

# Modulo V

## Sistema de Arquivos

Prof. Ismael H F Santos

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

1

## Ementa

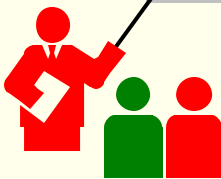
- File-System Interface
  - File Concept
  - Directory Structure
  - File Sharing
  - Protection
- File-System Structure
  - File-System Implementation
  - Allocation Methods
  - Efficiency and Performance
  - NFS – Network File System / Sun Microsystems

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

2

## SOP – CO023



Interface do  
Sistema de  
Arquivos

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

3

## Objectives

- To explain the function of file systems
- To describe the interfaces to file systems
- To discuss file-system design tradeoffs, including access methods, file sharing, file locking, and directory structures
- To explore file-system protection

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

4

## SOP – CO023



Conceito  
De  
Arquivo

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

5

## File Concept

- Contiguous logical address space
- Types:
  - Data
    - numeric
    - character
    - binary
  - Program

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

6

## File Structure

- None - sequence of words, bytes
- Simple record structure
  - Lines
  - Fixed length
  - Variable length
- Complex Structures
  - Formatted document
  - Relocatable load file
- Can simulate last two with first method by inserting appropriate control characters
- Who decides:
  - Operating system
  - Program

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

7

## File Attributes

- **Name** – only information kept in human-readable form
- **Identifier** – unique tag (number) identifies file within file system
- **Type** – needed for systems that support different types
- **Location** – pointer to file location on device
- **Size** – current file size

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

8

## File Attributes (cont.)

- **Protection** – controls who can do reading, writing, executing
- **Time, date, and user identification** – data for protection, security, and usage monitoring
- Information about files are kept in the directory structure, which is maintained on the disk

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

9

## File Operations

- File is an **abstract data type**
- **Create**
- **Write**
- **Read**
- **Reposition within file**
- **Delete**
- **Truncate**
- **Open( $F_i$ )** – search the directory structure on disk for entry  $F_i$  and move the content of entry to memory
- **Close ( $F_i$ )** – move the content of entry  $F_i$  in memory to directory structure on disk

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

10

## Open Files

- Several pieces of data are needed to manage open files:
  - **File pointer**: pointer to last read/write location, per process that has the file open
  - **File-open count**: counter of number of times a file is open – to allow removal of data from open-file table when last processes closes it
  - Disk location of the file: cache of data access information
  - Access rights: per-process access mode information

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

11

## Open File Locking

- Provided by some operating systems and file systems
- Mediates access to a file
- **Mandatory or advisory**:
  - **Mandatory** – access is denied depending on locks held and requested
  - **Advisory** – processes can find status of locks and decide what to do

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

12

## File Locking Example – Java API

```
import java.io.*;
import java.nio.channels.*;
public class LockingExample {
    public static final boolean EXCLUSIVE = false;
    public static final boolean SHARED = true;
    public static void main(String args[]) throws IOException {
        FileLock sharedLock = null;
        FileLock exclusiveLock = null;
        try {
            RandomAccessFile raf = new RandomAccessFile("file.txt", "rw");
            // get the channel for the file
            FileChannel ch = raf.getChannel();
            // this locks the first half of the file - exclusive
            exclusiveLock = ch.lock(0, raf.length()/2, EXCLUSIVE);
            /** Now modify the data . . . */
            // release the lock
            exclusiveLock.release();
        }
    }
}
```

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

13

## File Locking Example – Java API (cont)

```
// this locks the second half of the file - shared
sharedLock = ch.lock(raf.length()/2+1, raf.length(), SHARED);
/** Now read the data . . . */
// release the lock
exclusiveLock.release();
} catch (java.io.IOException ioe) {
    System.err.println(ioe);
} finally {
    if (exclusiveLock != null)
        exclusiveLock.release();
    if (sharedLock != null)
        sharedLock.release();
}
}
```

