# Big Data
# HDFS

**K V Subramaniam**

Computer Science and Engineering

**Lecture overview - HDFS**

- Need for Distributed file systems
- HDFS Introduction
- Architecture
- Operations
- Internals
- Fault Tolerance and replication
- Blocks

# Big Data: File systems and distributed file systems

## Data Growth – Why the need for HDFS?

As per RBI in May 2019,

#credit/debit card transactions~ 1.3 Billion

(https://rbidocs.rbi.org.in/rdocs/ATM/PDFs/ATM052019E96EC259708C4ED9AD9E0C6B5E8B6DD5.PDF)

If each transaction requires about 10K of data

# 13 TB of data



Lot of data and this is only for credit/debit card transactions

There are other transactions also

Suppose you want to look for fraudulent transactions
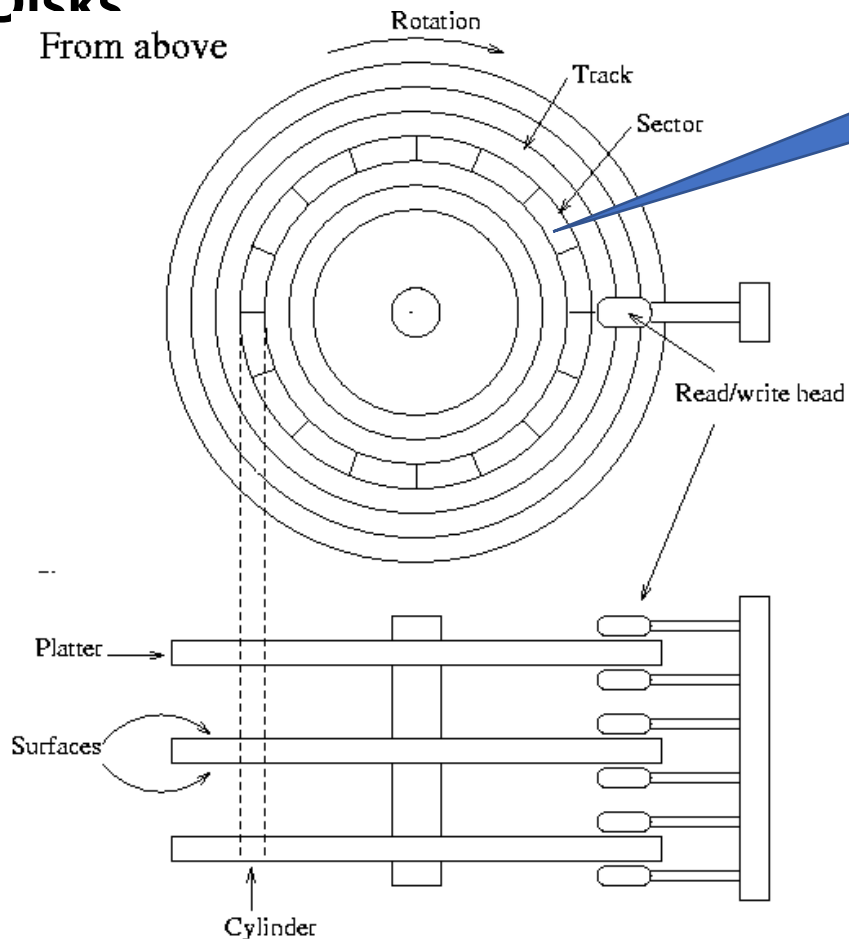
How to store and process this data?

This class

**Persistent Storage - Disks**



Block oriented device : storage divided into fixed size blocks

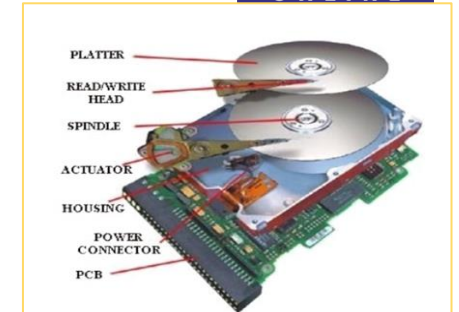Can we store the persistent data directly on these blocks?

Who will maintain the meta-data?

**Concerns**: Which block contains the data as files are not necessarily multiples of blocks size

A filesystem layer on top of disk manages blocks and maps data/metadata to blocks
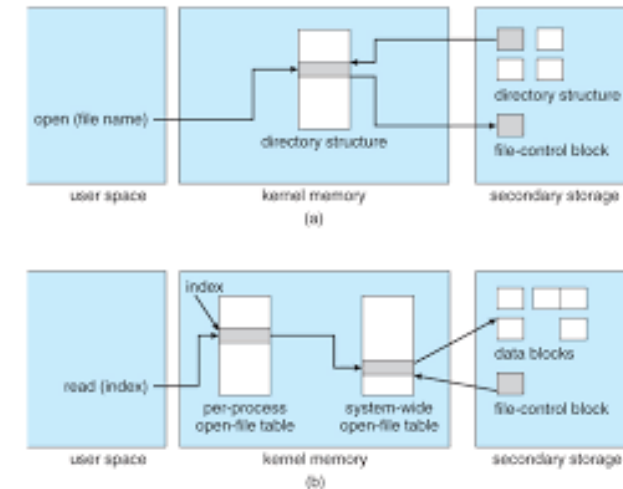
## HDFS − Hadoop distributed File system



**FILES**

- Named collection of related information
- on disks

**Desirable properties of files**

- Long-term existence
- Sharable between processes
- Access permissions

https://www.cs.uic.edu/~jbell/CourseNotes/OperatingSystems/12_FileSystemImplementation.html
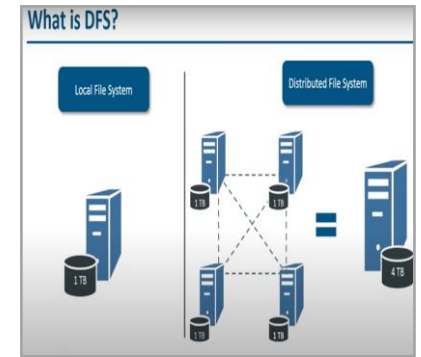
## Distributed File System

Consider case when data is so large that it cannot fit on a single disk.

