# Lesson 5 (b)

## File systems

Operating System Design

Bachelor in Informatics Engineering

# Recommended readings

## Base

1. Carretero 2007:
   1. Chapter 9

## Additional

1. Tanenbaum 2006(en):
   1. Chap.5
2. Stallings 2005:
   1. Three part
3. Silberschatz 2014:
   1. Chap. 10, 11 & 12

# Overview

1. Introduction

2. Main data structures in the secondary memory

3. Main data structures in the main memory

4. Block management

5. Complementary aspects

ARCOS @ UC3M

Alejandro Calderón Mateos

# Overview

1. **Introduction**

2. Main data structures in the secondary memory

3. Main data structures in the main memory

4. Block management
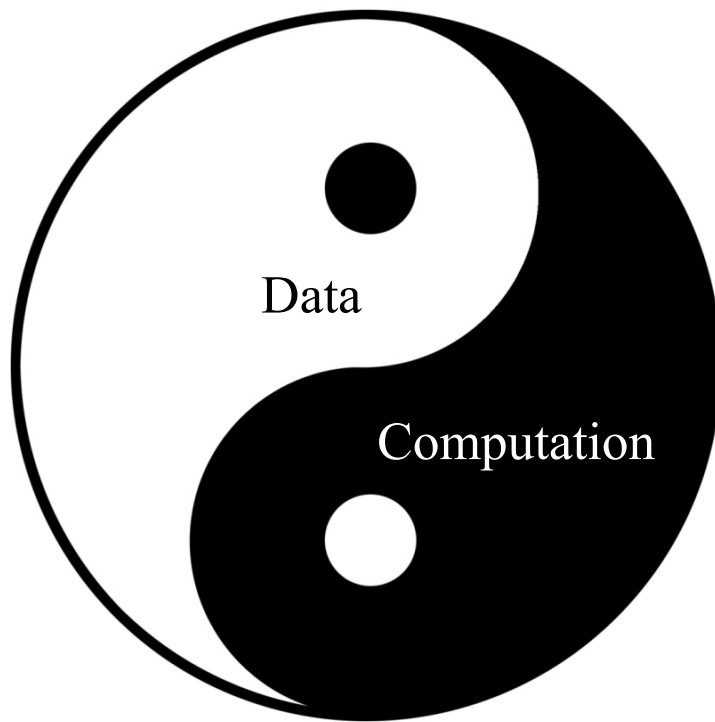
5. Complementary aspects

ARCOS @ UC3M
Alejandro Calderón Mateos

# Storage System Scope

# Storage System Scope



(1) Data abstraction

Process

(3) Another mgr

(2) Abstraction manager

ARCOS @ UC3M
Alejandro Calderón Mateos

# (**1/2**) The O.S. includes a basic and generic abstraction: file system

Alejandro Calderón Mateos

# (1/2) File system included

▸ The Operating System includes some implementations of a basic abstract representation for the storage systems: the file system.

# File system Characteristics



- ▸ Facilitate the secondary storage management.
    - ▸ Files, directories, etc.

- ▸ Independent from the physical device.

- ▸ Offer a unified logical view for users and applications.

sistemas operativos: una visión aplicada

ARCOS @ UC3M

Alejandro Calderón Mateos

# File system Architecture



Process (1)  Process (2)  …  Process (n)    User level

System level

Virtual File System

File organization modules

ext2    FAT    ...

Block server    Block cache

Device drivers

Device

…

sistemas operativos: una visión aplicada

ARCOS @ UC3M
Alejandro Calderón Mateos

# (**2/2**) The operating system supports the addition of other abstractions (& mgr.)

# File system Architecture



Process (n)

Process (1)    Process (2)    …    Py Mgr.

User level

System level

Virtual File System

File organization modules

ext2    FAT    …    xxxxx

Block server

Block cache

Device drivers

▸ A new file system implementation could be added.

▸ Other abstract representations could be implemented using the existing services on the Operating Systems.

sistemas operativos: una visión aplicada

ARCOS @ UC3M

Alejandro Calderón Mateos

# (1/2) Management structures



Process (1) | Process (2) | ... | Process (n) — User level

System level

Virtual File System

File organization modules

ext2 | FAT | ... | xxxxx

Block server | Block cache

Device drivers

Device

Main memory

Secondary memory

sistemas operativos: una visión aplicada

ARCOS @ UC3M

Alejandro Calderón Mateos

# (2/2) Management
# organization

Process (1)  Process (2)  …  Process (n)

Virtual File System

File organization modules

| ext2 | FAT | … | xxxxx |

Block server

Block cache

Device drivers

…

Device

**File System syscalls**

| Descriptors | | Usage of *namei* | | |
|---|---|---|---|---|
| open    pipe | | open    chown    unlink | | |
| creat    close | | creat    chmod    mknod | | |
| dup | | chdir    stat    mount | | |
| | | chroot    link    umount | | |

| Asig. i-n. | Attributes | I/O | File Sys. | View |
|---|---|---|---|---|
| creat | chown | read | mount | chdir |
| mknod | chmod | write | umount | chroot |
| link | stat | lseek | | |
| unlink | | | | |

**Low-level services for File Systems**

| namei | ialloc | alloc | bmap |
|---|---|---|---|
| iget    iput | ifree | free | |

**Block/cache management services**

| getblk | brelse | |
|---|---|---|
| bread | breada | bwrite |

sistemas operativos: una visión aplicada

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main requirements
## e.g.: Unix-like file system

- Processes have to use a secure interface, without direct access to the kernel data structures.

- Share the file offset position among processes from the same parent that open the file.

- Offer functionality for working with a file/directory in order to update the information that it contains.

- Go back and forth in the file system directory tree.

- Offer persistency of user data, seeking to minimize the impact on the performance and the space needed for the metadata.

- Keep track of the file systems registered in the kernel, and keep track of the mount point of these file systems.

ARCOS @ UC3M
Alejandro Calderón Mateos

# Getting the proper storage system for the requirements...
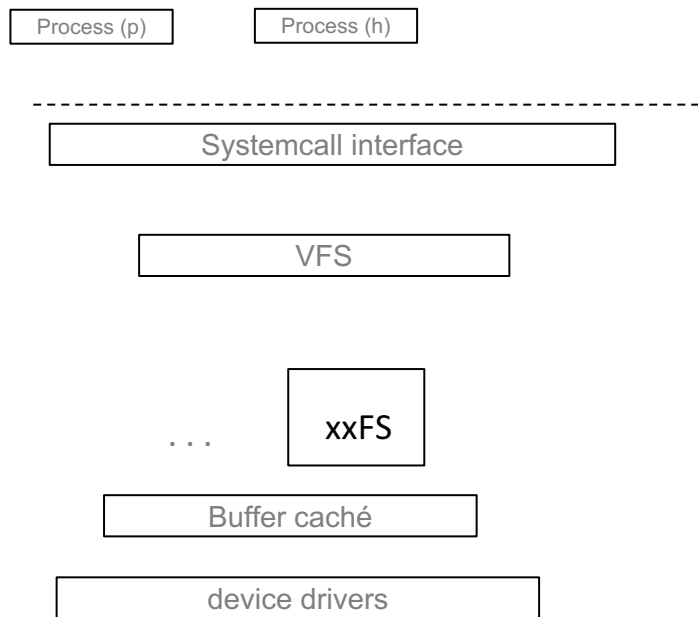
1. **To search** a file system that satisfies the requirements.

2. **To adapt** an existing file system in order to satisfy the requirements.

3. **To build** a file system that satisfies the requirements.

ARCOS @ UC3M

Alejandro Calderón Mateos

# File system organization
## main aspects: Linux

Process (p)    Process (h)

- - - - - - - - - - - - - - - - - - - - - - - - - -

Systemcall interface

VFS

. . .    xxFS

Buffer caché
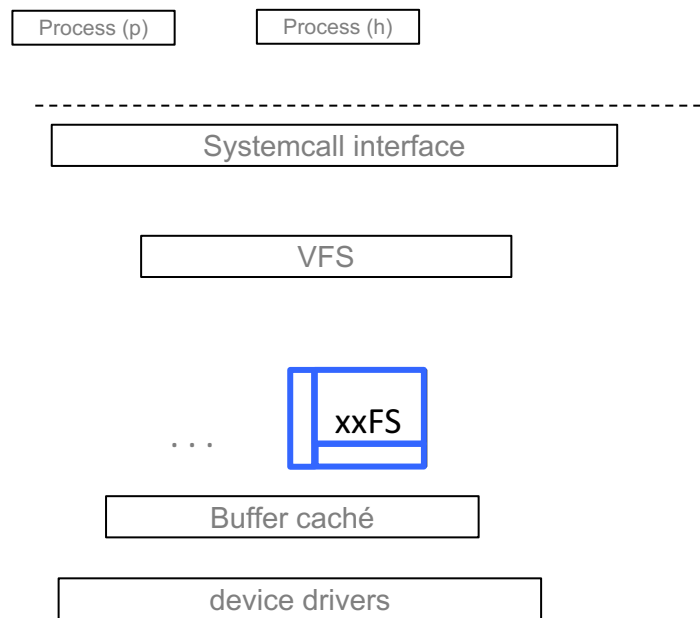
device drivers

▸ **Layered structure like UNIX.**

▸ **Main components:**

  ▸ System call interface

  ▸ VFS: *Virtual File System*

  ▸ xxFS: specific file system

  ▸ Buffer caché: block cache

  ▸ device drivers: *drivers*

http://publibn.boulder.ibm.com/doc_link/en_US/a_doc_lib/aixprggd/kernextc/virtual_fsys_over.htm#HDRA29C010F11

