# 4 – Linux File Systems

**Marian Marinov**
**CEO of 1H Ltd.**
**mm@1h.com**

**Stoyan Stoyanov**
**System Administrator**
**sto [ at ] softuni.bg**

# Agenda

- ➢ **File System Architecture**
- ➢ **Virtual File System Layer**
- ➢ **Directory Structure**
- ➢ **Mount operations**

# Agenda

- **File system types**
  - **Local**
  - **Log-structured - NAND**
  - **Pseudo**
  - **Network**
  - **Cluster**
  - **Distributed**

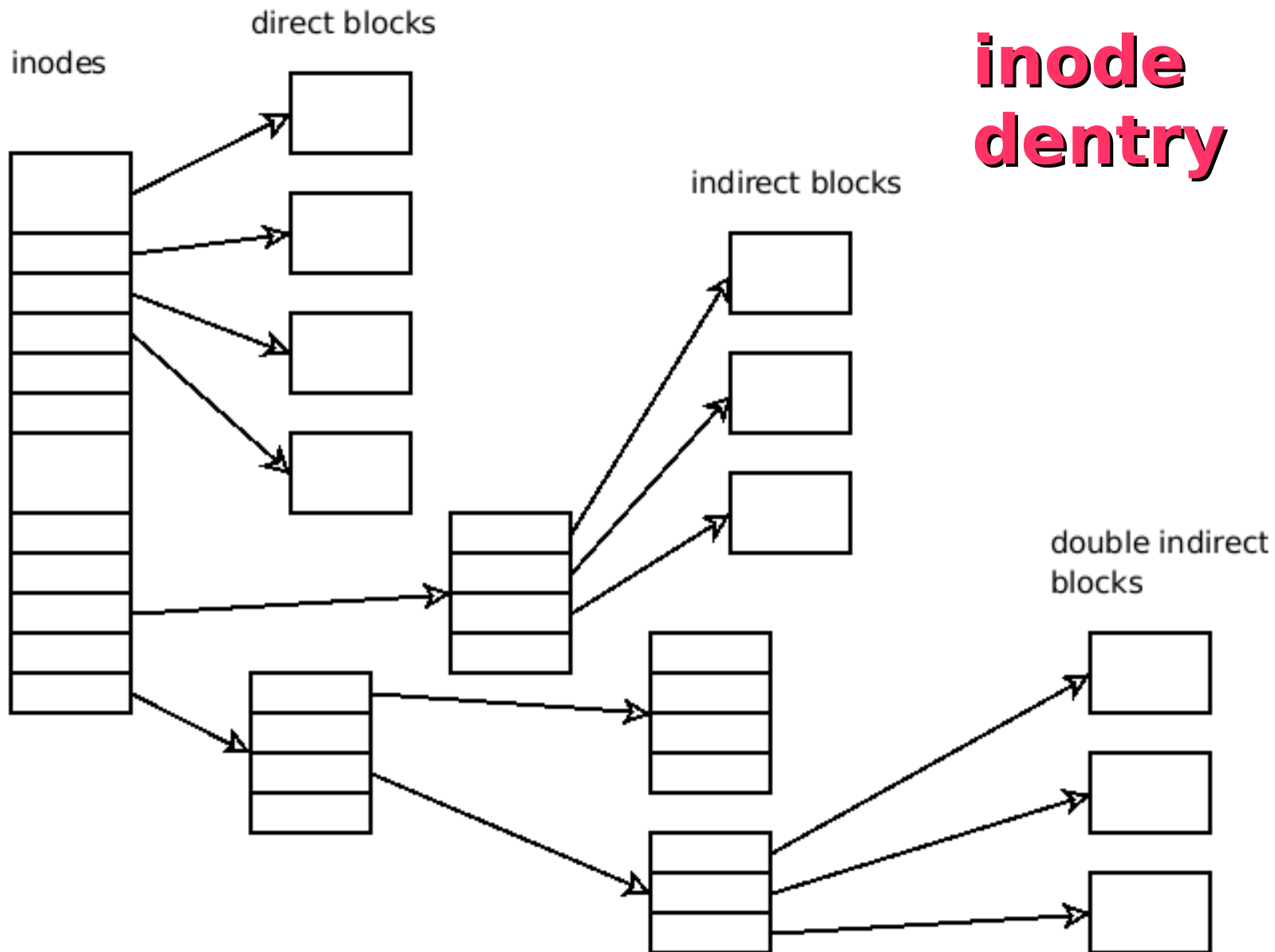# File System Architecture



**There is NO silver bullet!**

# File System Architecture
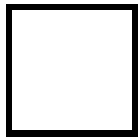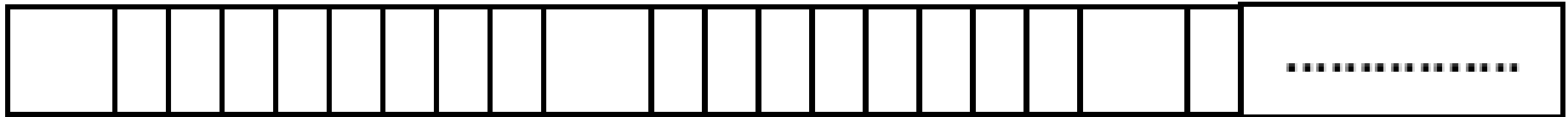


**There is NO DIRECTORY!**

# File System Architecture

## Super Blocks

inodes



super block

data block

# File System Architecture

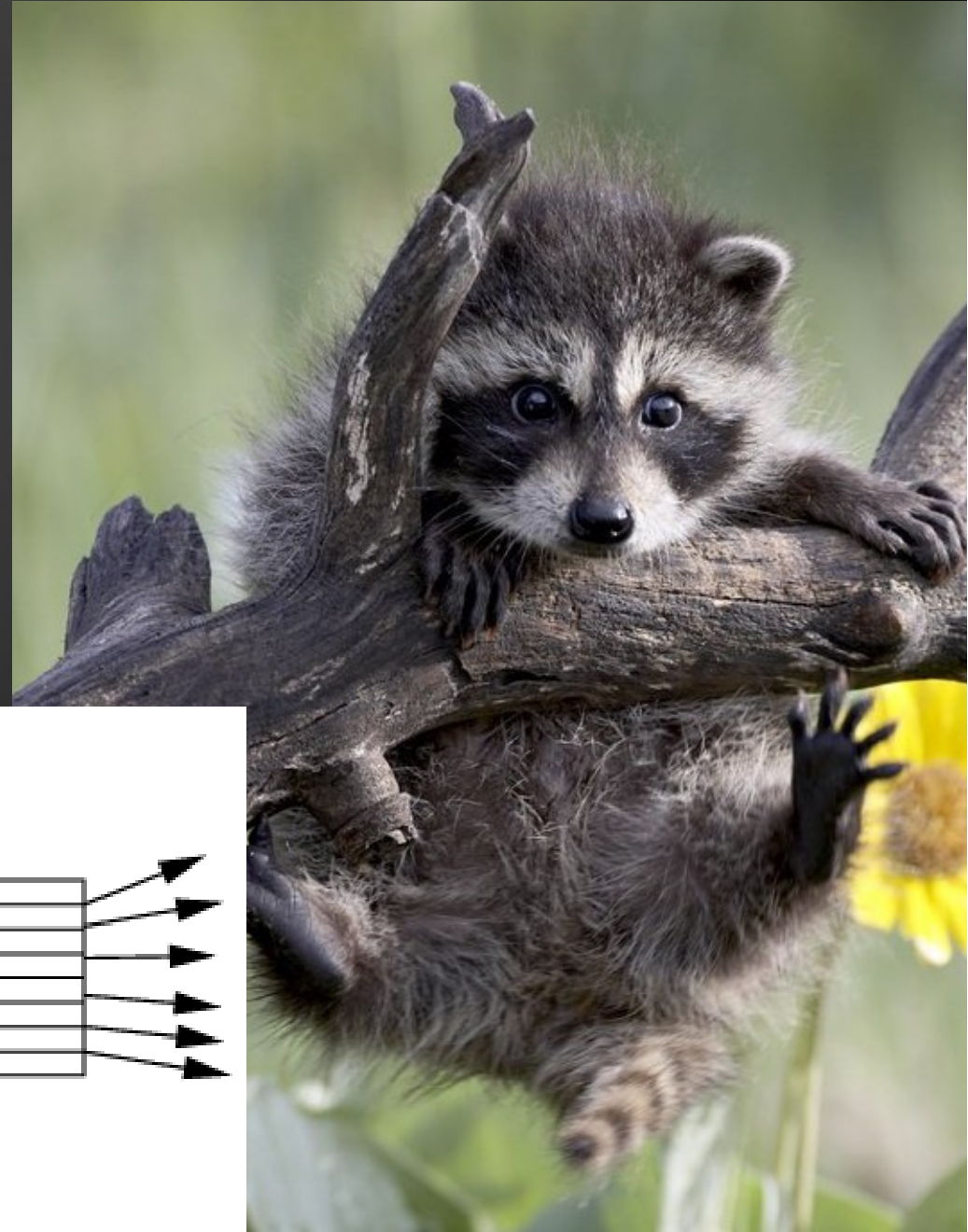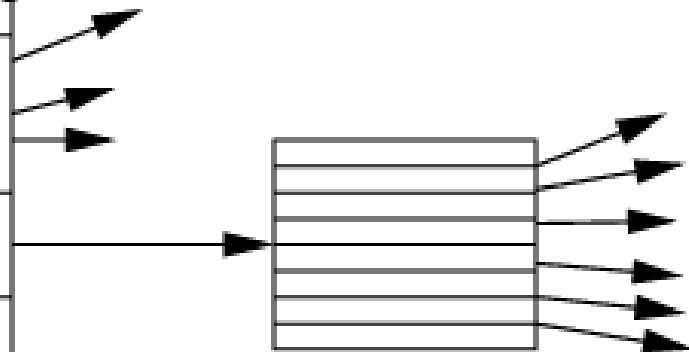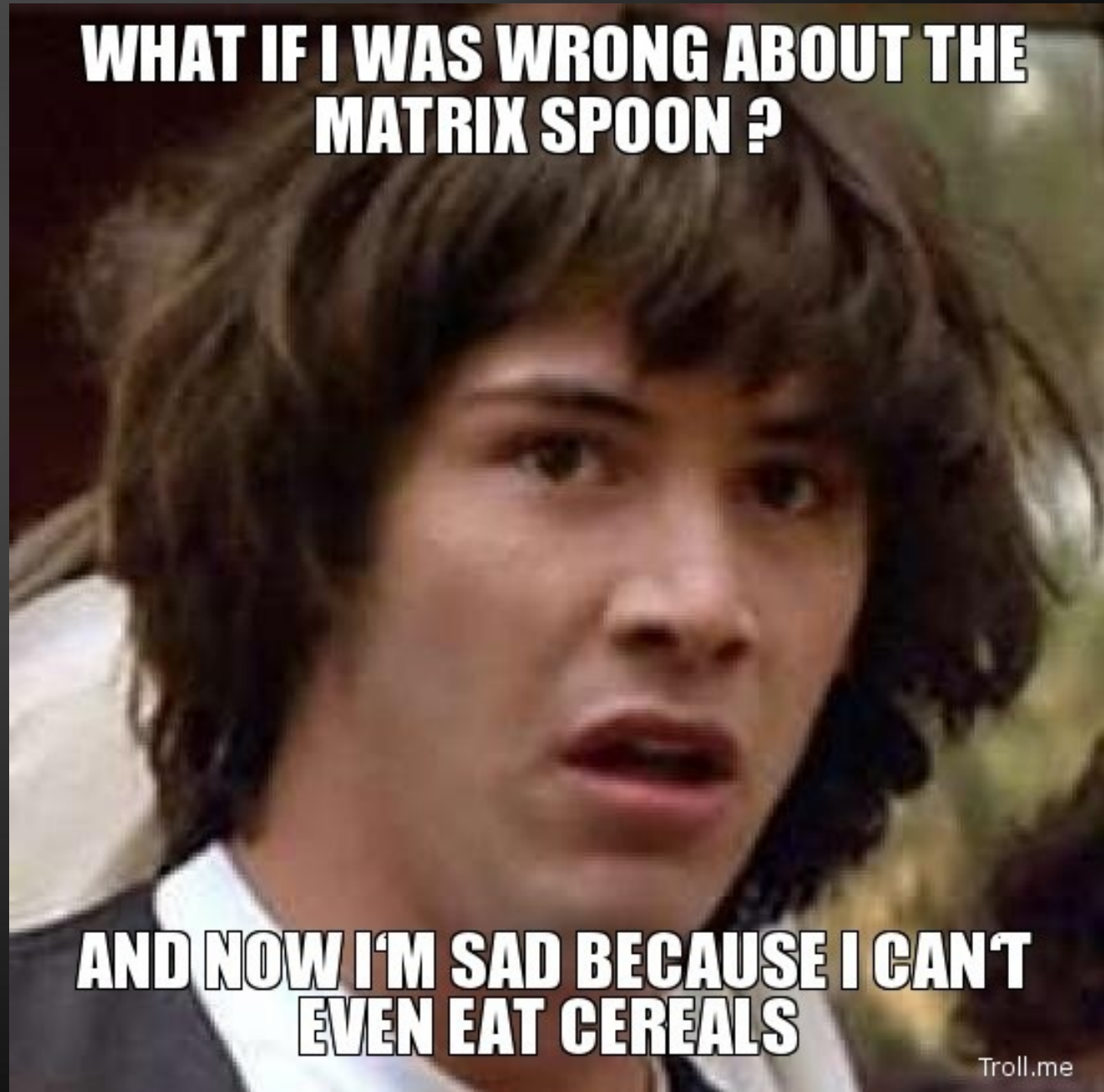| |
|---|
| File mode |
| Link count |
| Owner's id |
| Group id |
| File size |
| Last access time |
| Last mod time |
| Last inode access time |
| Addresses of first 10 blocks |
| Single indirect ptr |
| Double indirect ptr |
| Triple indirect ptr |