- A DFS manages files and folders across multiple computers.

DFS can organize and display the files as if they are stored on one computer.

It serves the same purpose as a traditional file system.

Designed to provide file storage and controlled access to files over local and wide area networks.


What is DFS?
Local File System    Distributed File System

Consider that you have 1TB of data?
Compare the time taken to read data in both the cases below
- single machine (*4 I/O channels each channel 100mb/s*)
- 10 machines (*each having 4 I/O channels each channel 100mb/s*).
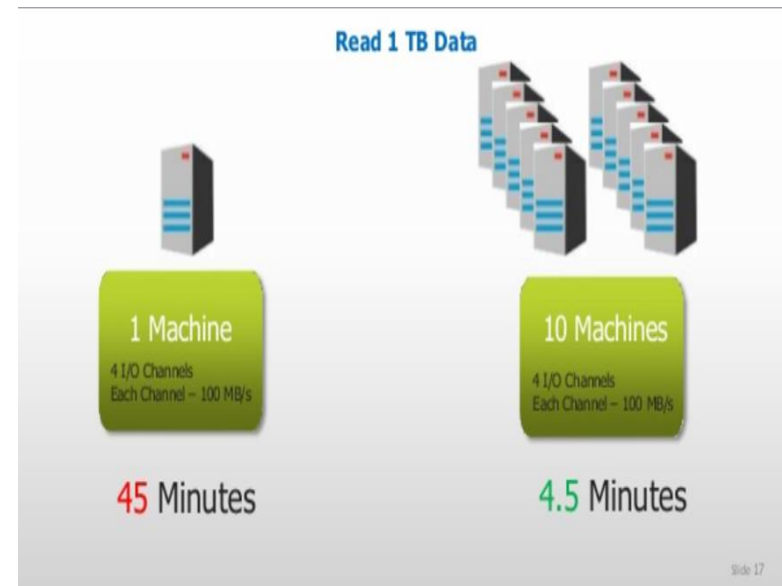
# Big Data: HDFS Introduction

Consider that you have 1TB of data?
Compare the time taken in both the cases below
- single machine (*4 I/O channels each channel 100mb/s*)
- 10 machines (*each having 4 I/O channels each channel 100mb/s*).

**HDFS – Inspired by GFS**

GFS – Google File System (2003)

Distributed File system on a cluster of machines

Developed by Doug Cutting and Mike Cafarella

Origin  - Apache Nutch

- Goal : web search engine on 1 Billion Pages

Open source

# HDFS – Hadoop distributed File system

" HDFS is a filesystem for storing <u>very large files</u> with <u>streaming data access</u> patterns, running on clusters of <u>commodity hardware</u>. "

**Very large**
- Files can be MB/GB/TB in size
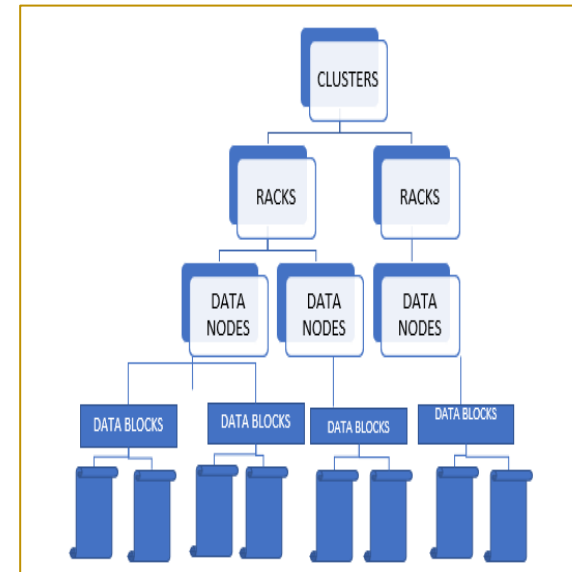- Hadoop clusters that are PB are currently operational

**Read Mostly data**
- most efficient data processing pattern is a write-once, read-many-times pattern.
- Each analysis will involve a large proportion of the dataset
- time to read the whole dataset is more important than the latency in reading the first record.

**Commodity hardware**
- Hadoop doesn't require expensive, high reliable hardware
- Designed to run on clusters of commodity hardware

# HDFS – Hadoop distributed File system

- If you want to store a file on disk what constitutes
  - Data
  - Metadata
- What are their access patterns?
  - How often do you think each one will be accessed during a normal file read
- How large are they (comparatively)? Why is this important?

