# COMP9313:
# Big Data Management

## Hadoop and HDFS

# Hadoop

- Apache Hadoop is an open-source software framework that
  - Stores big data in a distributed manner
  - Processes big data parallelly
  - Builds on large clusters of commodity hardware.
- Based on Google's papers on Google File System(2003) and MapReduce(2004).
- Hadoop is
  - Scalable to Petabytes or more easily (Volume)
  - Offering parallel data processing (Velocity)
  - Storing all kinds of data (Variety)

# Hadoop offers

- Redundant, Fault-tolerant data storage (HDFS)
- Parallel computation framework (MapReduce)
- Job coordination/scheduling  (YARN)

- Programmers no longer need to worry about
  - Where file is located?
  - How to handle failures & data lost?
  - How to divide computation?
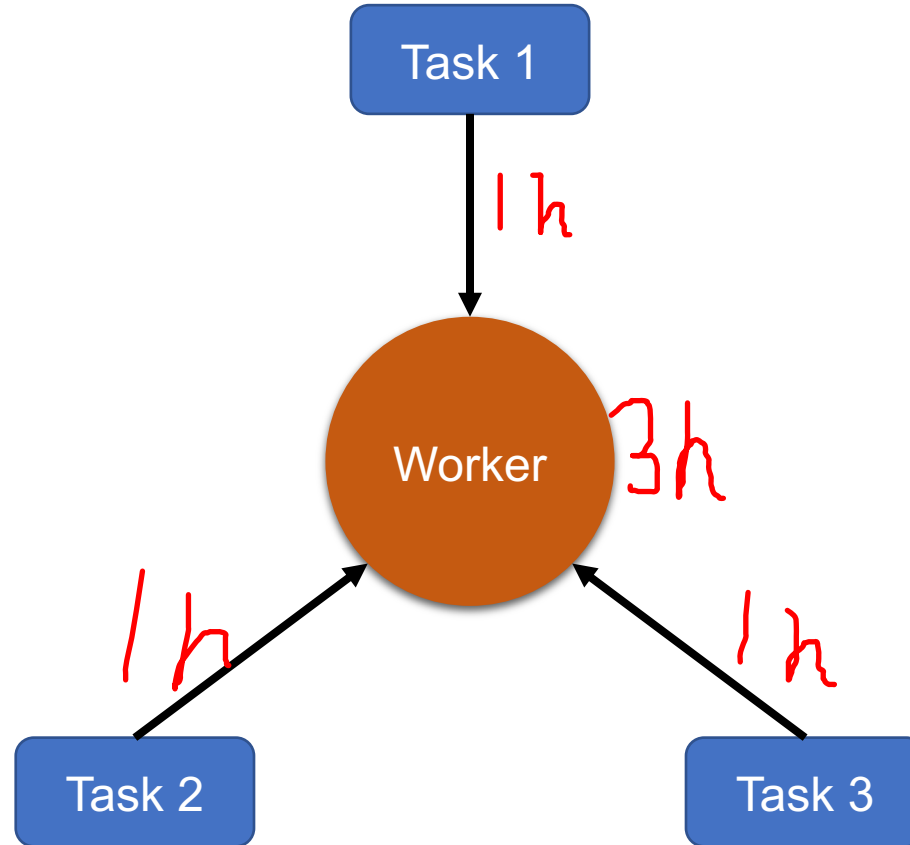  - How to program for scaling?

# Hadoop Ecosystem

- Hadoop

- Core of Hadoop
  - Hadoop distributed file system (HDFS)
  - MapReduce
  - YARN (Yet Another Resource Negotiator) (from Hadoop v2.0)

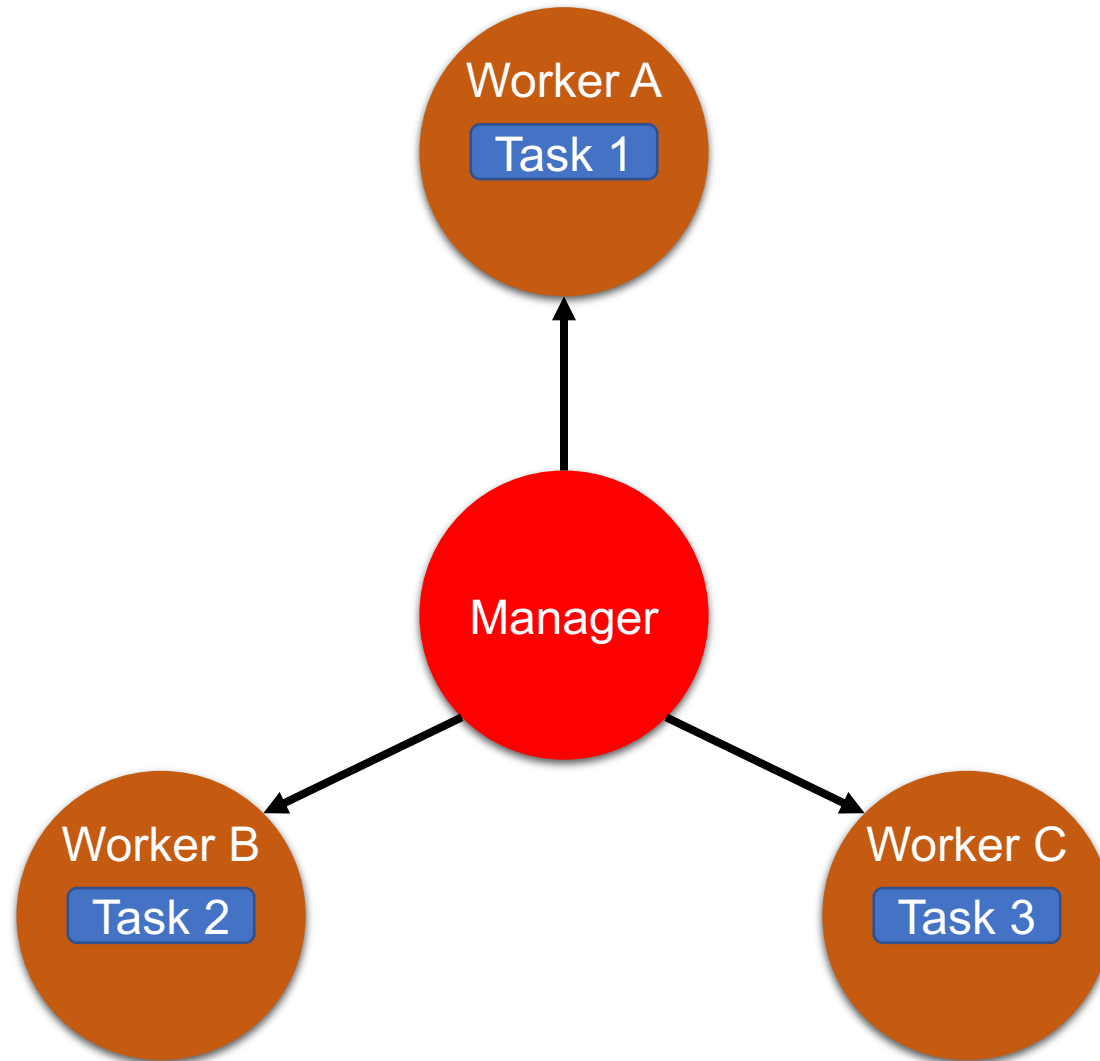- Additional software packages
  - Pig          Hadoop          Apache
  - Hive                  Hadoop                                          Hadoop Distributed File System
                          HDFS  HDFS                                                          low-cost
  - Spark                                          high throughput
                                          large data set                          HDFS          relax
  - HBase   POSIX                                  streaming access
             Hadoop                                          HDFS  MapReduce  HDFS
  - …                  MapReduce