ARCOS @ UC3M
Alejandro Calderón Mateos

# File system organization
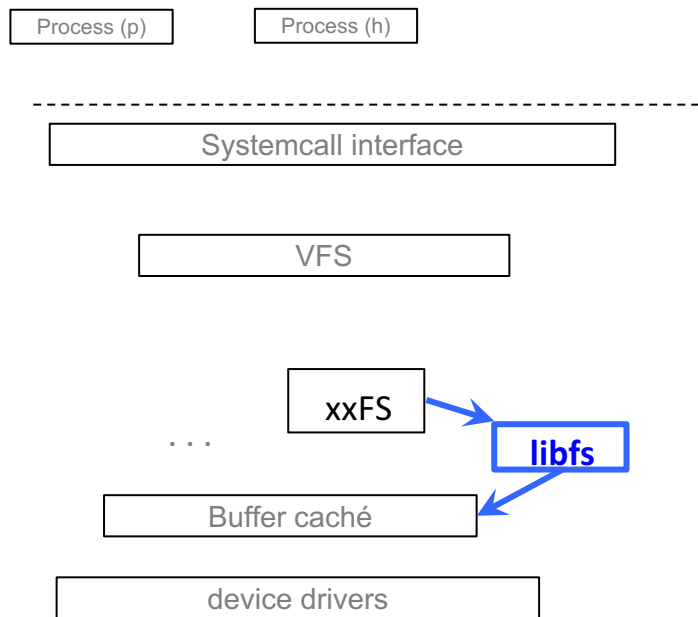## without *framework*, within kernel. E.g.: simplefs



- Interface:
  - **register**: to register the file system
  - …
  - **open**: to open a work session
  - **read**: read data
  - …
  - **namei**: convert from path to i-node
  - **iget**: read a i-node
  - **bmap**: compute an associated offset block
  - …

ARCOS @ UC3M
Alejandro Calderón Mateos

# File system organization
## with *framework*, within kernel: libfs

| | |
|---|---|
| Process (p) | Process (h) |

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Systemcall interface

VFS

xxFS

. . .  →  **libfs**

Buffer caché

device drivers
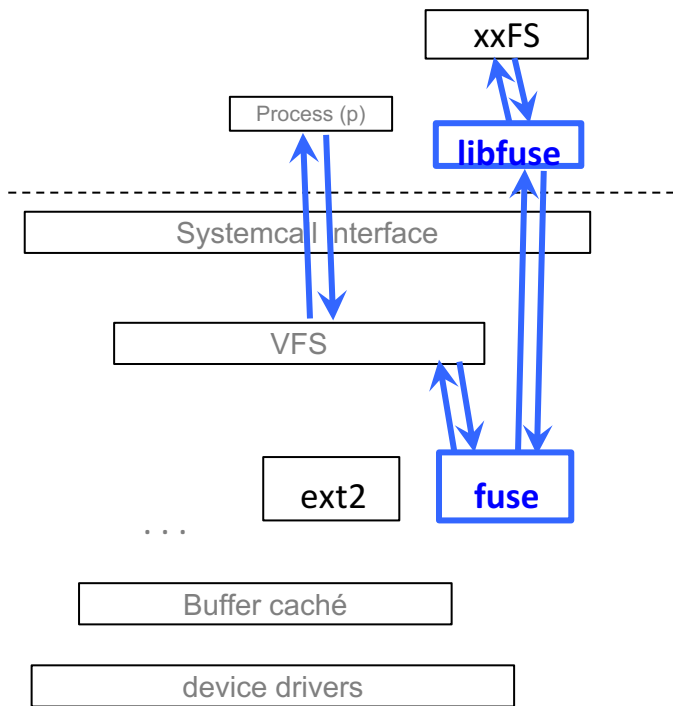
- ▸ Interface:
  libfs

  - ▸ **lfs_fill_super**: superblock

  - ▸ **lfs_create_file**: file creation

  - ▸ **lfs_make_inode**: default i-node

  - ▸ **lfs_open**: open a work session

  - ▸ **lfs_read_file**: read from file

  - ▸ **lfs_write_file**: write to file

  - ▸ ...

Alejandro Calderón Mateos

# File system organization
with *framework*, user space: fuse



Process (p)

xxFS

**libfuse**

Systemcall Interface

VFS

ext2        **fuse**

. . .

Buffer caché

device drivers

▸ Interface:
*File system in USer spacE*

struct fuse_operations {

  …

  int (*open) (const char *, struct fuse_file_info *);

  int (*read) (const char *, char *, size_t, off_t, struct fuse_file_info *);

  int (*write) (const char *, const char *, size_t, off_t,struct fuse_file_info *);

  int (*statfs) (const char *, struct statfs *);

  int (*flush) (const char *, struct fuse_file_info *);

  …

};

http://www.ibm.com/developerworks/linux/library/l-fuse/

# File system organization
## without *framework*, user space. E.g.: mtools



```
Process (p)        Process (h)              xxFS

- - - - - - - - - - - - - - - - - - - - - - - - - - - -
         Systemcall interface

              VFS

         ext2
    . . .
       Buffer caché

       device drivers
```

▸ To implement the file system interface in user space, and as library for other applications:

- ▸ **open**: to open a work session
- ▸ **read**: to read data
- ▸ …
- ▸ **namei**: to convert path into i-node
- ▸ **iget**: read i-node
- ▸ **bmap**: compute the associate block for a given offset
- ▸ …

http://www.gnu.org/software/mtools/manual/mtools.html

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main options for the file system organization

| | User space | Kernel space |
|---|---|---|
| With *Framework* | FUSE | libfs |
| Without *Framework* | E.g.: mtools | E.g.: simplefs |

E.g.:mtools

User level

System level

E.g.:simplefs

**FUSE**

**libfs**

ARCOS @ UC3M

Alejandro Calderón Mateos

# Overview

1. Introduction
2. **Main data structures on the secondary memory**
3. Main data structures in the main memory
4. Block management
5. Complementary aspects

ARCOS @ UC3M
Alejandro Calderón Mateos

# File system Structures

Entire disk

| label | partition | partition | partition |
|-------|-----------|-----------|-----------|

UNIX file system

| zone | zone | zone | zone | zone |
|------|------|------|------|------|

File system zone

| super block | inode bitmap | data bitmap | inode blocks | data blocks |
|-------------|--------------|-------------|--------------|-------------|

▸ UNIX/Linux

▸ FAT

ARCOS @ UC3M
Alejandro Calderón Mateos

# File system:
## Unix-like representation

Logical disk

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
|     |     |     |     |     |     |     |     |     |     |     |     |

# File system:
## Unix-like representation

Logical disk

Alejandro Calderón Mateos

# File system:
## Unix-like representation

Logical disk

| | 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Boot block | | | | | | | | | | | | |

Boot
sequence

Partition
table

ARCOS @ UC3M

Alejandro Calderón Mateos

# File system:
## Unix-like representation

Logical disk

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |

Boot block

Superblock

Boot sequence

Partition table

# allocation blocks

# i-nodes blocks

# data blocks

root i-node id.

...

# File system:
## Unix-like representation

Logical disk

| | 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Boot block | Superblock | Resources allocation | | | | | | | | | |