# File System Architecture
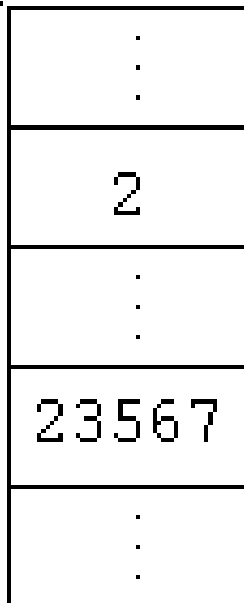
**Soft/symlinks
Hardlinks**

# File System Architecture

directory entry in /dirA

inode     name

| 12345 | name1 |

directory entry in /dirB

inode     name
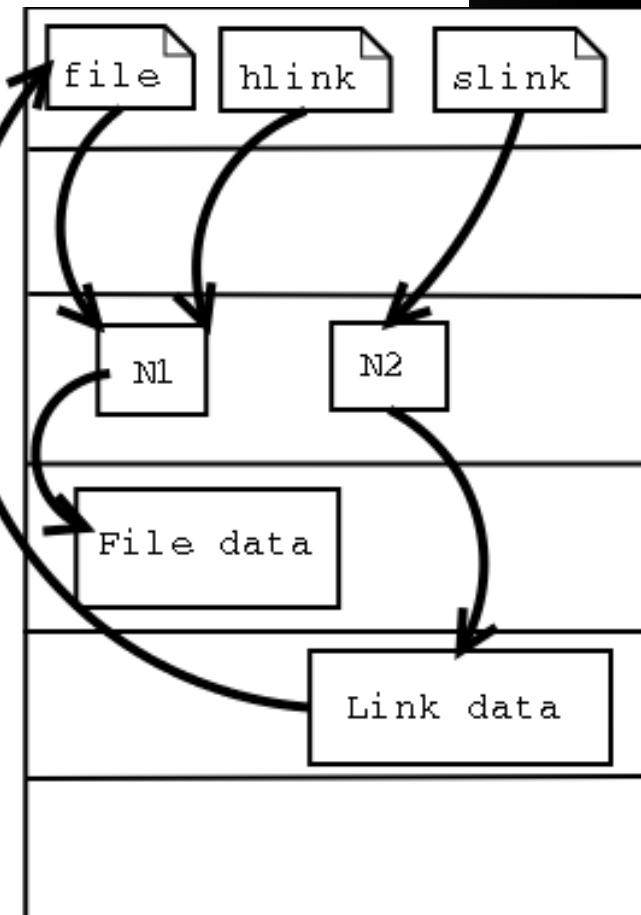
| 12345 | name2 |

inode 12345

| ⋮ |
| 2 |
| ⋮ |
| 23567 |
| ⋮ |

block 23567

"This is the text in the file."

file    hlink    slink

N1     N2

File data

Link data

Hard disk

# File System Architecture

In source destination          In -s source destination

**DATA is never erased... it gets OVERWRITTEN!**

# Local File System

# Network File System

# Cluster File System

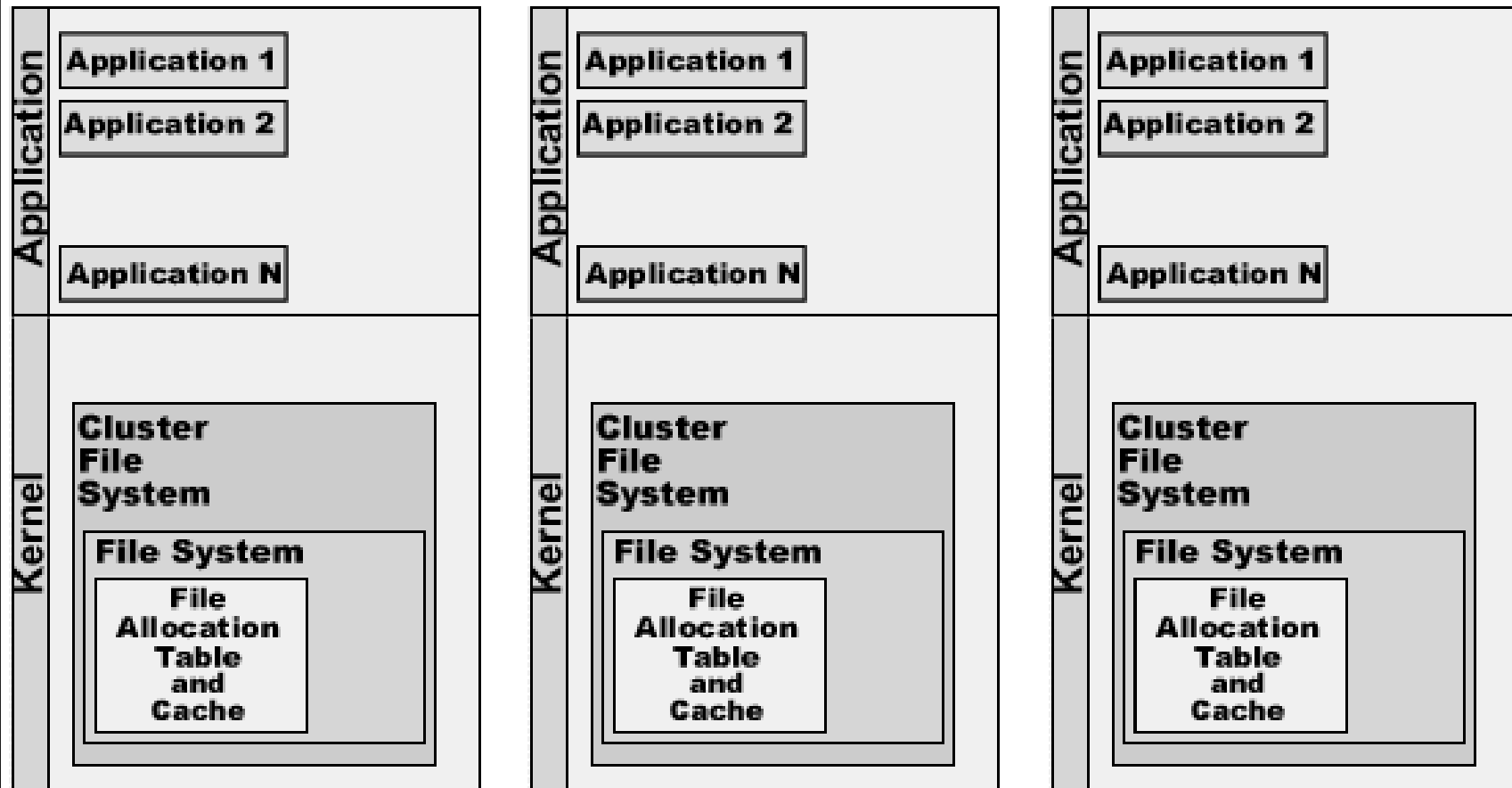# Cluster File System



**"One DLM to rule them all, one DLM to mind them, one DLM to sync them all, and in the cluster, bind them"**

# Distributed File System

# Virtual File System Layer

# Virtual File System Layer



**Introduced April 1992**

# VFS Caches

# FUSE

# Mounting

➢ **Attaching a device into the directory tree**

➢ **Mount point – a destination directory where a device is mounted**

➢ **Creates an entry in the kernel for each mounted device/dir**

➢ **/proc/mounts /etc/fstab, /etc/mtab**

# Mounting - CMD

➢ **cat /proc/partitions**

➢ **cat /proc/mounts**

➢ **mount**

➢ **umount**

➢ **/etc/fstab**

```
/dev/sdb2       /         ext4   defaults,noatime,nodiratime    0 0
/dev/sdb1       /boot   ext2   defaults,noatime,nodiratime    0 0
proc            /proc   proc   defaults                       0 0
tmpfs           /dev/shm  tmpfs  defaults                     0 0
/home/hackman  /fedora/home/hackman   none   rw,bind,auto  0 0
//10.2.0.11/share       /storage/beast          cifs
user=hackman,password=p1r@tk3,uid=1000,gid=1000,noauto 0 0
```
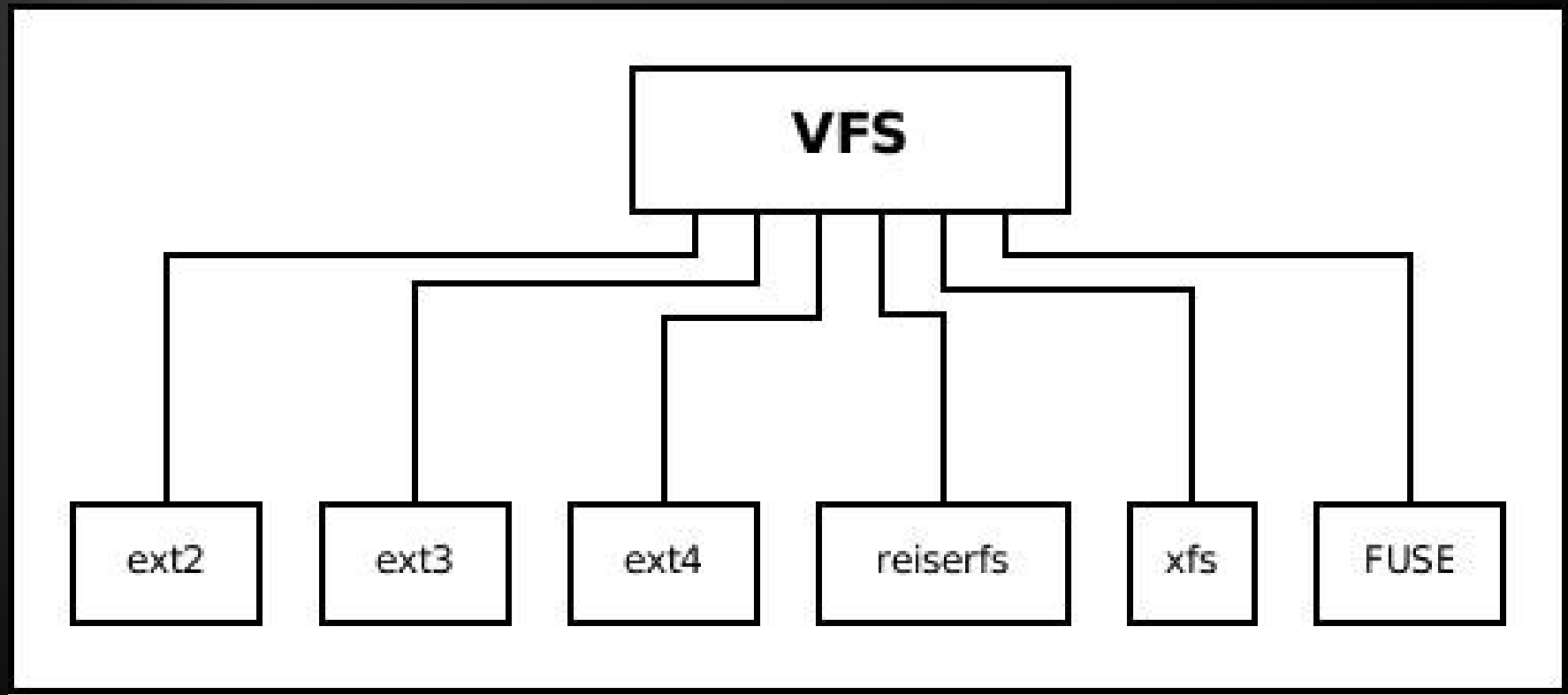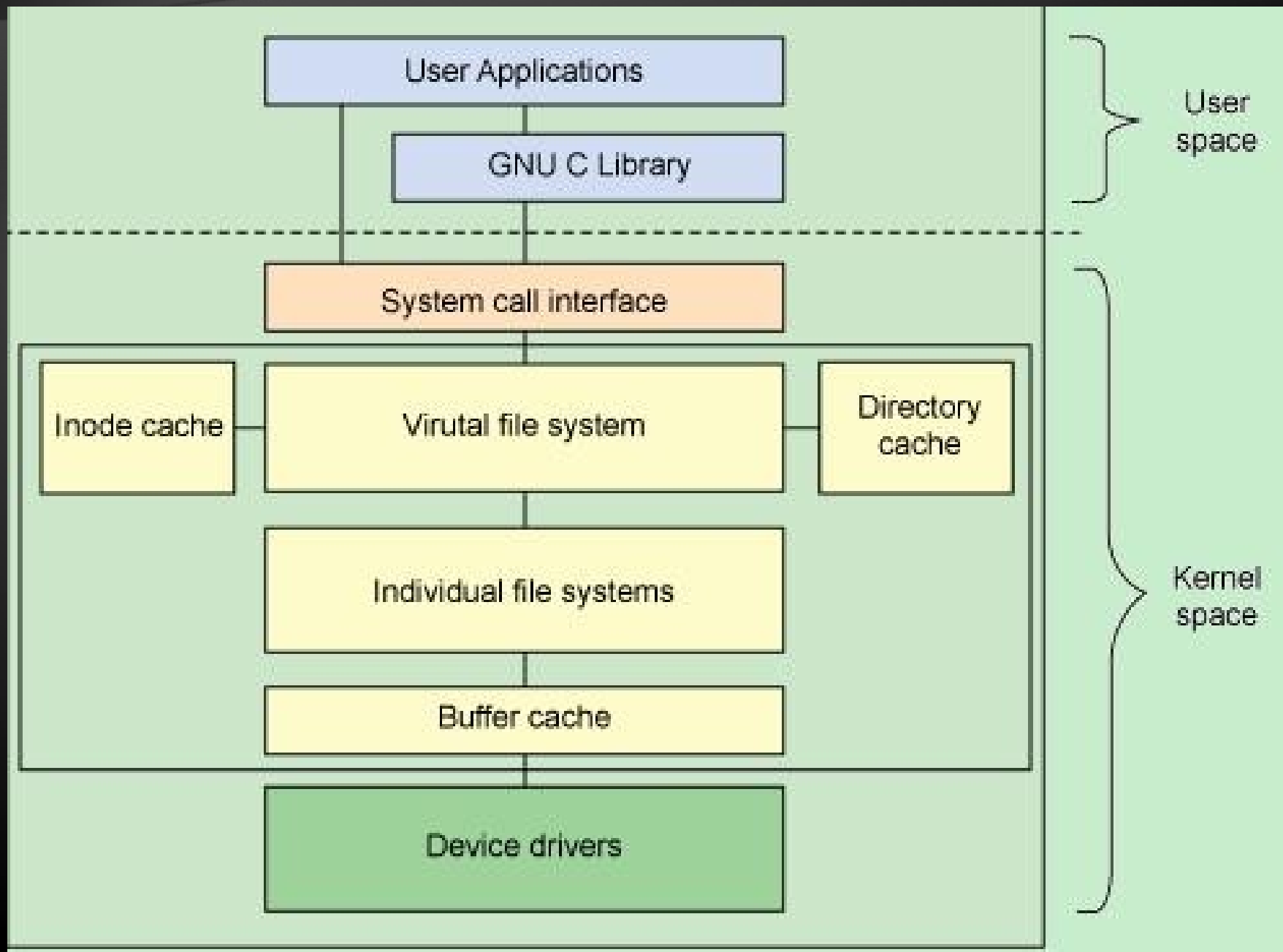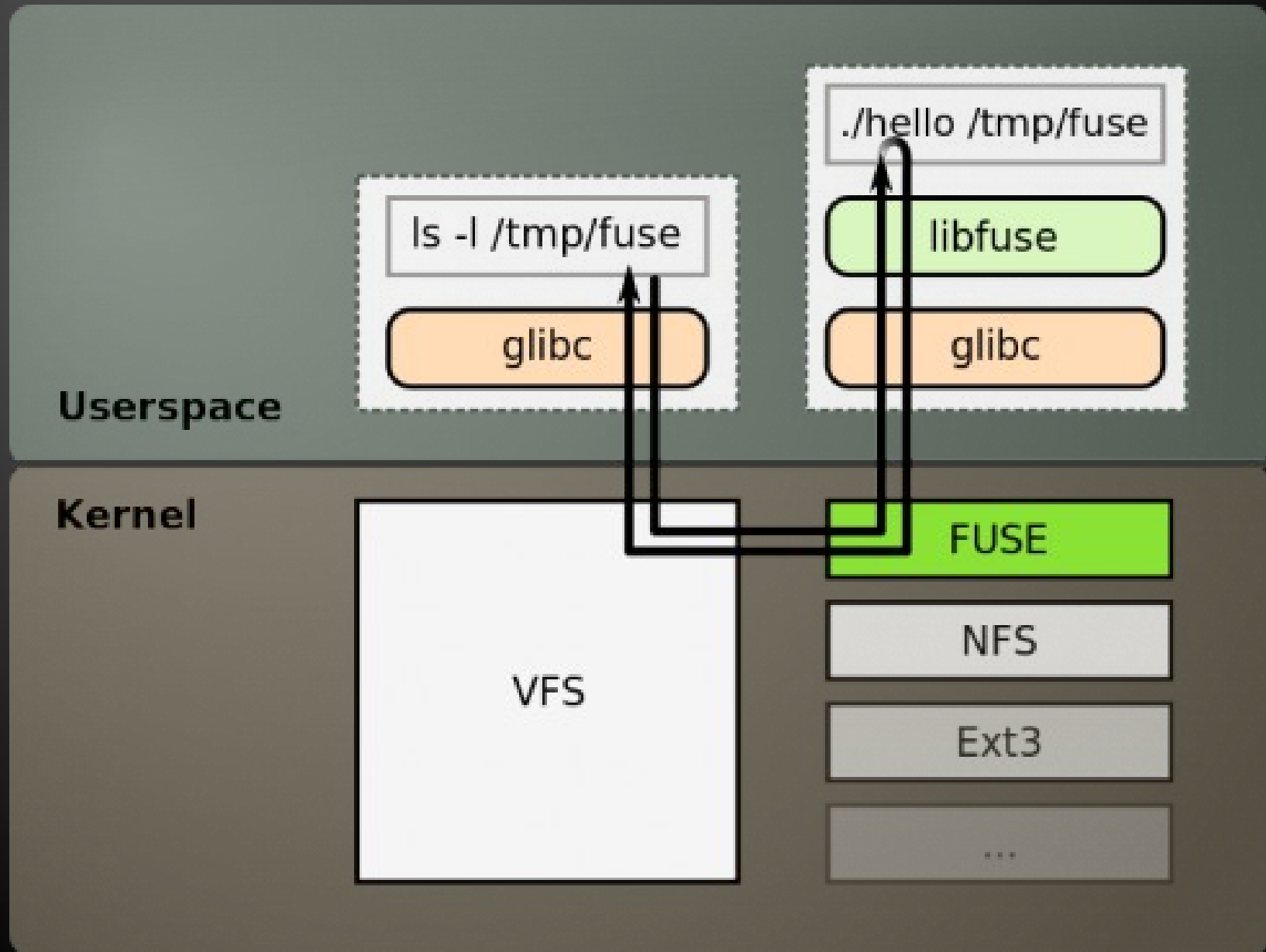
# Ext/2/3/4

