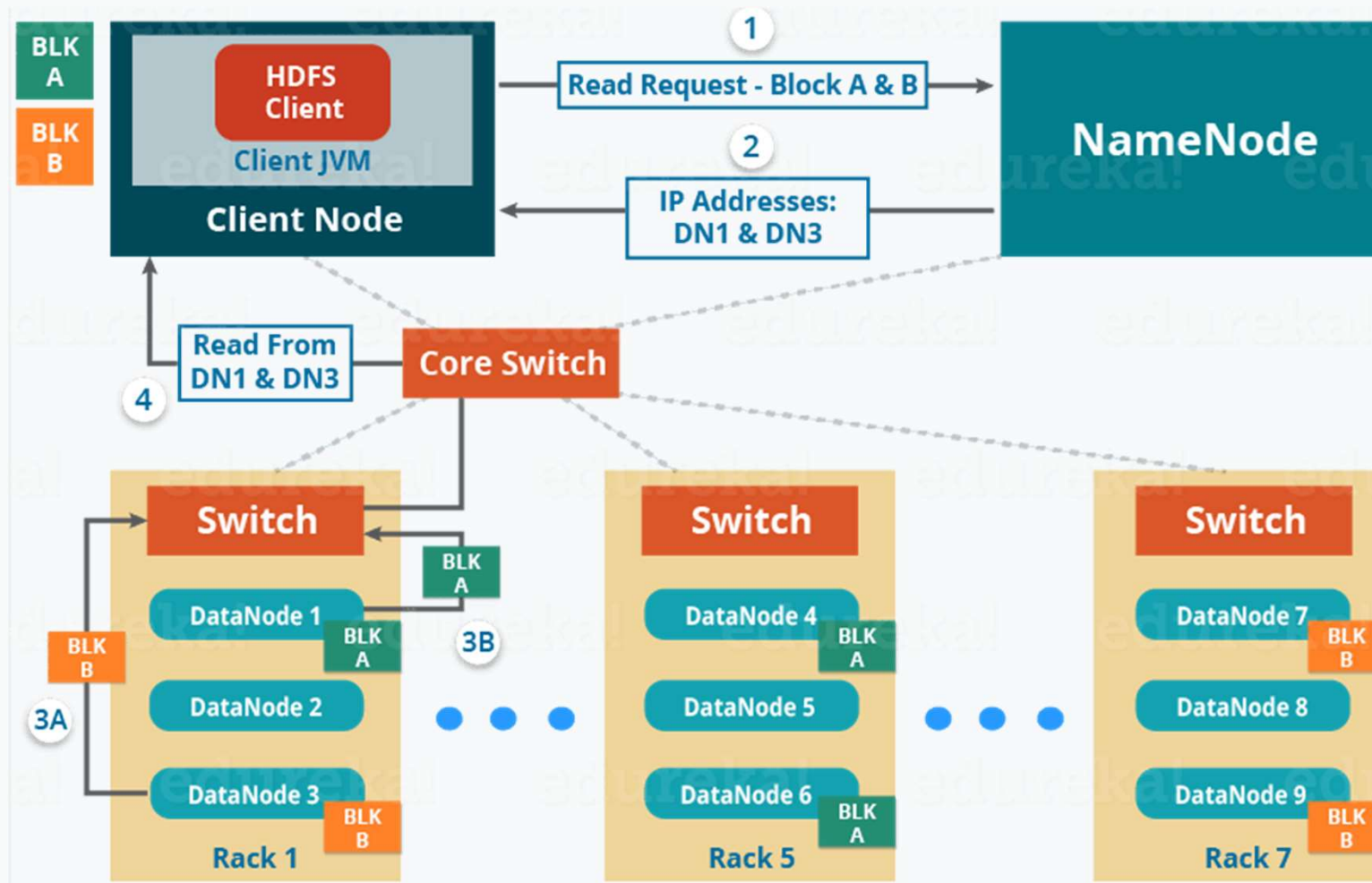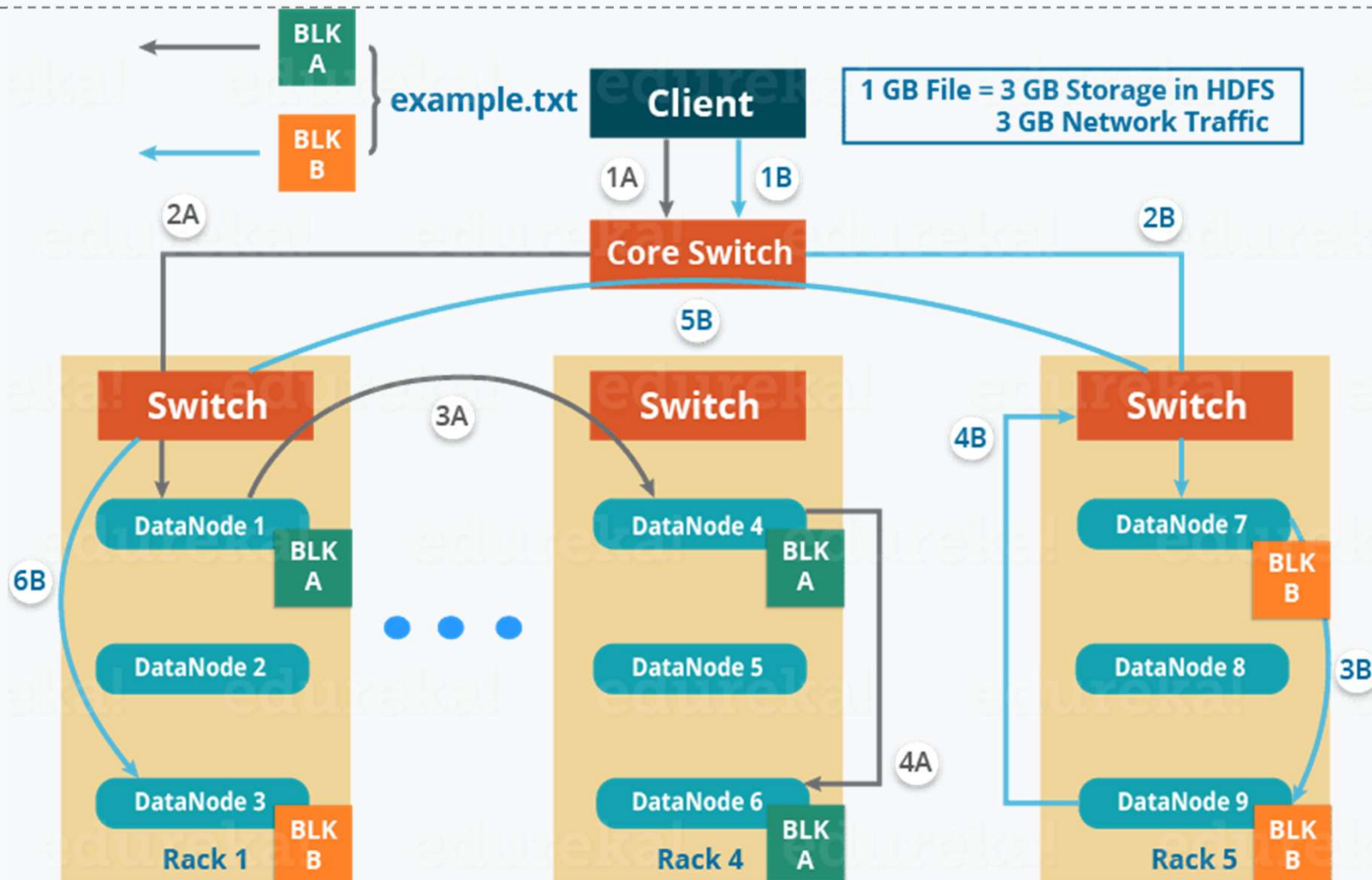