| Boot sequence |
|---|
| Partition table |

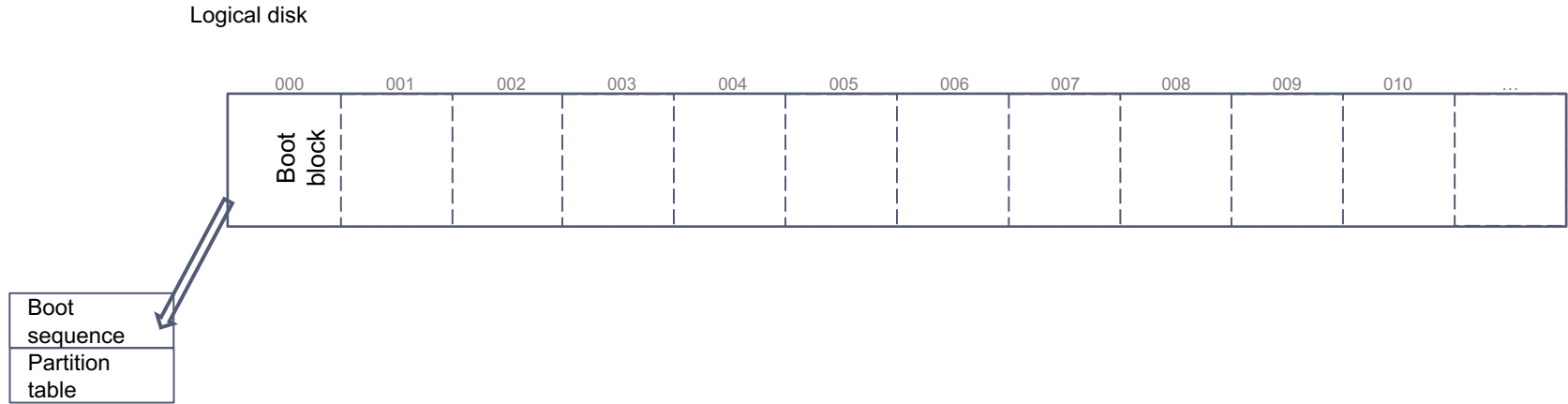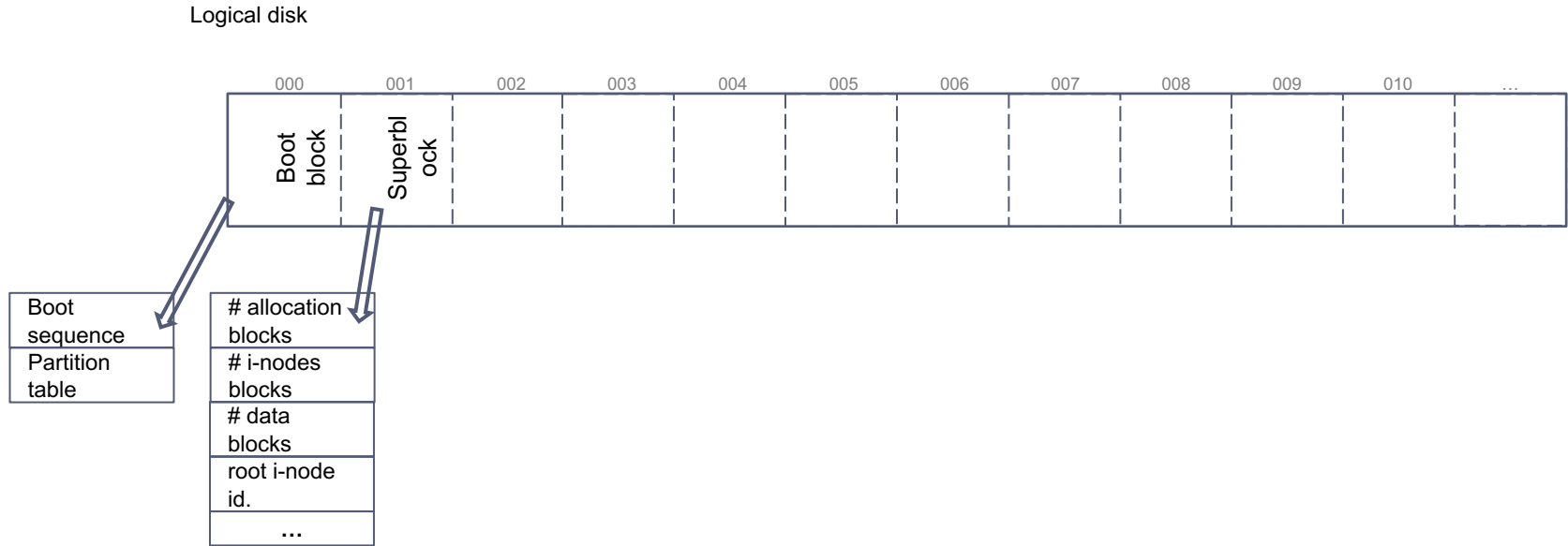| # allocation blocks |
|---|
| # i-nodes blocks |
| # data blocks |
| root i-node id. |
| ... |

0011…
11010…

# File system:
## Unix-like representation

Logical disk

# File system:
## Unix-like representation

Logical disk

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Boot block | Superblock | Resources allocation | | | | | | | | | |

| Boot sequence |
|---|
| Partition table |

| # allocation blocks |
|---|
| # i-nodes blocks |
| # data blocks |
| root i-node id. |
| ... |

0011... 11010...

‣ **Bit maps** or bit vectors:

   ‣ A vector with a bit per existing resource.
     If a resource is free then the bit value is 1, otherwise is 0.

      ‣ Easy to implement and simple to be used.

      ‣ Efficient if the device is not full or very fragmented.

‣ **Free resources list**:

   ‣ It keeps a linked list all the available resources and to use a pointer to the first element in the list.

      ‣ No very efficient method, unless for full devices or very fragmented devices.

‣ **Indexing**:

   ‣ Index table with the identification of the free resources.

ARCOS @ UC3M
Alejandro Calderón Mateos

# File system:
## Unix-like representation

Logical disk



Boot block | Superblock | Resources allocation | i-nodes

000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ...

| Boot sequence |
| Partition table |

| # allocation blocks |
| # i-nodes blocks |
| # data blocks |
| root i-node id. |
| ... |

0011…
11010…

| Entry attributes |
| Index blocks reference |

009
010

ARCOS @ UC3M
Alejandro Calderón Mateos

# File system:
## Unix-like representation

Logical disk



- ▸ Timestamp:
  - ▸ Creation, modification, last access, etc.
- ▸ Size:
  - ▸ Bytes or disk blocks used.
- ▸ Owner and protection:
  - ▸ Attributes, ACL, capacities, etc.
- ▸ Kind of file, link counter, etc.

# File system:
## Unix-like representation

Logical disk



| Boot sequence | Partition table |

| # allocation blocks |
| # i-nodes blocks |
| # data blocks |
| root i-node id. |
| ... |

0011…
11010…

Entry attributes
Index blocks reference

▸ **Contiguous allocation**:
  ▸ A contiguos list of blocks is allocated.
▸ **Non-contiguous allocation**:
  ▸ The first free block available is assigned.
  ▸ **Linked mechanism**:
    ▸ At the end of each block the identification of the following one is stored.
  ▸ **Indexed mechanism**:
    ▸ Blocks with the references of all file blocks.

ARCOS @ UC3M
Alejandro Calderón Mateos

# File systems:
## resources allocation alternatives
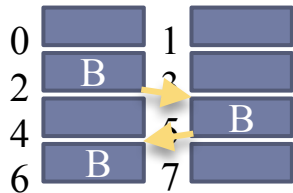


- **Contiguous** allocation:
  - The blocks of the files are contiguous.
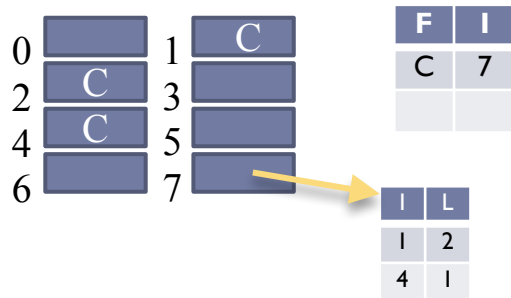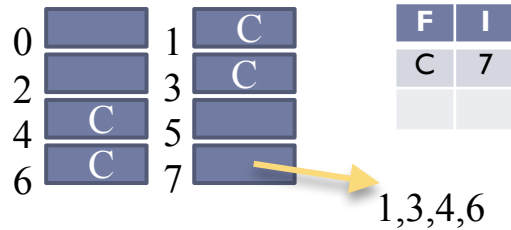  - It needs: first (I) and # of blocks (L)
  - To pack.

- **Non-contiguous** allocation :
  - Each block has the reference of the following one.
  - It needs: first (I) and # of blocks (L)
  - To defrag.

# File systems:
## resources allocation alternatives
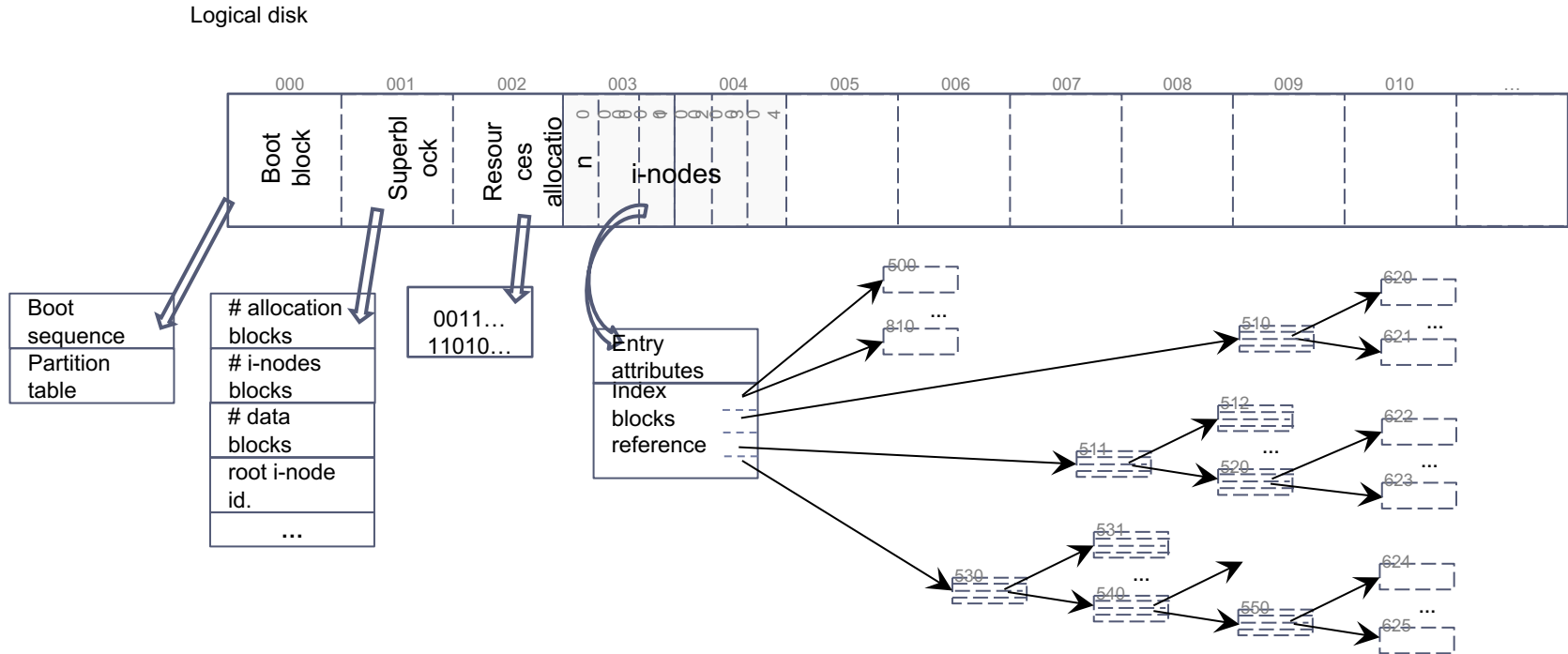


- **Indexed allocation (blocks):**
  - Some blocks are used to store the reference list of file data blocks.
  - It needs: id. Of the first index block.
  - To defrag.