- **First Linux FS – MinixFS**
- **And there we go….**
  - **Ext - April 1992, Linux 0.96c**
  - **Ext2 – January 1993**
  - **Ext3 – November 2001**
  - **Ext4 – October 2006**

# MinixFS

- ➢ **Max. partition size – 64MB**
- ➢ **Max. file name size – 14 chars**
- ➢ **Ownership – uid, gid**
- ➢ **Permissions – user, group, others**

# Ext

- Max. partition size – 2GB
- Max. file name size – 255 chars
- No support for time stamps
  - Access
  - Inode modification
  - Data modification

# Ext2

➢ **Max. partition size – 32TB**

➢ **Max. file name size – 255 chars**

➢ **Max. file size – 2TB**

➢ **Max. Number of files - $10^{18}$**
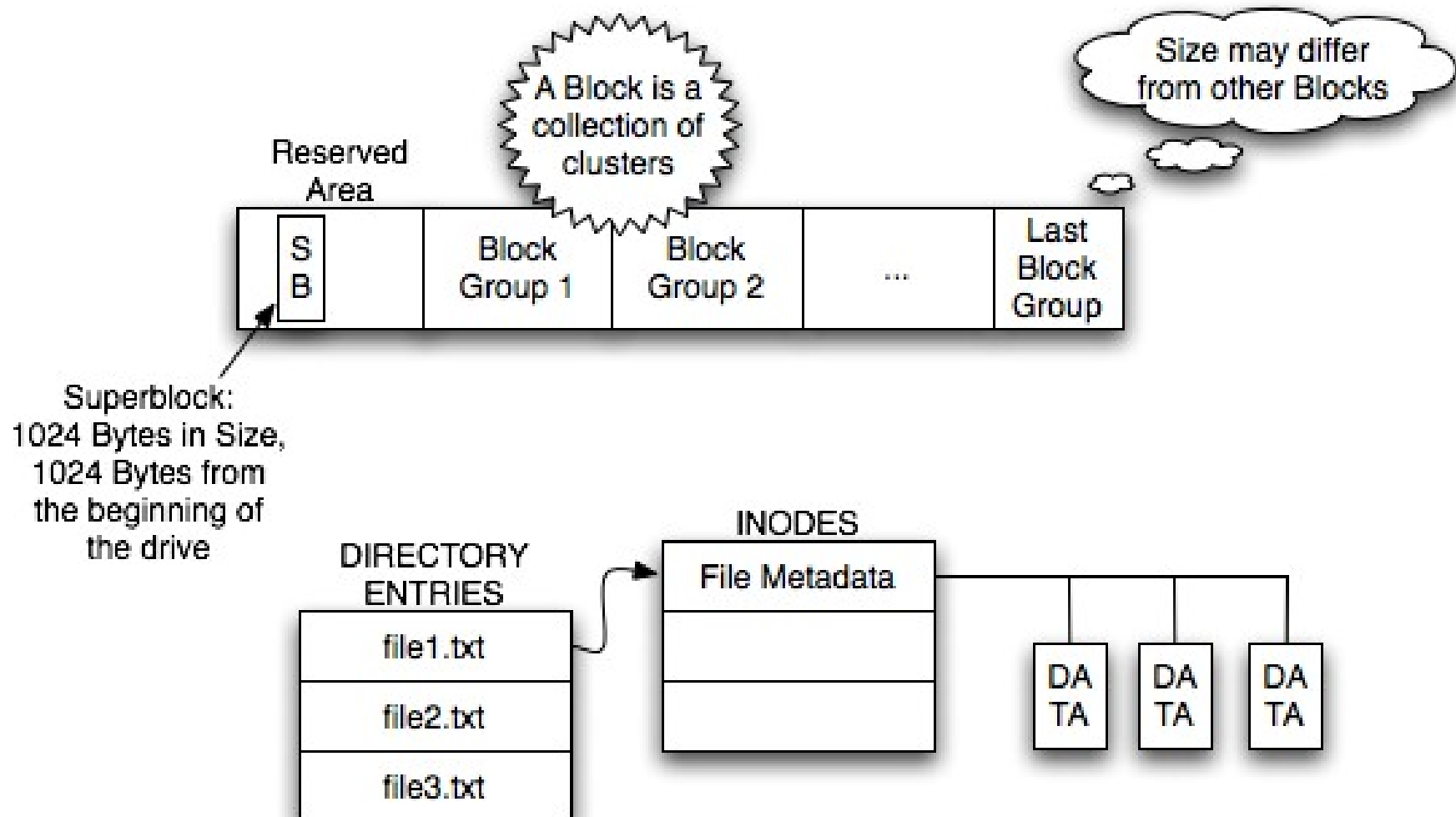
➢ **FS Perms**

➢ **Time stamps**

# Ext3

➢ **Max. partition size – 32TB**

➢ **Max. file name size – 255 chars**

➢ **Max. file size – 2TB**

➢ **Max. Number of files - $10^{18}$**

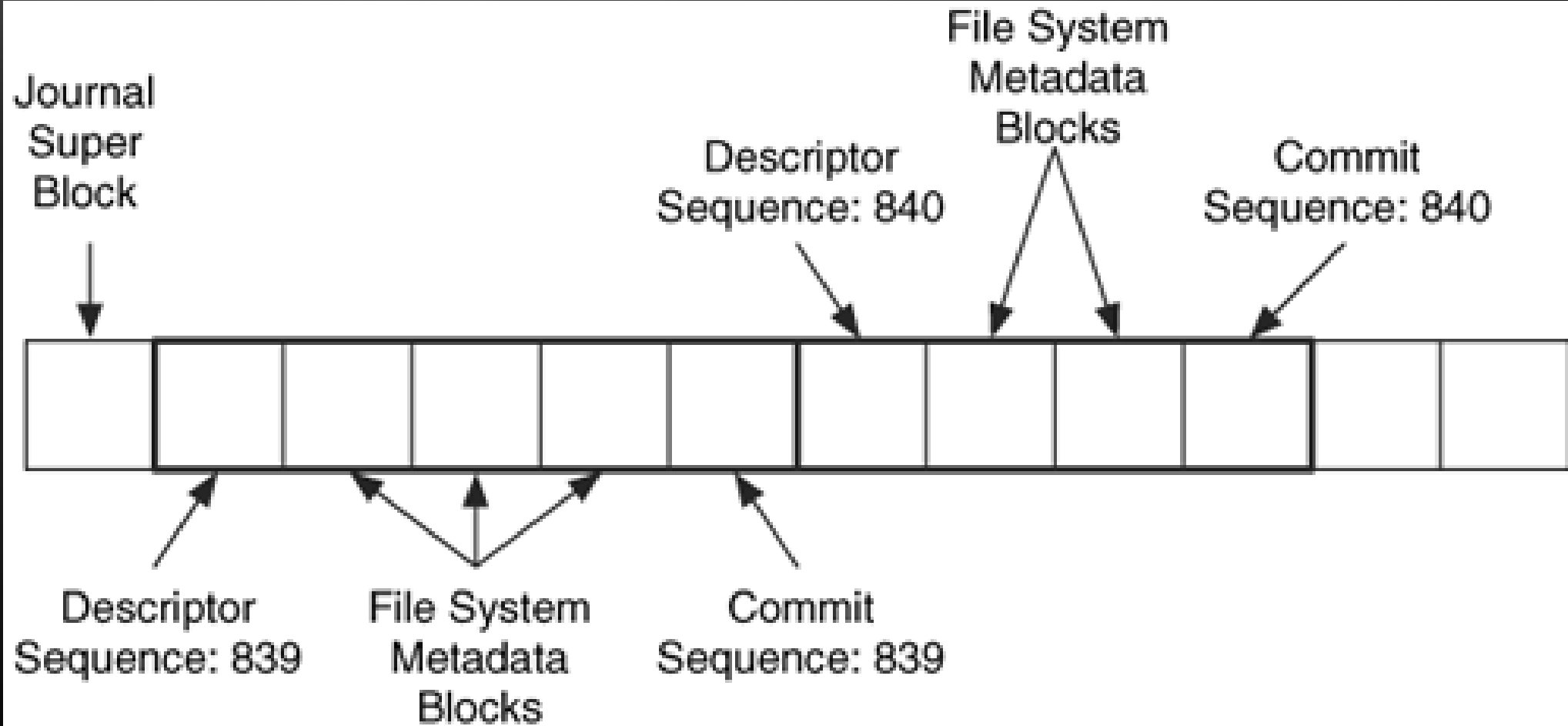➢ **Sub directory limit - 32,000**

➢ **Time stamps**

➢ **Has Jurnal**

# Ext4

- ➤ **Max. partition size – 1EB**
- ➤ **Max. file name size – 255 chars**
- ➤ **Max. file size – 2TB**
- ➤ **Max. Number of files - $10^{18}$**
- ➤ **Sub directory limit - 64,000**
- ➤ **File space pre-allocation**
- ➤ **File space delayed allocation**