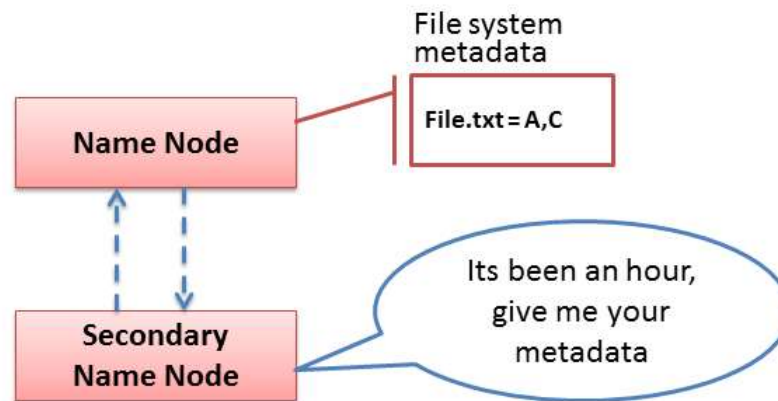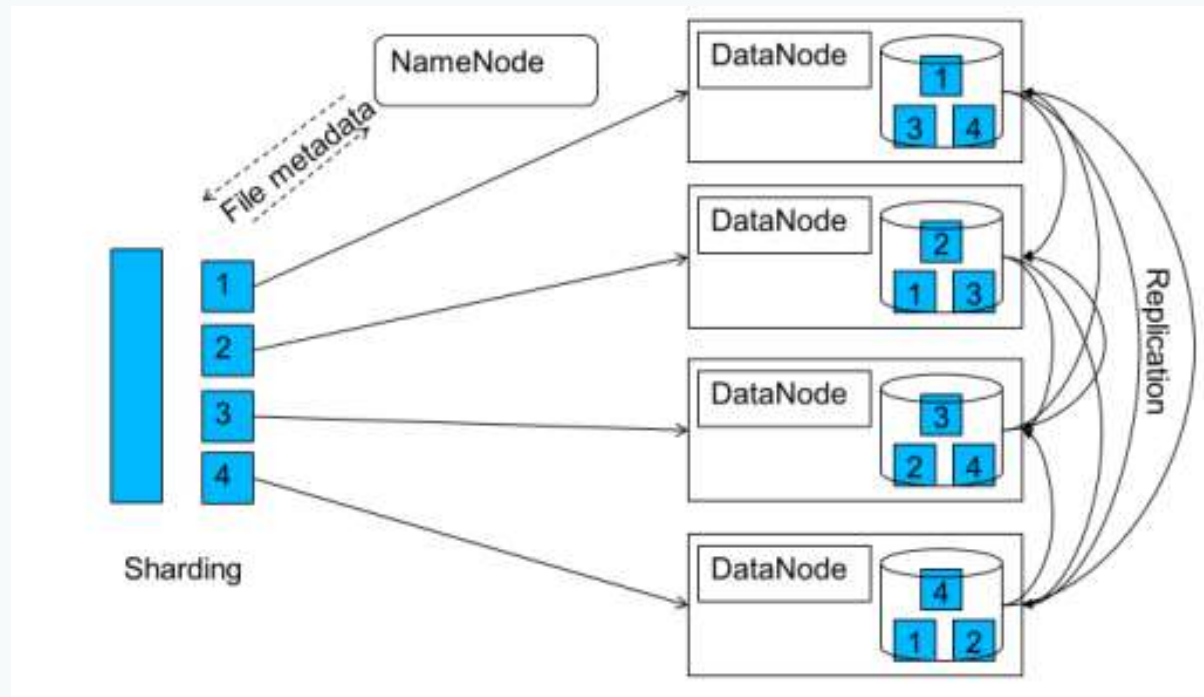