4

# The Master-Slave Architecture of Hadoop

Task 1

1h

Worker
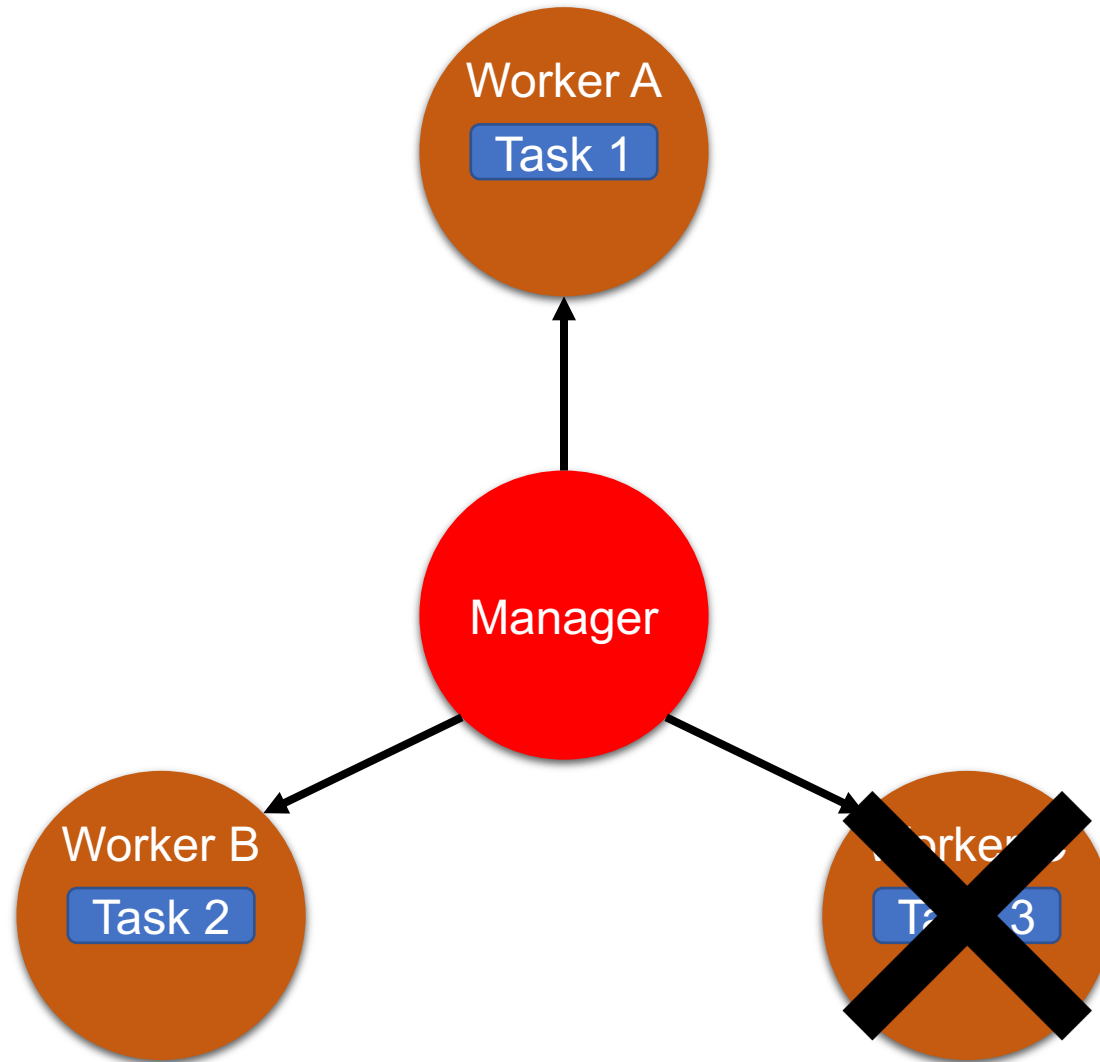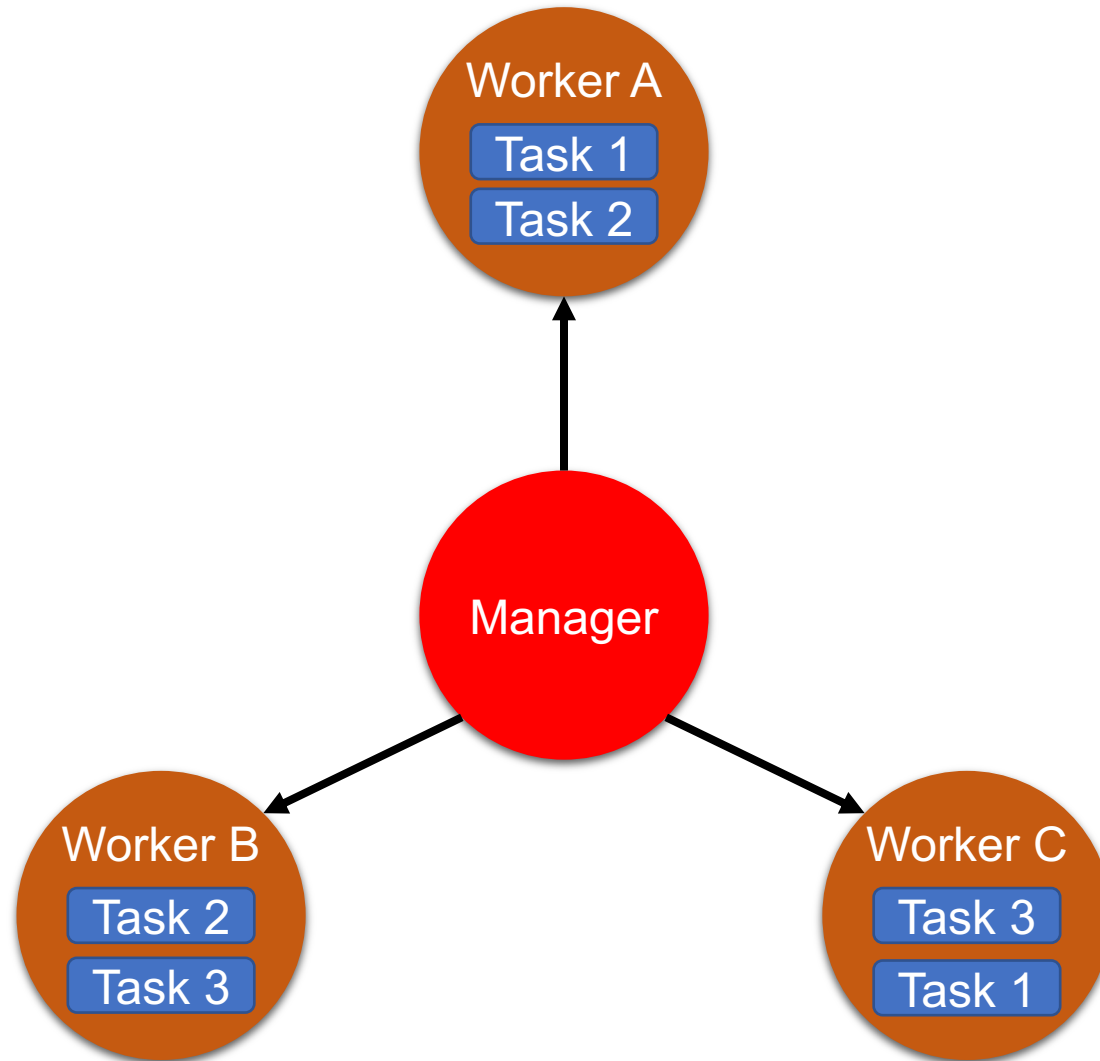
3h