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

14

## File Types – Name, Extension

| file type      | usual extension          | function  |
|----------------|--------------------------|---|
| executable     | exe, com, bin or none    | ready-to-run machine-language program   |
| object         | obj, o                   | compiled machine language, not linked   |
| source code    | c, cc, java, pas, asm, a | source code in various languages  |
| batch          | bat, sh                  | commands to the command interpreter   |
| text           | txt, doc                 | textual data, documents   |
| word processor | wp, tex, rtf, doc        | various word-processor formats  |
| library        | lib, a, so, dll          | libraries of routines for programmers   |
| print or view  | ps, pdf, jpg             | ASCII or binary file in a format for printing or viewing                            |
| archive        | arc, zip, tar            | related files grouped into one file, sometimes compressed, for archiving or storage |
| multimedia     | mpeg, mov, rm, mp3, avi  | binary file containing audio or A/V information                                     |

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

15

## Access Methods

### Sequential Access

read next  
write next  
reset  
no read after last write (rewrite)

### Direct Access

read *n*  
write *n*  
position to *n*  
read next  
write next  
rewrite *n*

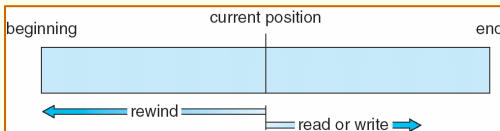
*n* = relative block number

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

16

## Sequential-access File



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

17

## Simulation of Sequential Access on a Direct-access File

| sequential access | implementation for direct access                |
|-------------------|---|
| reset             | <i>cp</i> = 0;                                  |
| read next         | read <i>cp</i> ;<br><i>cp</i> = <i>cp</i> + 1;  |
| write next        | write <i>cp</i> ;<br><i>cp</i> = <i>cp</i> + 1; |

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

18

## SOP – CO023

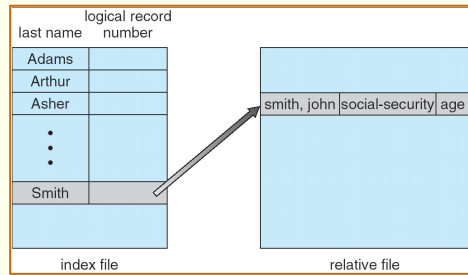


April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

19

## Example of Index and Relative Files



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

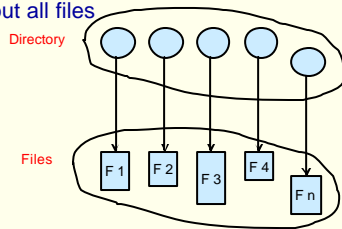
20

## Directory Structure

- A collection of nodes containing information about all files

Both the directory structure and the files reside on disk.

Backups of these two structures are kept on tapes

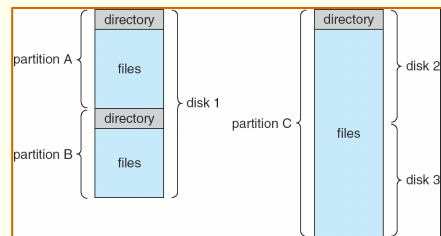


April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

21

## A Typical File-system Organization



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

22

## Operations Performed on Directory

- Search for a file
- Create a file
- Delete a file
- List a directory
- Rename a file
- Traverse the file system

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

23

## Organize the Directory (Logically) to Obtain

- **Efficiency** – locating a file quickly
- **Naming** – convenient to users
  - Two users can have same name for different files
  - The same file can have several different names
- **Grouping** – logical grouping of files by properties, (e.g., all Java programs, all games, ...)