- **Indexed allocation (extends):**
  - Some blocks are used to store the reference list of continuous file data blocks sequences.
  - It needs: id. of the first index block.
  - To defrag.

ARCOS @ UC3M
Alejandro Calderón Mateos

# File system:
## Unix-like representation



Logical disk

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |

Boot block | Superblock | Resources allocation | i-nodes

**Boot**
Boot sequence
Partition table

**Superblock**
# allocation blocks
# i-nodes blocks
# data blocks
root i-node id.
...

**Resources allocation**
0011…
11010…

**Entry**
Entry attributes
Index blocks reference

Alejandro Calderón Mateos

# How elements are represented



- Files



- Directories



- Links

# How elements are represented

▸ Files

▸ Directories

▸ Links

Alejandro Calderón Mateos

# File system:
## Unix-like representation: files

ARCOS @ UC3M
Alejandro Calderón Mateos

# How elements are represented



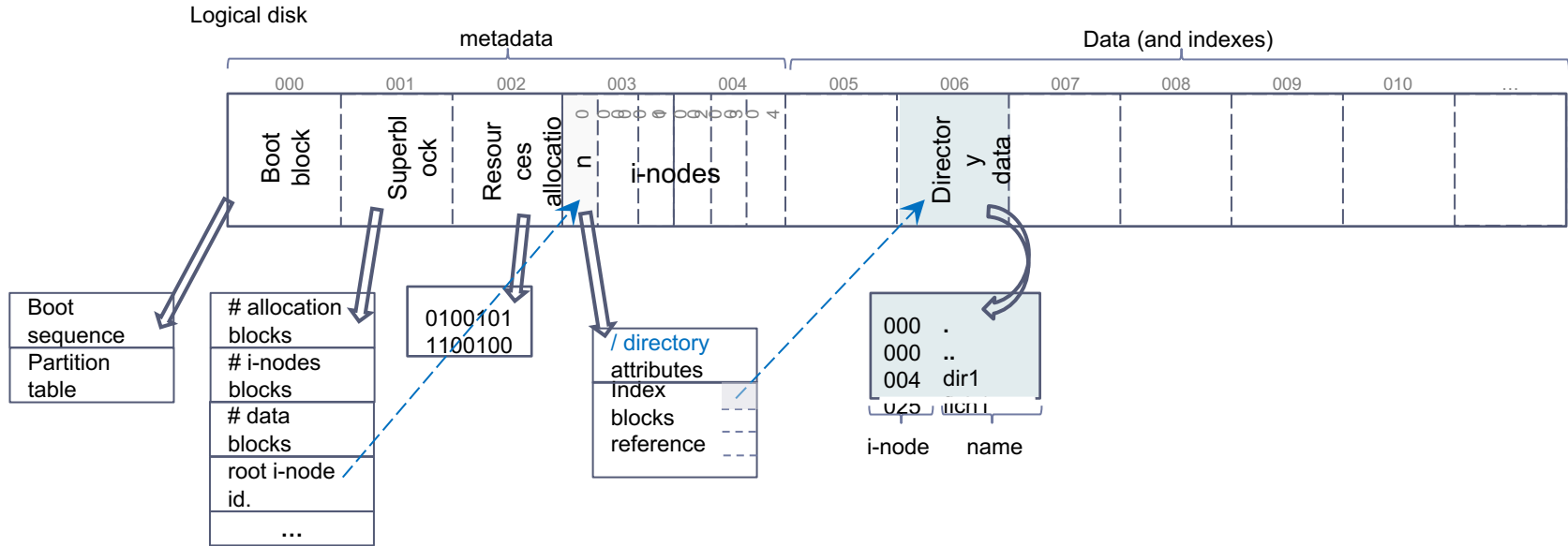▸ Files     ▸ Directories     ▸ Links

# File system:
## Unix-like representation: directories

# File system:
## Unix-like representation: directories

Alejandro Calderón Mateos

# File system:
## Unix-like representation: directories



ls –l /dir1/dir2/fich5.txt
- / + dir1 + dir2 + fich5.txt
- 4 i-nodes + 3 data blocks

Alejandro Calderón Mateos

# How elements are represented

- Files

- Directories

- Links

Alejandro Calderón Mateos

# File system:
## Unix-like representation: Symbolic link (soft link)

Logical disk

metadata | Data (and indexes)

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |

Boot block

Superblock

Resources allocation

i-nodes (0 1 2 3 4 ... 0 1 2 3 4)

Directory data

**Boot sequence** / **Partition table**

**# allocation blocks** / **# i-nodes blocks** / **# data blocks** / **root i-node id.** / ...

0100101
1100100

/dir1
attributes
Index
blocks
reference

**data(004)**

| 004 | . |
| 000 | .. |
| 054 | dir2 |
| 055 | fich2.c |

i-node    name

ln -s  /dir1  /dir1/soft

Alejandro Calderón Mateos

# File system:
## Unix-like representation: Symbolic link (soft link)



Logical disk

metadata

Data (and indexes)

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |

Boot block

Superblock

Resources allocation

i-nodes

Directory data

Boot sequence
Partition table

# allocation blocks
# i-nodes blocks
# data blocks
root i-node id.
...

0100101
1100100

**i-node 064**

Attributes

/dir1
attributes
Index
blocks
reference

data(**004**)

| 004 | . |
| 000 | .. |
| 054 | dir2 |
| 055 | fich2.c |

i-node     name

ln -s  /dir1  /dir1/soft

# File system:
## Unix-like representation: Symbolic link (soft link)



Logical disk

metadata | Data (and indexes)

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |

Boot block | Superblock | Resources allocation | i-nodes | | | | | Directory data | | |

i-nodes index: 0 1 2 3 4 ... 1 2 3 4

Boot sequence
Partition table

# allocation blocks
# i-nodes blocks
# data blocks
root i-node id.
...

0100101
1100100

Atributos del

/dir1
attributes
Index blocks
reference

**i-node 064**

data(**004**)

| 004 | . |
| 000 | .. |
| 054 | dir2 |
| 055 | fich2.c |
| **064** | **soft** |

i-node     name

ln -s  /dir1  /dir1/soft

# File system:
## Unix-like representation: Symbolic link (soft link)



Logical disk

metadata | Data (and indexes)

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |

Boot block | Superblock | Resources allocation | i-nodes | | | | | Directory data | | | Soft link data |

Boot sequence
Partition table

# allocation blocks
# i-nodes blocks
# data blocks
root i-node id.
...

0100101
1100100

Atributos del

/dir1
attributes
Index
blocks
reference

i-node 064

data(**004**)

| 004 | . |
| 000 | .. |
| 054 | dir2 |
| 055 | fich2.c |
| **064** | **soft** |

i-node    name

data(**064**)

/dir1

ln -s  /dir1  /dir1/soft

# File system:
## Unix-like representation: hard link



Logical disk

metadata | Data (and indexes)

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |

Boot block | Superblock | Resources allocation | i-nodes | | | | | Directory data | | |

Boot sequence
Partition table

# allocation blocks
# i-nodes blocks
# data blocks
root i-node id.
...

0100101
1100100

**i-node 004**
/dir1
attributes
Index
blocks
reference

data(**004**)
| 004 | . |
| 000 | .. |
| 054 | dir2 |
| 055 | fich2.c |
| 064 | soft |

i-node | name

ln    /dir1  /dir1/hard

# File system:
## Unix-like representation: hard link

Logical disk

metadata | Data (and indexes)

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |

Boot block

Superblock

Resources allocation

i-nodes

Directory data

Boot sequence
Partition table

# allocation blocks
# i-nodes blocks
# data blocks
root i-node id.
...

0100101
1100100

**i-node 004**
/dir1
attributes
Index
blocks
reference

**data(004)**

| 004 | . |
| 000 | .. |
| 054 | dir2 |
| 055 | fich2.c |
| 064 | soft |
| **004** | **hard** |

i-node | name

ln    /dir1  /dir1/hard

# File system:
## Unix-like representation: hard link



Logical disk

metadata — Data (and indexes)

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |

Boot block | Superblock | Resources allocation | i-nodes | | | | | Directory data

Boot sequence
Partition table

# allocation blocks
# i-nodes blocks
# data blocks
root i-node id.
...

0100101
1100100

**i-node 004**
/dir1       ③
attributes
Index blocks
reference

data(**004**)
004    .
000    ..
054    dir2
055    fich2.c
064    soft
**004    hard**
i-node    name

ln    /dir1  /dir1/hard

# File system:
## hard link vs soft link



Logical disk

metadata — Data (and indexes)