# Ext/2/3

# Ext3
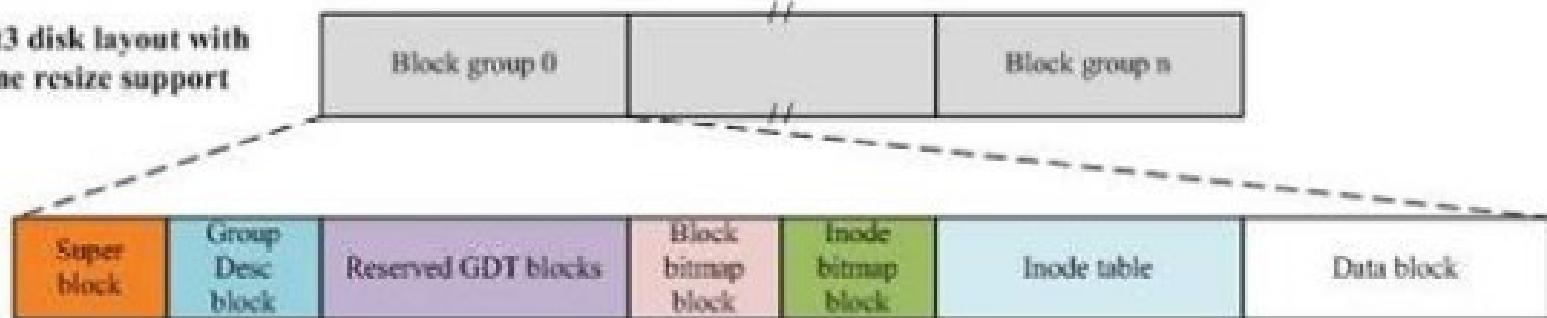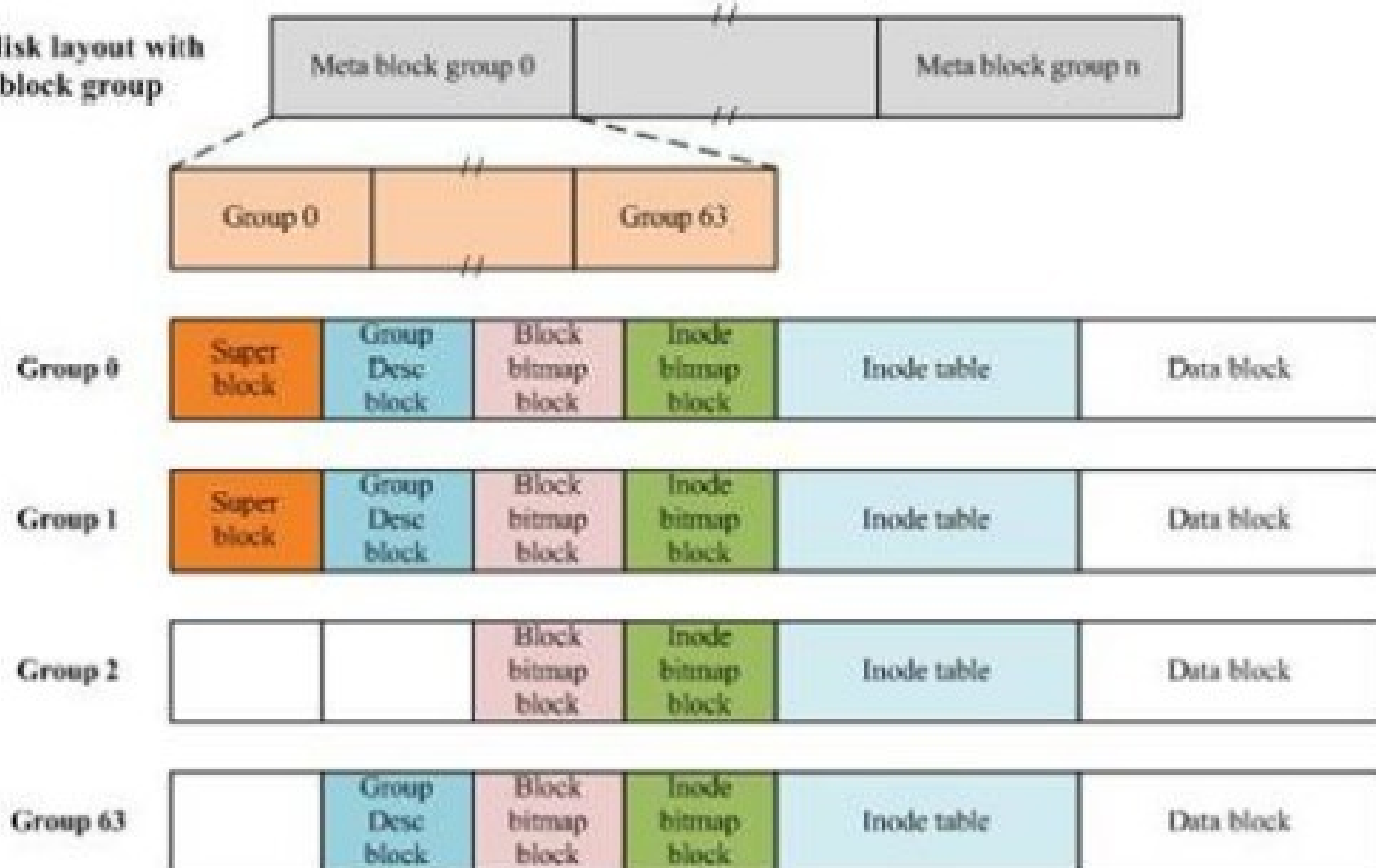
# Ext4

# ReiserFS

➢ **Introduced 2001**

➢ **Metadata-only journaling**

➢ **Online resizing (growth only)**

➢ **Tail packing, a scheme to reduce internal fragmentation.**

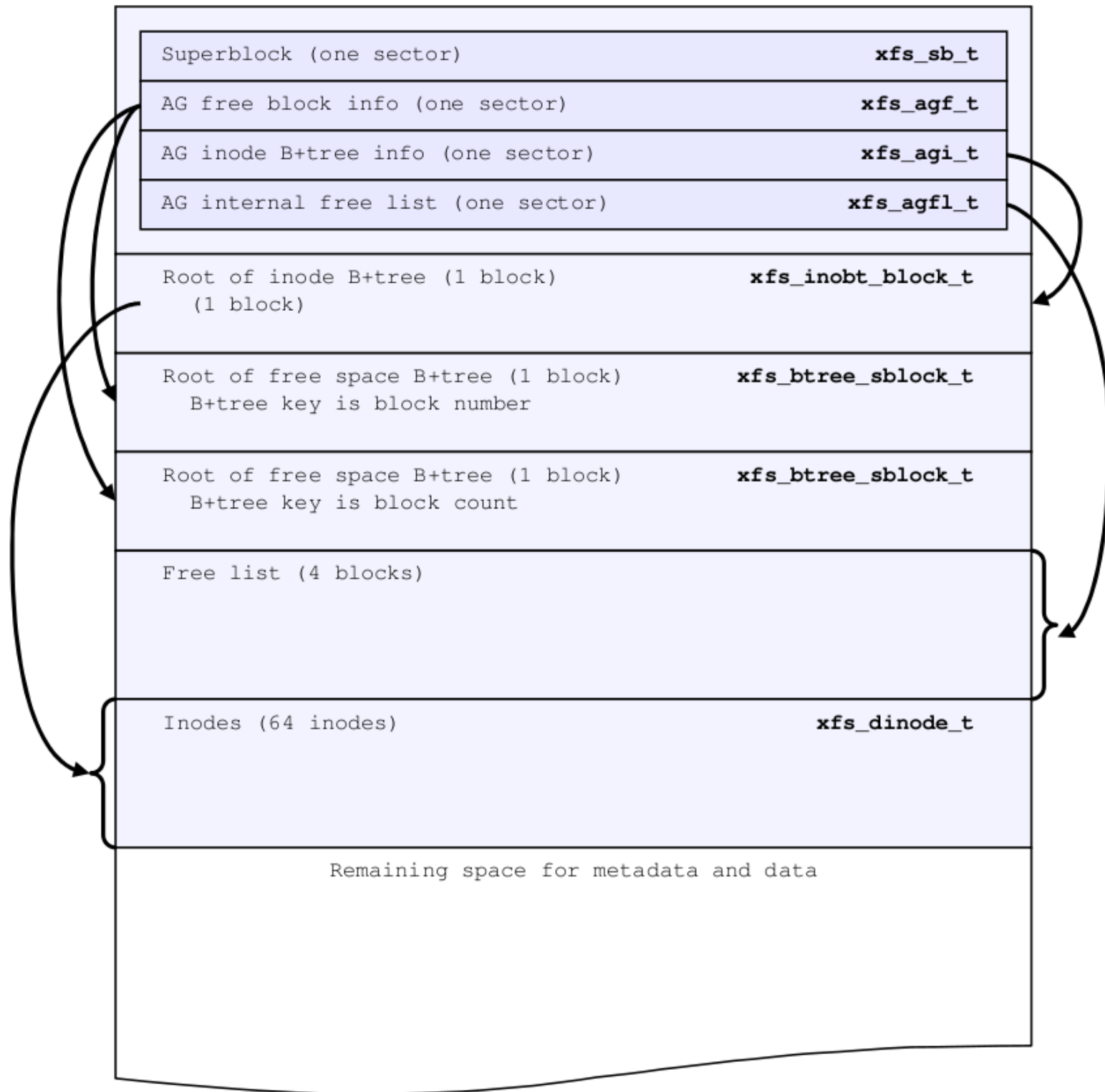➢ **Max. file size - 1EB**

➢ **Max. number of files - $2^{32}$**

# ReiserFS

# ReiserFS

# XFS

- **Introduced 2001**
- **Max file size 8 EB**
- **Max volume size 16 EB**
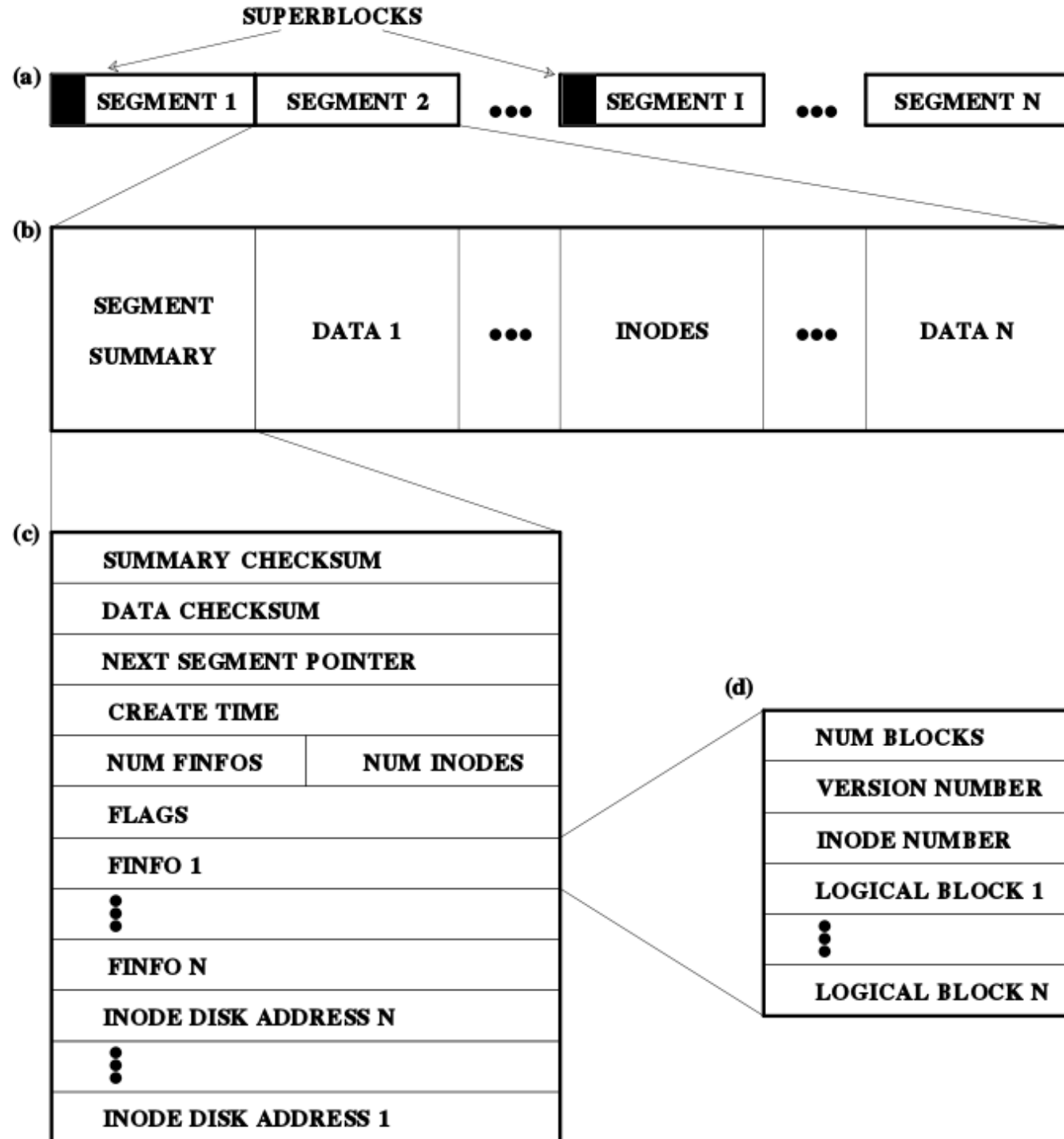- **Online resize(growth only)**
- **Online defragmentation**

# XFS

**Equally sized chunks**

**Allocation groups – AG**

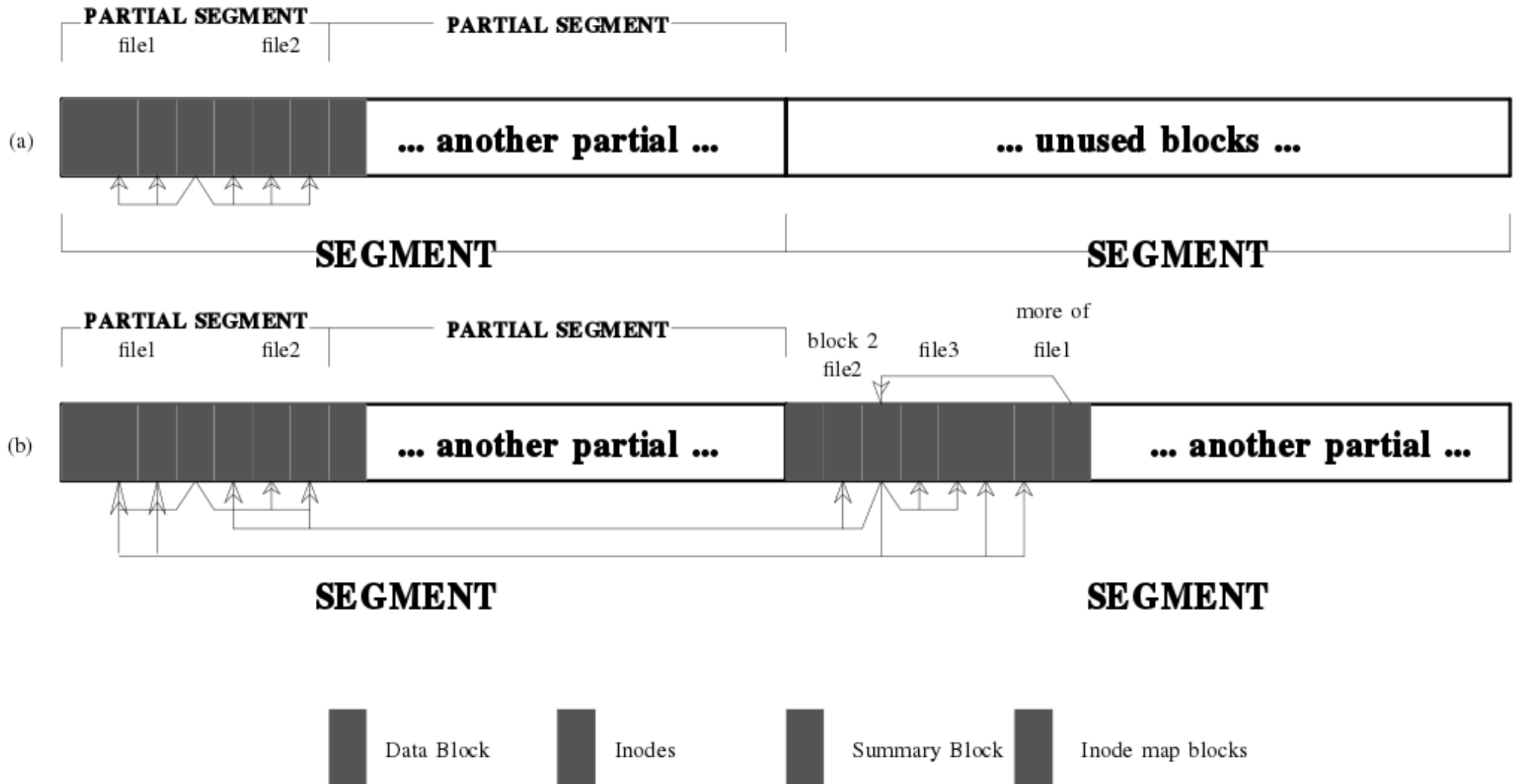| | |
|---|---|
| Superblock (one sector) | **xfs_sb_t** |
| AG free block info (one sector) | **xfs_agf_t** |
| AG inode B+tree info (one sector) | **xfs_agi_t** |
| AG internal free list (one sector) | **xfs_agfl_t** |
| Root of inode B+tree (1 block)<br>  (1 block) | **xfs_inobt_block_t** |
| Root of free space B+tree (1 block)<br>  B+tree key is block number | **xfs_btree_sblock_t** |
| Root of free space B+tree (1 block)<br>  B+tree key is block count | **xfs_btree_sblock_t** |
| Free list (4 blocks) | |
| Inodes (64 inodes) | **xfs_dinode_t** |
| Remaining space for metadata and data | |

**Log-structured File Systems Architecture**

# Log-structured File Systems Architecture

# Log-structured NAND File Systems

| | System requirement | JFFS2 | YAFFS2 | LogFS | UBIFS |
|---|---|---|---|---|---|
| 1 | Boot time | Poor | Good | Excellent | Good |
| 2 | I/O performance | Good | Good | Fair | Excellent |
| 3 | Resource usage | Fair | Excellent | Good | Fair |
| 4 | NAND device life expectancy | Good | Fair | N/A | Excellent |
| 5 | Tolerance for unexpected power-off | Good | Good | Poor | Good |
| 6 | Integrated in mainline | Yes | No | Yes | Yes |