1h

1h

Task 2

Task 3

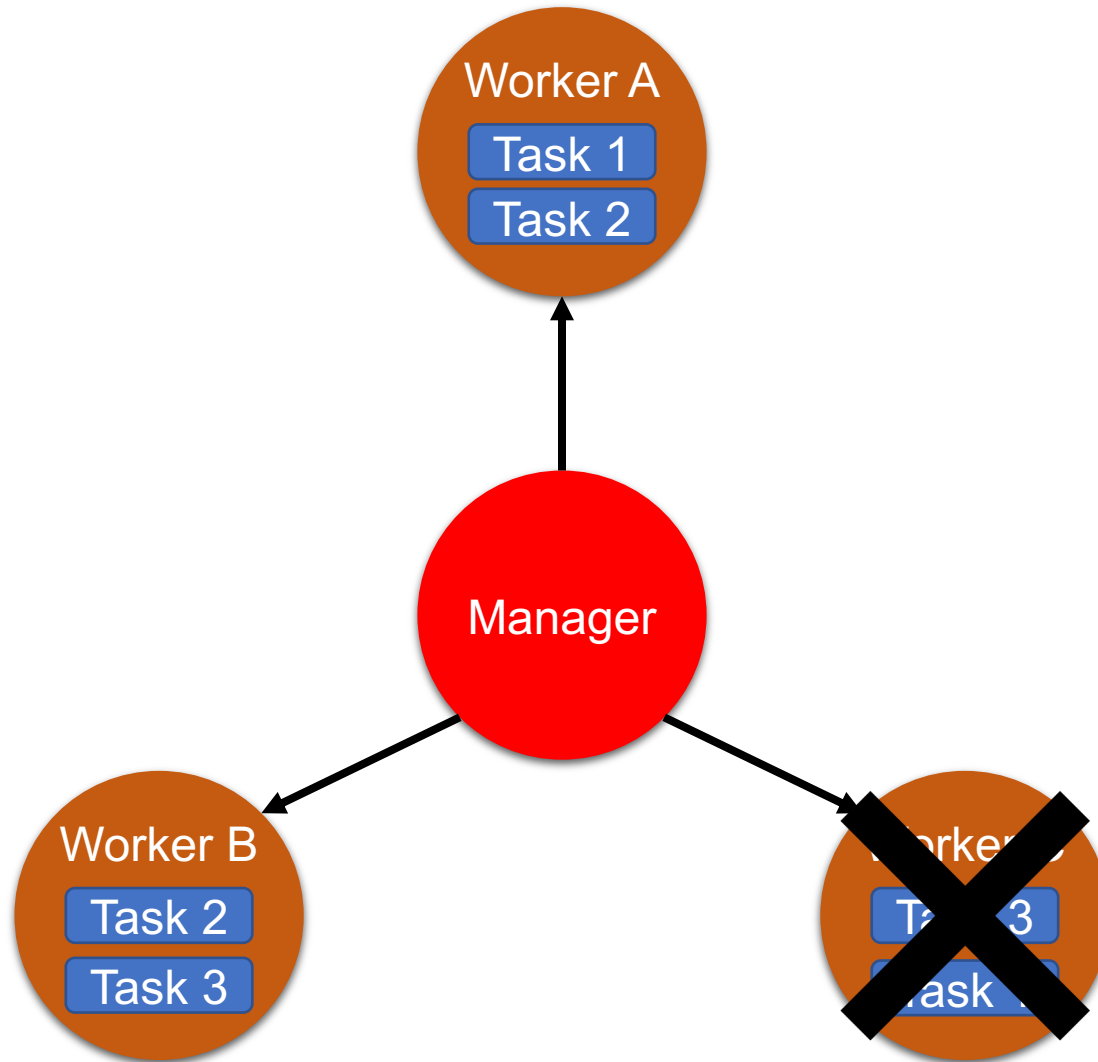# The Master-Slave Architecture of Hadoop

# The Master-Slave Architecture of Hadoop

# The Master-Slave Architecture of Hadoop

# The Master-Slave Architecture of Hadoop

# Hadoop Distributed File Systems (HDFS)

- HDFS is a file system that  <span style="color:red">HDFS</span>
  - follows master-slave architecture
  - allows us to store data over multiple nodes (machines) ,
  - allows multiple users to access data.
  - just like file systems in your PC