# Big Data: HDFS Architecture

- If you want to store a file on disk what constitutes
  - Data
  - Metadata
- What are their access patterns?
  - How often do you think each one will be accessed during a normal file read
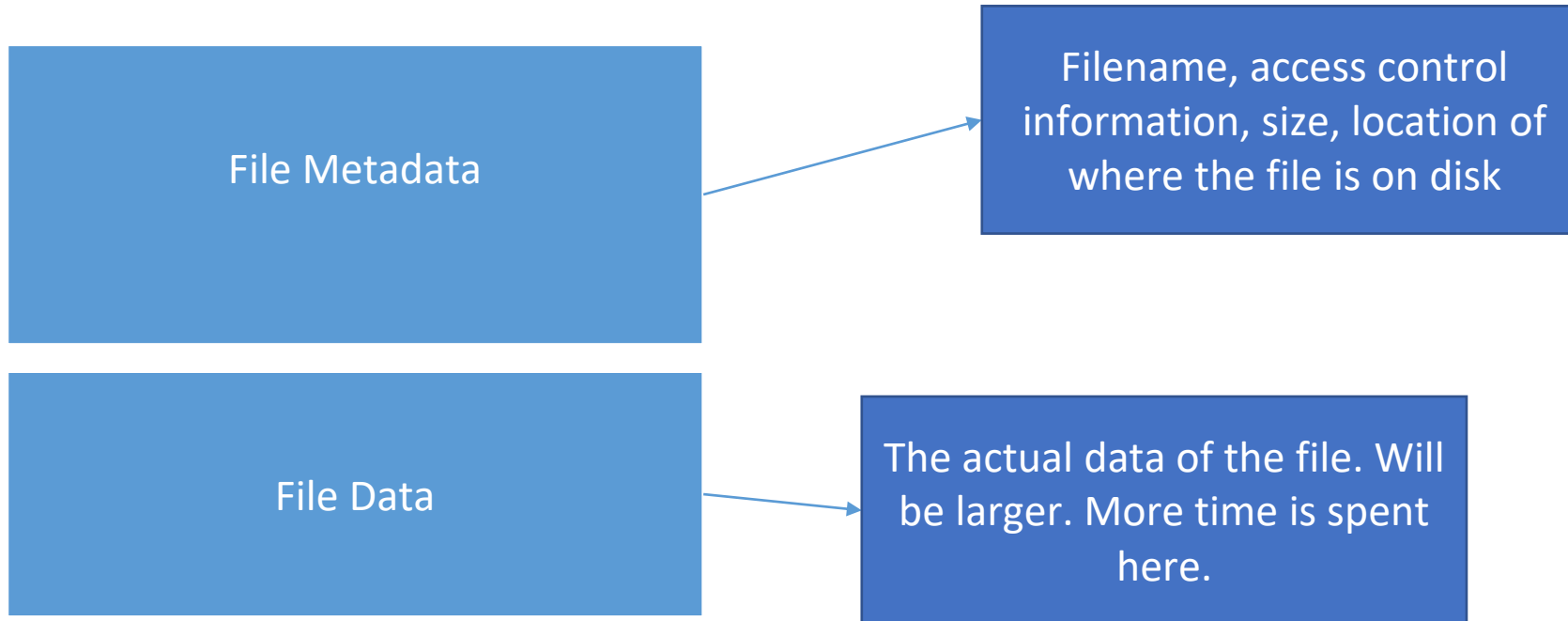- How large are they (comparatively)? Why is this important?

Data: much larger in size. Occupies multiple blocks

Metadata – smaller compared to data
Only information on filenames and blocks it occupies

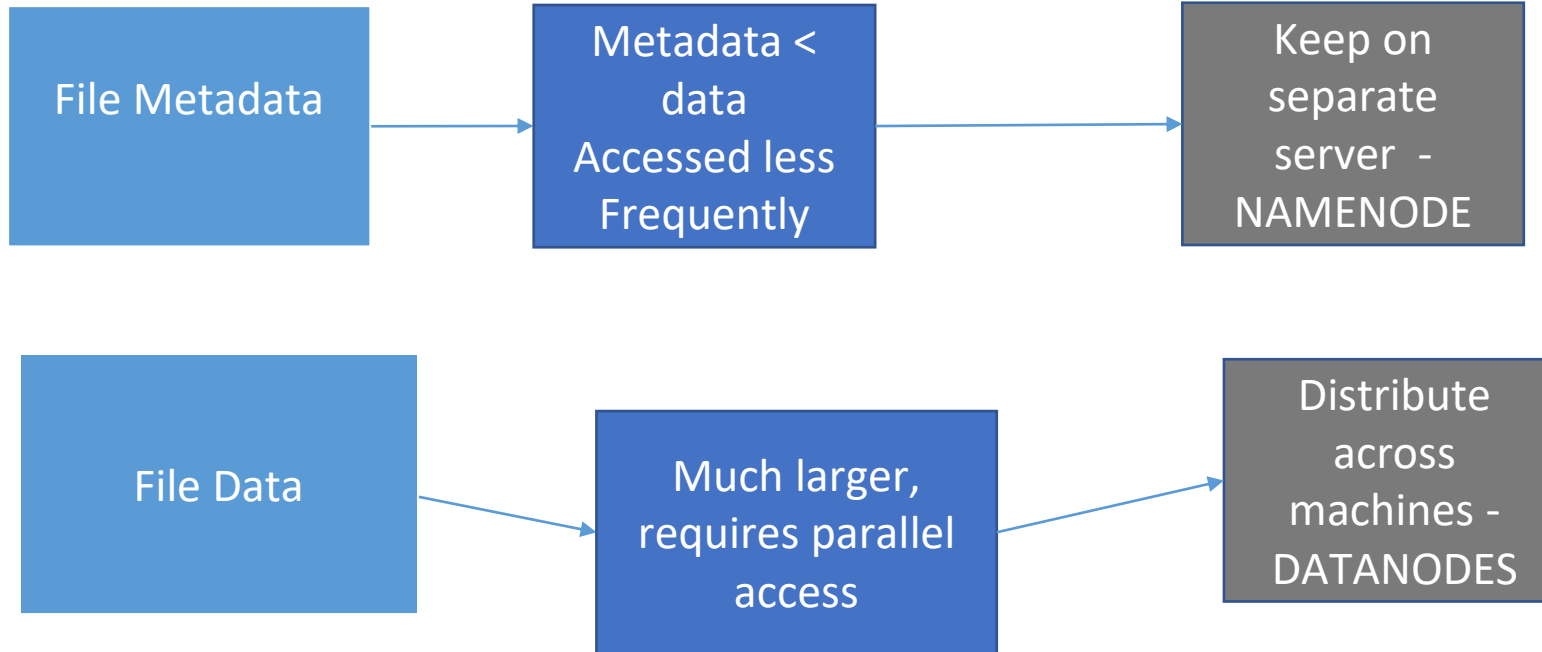Since data is much larger. Most time spent in fetching data.

File Metadata

Filename, access control information, size, location of where the file is on disk

File Data

The actual data of the file. Will be larger. More time is spent here.