**NILFS2**
**F2FS**

# Pseudo File Systems

➢ **procfs**

➢ **sysfs**

➢ **debugfs**

➢ **configfs**

➢ **tmpfs**

➢ **others**

# Pseudo File Systems

➢ **debugfs is designed to provide Kernel Devs with simple way to push data into User space**

➢ **configfs is for creating, managing and destroying kernel objects from user-space**

➢ **sysfs is for viewing and manipulating objects from user-space which are created and destroyed by kernel space**
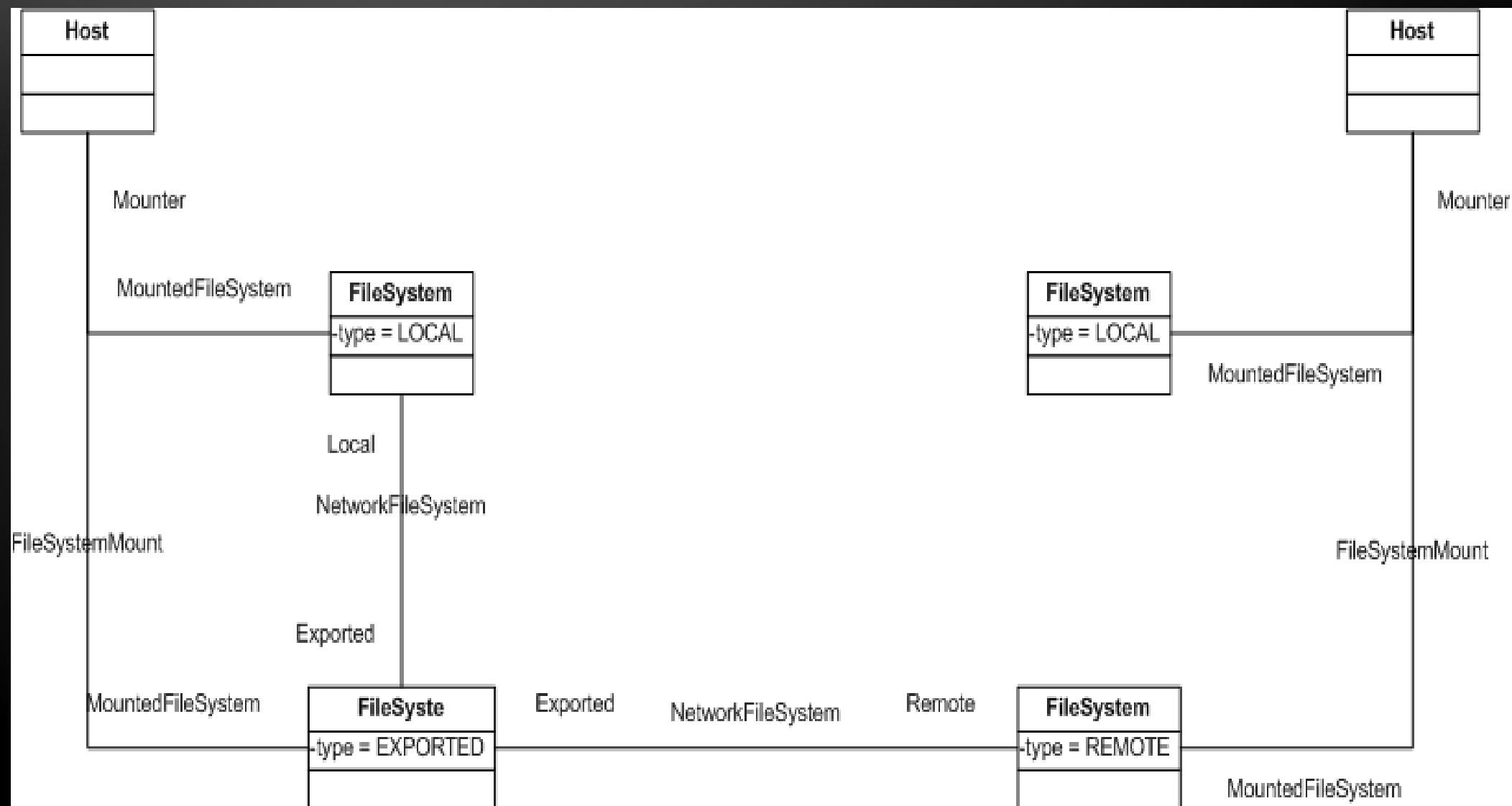
# Pseudo File Systems

➢ **procfs is the first FS to provide easy access to kernel-space from user-space**
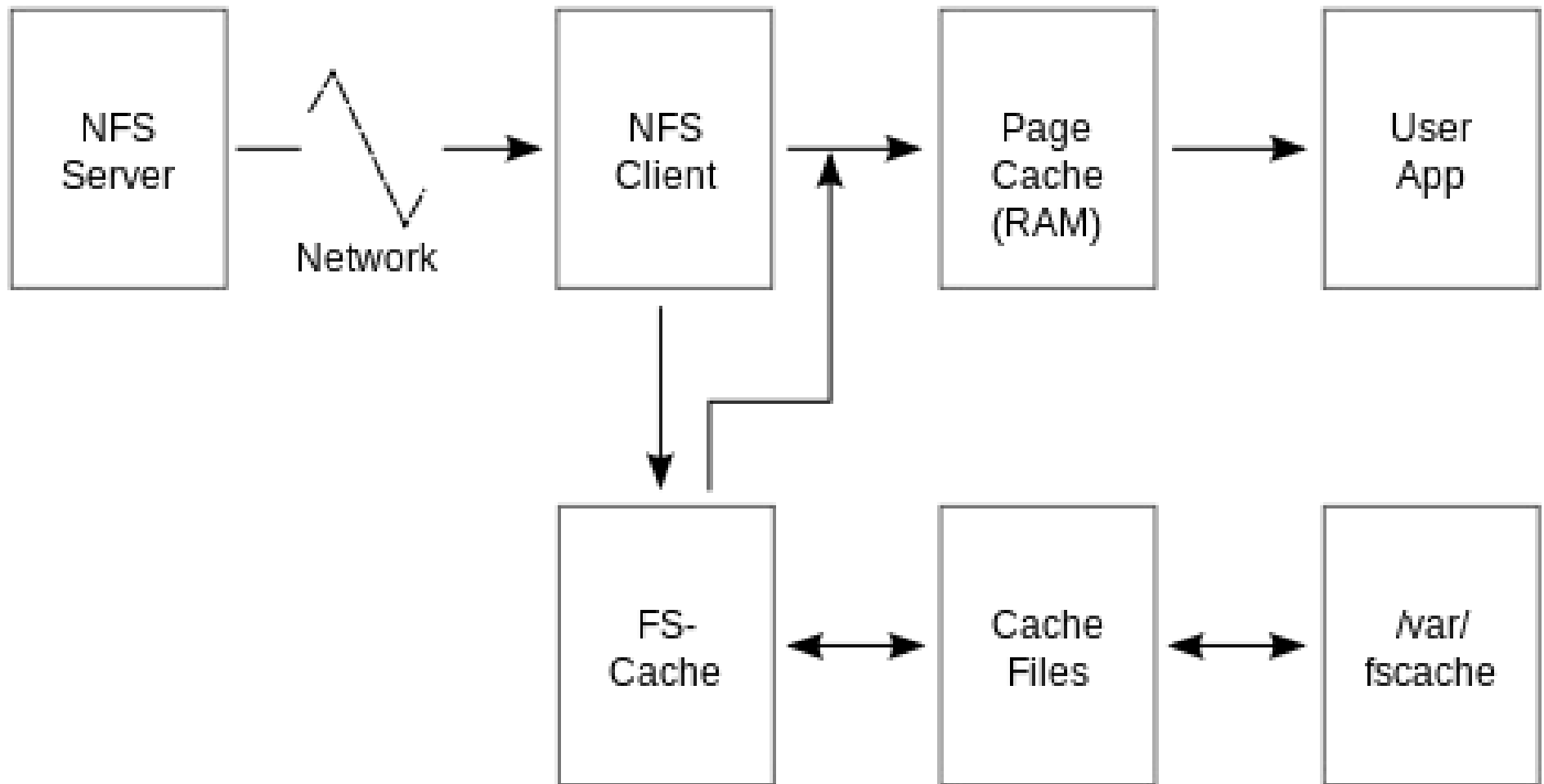
➢ **tmpfs is a very fast in-memory file system**