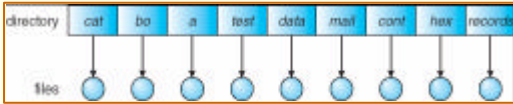
April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

24

## Single-Level Directory

- A single directory for all users



Naming problem

Grouping problem

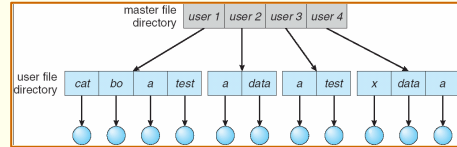
April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

25

## Two-Level Directory

- Separate directory for each user



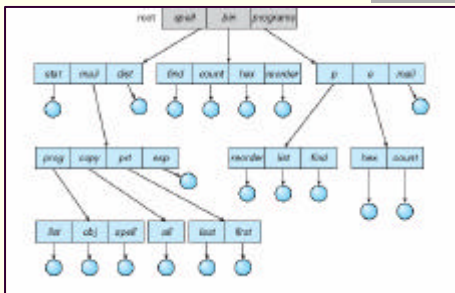
- Path name
- Can have the same file name for different user
- Efficient searching
- No grouping capability

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

26

## Tree-Structured Directories



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

27

## Tree-Structured Directories (Cont)

- Efficient searching
- Grouping Capability
- Current directory (working directory)
  - cd /spell/mail/prog
  - type list

April 05

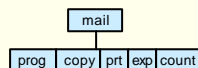
Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

28

## Tree-Structured Directories (Cont)

- **Absolute** or **relative** path name
- Creating a new file is done in current directory
- Delete a file: **rm <file-name>**
- Creating a new subdirectory is done in current directory: **mkdir <dir-name>**

Example: if in current directory **/mail**  
**mkdir count**



Deleting "mail" ⇒ deleting the entire subtree rooted by "mail"

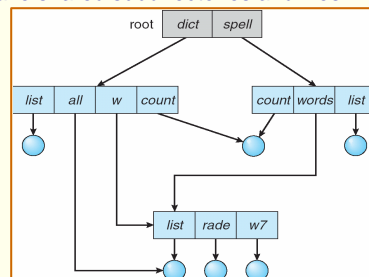
April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

29

## Acyclic-Graph Directories

- Have shared subdirectories and files



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

30

## Acyclic-Graph Directories (Cont.)

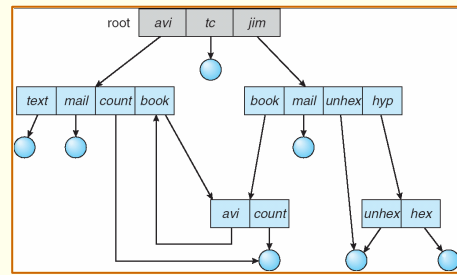
- Two different names (aliasing)
  - If *dict* deletes *list* ⇒ dangling pointer
- Solutions:
- Backpointers, so we can delete all pointers  
Variable size records a problem
  - Backpointers using a daisy chain organization
  - Entry -hold-count solution
- New directory entry type
    - **Link** – another name (pointer) to an existing file
    - **Resolve the link** – follow pointer to locate the file

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

31

## General Graph Directory



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

32

## General Graph Directory (Cont.)

- How do we guarantee no cycles?
  - Allow only links to file not subdirectories
  - Garbage collection
  - Every time a new link is added use a cycle detection algorithm to determine whether it is OK

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

33

## File System Mounting

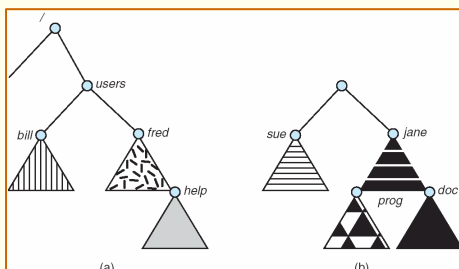
- A file system must be **mounted** before it can be accessed
- A unmounted file system (i.e. Fig. 11-11(b)) is mounted at a **mount point**

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

34

## (a) Existing. (b) Unmounted Partition



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

35

## SOP – CO023

Compartilhamento  
De  
Arquivos

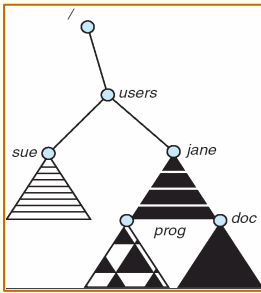


April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

36

## Mount Point



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

37

## File Sharing

- Sharing of files on multi-user systems is desirable
- Sharing may be done through a **protection** scheme
- On distributed systems, files may be shared across a network
- **Network File System (NFS)** is a common distributed file-sharing method