- HDFS supports
  - distributed storage
  - distributed computation
  - horizontal scalability

# Vertical Scaling vs. Horizontal Scaling



Vertical Scaling

Horizontal Scaling

# HDFS Architecture



HDFS Client

NameNode

Secondary NameNode

file 1

file 2

block

Rack 1

Rack 2

DataNode

DataNode

DataNode

DataNode

DataNode

Local Disks

# NameNode

NameNode          DataNodes

- NameNode maintains and manages the blocks in the DataNodes (slave nodes).
  - Master node

- Functions:
  - records the metadata of all the files
    - FsImage: file system namespace   NameNode          FsImage
    - EditLogs: all the recent modifications
  - records each change to the metadata
  - regularly checks the status of datanodes
  - keeps a record of all the blocks in HDFS
  - if the DataNode failure, handle data recoveray

DataNode                                        FsImage

EditLog

                   NameNode

# DataNode

- A commodity hardware stores the data
  - Slave node

- Functions
  - stores actual data
  - perform the read and write requests
  - report the health to NameNode (heartbeat)

DataNode HDFS Hadoop
NameNode DataNode DataNode
Hadoop
[5]
DataNode HDFS NameNode
NameNode DataNode
heartbeat NameNode
DataNode
NameNode

# NameNode vs. DataNode

|  | NameNode | DataNode |
|---|---|---|
| Quantity | One | Multiple |
| Role | Master | Slave |
| Stores | Metadata of files | Blocks |
| Hardware Requirements | High Capacity Memory | High Volume Hard Drive |
| Failure rate | Lower | Higher |
| Solution to Failure | Secondary NameNode | Replications |

# If NameNode failed…

- All the files on HDFS will be lost
  - there's no way to reconstruct the files from the blocks in DataNodes without the metadata in NameNode

NameNode

- In order to make NameNode resilient to failure
  - back up metadata in NameNode (with a remote NFS mount)
  - Secondary NameNode

# Secondary NameNode

NameNode

- Take checkpoints of the file system metadata present on NameNode

NameNode

  - It is not a backup NameNode!

- Functions:
  - Stores a copy of FsImage file and Editlogs
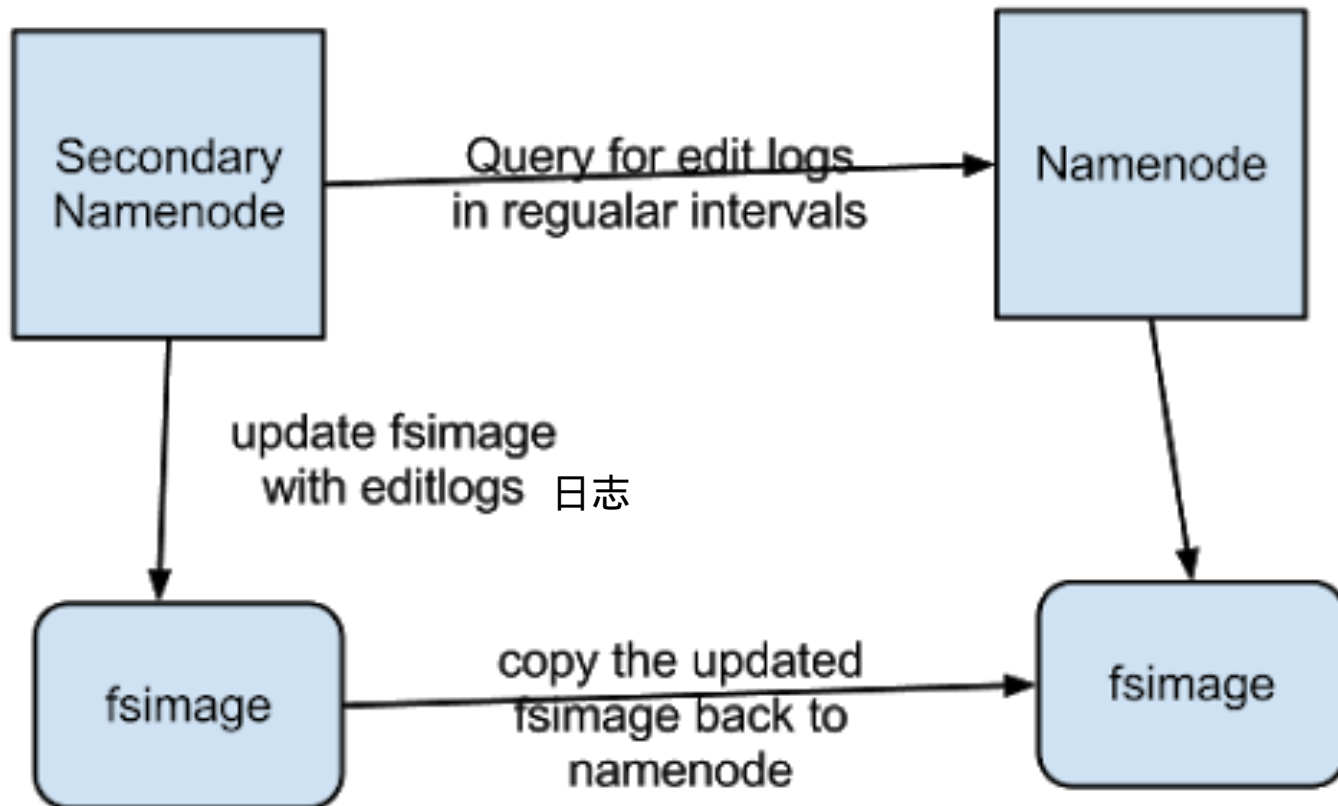  - Periodically applies Editlogs to FsImage and refreshes the Editlogs.