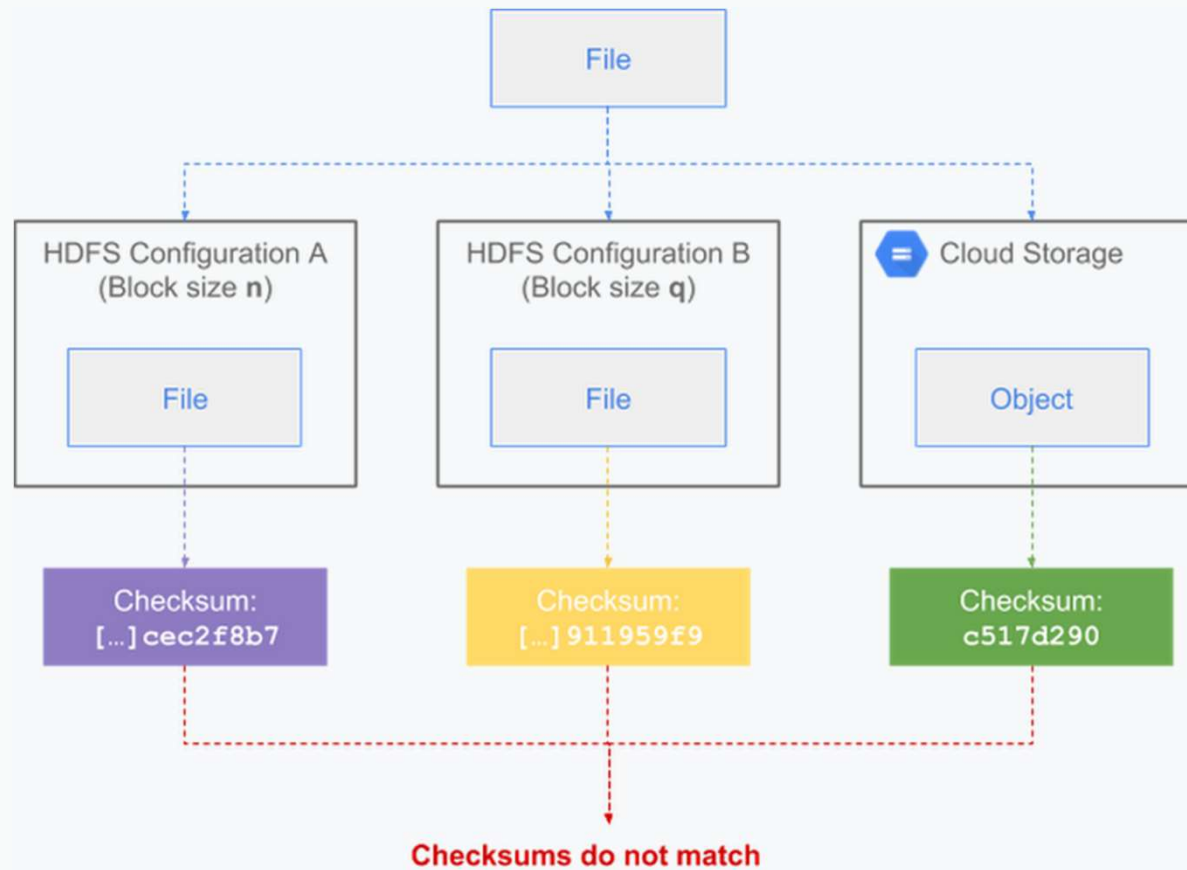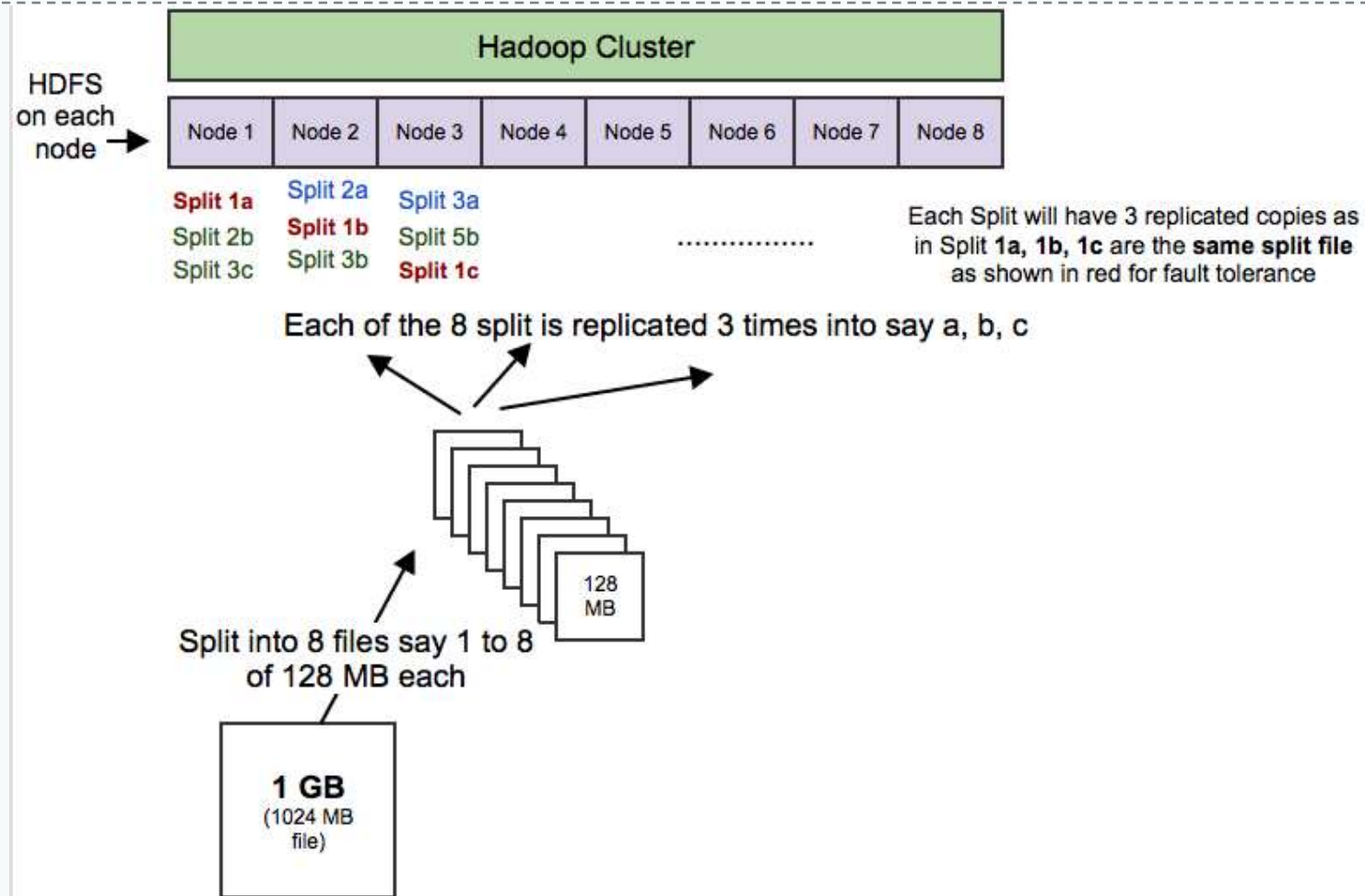
**Checksums do not match**

# Fault Tolerance

# Thank You