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

38

## File Sharing – Multiple Users

- **User IDs** identify users, allowing permissions and protections to be per-user
- **Group IDs** allow users to be in groups, permitting group access rights

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

39

## File Sharing – Remote File Systems

- Uses networking to allow file system access between systems
  - Manually via programs like FTP
  - Automatically, seamlessly using **distributed file systems**
  - Semi automatically via the **world wide web**
- **Client-server** model allows clients to mount remote file systems from servers
  - Server can serve multiple clients
  - Client and user-on-client identification is insecure or complicated

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

40

## File Sharing – Remote File Systems

- **NFS** is standard UNIX client-server file sharing protocol
- **CIFS** is standard Windows protocol
- Standard operating system file calls are translated into remote calls
- Distributed Information Systems (**distributed naming services**) such as **LDAP**, **DNS**, **NIS**, **Active Directory** implement unified access to information needed for remote computing

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

41

## File Sharing – Failure Modes

- Remote file systems add **new failure modes**, due to network failure, server failure
- Recovery from failure can involve state information about status of each remote request
- Stateless protocols such as NFS include all information in each request, allowing easy recovery but less security

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

42

## File Sharing – Consistency Semantics

- Consistency semantics specify how multiple users are to access a shared file simultaneously
  - Similar to process synchronization algorithms
    - Tend to be less complex due to disk I/O and network latency (for remote file systems)
  - Andrew File System (AFS) implemented complex remote file sharing semantics

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

43

## File Sharing – Consistency Semantics

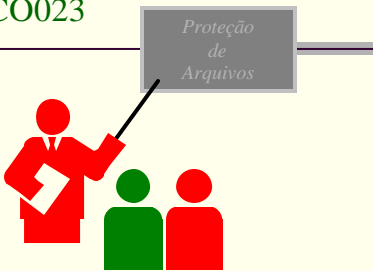
- Unix file system (UFS) implements:
  - Writes to an open file visible immediately to other users of the same open file
  - Sharing file pointer to allow multiple users to read and write concurrently
- AFS has session semantics
  - Writes only visible to sessions starting after the file is closed

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

44

## SOP – CO023



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

45

## Protection

- File owner/creator should be able to control:
  - what can be done
  - by whom
- Types of access
  - Read
  - Write
  - Execute
  - Append
  - Delete
  - List

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

46

## Access Lists and Groups

- Mode of access: read, write, execute
- Three classes of users
 

|    |                      |   |            |
|----|----------------------|---|------------|
|    |                      |   | <b>RWX</b> |
| a) | <b>owner access</b>  | 7 | ⇒ 111      |
| b) | <b>group access</b>  | 6 | ⇒ 110      |
| c) | <b>public access</b> | 1 | ⇒ 001      |
- Ask manager to create a group (unique name), say **Group1**, and add some users to the group.
- For a particular file (say **game**) or subdirectory, define an appropriate access.
 

```

                owner  group  public
                 |    |    |
                chmod 761 game
                chgrp Group1 game
                Attach a group to a file
            
```

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

47

## A Sample UNIX Directory Listing

```

-rw-rw-r-- 1 pbg staff 31200 Sep 3 08:30 intro.ps
drwx----- 5 pbg staff 512 Jul 8 09:33 private/
drwxrwxr-x 2 pbg staff 512 Jul 8 09:35 doc/
drwxrwx--- 2 pbg student 512 Aug 3 14:13 student-proj/
-rw-r--r-- 1 pbg staff 9423 Feb 24 2003 program.c
-rwxr-xr-x 1 pbg staff 20471 Feb 24 2003 program
drwx--x--x 4 pbg faculty 512 Jul 31 10:31 lib/
drwx----- 3 pbg staff 1024 Aug 29 06:52 mail/
drwxrwxrwx 3 pbg staff 512 Jul 8 09:35 test/
    
```