    FsImage

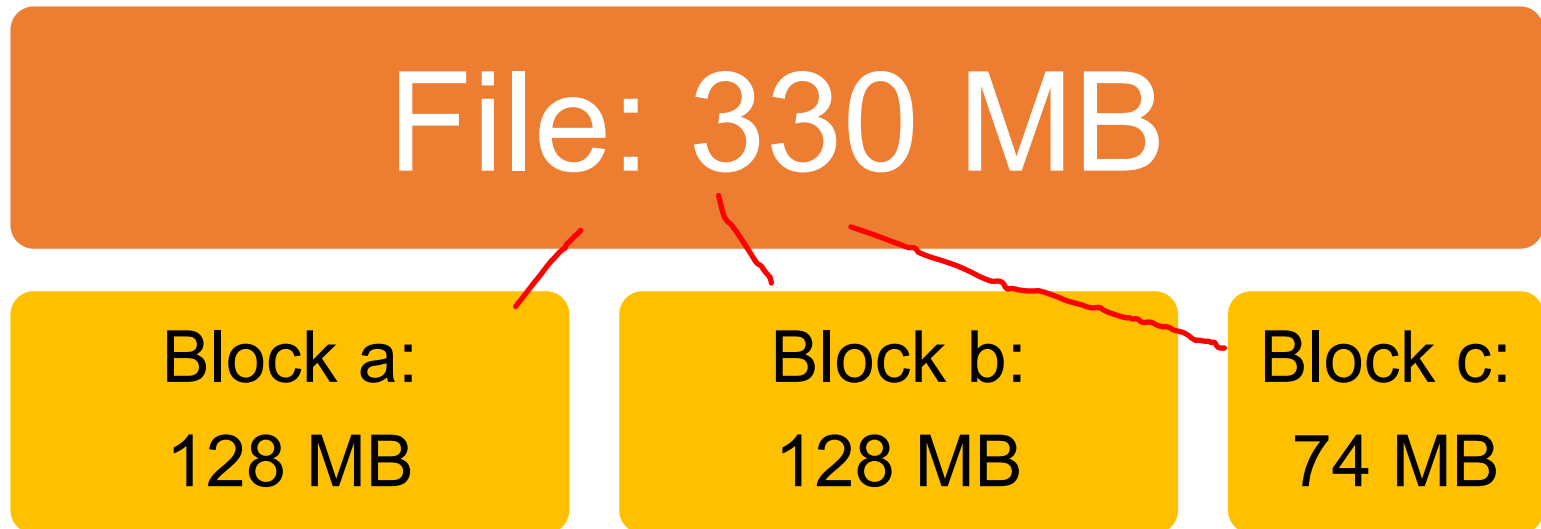  - If NameNode is failed, File System metadata can be recovered from the last saved FsImage on the Secondary NameNode.

# NameNode vs. Secondary NameNode

# Blocks

- Block is a sequence of bytes that stores data
  - Data stores as a set of blocks in HDFS
  - Default block size is 128MB (Hadoop 2.x and 3.x)
  - A file is spitted into multiple blocks

| File: 330 MB | | |
|---|---|---|
| Block a: 128 MB | Block b: 128 MB | Block c: 74 MB |

# Why Large Block Size?

- HDFS stores huge datasets
- If block size is small (e.g., 4KB in Linux), then the number of blocks is large:
  - too much metadata for NameNode
  - too many seeks affect the read speed
  - harm the performance of MapReduce too

MapReduce

- We don't recommend using HDFS for small files due to similar reasons.
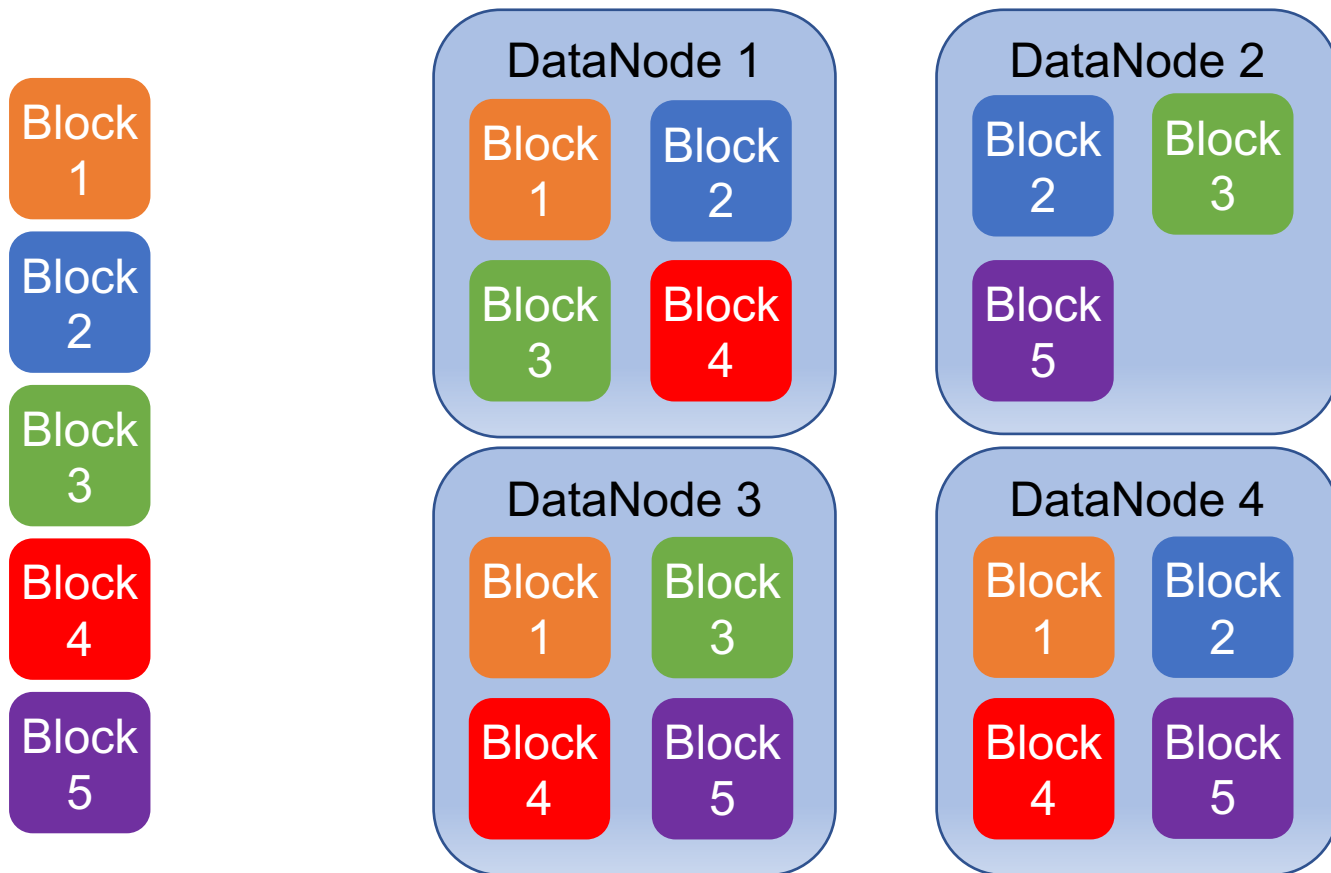  - Even a 4KB file will occupy a whole block.