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |

Boot block | Superblock | Resources allocation | i-nodes | | | | | Directory data | | | Soft link data

Boot sequence
Partition table

# allocation blocks
# i-nodes blocks
# data blocks
root i-node id.
...

0100101
1100100

**i-node 004**
/dir1    3
attributes
Index
blocks
reference

**i-node 004**

| 004 | . |
| 000 | .. |
| 054 | dir2 |
| 055 | fich2.c |
| 064 | soft |
| 004 | hard |
| i-node | name |

**i-node 064**
/dir1

ln -s  /dir1  /dir1/soft
ln     /dir1  /dir1/hard

# File system structures



Entire disk

| label | partition | partition | partition |

UNIX file system

| zone | zone | zone | zone | zone |

File system zone

| super block | inode bitmap | data bitmap | inode blocks | data blocks |

▶ UNIX/Linux

▶ FAT

Alejandro Calderón Mateos

# File sytem structures:
## FAT

| Boot block | FAT$_1$ | FAT$_2$ | Root directory | Data block |
|---|---|---|---|---|

sistemas operativos: una visión aplicada

ARCOS @ UC3M

Alejandro Calderón Mateos

# Files and directories representation:
## FAT



**Root directory**

| name | Atrib. | KB | Agrup. |
|------|--------|-----|--------|
| dir1 | dir | 5 | 27 |
| fich1.txt | | 12 | 45 |

**dir1 directory**

| name | Atrib. | KB | Agrup. |
|------|--------|-----|--------|
| index.html | | 24 | 74 |
| prueba.zip | | 16 | 91 |

**FAT**

| | |
|-----|--------|
| 27 | <eof> |
| 45 | 58 |
| 51 | <eof> |
| 58 | <eof> |
| 74 | 75 |
| 75 | 76 |
| 76 | <eof> |
| 91 | 51 |

FAT 12 bits ($2^{12}$ blocks)
FAT 16 bits ($2^{16}$ blocks)
FAT 32 bits ($2^{32}$ blocks)

ARCOS @ UC3M
Alejandro Calderón Mateos

# Overview

1. Introduction

2. Main data structures on the secondary memory

3. **Main data structures in the main memory**

4. Block management

5. Complementary aspects

ARCOS @ UC3M

Alejandro Calderón Mateos

# Initial design…



**Module**

insmod ----→ **init_module()** ----→

rmmod ----→ **cleanup_module()** ----→

**open()**

**close()**

**read()**

**write()**

**ioctl()**

**...**

**Kernel core**

**register_capability()**

**unregister_capability()**

**File Systems**

**Mounted F.S.**

**Super-block**

**Resource assignment**

**i-nodes used**

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Boot block | Super-block | Resource allocation | i-nodes | | | | | | | | |

ARCOS @ UC3M
Alejandro Calderón Mateos

# Initial design...

Process (p)     Process (h)

**File Systems**

**Mounted
F.S.**

**Super-block**     **Resource
assignment**     **i-nodes
used**

| 000 | 001 | 002 | 003 | 004 | 005 | 006 | 007 | 008 | 009 | 010 | ... |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Boot block | Super-block | Resource allocation | i-nodes |  |  |  |  |  |  |  |  |

59

# Initial design...



open("/f1") -> 0x100
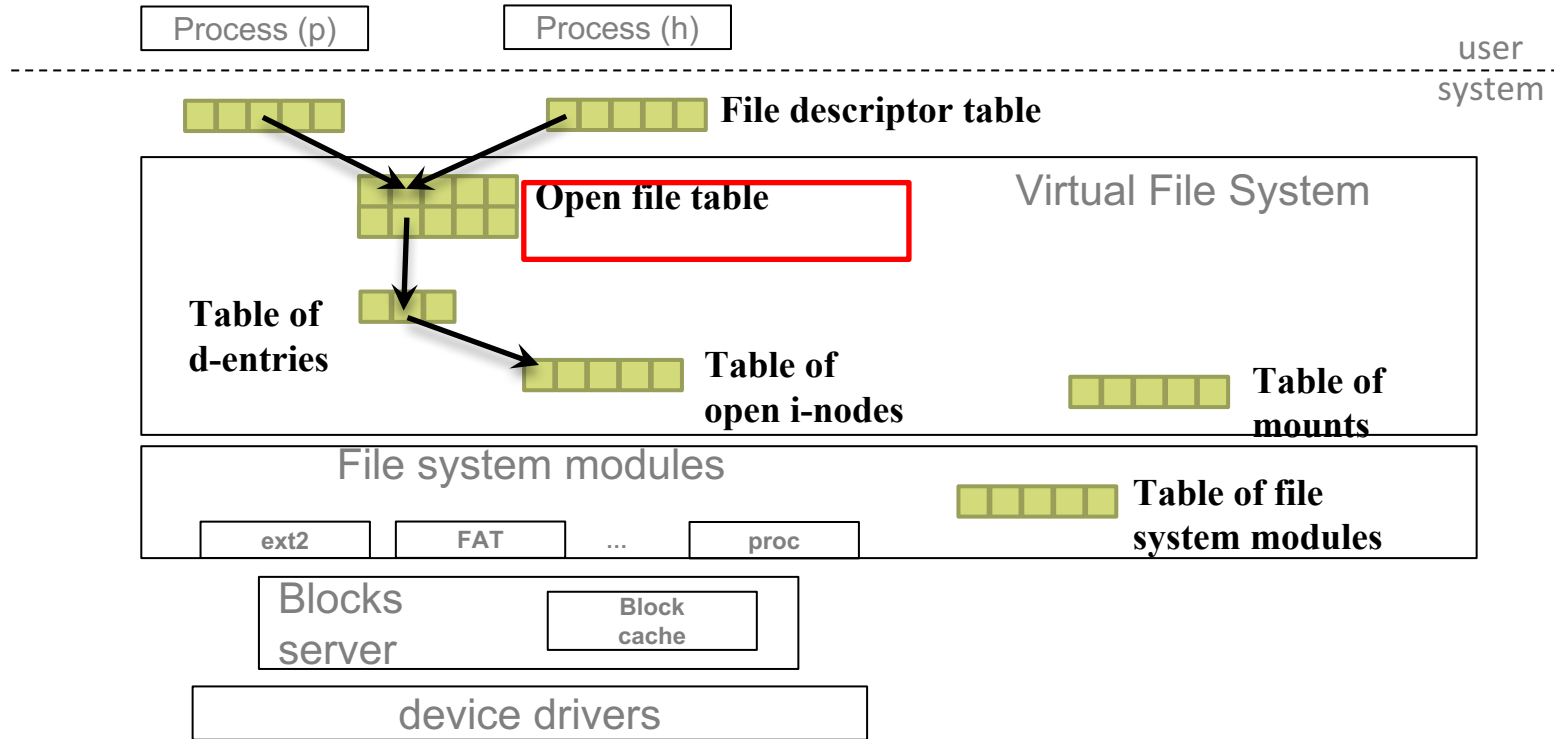...
read(0x150, buffer, 10)

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main goals
## (for a Unix-like file system)

▸ **Processes have to use a secure interface**, without direct access to the kernel data structures.

▸ **Share the file offset position** among processes from the same parent that open the file.

▸ Offer functionality for working with a file/directory in order to update the information that it contains.

▸ Go back and forth in the file system directory tree.

▸ Offer persistency of user data, seeking to minimize the impact on the performance and the space needed for the metadata.

▸ Keep track of the file systems registered in the kernel, and keep track of the mount point of these file systems.

# Main management structures

Process (p)    Process (h)

**File descriptor table**

**Open file table**    Virtual File System

**Table of
d-entries**

**Table of
open i-nodes**    **Table of
mount points**

File system modules

**Table of file
system modules**

| ext2 | FAT | ... | proc |

Blocks
server    Block
cache

device drivers

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures

ARCOS @ UC3M

Alejandro Calderón Mateos

## File descriptor table: Linux

```
struct fs_struct {

    atomic_t    count;           /* structure's usage count */

    spinlock_t  file_lock;        /* lock protecting this structure */

    int         max_fds;          /* maximum number of file objects */

    int         max_fdset;        /* maximum number of file descriptors */

    int         next_fd;          /* next file descriptor number */

    struct file **fd;             /* array of all file objects */

    fd_set      *close_on_exec;   /* file descriptors to close on exec() */

    fd_set      *open_fds;        /* pointer to open file descriptors */

    fd_set      close_on_exec_init; /* initial files to close on exec() */

    fd_set      open_fds_init;    /* initial set of file descriptors */

    struct file *fd_array[NR_OPEN_DEFAULT]; /* array of file objects */

};
```

http://www.makelinux.net/books/lkd2/ch12lev1sec10    ARCOS @ UC3M

Alejandro Calderón Mateos

# Main management structures
## Descriptors table (open files): Linux



**Module**

insmod - - - -> init_module() - - - -> **Kernel core** register_capability()

rmmod - - - -> cleanup_module() - - - -> unregister_capability()

open()

close()

read()

write()

ioctl()

...