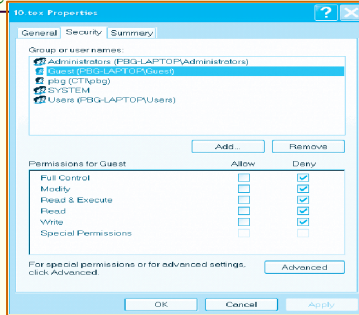
April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

48



## Windows XP Access-control List Management



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

49

## SOP – CO023

Implementação  
Sistema de  
Arquivos



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

50

## Objectives

- To describe the details of implementing local file systems and directory structures
- To describe the implementation of remote file systems
- To discuss block allocation and free-block algorithms and trade-offs

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

51

## SOP – CO023

Estrutura do  
Sistema de  
Arquivos



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

52

## File-System Structure

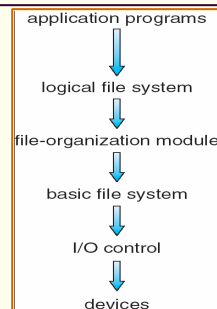
- File structure
  - Logical storage unit
  - Collection of related information
- File system resides on secondary storage (disks)
- File system organized into layers
- **File control block** – storage structure consisting of information about a file

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

53

## Layered File System



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

54

## A Typical File Control Block

|  |
|--|
| file permissions                                 |
| file dates (create, access, write)               |
| file owner, group, ACL                           |
| file size  |
| file data blocks or pointers to file data blocks |

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

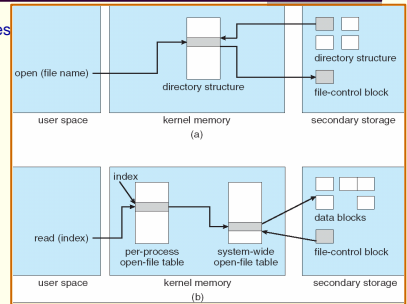
55

## In-Memory File System Structures

- The following figure illustrates the necessary file system structures provided by the operating systems.

- Figure (a) refers to opening a file.

- Figure (b) refers to reading a file.



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

56

## Virtual File Systems

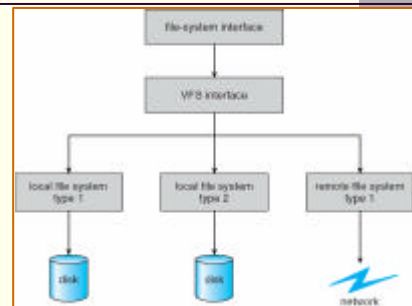
- Virtual File Systems (VFS)** provide an object-oriented way of implementing file systems.
- VFS allows the same system call interface (the API) to be used for different types of file systems.
- The API is to the VFS interface, rather than any specific type of file system.

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

57

## Schematic View of Virtual File System



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

58

## SOP – CO023



April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

59

## Directory Implementation

- Linear list** of file names with pointer to the data blocks.
  - simple to program
  - time-consuming to execute
- Hash Table** – linear list with hash data structure.
  - decreases directory search time
  - collisions** – situations where two file names hash to the same location
  - fixed size

April 05

Prof. Ismael H. F. Santos - ismael@tecgraf.puc-rio.br

60

## Allocation Methods

- An allocation method refers to how disk blocks are allocated for files:
- **Contiguous allocation**
- **Linked allocation**
- **Indexed allocation**

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

61

## Contiguous Allocation of Disk Space

- Each file occupies a set of contiguous blocks on the disk
- **Simple** – only starting location (block #) and length (number of blocks) are required
- **Random access**
- **Wasteful of space** (dynamic storage-allocation problem).
- **Files cannot grow !**

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

62

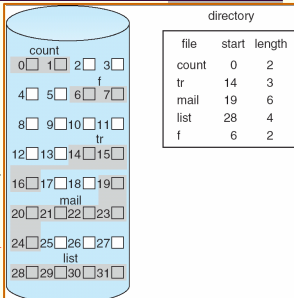
## Contiguous Allocation of Disk Space

- Mapping from logical to physical



Block to be accessed =  $Q +$   
starting address

Displacement into block =  $R$



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

63

## Extent-Based Systems

- Many newer file systems (i.e. Veritas File System) use a modified contiguous allocation scheme
- **Extent-based file systems allocate disk blocks in extents**
- **An extent is a contiguous block of disks**
  - Extents are allocated for file allocation
  - A file consists of one or more extents.