# If DataNode Failed…

- Commodity hardware fails
  - If NameNode hasn't heard from a DataNode for 10mins, The DataNode is considered dead…
- HDFS guarantees data reliability by generating multiple replications of data
  - each block has 3 replications by default
  - replications will be stored on different DataNodes
  - if blocks were lost due to the failure of a DataNode, they can be recovered from other replications
  - the total consumed space is 3 times the data size

- It also helps to maintain data integrity

# Replication Management

- Each block is replicated 3 times and stored on different DataNodes

# Why default replication factor = 3?

- If 1 replicate
  - DataNode fails, block lost
- Assume
  - # of nodes N = 4000
  - # of blocks R = 1,000,000
  - Node failure rate FPD = 1 per day
    B=
- If one node fails, then R/N = 250 blocks are lost
  - E(# of losing blocks in one day) = 250
- Let the number of losing blocks follows Poisson distribution, then
  - Pr[# of losing blocks in one day >= 250] = 0.508

# Why default replication factor = 3?

- Assume
  - # of nodes N = 4000
  - Capacity of each node GB = 4000 Gigabytes
  - # of block replicas R = 1,000,000 * 3
  - Node failure rate FPD = 1 per day
  - Replication speed = 1.35 MB per second per node
- If one node fails, B = R/N = 750 replicas/blocks are unavailable
- There are on average S = 2B/(N-1) = 0.38 replicas per node for the blocks in the failed node
- So if second node fails, 0.38 blocks now have only a single replica

# Why default replication factor = 3?

- If the third node fails,
  - The probability that it has the only remaining replica of a particular block is
    - $Pr[last] = 1/(N-2) = 0.000250$
  - The probability that it has none of those replicas is
    - $Pr[none] = (1-Pr[last])^S = 0.999906$
  - The probability of losing the last replica of a block is
    - $Pr[lose] = 1- Pr[none] = 9.3828E-05$

- Recall:
  - N is # of nodes
  - S is the # of replicas per node for the blocks in the first failed node

# Why default replication factor = 3?

- Assume # of node failures follows Poisson distribution with rate
  - $\omega$=FPD/(24*3600)=1.1574E-05 per second
- Re-replication is a fully parallel operation on the remaining nodes
  - Recovery (re-create the lost replicas) time is
    - 1000 * GB / MPS / (N-1) = 740.93 seconds
    - Recovery rate $\mu$= 1/ 740.93 per second
  - E(# of failed nodes in 1 sec) =$\omega/\mu$ = 0.008576
- At any second, the probability of k failed nodes follows Poisson distribution
  - Pr[0 failed node] = 0.991461
  - Pr[1 failed node] = 0.008502
  - Pr[2 or more failed nodes] = 1- Pr(0) - Pr(1) = 0.00003656
- Thus, the rate of third failure is
  - Pr[2 or more failed nodes] *$\omega$= 4.2315E-10 per sec
- The rate of losing a data block is
  - $\lambda$=Pr[2 or more failed nodes] *$\omega$* Pr[lose]  = 3.9703E-14

# Why default replication factor = 3?

- Recall that in one second, the rate of losing a data block is
  - $\lambda = 3.9703\text{E-}14$ per second
- According to exponential distribution, we have:
  - Pr[losing a block in one year] = $1-e^{-\lambda t}$ = 0.00000125
    - $t = 365*24*3600$

- So replication factor = 3 is good enough.

# What about Simultaneous Failure?

- If one node fails, we've lost B (first) replicas

- If two nodes fail, we've lost some second replicas and more first replicas

- If three nodes fail, we've lost some third replicas, some second replicas and some first replicas

- …

# What about Simultaneous Failure?

- Assume k of N nodes have failed simultaneously, let there be
    - L1(k,N) blocks have lost one replica
    - L2(k,N) blocks have lost two replicas
    - L3(k,N) blocks have lost three replicas
    - B is # of unavailable blocks if one node fails
- k=0:
    - L1(0,N) = L2(0,N) = L3(0,N) = 0
- k=1:
    - L1(1,N) = B
    - L2(1,N) = L3(1,N) = 0
- k=2:
    - L1(2,N) = 2B-2*L2(2,N)
    - L2(2,N) = 2*L1(1,N)/(N-1)
    - L3(2,N) = 0
- k=3:
    - L1(3,N) = 3B-2*L2(3,N)-3*L3(3,N)
    - L2(3,N) = 2*L1(2,N)/(N-2)+L2(2,N)-L3(3,N)
    - L3(3,N) = L2(2,N)/(N-2)

B= R/N

# What about Simultaneous Failure?

7    9

- In general
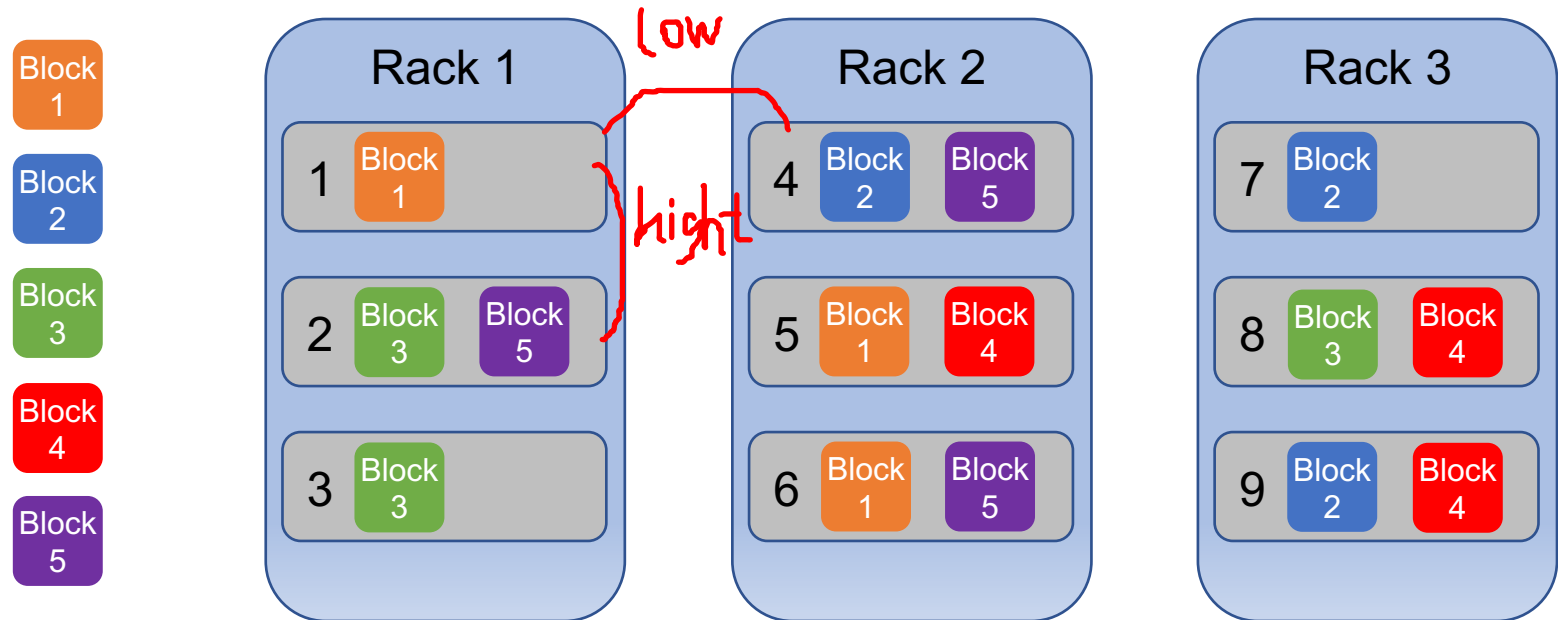  - L1(k,N) = k*B-2*L2(k,N)-3*L3(k,N)
  - L2(k,N) = 2*L1(k-1,N)/(N-k+1)+L2(k-1,N)- ~~L3(k,N)~~  L2(k-1, N) /( N-k+1)
  - L3(k,N) = L2(k-1,N)/(N-k+1)+L3(k-1,N)
- Let N = 4000, B = 750, we have