## HDFS − Architecture Motivation

**Solution**

```
┌──────────────┐      ┌──────────────┐      ┌──────────────┐
│              │      │  Metadata <  │      │   Keep on    │
│ File Metadata│ ───▶ │     data     │ ───▶ │  separate    │
│              │      │ Accessed less│      │  server  -   │
│              │      │  Frequently  │      │  NAMENODE    │
└──────────────┘      └──────────────┘      └──────────────┘


┌──────────────┐      ┌──────────────┐      ┌──────────────┐
│              │      │              │      │  Distribute  │
│              │      │ Much larger, │      │   across     │
│  File Data   │ ───▶ │requires paral│ ───▶ │  machines -  │
│              │      │ lel access   │      │  DATANODES   │
│              │      │              │      │              │
└──────────────┘      └──────────────┘      └──────────────┘
```

# BIG DATA

## HDFS – Master Slave Architecture
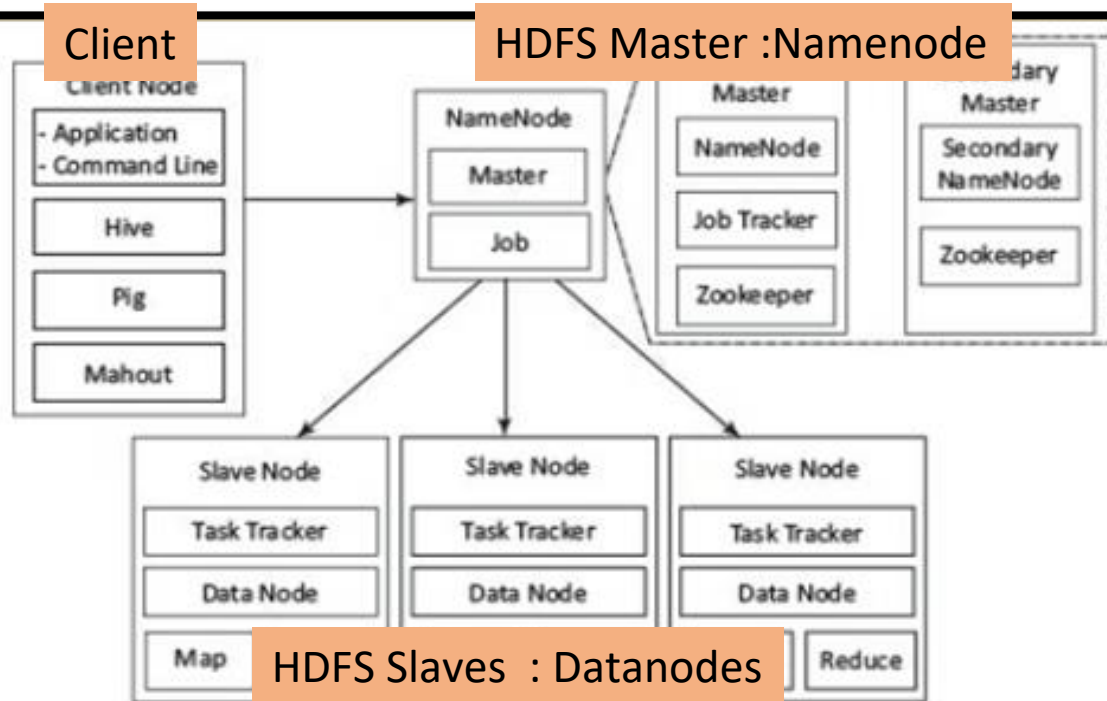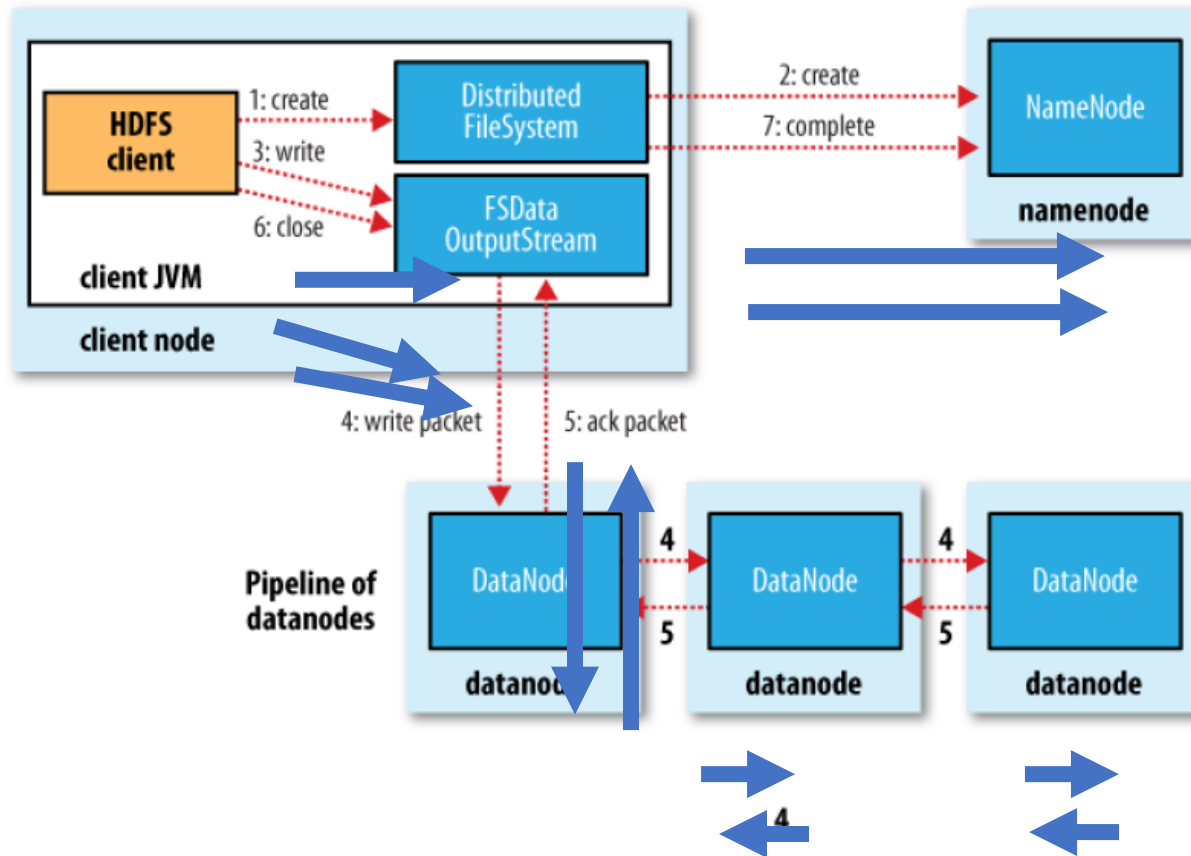


Client

HDFS Master :Namenode

HDFS Slaves : Datanodes

**Figure 2.4** The client, master NameNode, MasterNodes and slave nodes

# Big Data: HDFS Operations
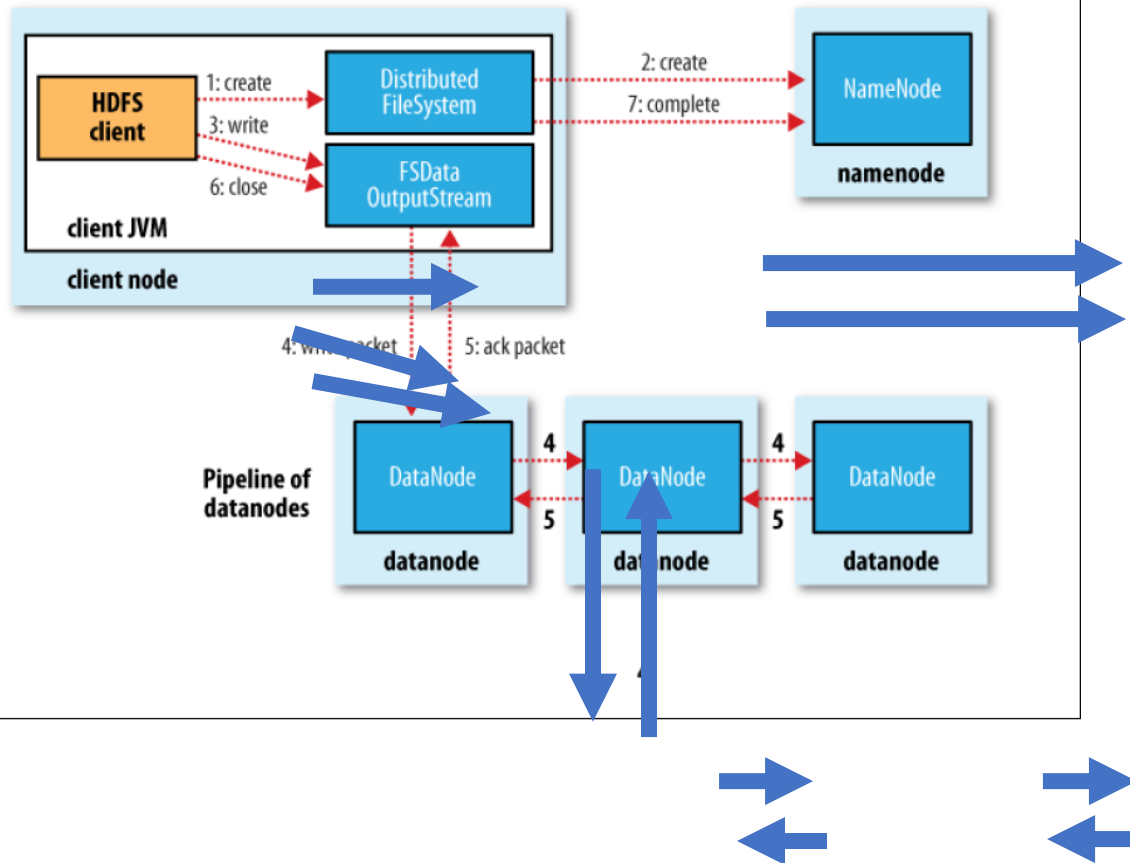
## HDFS – Writing a file
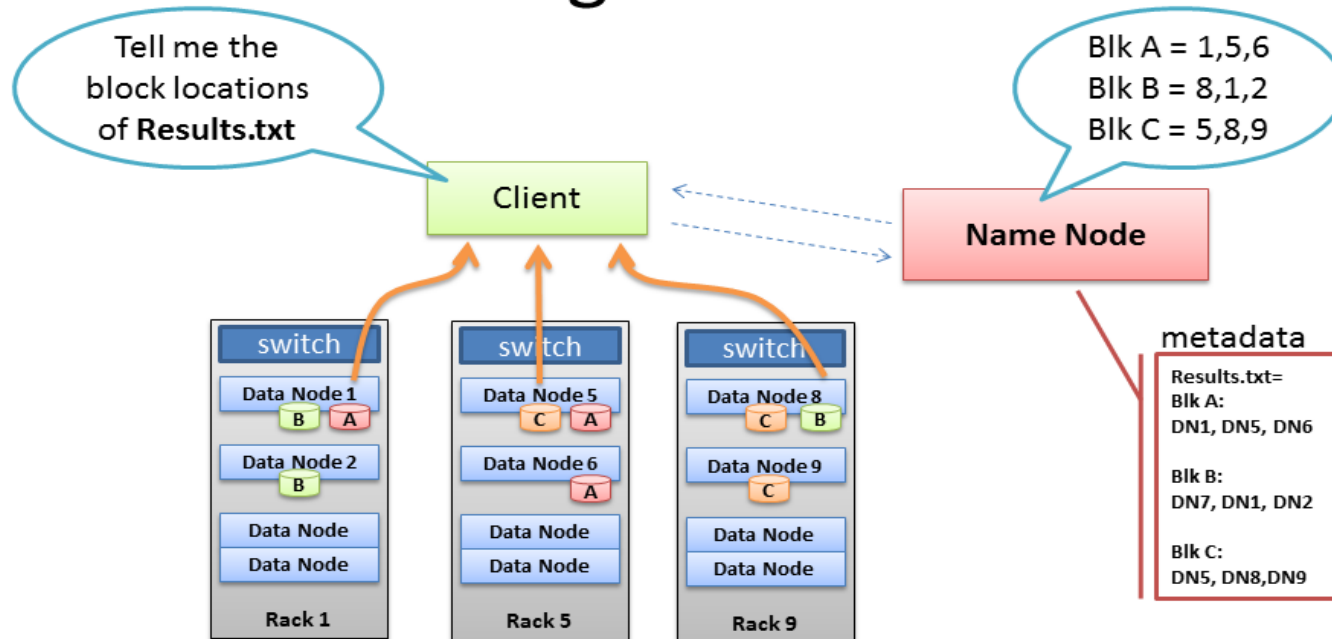