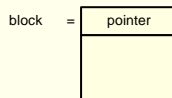
April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

64

## Linked Allocation

- Each file is a linked list of disk blocks: blocks may be scattered anywhere on the disk.



- **Simple** – need only starting address
- **Free-space management system** – no waste of space
- **No random access**
- **Mapping**

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

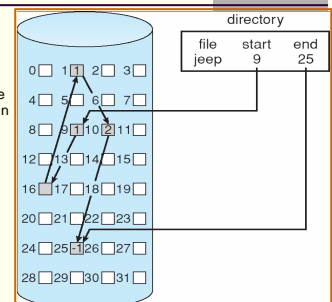
65

## Linked Allocation



Block to be accessed is the  $Q$ th block in the linked chain of blocks representing the file.

Displacement into block =  $R + 1$



April 05

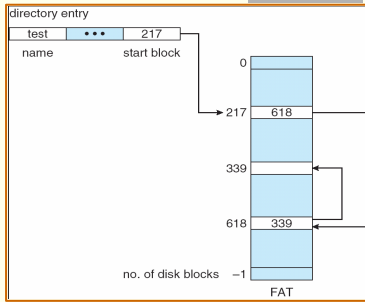
Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

66

# File-Allocation Table (FAT)

File-allocation table

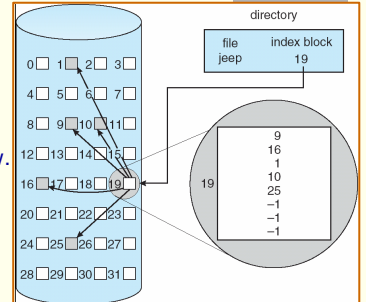
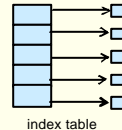
disk-space allocation used by MS-DOS and OS/2.



# Indexed Allocation

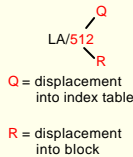
- Brings all pointers together into the **index block**.

- Logical view.



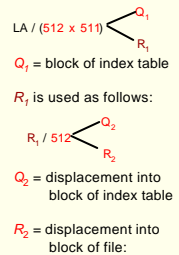
# Indexed Allocation (Cont.)

- Need index table
- Random access
- Dynamic access without external fragmentation, but have overhead of index block.
- Mapping from logical to physical in a file of **maximum size of 256K words and block size of 512 words**. We need only 1 block for index table.



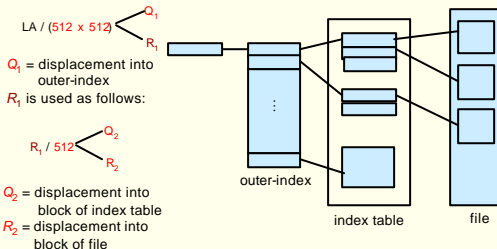
# Indexed Allocation – Mapping (Cont.)

- Mapping from logical to physical in a file of unbounded length (block size of 512 words).
- Linked scheme – Link blocks of index table (no limit on size).