| Failed Nodes | 1st replicas lost | 2nd replicas lost | 3rd replicas lost |
|---|---|---|---|
| 50 | 36,629 | 433 | 2 |
| 100 | 72,002 | 1,479 | 13 |
| 150 | 107,374 | 2,504 | 39 |
| 200 | 143,963 | 2,905 | 76 |

# Rack Awareness Algorithm

- If the replication factor is 3:
  - 1st replica will be stored on the local DataNode
  - 2nd on a different rack from the first.
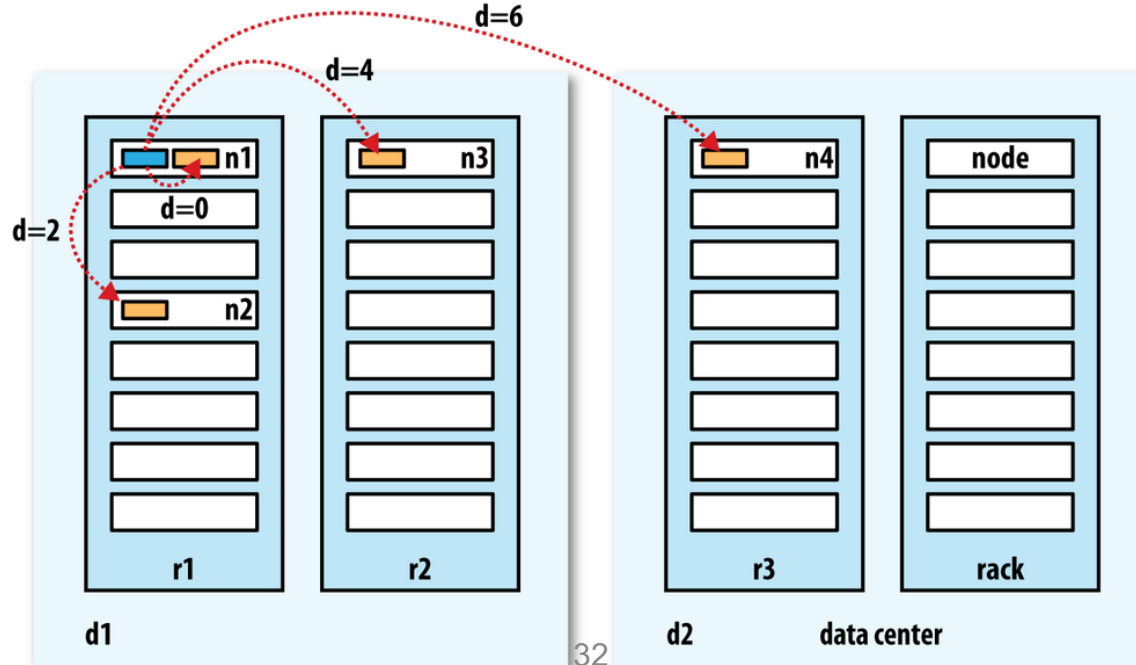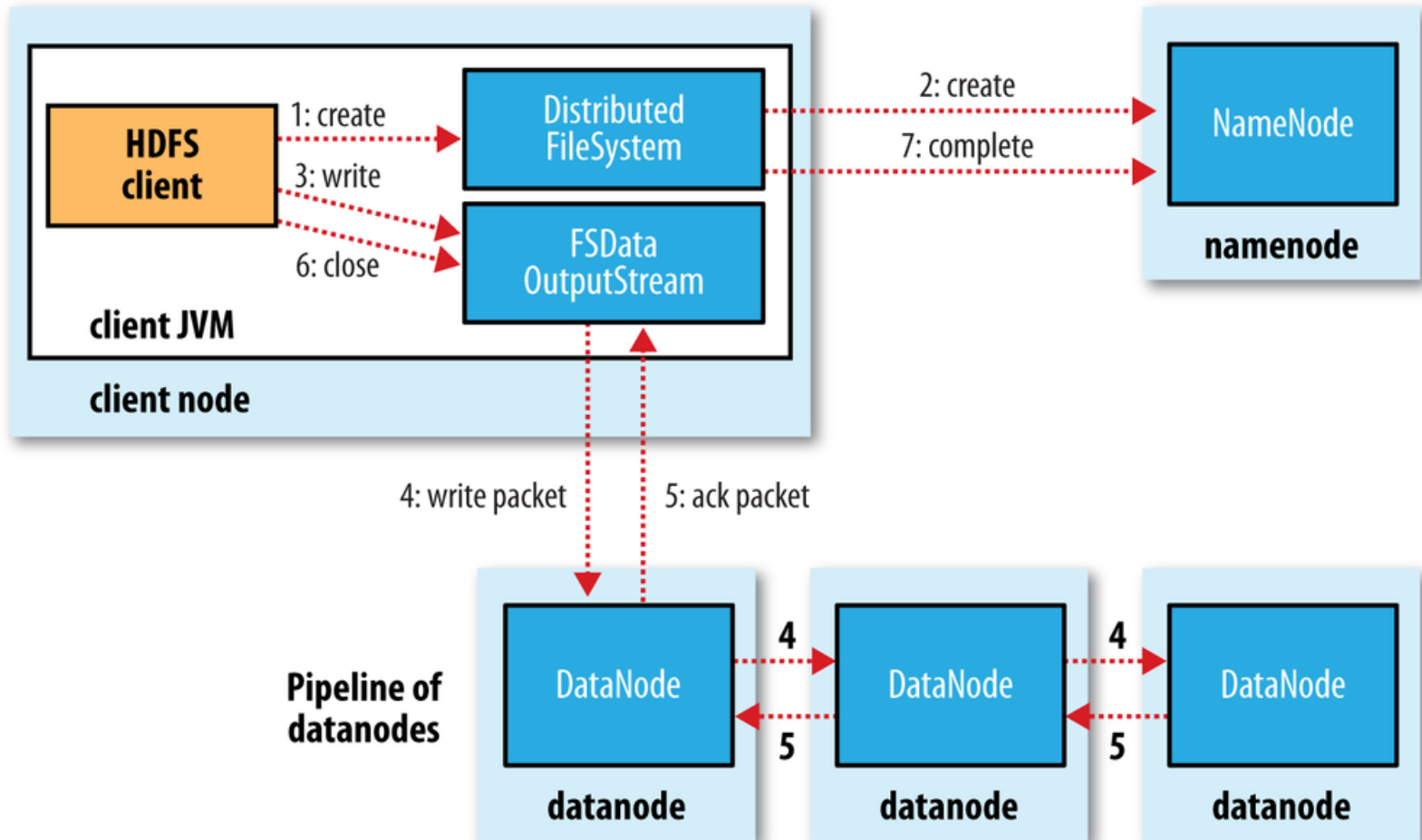  - 3rd on the same rack as 2nd, but on a different node.