Image source: Hadoop Definitive Guide, Tom White, O'Reilly

# HDFS – Writing a file



- The write operation is outlined in the figure

- Can you outline the steps in the read operation?

Image source: Hadoop Definitive Guide, Tom White, O'Reilly

# HDFS: File reading

# Big Data: HDFS Internals

# HDFS − Hadoop distributed File system

## HDFS Architecture

## Namenode features



Primary namenode

Master Namenode

Write latest changes into logs

Get/Read initial information from fsimage and store in main memory

fsimage

edit logs

Disk Storage

- Manages file system namespace
- Regulate access to files by client
- Opening closing and renaming files
- Manages block creation, deletion and replication
- Determines mapping of blocks to data nodes
- Handles block failure
- Transaction Log
- Contains the metadata in memory

Image source: https://hrouhani.org/hdfs-file-system/

## Namenode organization

| FSImage | Edit Log |
|---|---|
| • Serialized form of filesystem tree | • Every write first writes to edit log |
| • Not updated on every write | • Flushed and synced after every transaction |
|     • To avoid recopy of data | • Append only operation |
| • Stores | |
|     • Filename | |
|     • access time | |
|     • #blocks | |
|     • blocks consisting the file | |

Since no modify operation is done on either file, it can be done really fast.

## Namenode memory requirements

- Stored in memory
- What takes up space?
  - #blocks/file
  - Filename length
  - #directories
- Limited amount of memory
- Rule of thumb – 1000MB per million blocks
- Solution:
  - Limit the responsibility of each node

- Example calculation
  - 200 node cluster
  - 24TB/node
  - 128MB block size
  - Replication factor = 3
- How much space is required?
  - #blocks = 200*24*2^20/(128*3)
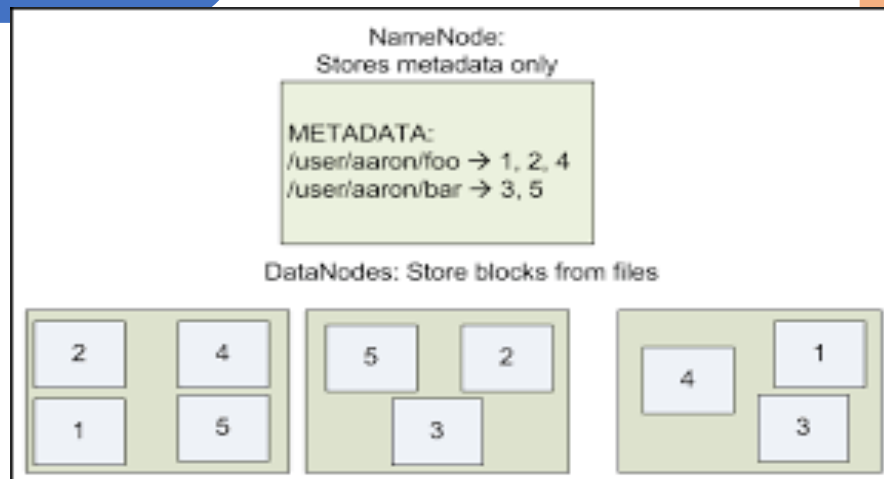  - ~12Million blocks
  - ~12,000 MB memory.