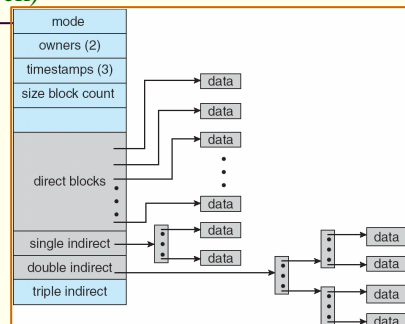


# Indexed Allocation – Mapping (Cont.)

- Two-level index (maximum file size is  $512^3$ )

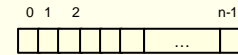


# Combined Scheme: UNIX (4K bytes per block)



## Free-Space Management

- Bit vector ( $n$  blocks)



$$\text{bit}[j] = \begin{cases} 0 \Rightarrow \text{block}[j] \text{ free} \\ 1 \Rightarrow \text{block}[j] \text{ occupied} \end{cases}$$

Block number calculation

(number of bits per word) \*  
(number of 0-value words) +  
offset of first 1 bit

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

73

## Free-Space Management (Cont.)

- Bit map requires extra space

- Example:

block size =  $2^{12}$  bytes

disk size =  $2^{30}$  bytes (1 gigabyte)

$n = 2^{30}/2^{12} = 2^{18}$  bits (or 32K bytes)

- Easy to get contiguous files

- Linked list (free list)

- Cannot get contiguous space easily

- No waste of space

- Grouping

- Counting

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

74

## Free-Space Management (Cont.)

- Need to protect:

- Pointer to free list

- Bit map

- Must be kept on disk

- Copy in memory and disk may differ

- Cannot allow for block[ $j$ ] to have a situation where bit[ $j$ ] = 1 in memory and bit[ $j$ ] = 0 on disk

- Solution:

- Set bit[ $j$ ] = 1 in disk

- Allocate block[ $j$ ]

- Set bit[ $j$ ] = 1 in memory

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

75

## Directory Implementation

- Linear list of file names with pointer to the data blocks

- simple to program

- time-consuming to execute

- Hash Table – linear list with hash data structure

- decreases directory search time

- **collisions** – situations where two file names hash to the same location

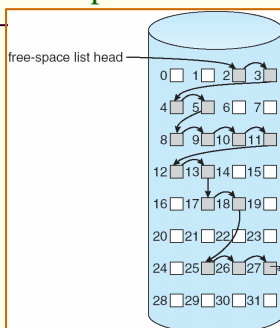
- fixed size

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

76

## Linked Free Space List on Disk



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

77

## SOP – CO023

Eficiência  
e  
Performance



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

78

## Efficiency and Performance

- **Efficiency dependent on:**
  - disk allocation and directory algorithms
  - types of data kept in file's directory entry
- **Performance**
  - **disk cache** – separate section of main memory for frequently used blocks
  - **free-behind and read-ahead** – techniques to optimize sequential access
  - improve PC performance by dedicating section of memory as **virtual disk**, or RAM disk

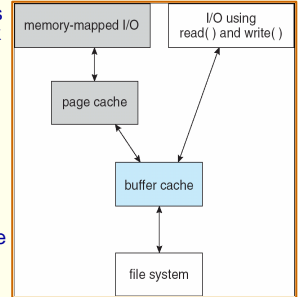
April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

79

## Page Cache

- A **page cache** caches pages rather than disk blocks using virtual memory techniques
- **Memory-mapped I/O** uses a page cache
- Routine I/O through the file system uses the buffer (disk) cache



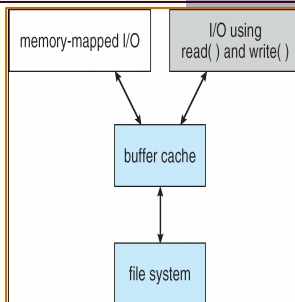
April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

80

## Unified Buffer Cache

- A **unified buffer cache** uses the same page cache to cache both memory-mapped pages and ordinary file system I/O



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

81

## Recovery

- **Consistency checking** – compares data in directory structure with data blocks on disk, and tries to fix inconsistencies
- Use system programs to **back up** data from disk to another storage device (floppy disk, magnetic tape, other magnetic disk, optical)
- **Recover** lost file or disk by **restoring** data from backup

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

82

## Log Structured File Systems

- **Log structured** (or journaling) file systems record each update to the file system as a **transaction**
- All transactions are written to a **log**
  - A transaction is considered **committed** once it is written to the log
  - However, the file system may not yet be updated
- The transactions in the log are asynchronously written to the file system
  - When the file system is modified, the transaction is removed from the log
- If the file system crashes, all remaining transactions in the log must still be performed