# Network File Systems

➤ Network File System – NFS v3/v4
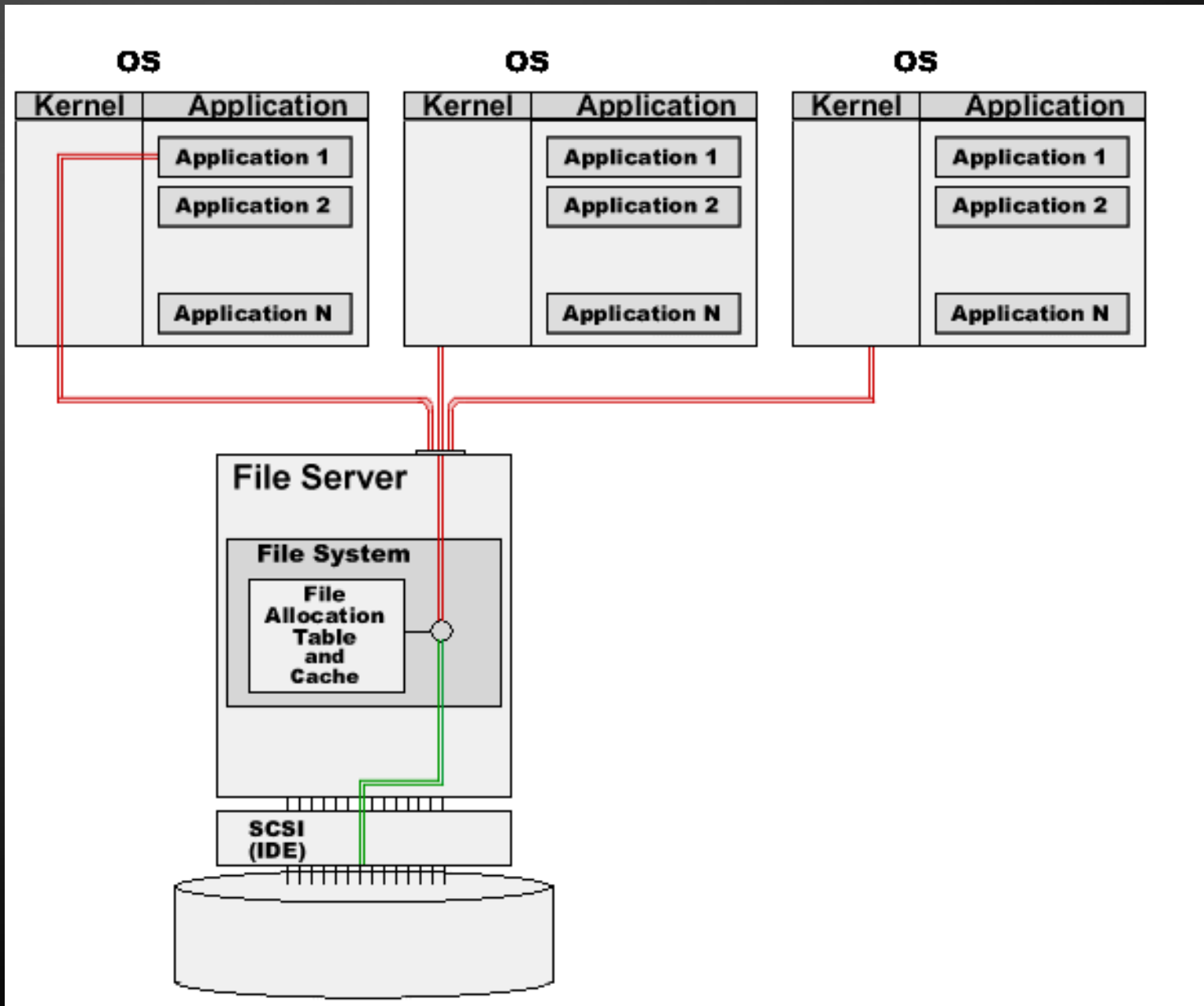
➤Common Internet File System - CIFS
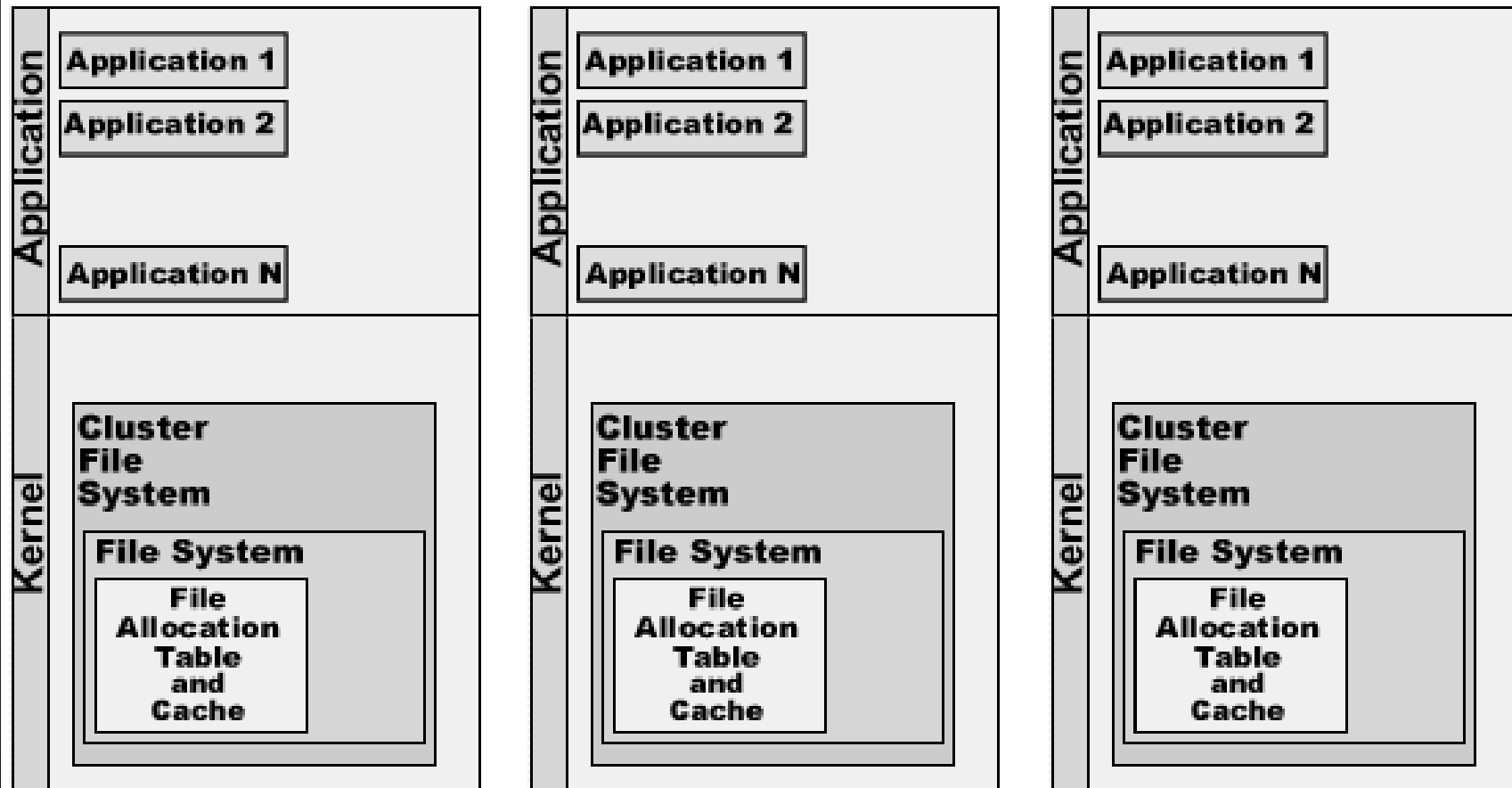
# Network File Systems
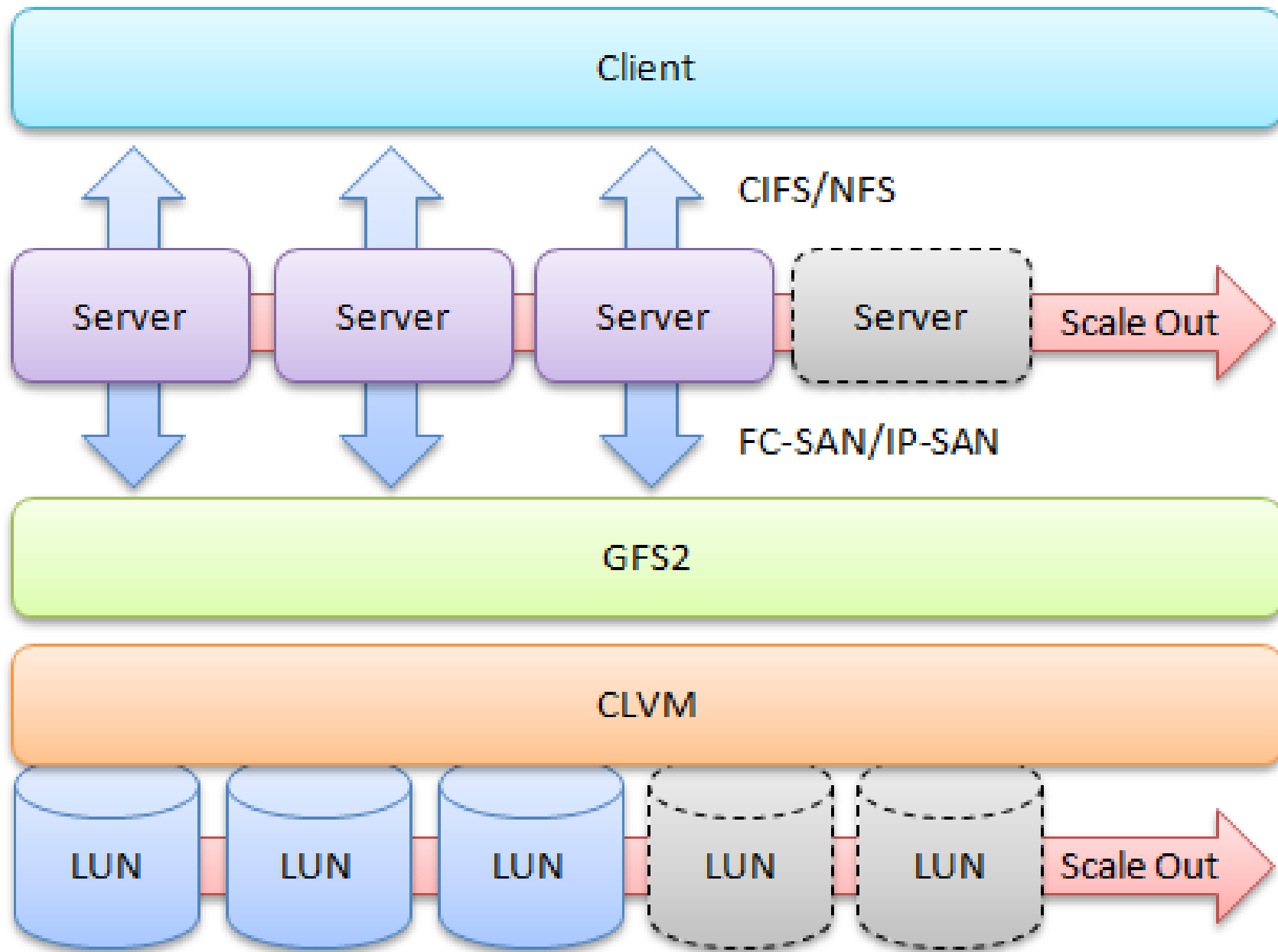
# Network File Systems
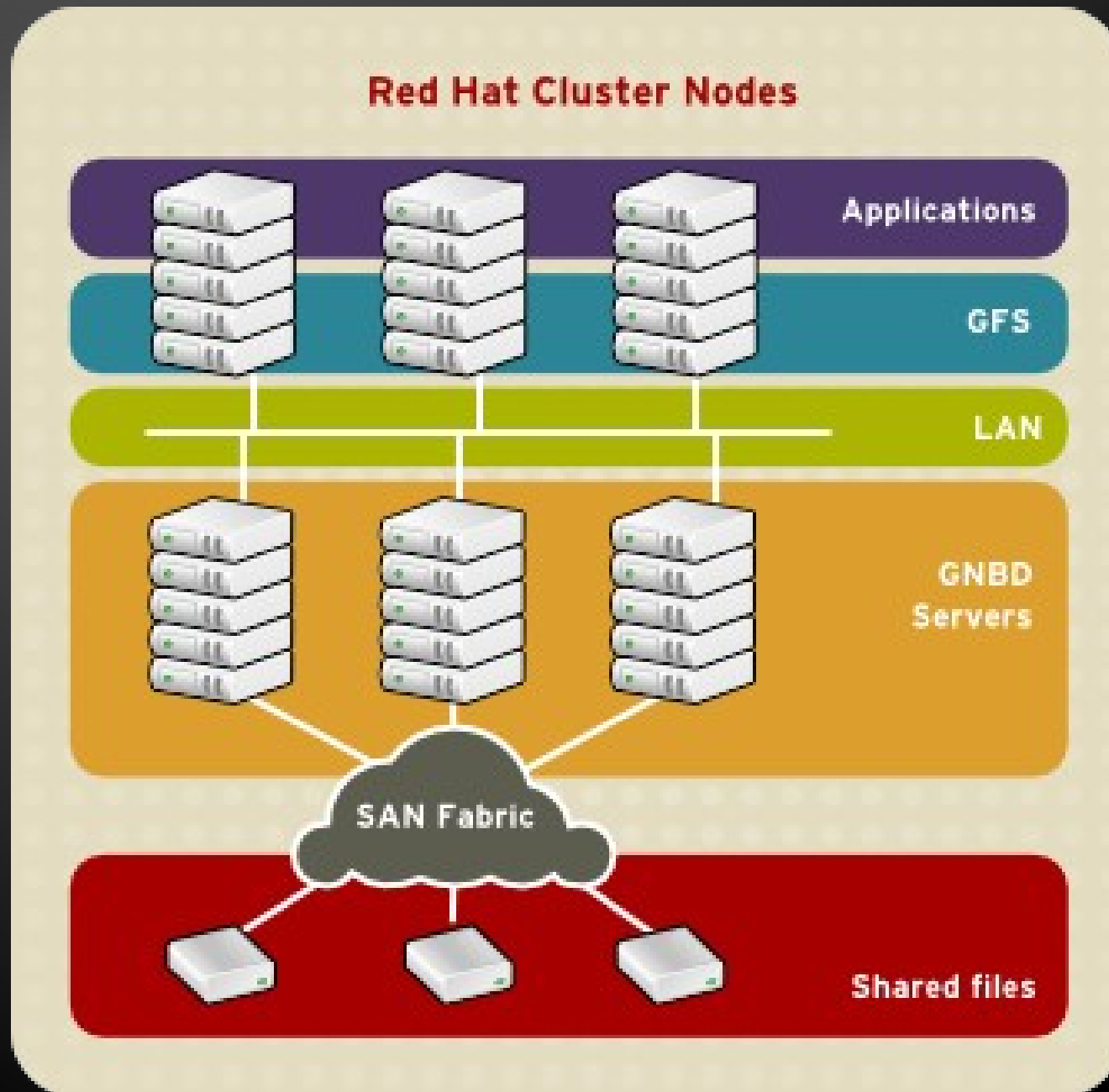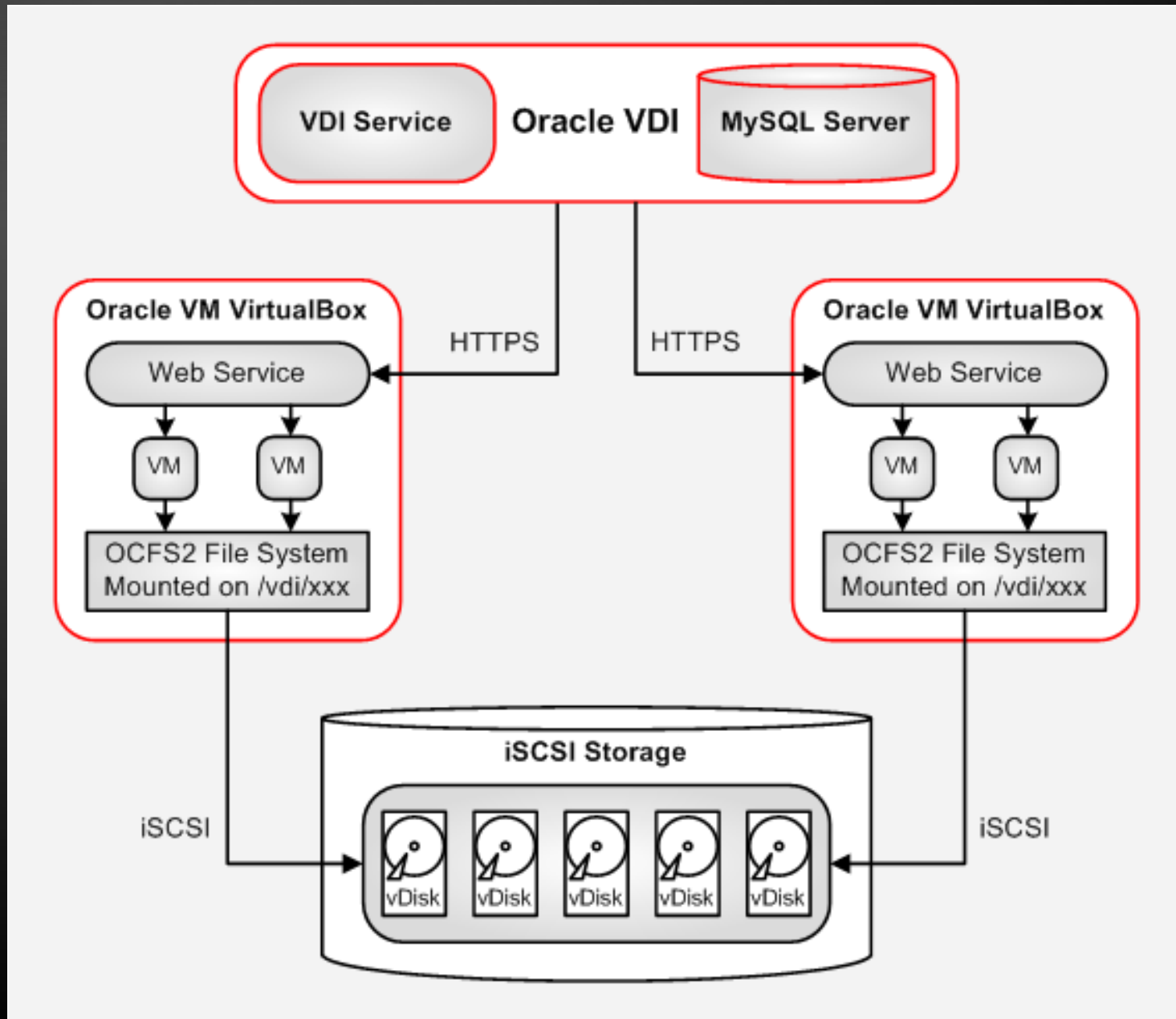
# Network File System