Function call ———>

Function pointer ———>

Data pointer ———>

task_struct

files_struct  fs

... file

files

fd_array

f_op

kfd = current->files.fd_array[**fd**];

http://www.xml.com/ldd/chapter/book/ch02.html

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main goals
## (for a Unix-like file system)

- The processes have to use a secure interface, without direct access to the kernel representation.

- To share the file offset among process from the same parent that open the file.

- To have a working session with the file/directory in order to update the information that it contains.

- Go back and forth in the file system directory tree.

- Offer persistency of user data, seeking to minimize the impact on the performance and the space needed for the metadata.

- Keep track of the file system registered in the kernel, and keep track of the mount points of these file systems.

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures

Process (p)          Process (h)

**File descriptor table**

**Open file table**            Virtual File System

**Table of
d-entries**

**Table of
open i-nodes**

**Table of
mounts**

File system modules

| ext2 | FAT | ... | proc |

**Table of file
system modules**

Blocks
server

Block
cache

device drivers

# Main management structures
## Seek pointers table

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures
## Seek pointers table: Linux

**FD table P1**

| | |
|---|---|
| 0 | 23 |
| 1 | 4563 |
| 2 | 56 |
| 3 | 3 |
| 4 | 678 |

fd

**FD table P2**

| | |
|---|---|
| 0 | 230 |
| 1 | 563 |
| 2 | 98 |
| 3 | 3 |
| 4 | 247 |

fd

**FD table P3**

| | |
|---|---|
| 0 | 2300 |
| 1 | 53 |
| 2 | 4 |
| 3 | 3465 |
| 4 | 347 |

fd

**i-node table**

**d-entry table**

| i-node | Position |
|---|---|
| 92 | 345 |
| 92 | 5678 |
| | |

**Open file table**

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures
## File table: Linux

struct **file** {

    struct dentry             *f_dentry;

    struct vfsmount       *f_vfsmnt;

    struct file_operations  *f_op; →     struct **file_operations**  {

    mode_t              f_mode;                       int      (*open)   (struct inode *, struct file *);

    loff_t              f_pos;                      ssize_t (*read)   (struct file *, char *, size_t, loff_t *);

    struct fown_struct    f_owner;            ssize_t (*write)  (struct file *, const char *, size_t, loff_t *);

    unsigned int         f_uid, f_gid;       loff_t  (*llseek)  (struct file *, loff_t, int);

    unsigned long       f_version;         int      (*ioctl)   (struct inode *, struct file *,
                                                  unsigned int, ulong);

    ...                                 int      (*readdir) (struct file *, void *, filldir_t);

} ;                                   int      (*mmap)  (struct file *, struct vm_area_struct *);

...

} ;

http://www.win.tue.nl/~aeb/linux/lk/lk-8.html

# Main management structures
## File table: Linux



**Module**

insmod  ----→  init_module()  ----→

rmmod  ----→  cleanup_module()  ----→

open()

close()

read()

write()

ioctl()

...

**Kernel core**

register_capability()

unregister_capability()

... 

**files_operations**

read

write

...

**f_op**

**files_struct**

files

fd_array

**task_struct**

fs

file

lseek()

read()

...()

Function call  ---→

Function pointer  ---→

Data pointer  ---→

lseek(fd, newPos)
 kfd = current->files.fd_array[**fd**] ;
 kfd->f_pos = **newPos**;

http://www.xml.com/ldd/chapter/book/ch02.html

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures
## File table: Linux



**Module**

insmod - - - -> init_module() - - - ->
rmmod - - - -> cleanup_module() - - - ->

open()
close()
read()
write()
ioctl()
...

**Kernel core**

register_capability()
unregister_capability()

...

task_struct
files_struct        fs
files               file
files_operations
read                fd_array
write        f_op
...

lseek()
read()
...()

Function call - - - ->
Function pointer ———>
Data pointer ———>

read(fd, buf, size)
   kfd = current->files.fd_array[**fd**] ;
   kfd->f_op->**read(kfd,buf,size**);

http://www.xml.com/ldd/chapter/book/ch02.html
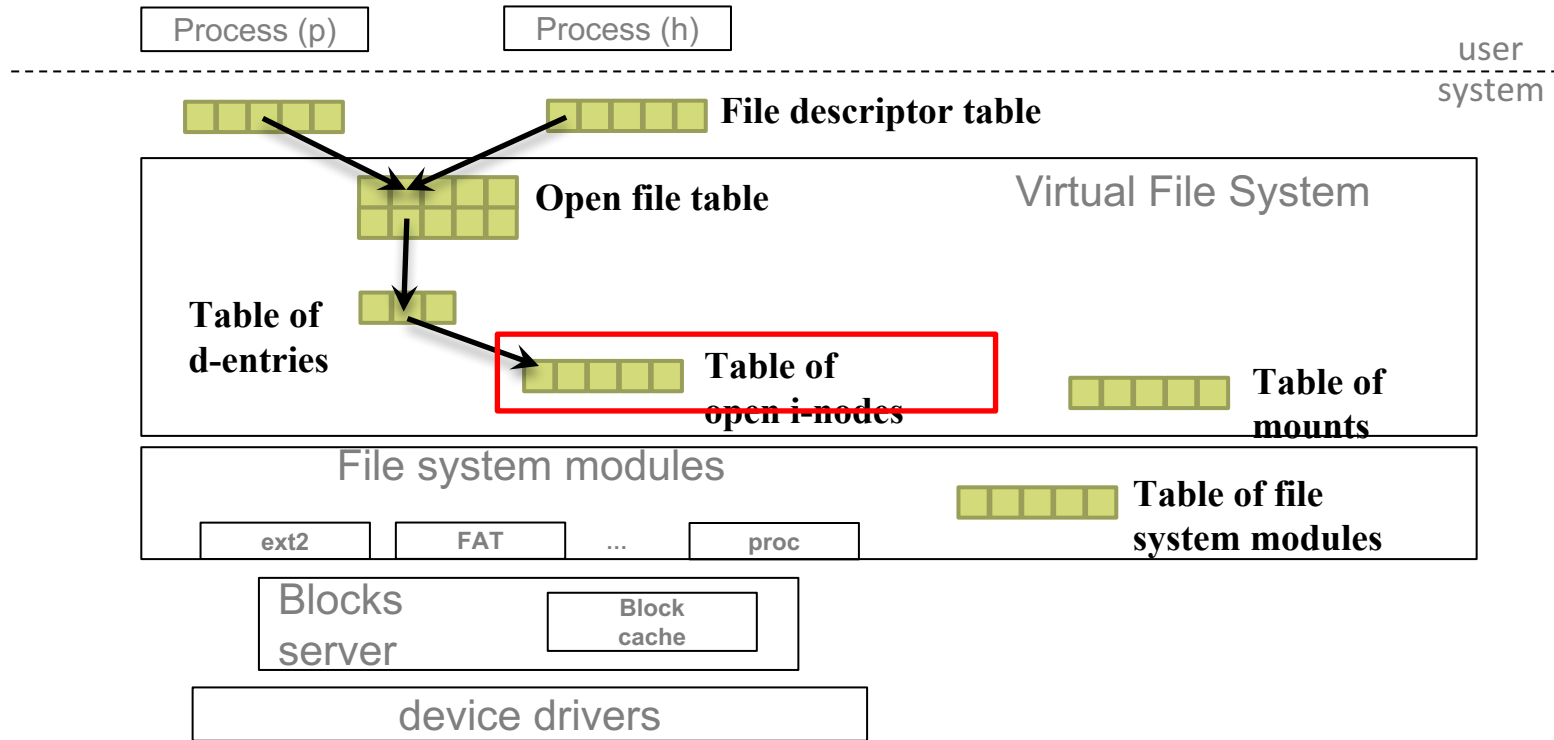
ARCOS @ UC3M
Alejandro Calderón Mateos

# Main goals
## (for a Unix-like file system)

- ▸ The processes have to use a secure interface, without direct access to the kernel representation.

- ▸ To share the file offset among process from the same parent that open the file.

- ▸ To have a working session with the file/directory in order to update the information that it contains.

- ▸ Go back and forth in the file system directory tree.

- ▸ Offer persistency of user data, seeking to minimize the impact on the performance and the space needed for the metadata.

- ▸ Keep track of the file system registered in the kernel, and keep track of the mount points of these file systems.

Alejandro Calderón Mateos

# Main management structures

Process (p)

Process (h)

File descriptor table

Open fie table

Virtual File System

Table of
d-entries

Table of
open i-nodes

Table of
mounts

File system modules

Table of file
system modules

ext2        FAT        ...        proc

Blocks
server

Block
cache

device drivers

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures
## Table of d-entries (directory entries): Linux

struct **dentry** {

    <span style="color:red">struct  inode          *d_inode;</span>

    struct  dentry       *d_parent;

    <span style="color:red">struct  qstr          d_name;</span>

    struct  dentry_operations  *d_op;

    struct  super_block       *d_sb;

    struct  list_head        d_subdirs;

    ...

}

struct **dentry_operations** {

    int (*d_revalidate) (struct dentry *, int);

    int (*d_hash)      (struct dentry *,  struct qstr *);

    int (*d_compare)  (struct dentry *, struct qstr *,
                        struct qstr *);

    int (*d_delete)    (struct dentry *);

    void (*d_release)  (struct dentry *);

    void (*d_iput)     (struct dentry *,
                    struct inode *);

}

Alejandro Calderón Mateos

# Main management structures
## Table of d-entries (directory entries): Linux

http://www.coins.tsukuba.ac.jp/~yas/classes/os2-2012/2013-02-19/index.html

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main goals
## (for a Unix-like file system)

▸ The processes have to use a secure interface,
without direct access to the kernel representation.