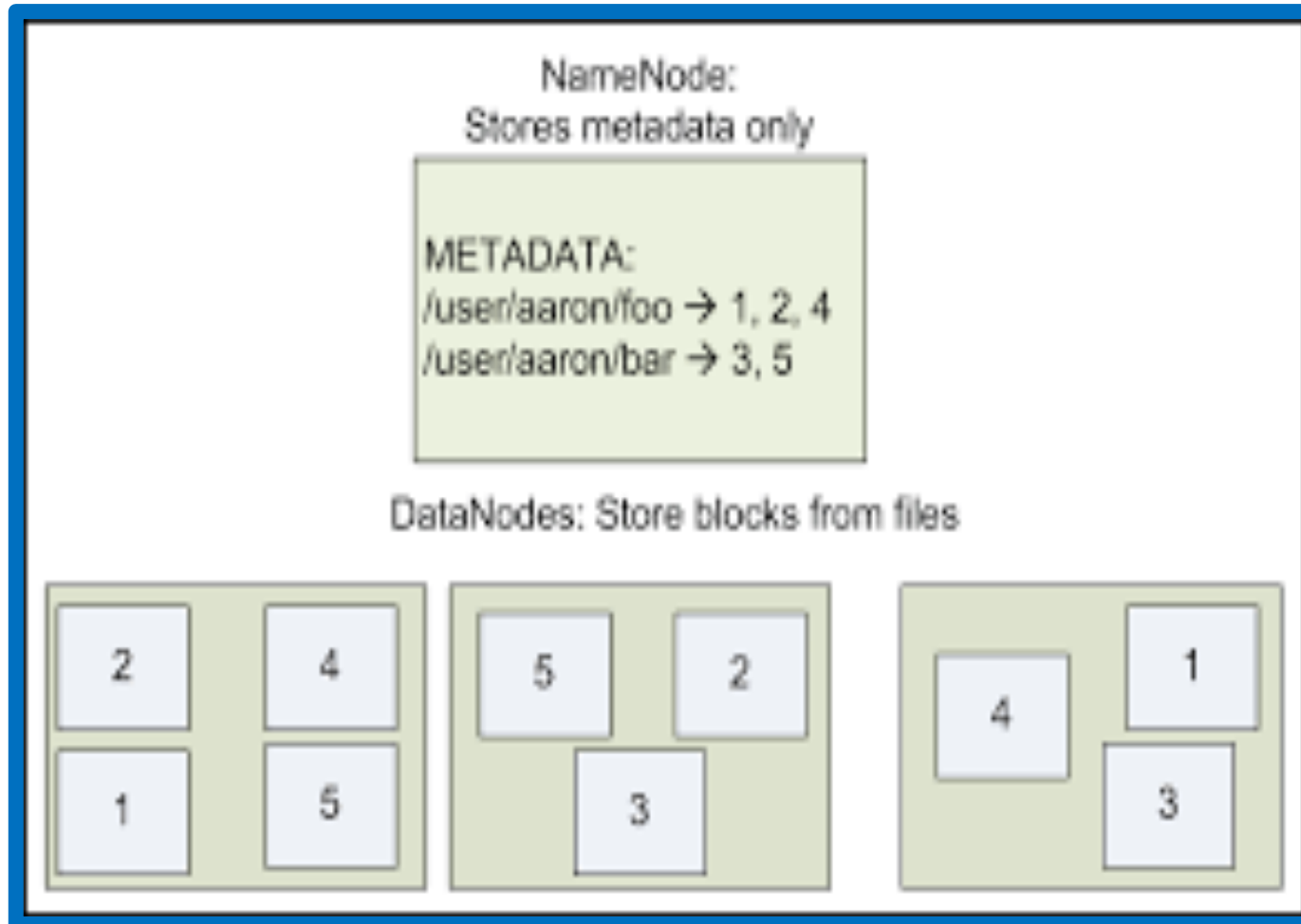Have lots of disk storage and moderate amounts of processing capabilities and DRAM

Responsible to store the data and process the computation tasks

Periodically sends a report to the name node.



NameNode:
Stores metadata only

METADATA:
/user/aaron/foo → 1, 2, 4
/user/aaron/bar → 3, 5

DataNodes: Store blocks from files

| 2 | 4 |
| 1 | 5 |

| 5 | 2 |
| | 3 |

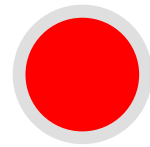| | 1 |
| 4 | |
| | 3 |

## HDFS − Hadoop distributed File system

### Data Blocks

Each file split into datablocks – 128MB for HDFS v2, 64MB for HDFS v1

Each block is stored on one or more nodes
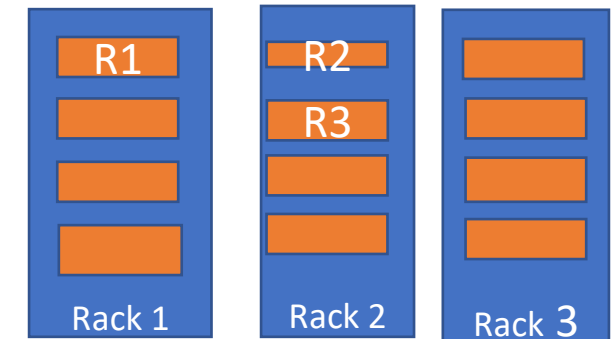
Each copy of the block is called replica