# Cluster File System

# Cluster File Systems

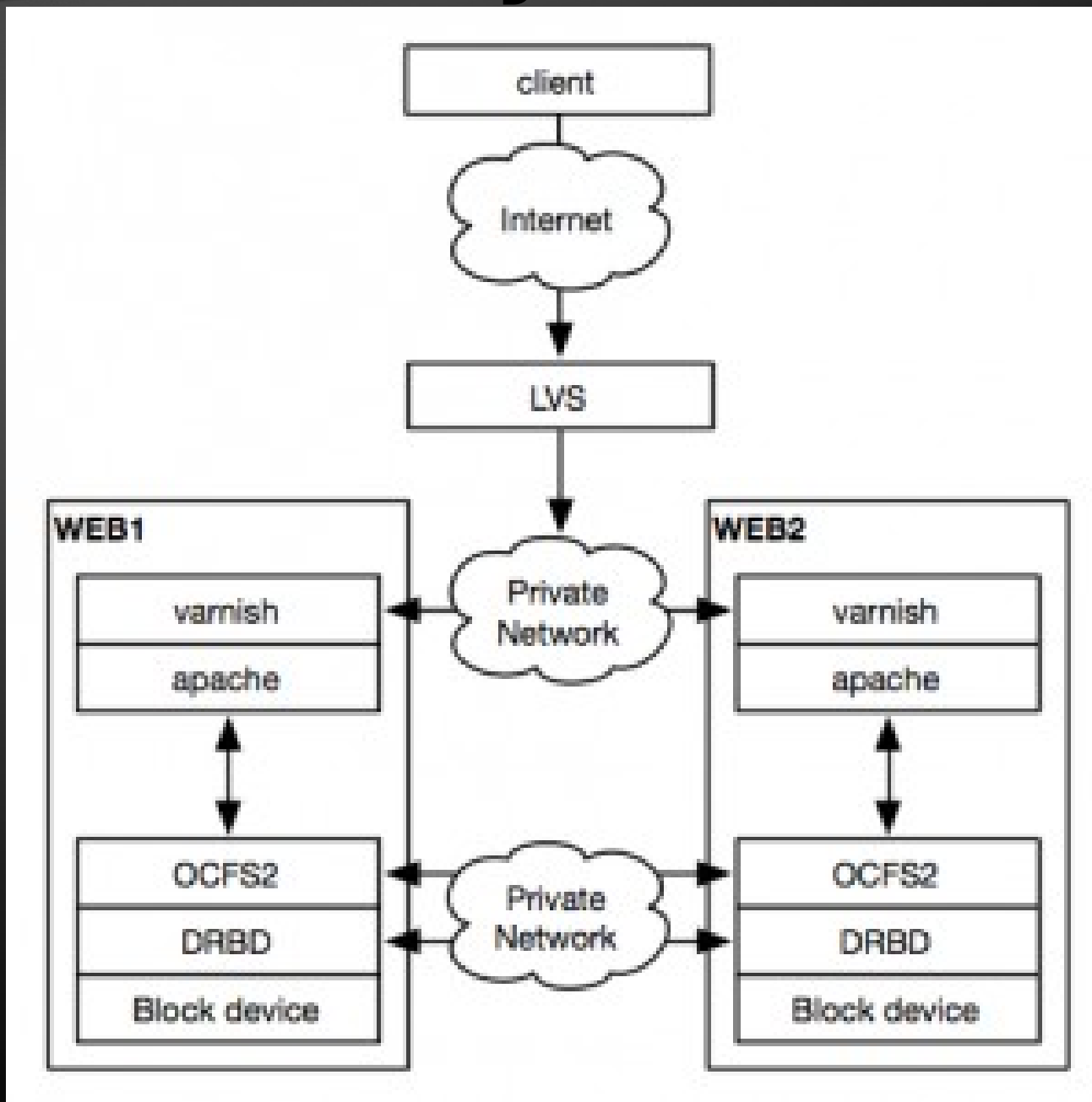- GFS, GFS2
- OCFS2

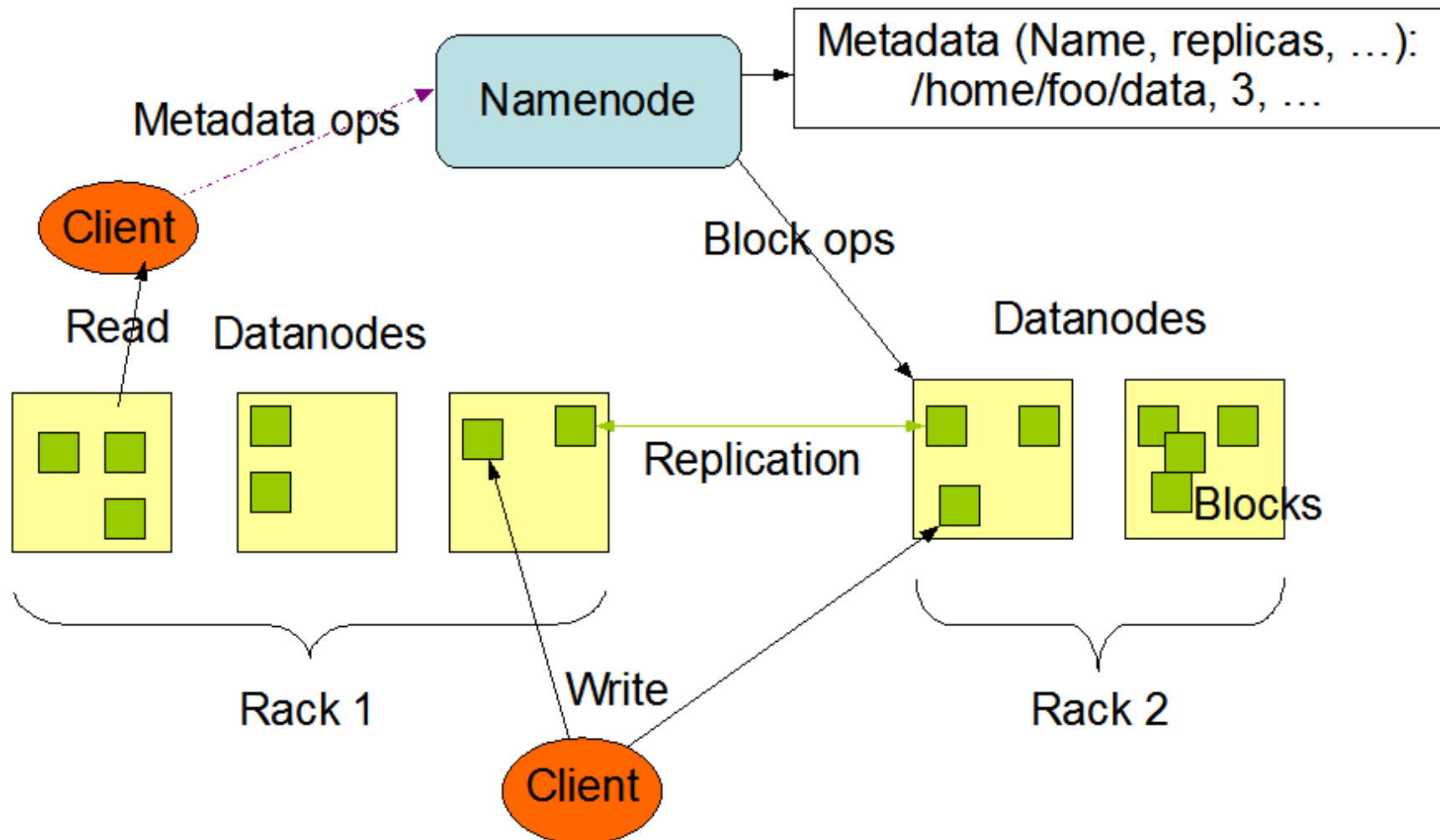# Cluster File Systems - GFS2

# Cluster File Systems - OCFS2

# Distributed File Systems

➢ **Hadoop**

➢ **Lustre**

➢ **GlusterFS**

➢ **GFarm**

➢ **FhgFS**
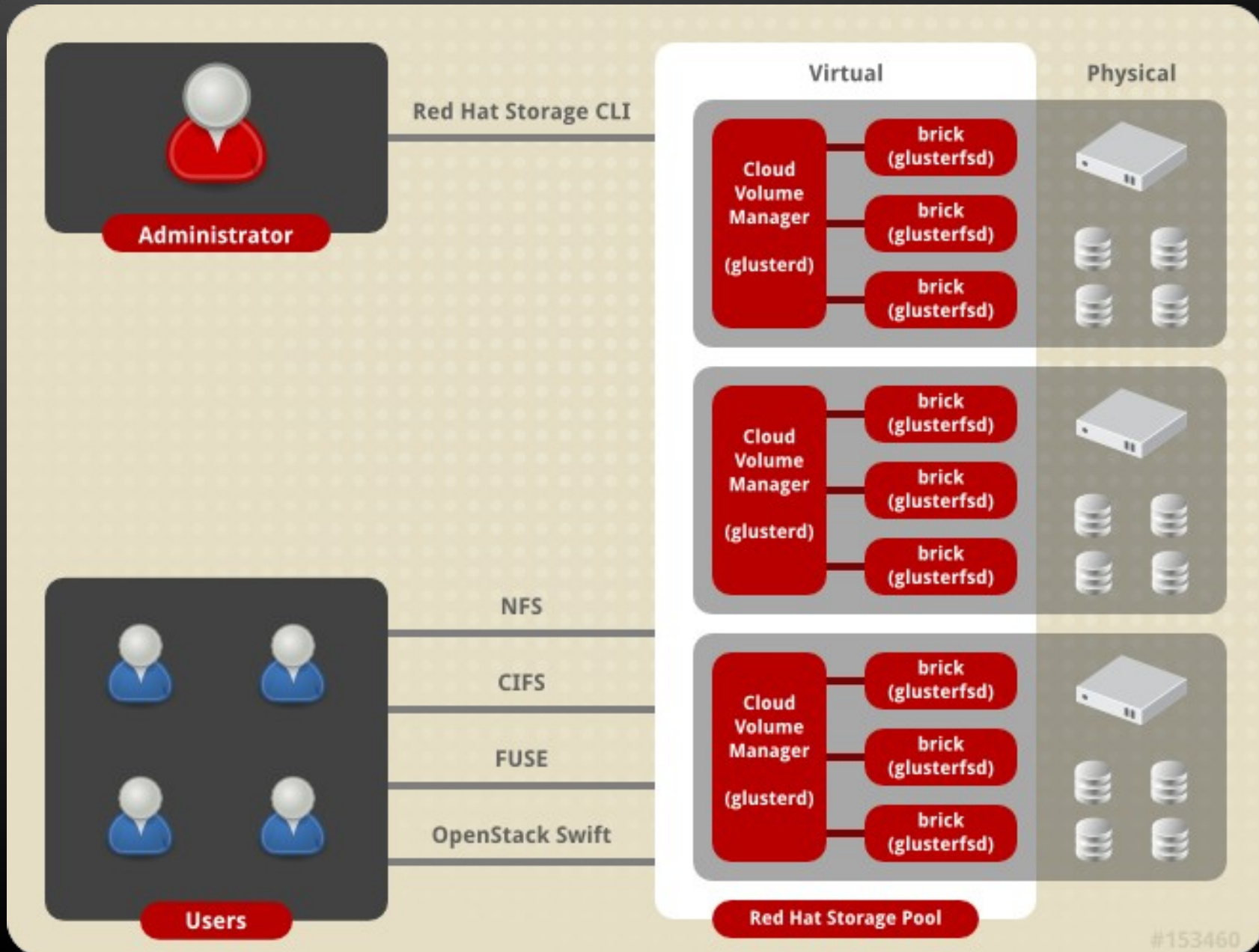
➢ **PohmelFS**
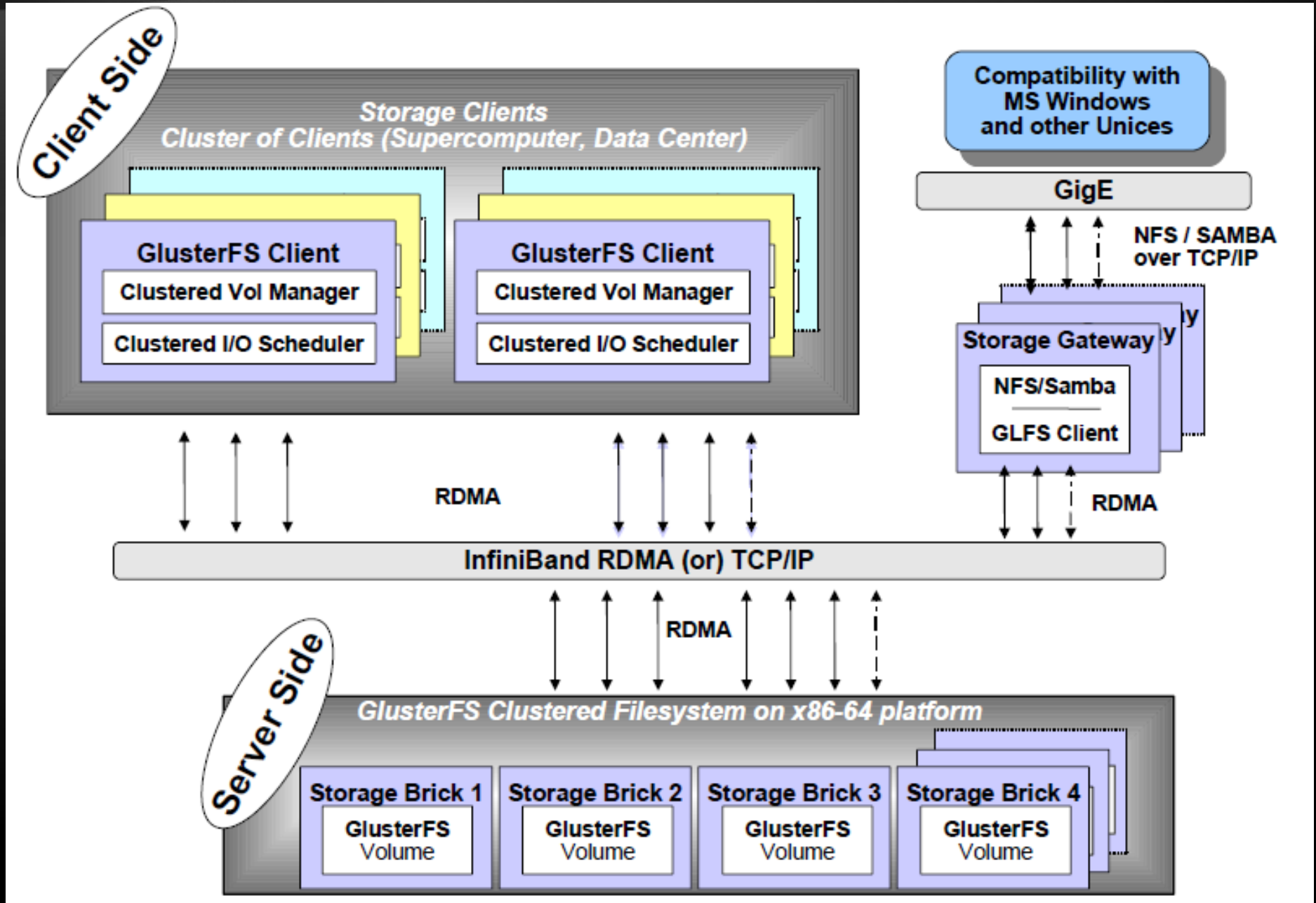
➢ **Ceph**