▸ To share the file offset among process from the same parent that open the file.

▸ To have a working session with the file/directory in order to update the information that
it contains.

▸ Go back and forth in the file system directory tree.

▸ Offer persistency of user data, seeking to minimize the impact on the performance and
the space needed for the metadata.

▸ Keep track of the file system registered in the kernel, and keep track of the mount points
of these file systems.

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures

Process (p)

Process (h)

**File descriptor table**

**Open file table**

Virtual File System

**Table of
d-entries**

**Table of
open i-nodes**

**Table of
mounts**

File system modules

**Table of file
system modules**

| ext2 | FAT | ... | proc |

Blocks
server

Block
cache

device drivers

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures
## Table of i-nodes: Linux

```
struct inode {
        unsigned long       i_ino;
        umode_t             i_mode;
        uid_t               i_uid;
        gid_t               i_gid;
        kdev_t              i_rdev;
        loff_t              i_size;
        struct timespec     i_atime;
        struct timespec     i_ctime;
        struct timespec     i_mtime;
        struct super_block          *i_sb;
        struct inode_operations  *i_op;
        struct address_space        *i_mapping;
        struct list_head            i_dentry;
        ...
} ;
```

http://www.win.tue.nl/~aeb/linux/lk/lk-8.html

ARCOS @ UC3M

Alejandro Calderón Mateos

# Main management structures
## Table of i-nodes: Linux

```
struct inode_operations {
        int (*create) (struct inode *,
                struct dentry *, int);
        int (*unlink) (struct inode *,
                struct dentry *);
        int (*mkdir) (struct inode *,
                struct dentry *, int);
        int (*rmdir) (struct inode *,
                struct dentry *);
        int (*mknod) (struct inode *,
                struct dentry *,
                int, dev_t);
        int (*rename) (struct inode *,
                struct dentry *,
                struct inode *,
                struct dentry *);
        void (*truncate) (struct inode *);
        struct dentry * (*lookup) (struct inode *,
                        struct dentry *);

        int (*permission) (struct inode *, int);
        int (*setattr) (struct dentry *,
                struct iattr *);
        int (*getattr) (struct vfsmount *mnt,
                struct dentry *,
                struct kstat *);
        int (*setxattr) (struct dentry *,
                const char *,
                const void *,
                size_t, int);
        ssize_t (*getxattr) (struct dentry *,
                const char *,
                void *, size_t);
        ssize_t (*listxattr) (struct dentry *,
                char *, size_t);
        int (*removexattr) (struct dentry *,
                const char *);

        int    (*link) (struct dentry *,
                struct inode *,
                struct dentry *);
        int (*symlink) (struct inode *,
                struct dentry *,
                const char *);
        int (*readlink) (struct dentry *,
                char *, int);
        int (*follow_link) (struct dentry *,
                struct nameidata *);
};
```

http://www.win.tue.nl/~aeb/linux/lk/lk-8.html

ARCOS @ UC3M
Alejandro Calderón Mateos

## Table of i-nodes: Linux

http://www.coins.tsukuba.ac.jp/~yas/classes/os2-2012/2013-02-19/index.html

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main goals
## (for a Unix-like file system)

- The processes have to use a secure interface, without direct access to the kernel representation.

- To share the file offset among process from the same parent that open the file.

- To have a working session with the file/directory in order to update the information that it contains.

- Go back and forth in the file system directory tree.

- Offer persistency of user data, seeking to minimize the impact on the performance and the space needed for the metadata.

- Keep track of the file systems registered in the kernel, and keep track of the mount points of these file systems.

Alejandro Calderón Mateos

# Main management structures

Alejandro Calderón Mateos

# Main management structures
## File system table: Linux

file_systems

```
struct  file_system_type  {
        const char  *name;
        int             fs_flags;
        struct dentry *(*mount) (struct file_system_type *,
                                    int, const char *, void *);
        void            (*kill_sb) (struct super_block *);
        struct  module          *owner;
        struct  file_system_type    *next;
        struct  list_head           fs_supers;
        struct  lock_class_key       s_lock_key;

        …

}
```

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures
## File system table: Linux

http://www.coins.tsukuba.ac.jp/~yas/classes/os2-2012/2013-02-19/index.html

ARCOS @ UC3M
Alejandro Calderón Mateos

## Table of mounts: Linux

```
struct vfsmount {
        struct vfsmount   *mnt_parent;  /* fs we are mounted on */
        struct dentry      *mnt_mountpoint;  /* dentry of mountpoint */
        struct dentry      *mnt_root;    /* root of the mounted tree */
        struct super_block *mnt_sb;      /* pointer to superblock */
        struct list_head    mnt_hash;
        struct list_head    mnt_mounts;  /* list of children, anchored here */
        struct list_head    mnt_child;    /* and going through their mnt_child */
        struct list_head    mnt_list;
        atomic_t            mnt_count;
        int                  mnt_flags;
        char                *mnt_devname;  /* Device name, e.g. /dev/hda1 */
};
```

current->namespace->list

http://www.win.tue.nl/~aeb/linux/lk/lk-8.html

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures
## Superblock table: Linux

```
struct super_block {

        dev_t                        s_dev;
        unsigned long                s_blocksize;
        struct file_system_type    *s_type;
        struct super_operations    *s_op;
        struct dentry              *s_root;
        ...
} ;
```

current->namespace->list->mnt_sb

http://www.win.tue.nl/~aeb/linux/lk/lk-8.html

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures
## Superblock table: Linux

```
struct super_operations  {

    struct inode *(*alloc_inode)(struct super_block *sb);

    void (*destroy_inode)(struct inode *);

    void (*read_inode) (struct inode *);

    void (*dirty_inode) (struct inode *);

    void (*write_inode) (struct inode *, int);

    void (*put_inode) (struct inode *);

    void (*drop_inode) (struct inode *);

    void (*delete_inode) (struct inode *);

    void (*clear_inode) (struct inode *);

    void (*put_super) (struct super_block *);

    void (*write_super) (struct super_block *);

    int   (*sync_fs)(struct super_block *sb, int wait);

    void (*write_super_lockfs) (struct super_block *);

    void (*unlockfs) (struct super_block *);

    int   (*statfs) (struct super_block *, struct statfs *);

    int   (*remount_fs) (struct super_block *, int *, char *);

    void (*umount_begin) (struct super_block *);

    int (*show_options)(struct seq_file *, struct vfsmount *);

};
```

http://www.win.tue.nl/~aeb/linux/lk/lk-8.html

ARCOS @ UC3M
Alejandro Calderón Mateos

http://www.coins.tsukuba.ac.jp/~yas/classes/os2-2012/2013-02-19/index.html

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures
## summary



Process (p)  Process (h)

**File descriptor table**

**Open file table**  Virtual File System

**Table of
d-entries**

**Table of
open i-nodes**

**Table of
mounts**

File system modules

**Table of file
system modules**

| ext2 | FAT | ... | proc |

Blocks
server

Block
cache

device drivers

# Main management structures
## summary

http://www.coins.tsukuba.ac.jp/~yas/classes/os2-2012/2013-02-19/index.html

ARCOS @ UC3M

Alejandro Calderón Mateos

# Main management structures
## summary (usage)

http://www.xml.com/ldd/chapter/book/ch02.html

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main goals
## (for a Unix-like file system)

✓

   ▸ The processes have to use a secure interface,
without direct access to the kernel representation.

✓

   ▸ To share the file offset among process from the same parent that open the file.

✓

   ▸ To have a working session with the file/directory in order to update the information that it contains.

✓

   ▸ Go back and forth in the file system directory tree.

✓

   ▸ Offer persistency of user data, seeking to minimize the impact on the performance and the space needed for the metadata.

   ▸ Keep track of the  file system registered in the kernel, and keep track of the  mount points of these file systems.

✓

Alejandro Calderón Mateos

# Overview

1. Introduction

2. Main data structures on the secondary memory

3. Main data structures in the main memory

4. **Block management**

5. Complementary aspects

ARCOS @ UC3M

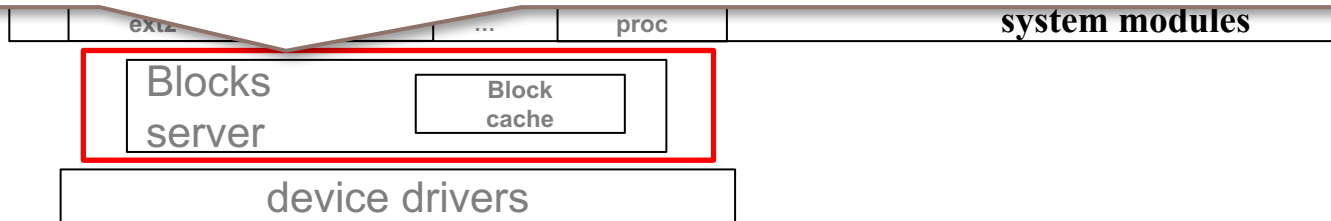Alejandro Calderón Mateos

# Main management structures



Process (p)    Process (h)

user
system

File descriptor table

Open file table    Virtual File System

Table of
d-entries

Table of
open i-nodes

Table of
mounts

File system modules

ext2    FAT    ...    proc

Table of file
system modules

Blocks
server

Block
cache

device drivers

ARCOS @ UC3M
Alejandro Calderón Mateos

# Main management structures

Process (p)          Process (h)

system

▸ **getblk**: find/reserve in cache a v-node block with its offset and size.