# Why Rack Awareness?

- Reduce latency
  - Write: to 2 racks instead of 3 per block
  - Read: blocks from multiple racks

- Fault tolerance
  - Never put your eggs in the same basket
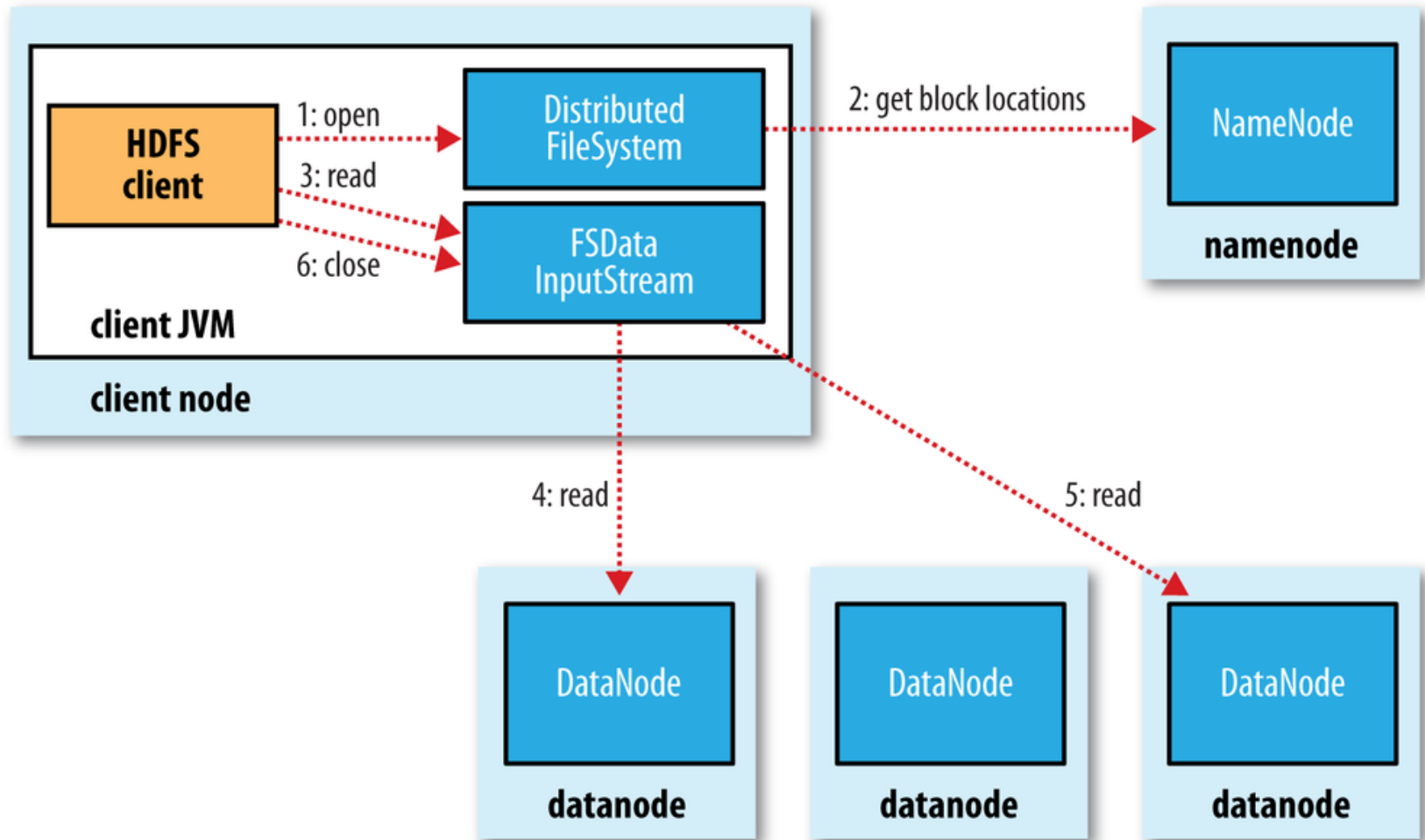
# Write in HDFS

- Create file – Write file – Close file

# Write in HDFS

- There is only single writer allowed at any time
- The blocks are writing simultaneously
- For one block, the replications are replicating sequentially
- The choose of DataNodes is random, based on replication management policy, rack awareness, …

  DataNodesis

# Read in HDFS

# Read in HDFS

- Multiple readers are allowed to read at the same time
- The blocks are reading simultaneously
- Always choose the closest DataNodes to the client (based on the network topology)
- Handling errors and corrupted blocks
  - avoid visiting the dataNode again
  - report to NameNode