➢ **PVFS2**

# Hadoop

# Hadoop

➢ **Large block FS – 64MB**

➢ **Write mostly FS**

➢ **Writes smaller then one block wait**

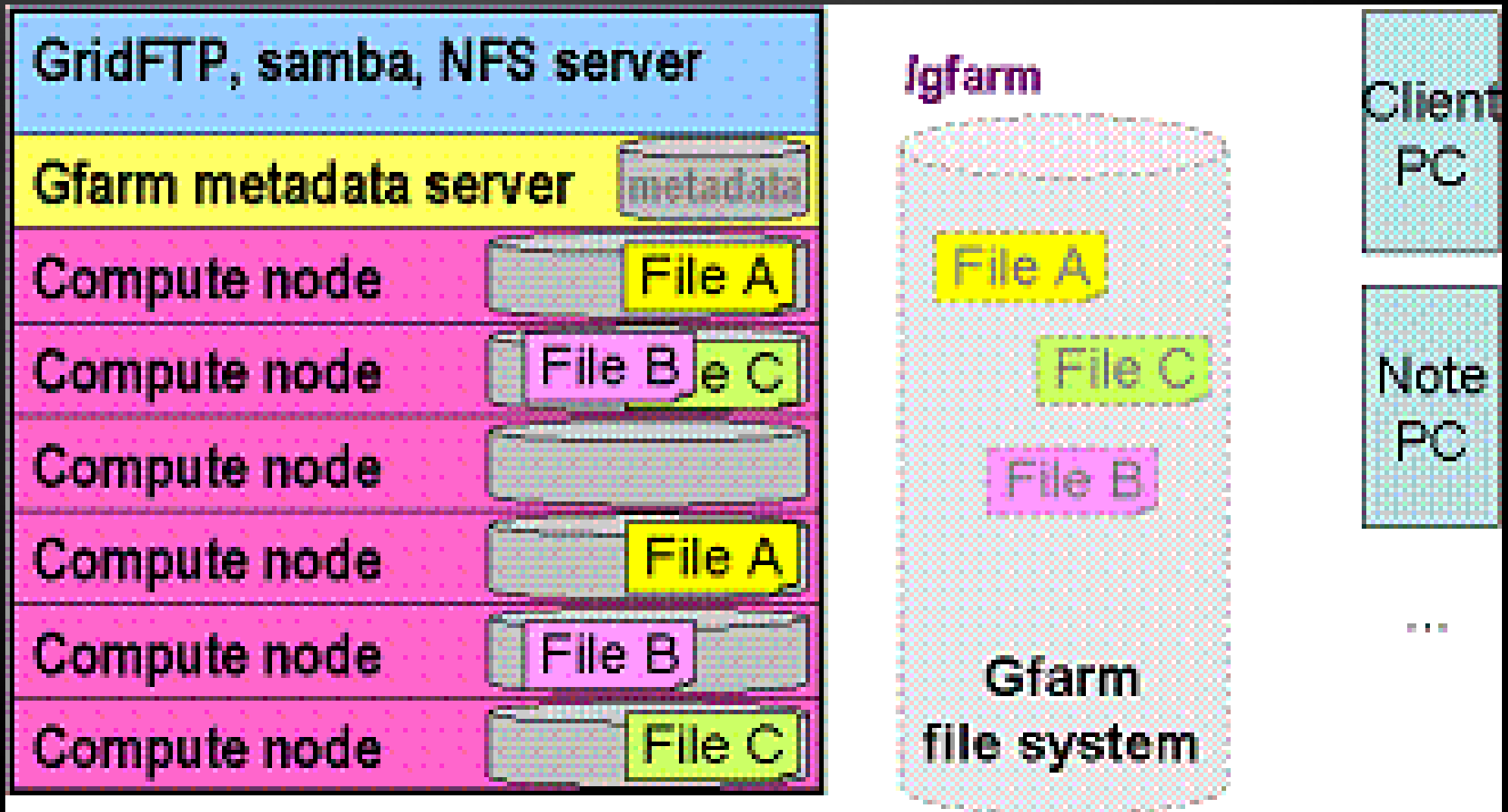➢ **Adding/removing nodes requires restart of the cluster**
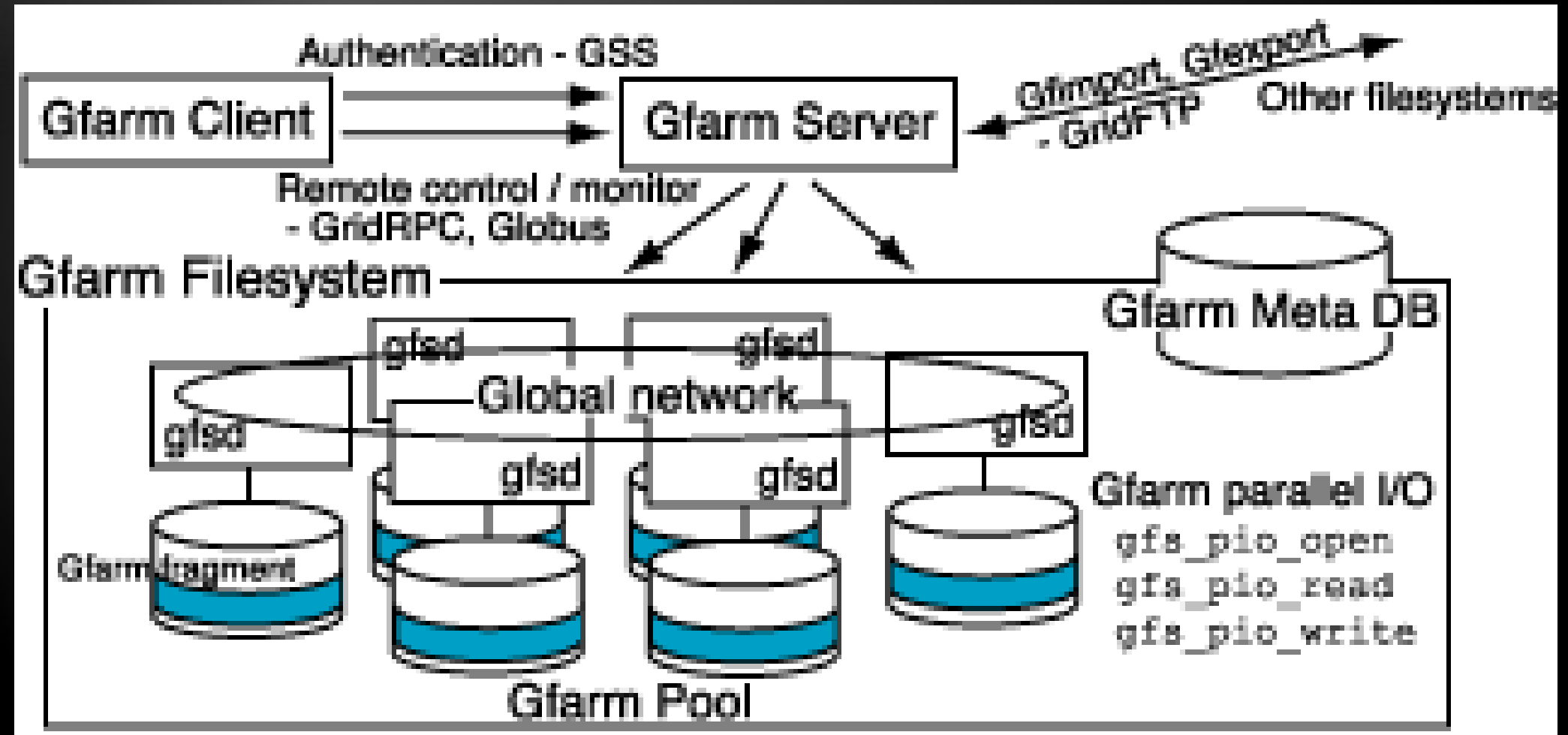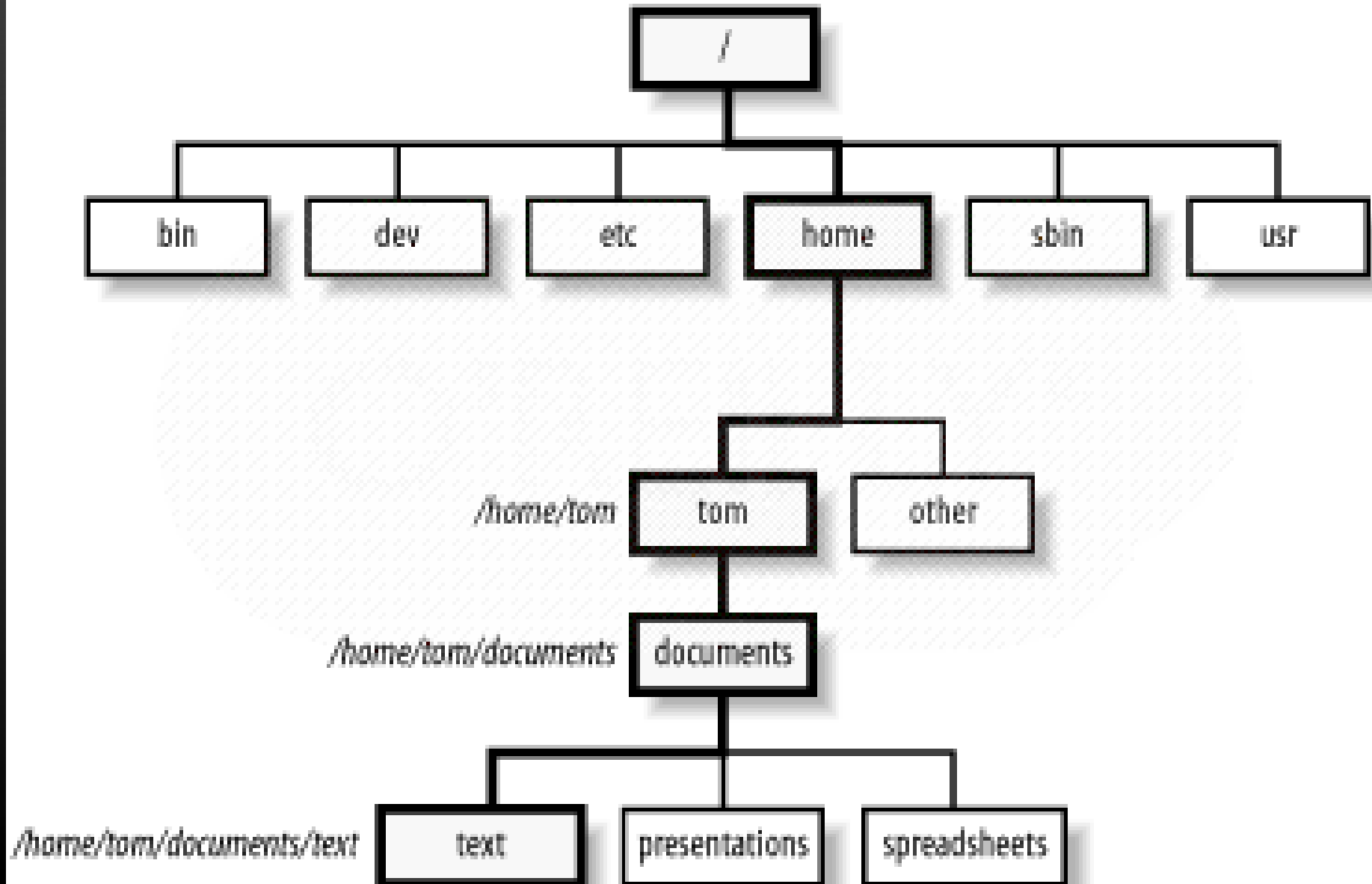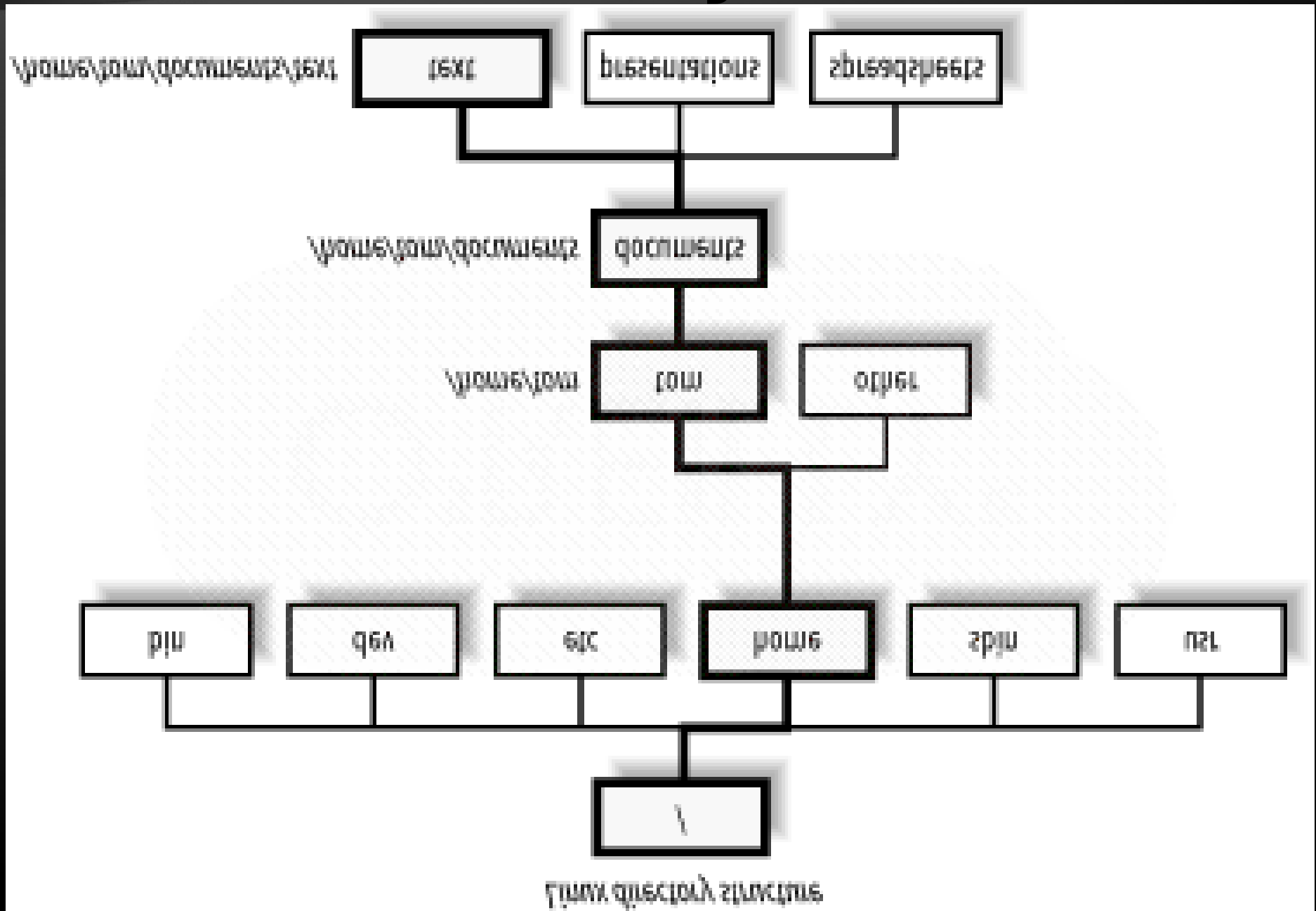
# GlusterFS

# GlusterFS

# GFarm

# GFarm

Linux directory structure

# FILESYSTEM HIERARCHY STANDARD ( FHS )

**ROOT DIRECTORY OF THE ENTIRE FILE SYSTEM HIERARCHY**

/

*PRIMARY HIERARCHY*

| Directory | Description |
|---|---|
| /bin/ | ESSENTIAL USER COMMAND BINARIES |
| /boot/ | STATIC FILES OF THE BOOT LOADER |
| /dev/ | DEVICE FILES |
| /etc/ | HOST-SPECIFIC SYSTEM CONFIGURATION<br>REQUIRED DIRECTORIES: OPT, X11, SGML, XML |
| /home/ | USER HOME DIRECTORIES |
| /lib/ | ESSENTIAL SHARED LIBRARIES AND KERNEL MODULES |
| /media/ | MOUNT POINT FOR REMOVABLE MEDIA |
| /mnt/ | MOUNT POINT FOR A TEMPORARILY MOUNTED FILESYSTEMS |
| /opt/ | ADD-ON APPLICATION SOFTWARE PACKAGES |
| /sbin/ | SYSTEM BINARIES |
| /srv/ | DATA FOR SERVICES PROVIDED BY THIS SYSTEM |
| /tmp/ | TEMPORARY FILES |
| /usr/ | (MULTI-)USER UTILITIES AND APPLICATIONS<br>SECONDARY HIERARCHY<br>REQUIRED DIRECTORIES: BIN, INCLUDE, LIB, LOCAL, SBIN, SHARE |
| /var/ | VARIABLE FILES |
| /root/ | HOME DIRECTORY FOR THE ROOT USER |
| /proc/ | VIRTUAL FILESYSTEM DOCUMENTING KERNEL AND PROCESS STATUS AS TEXT FILES |

/home/student/ → /home/student/dir

/home/linuxgym

/usr/local → /usr/local/bin

/usr/local → /usr/local/games

LINUXCONFIG.ORG

Wikipedia - Comparison of file systems

Ext2 and OCFS2 on-disk layout

Ext2 on-disk layout

XFS on-disk structure

ReiserFS on-disk structure

RFSTool for Windows

XFS Scalability

BtrFS on-disk structure

NILFS2 the new kid on the block

**Usenix paper on Log-Structured File Systems**

# Questions?

THAT TIME AGAIN?
Yes, then we'll have a beer