▸ **brelse**: to free a buffer and to insert it into the free list.

▸ **bwrite**: to write a cache block to the disk.

▸ **bread**: to read a disk block and store it in cache.

▸ **breada**: to read a block (and the following one) from disk to cache.

| ext2 | ... | proc | **system modules** |

Blocks
server    Block cache

device drivers

# Block server

▶ It is responsible for:

  ▶ Issuing commands to read and write device drivers blocks (by using the specific device routines)

  ▶ Optimizing the I/O requests.

    ▶ E.g.: Block cache.

  ▶ Offering a logical device namespace.

    ▶ E.g.: /dev/hda3 (third partition of the first disk)

ARCOS @ UC3M
Alejandro Calderón Mateos

# Block server

- **General behavior:**

  - If the block is in the cache

    - Copy the content (and update the block usage metadata)

  - If it is not in the cache

    - Read the block from the device and store it in cache

    - Copy the content (and to update the block metadata)

    - If the block has been modified (*dirty*)

      - Cache write policy

    - If the cache is full, it is necessary get some free slots

      - Cache replacement policy

ARCOS @ UC3M
Alejandro Calderón Mateos

# Block server

- ▶ **General behavior:**

  - ▶ If the block is in the cache

    - o **Read-ahead**:
      - o Read the following blocks into the cache (in order to improve the performance on sequential accesses)

    - ▶ To read the block from the device and store it in cache

    - ▶ To copy the content (and to update the block metadata)

    - ▶ If the block has been modified (*dirty*)

      - □ Cache write policy

    - ▶ If the cache is full, it is necessary get some free slots

      - □ Cache replacement policy

Alejandro Calderón Mateos

# Block server

o **write-through**:
  o Each time a block is modified it is also flushed to disk (lower performance)

o **write-back**:
  o The blocks are flushed to disk only when the block has to be evicted from the cache and it was dirty (better performance but reliability problems)

o **delayed-write**:
  o The modified blocks are saved to disk periodically (e.g., every 30 seconds in Unix) (trade-off for the former options)

o **write-on-close**:
  o When the file descriptor is closed, all file blocks are flushed to disk.

▸ If the block    been modified (*dirty*)

  □ Cache write policy

▸ If the cache is full, it is necessary get some free slots

  □ Cache replacement policy

# Block server

▸ **General behavior:**

  ▸ If the block is in the cache

    ▸ To copy the content (and to update the block usage metadata)

  ▸ If it is not in the cache

    ▸ To read the block from the device into the cache

      o **FIFO** *(First in First Out)*
      o **Clock algorithm** *(Second opportunity)*
      o **MRU** *(Most Recently Used)*
      o **LRU** *(Least Recently Used)*

    ▸ If the cache is full, it is necessary get some free slots

      □ Cache replacement policy

Alejandro Calderón Mateos

# Overview

1. Introduction

2. Main data structures on the secondary memory

3. Main data structures in the main memory

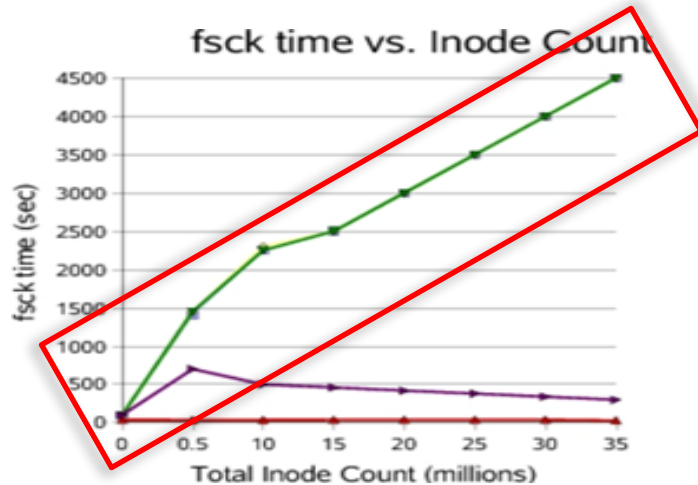4. Block management

5. **Complementary aspects**

ARCOS @ UC3M
Alejandro Calderón Mateos

# Advanced features

▸ **Journaling**

▸ Snapshots

▸ Dynamic file system expansion

# Without *Journaling*



fsck time vs. Inode Count
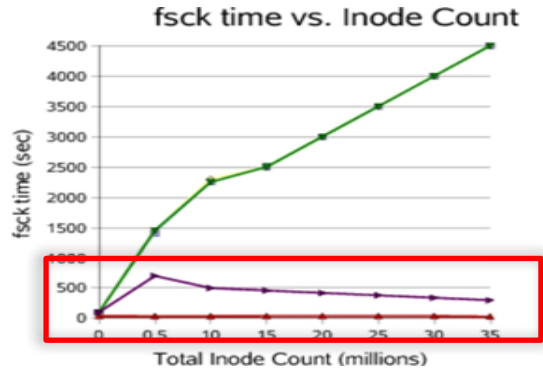
▸ If the computer is shut down abruptly, the file system might remain be inconsistent.

▸ In order to repair the file system, all metadata has to be reviewed:

  ▸ The required time depends of the file system size (all the metadata has to be reviewed, the more metadata to be reviewed the more time is needed).

Storage Networks Explained: Basics and Application of Fibre Channel SAN, NAS… (Página 139)

ARCOS @ UC3M

Alejandro Calderón Mateos

# With *Journaling*



3. Files "punch out" after work is done

1. Files "punch in" for disk write

**Journal**

2. Files do their work



fsck time vs. Inode Count

▸ The file system writes the changes in a log before changing the file.

▸ If the computer is shut down abruptly, the file system checks has to review the log for the pending changes, and do these changes (commit):

  ▸ The time needed depends of the number of pending changes in the log, and does not depend on the file system size.
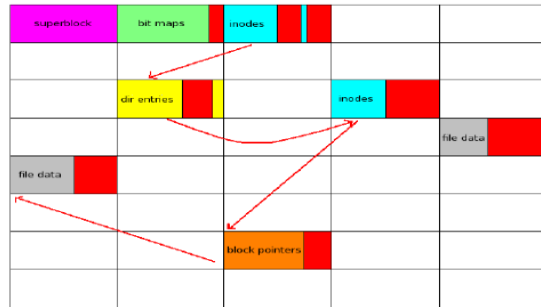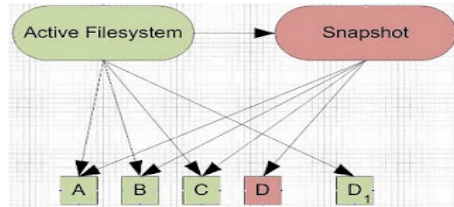
  ▸ From hours to seconds…

http://www.howtogeek.com/howto/33552/htg-explains-which-linux-file-system-should-you-choose/

ARCOS @ UC3M
Alejandro Calderón Mateos

# Advanced features



- Journaling

- Snapshots

- Dynamic file system expansion

ARCOS @ UC3M

Alejandro Calderón Mateos

# *Snapshot*





▸ A Snapshot represents the state of the file system at a point of time:

   ▸ In a few seconds is done.

   ▸ It is possible to access to all the file system snapshots on this disk.
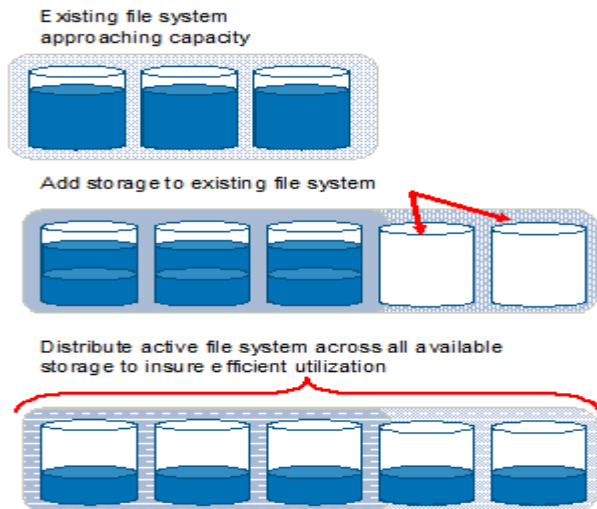
▸ E.g.: system updates, backups, etc.

ARCOS @ UC3M
Alejandro Calderón Mateos

# Advanced features



▸ Journaling

▸ Snapshots

▸ Dynamic file system expansion

ARCOS @ UC3M
Alejandro Calderón Mateos

# Dynamic file system expansion



Existing file system approaching capacity

Add storage to existing file system

Distribute active file system across all available storage to insure efficient utilization

▸ It is important to design the file system in a way that it could be resized (add more space, remove space, etc.) without losing information.

  ▸ Dynamic and flexible structures

  ▸ Metadata is distributed along the disk

http://www.storagenetworks.com/documents/productdatasheets/x9000_nas_proposal.htm

ARCOS @ UC3M
Alejandro Calderón Mateos

ARCOS Group

Universidad Carlos III de Madrid

# Lesson 5 (b)

## File systems

Operating System Design

Bachelor in Informatics Engineering