### Block placement policy

First replica is placed on the local node

Second replica is placed in a different rack

Third replica is placed in the same rack as the second replica
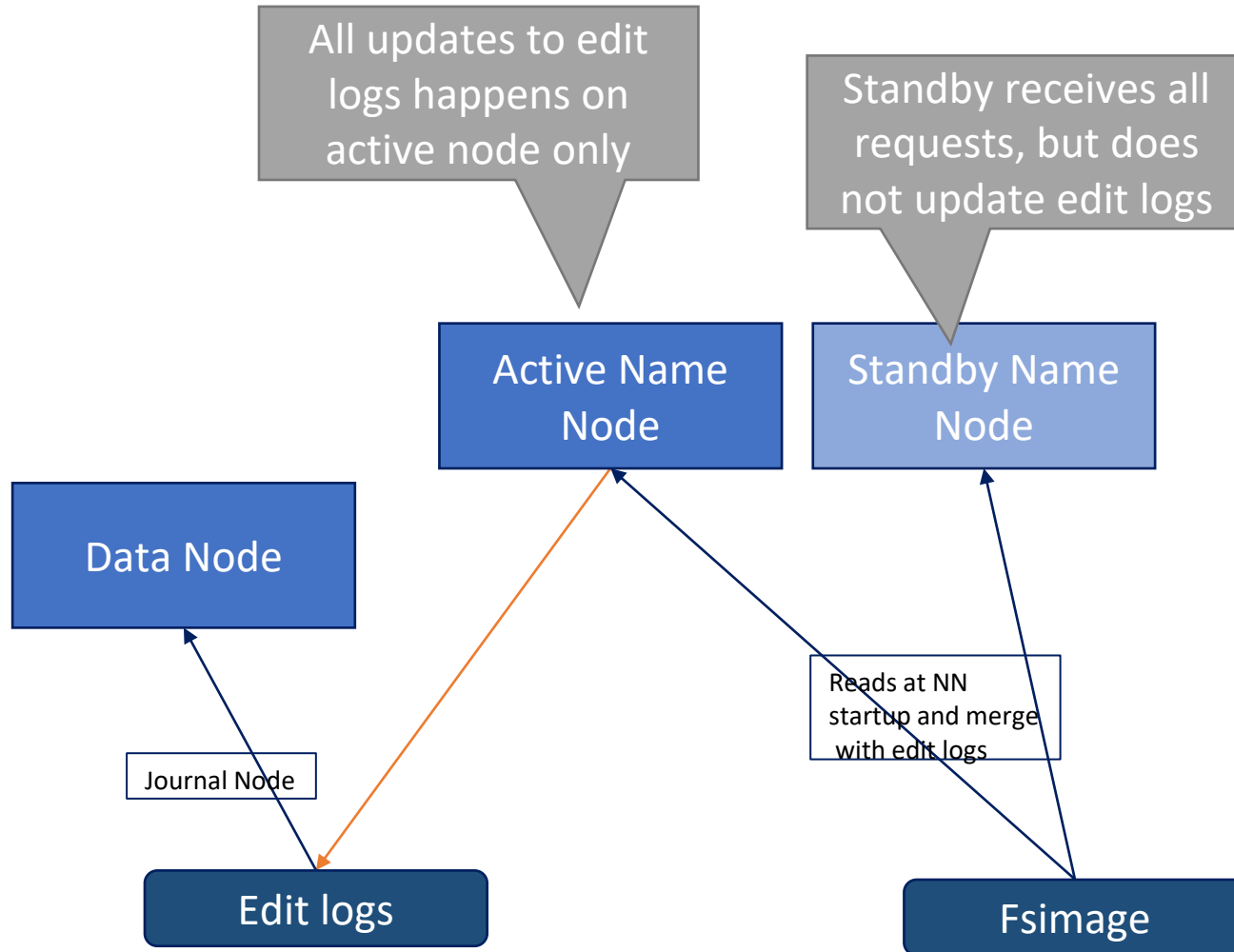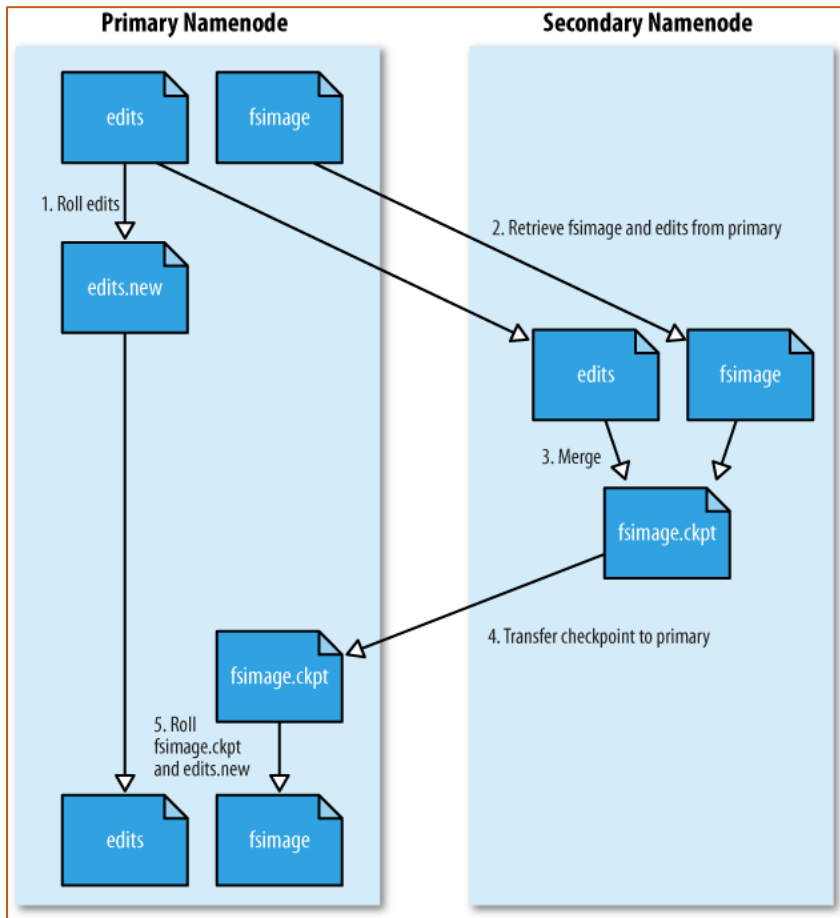
**Block placement policy**



| Rack 1 | Rack 2 | Rack 3 |
| --- | --- | --- |
| R1 | R2 | |
| | R3 | |

# Big Data: HDFS Fault Tolerance and Replication

Used to combine fsimage with edits

Offload the operation to secondary

Primary can focus on serving request

On completion of merging, new fsimage will be used.

## Secondary Namenodes

**Secondary node**
- keeps a copy of NameNode meta data.

**Stored meta data**
- merges the metadata files to obtain an updated Fsimage

**Checkpointing**
- Primary and secondary sync up the data

**Why?**
- To free namenode from task of merging

# Big Data: HDFS Blocks

**Benefits of block abstraction.**

**% hadoop fsck -files –blocks**

A file can be larger than any single disk in the network.

will list the blocks that make up each file in the filesystem

Simplifies the storage subsystem

Blocks fit well with replication for providing fault tolerance and availability.

# HDFS block size – why 128MB?

- Time to read data from disk
  - =  seek time + rotational latency + transfer time
  - Seek time – time to position head onto track (variable) - mechanical
  - Rotational latency – time to spin disk to get sector under head (fixed)
  - Transfer time – time to read data (fixed)

- Improving performance of read      reduce seek times.

- Large block size      transfer time >> seek time

**Example**
**If the seek time is around 10 ms, and the transfer rate is 100 MB/s, then to make the seek time 1% of the transfer time, we need to make the block size around 100 MB.**

   **Hadoop v1 default – 64MB**
   **Hadoop v2 default – 128MB**
   **Review: why has this increased?**

# THANK YOU

**K V Subramaniam**

Department of Computer Science and Engineering

**subramaniamkv@pes.edu**

+91 80 6666 3333 Extn 877

Content Creators
Prof Resma K S
Prof Prafullata K A
Prof Usha Devi
Prof Rachana