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

83

## The Sun Network File System (NFS)

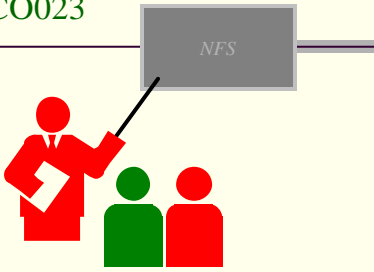
- An implementation and a specification of a software system for accessing remote files across LANs (or WANs)
- The implementation is part of the Solaris and SunOS operating systems running on Sun workstations using an unreliable datagram protocol (UDP/IP protocol and Ethernet)

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

84

## SOP – CO023



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

85

## NFS (Cont.)

- Interconnected workstations viewed as a set of independent machines with independent file systems, which allows sharing among these file systems in a transparent manner
  - A remote directory is mounted over a local file system directory
    - The mounted directory looks like an integral subtree of the local file system, replacing the subtree descending from the local directory

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

86

## NFS (Cont.)

- Specification of the remote directory for the mount operation is nontransparent; the host name of the remote directory has to be provided
  - Files in the remote directory can then be accessed in a transparent manner
- Subject to access-rights accreditation, potentially any file system (or directory within a file system), can be mounted remotely on top of any local directory

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

87

## NFS (Cont.)

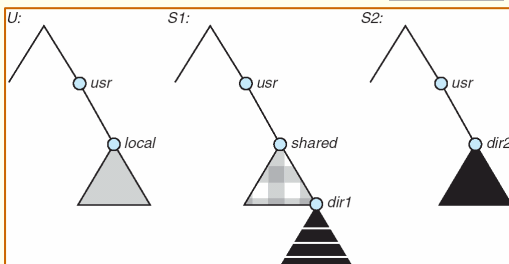
- NFS is designed to operate in a heterogeneous environment of different machines, operating systems, and network architectures; the NFS specifications independent of these media
- This independence is achieved through the use of RPC primitives built on top of an External Data Representation (XDR) protocol used between two implementation-independent interfaces
- The NFS specification distinguishes between the services provided by a mount mechanism and the actual remote-file-access services

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

88

## Three Independent File Systems

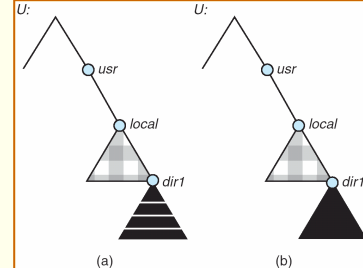


April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

89

## Mounting in NFS



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

90

## NFS Mount Protocol

- Establishes initial logical connection between server and client
- Mount operation includes name of remote directory to be mounted and name of server machine storing it
  - **Mount request** is mapped to corresponding RPC and forwarded to mount server running on server machine
  - **Export list** – specifies local file systems that server exports for mounting, along with names of machines that are permitted to mount them

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

91

## NFS Mount Protocol

- Following a mount request that conforms to its export list, the server returns a file handle – a key for further accesses
- **File handle** – a file-system identifier, and an inode number to identify the mounted directory within the exported file system
- The mount operation changes only the user's view and does not affect the server side

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

92

## NFS Protocol

- Provides a set of remote procedure calls for remote file operations. The procedures support the following operations:
  - searching for a file within a directory
  - reading a set of directory entries
  - manipulating links and directories
  - accessing file attributes
  - reading and writing files

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

93

## NFS Protocol

- NFS servers are **stateless**; each request has to provide a full set of arguments (NFS V4 is just coming available – very different, stateful)
- Modified data must be committed to the server's disk before results are returned to the client (lose advantages of caching)
- The NFS protocol does not provide concurrency-control mechanisms

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

94

## Three Major Layers of NFS Architecture

- **UNIX file-system interface** (based on the **open**, **read**, **write**, and **close** calls, and **file descriptors**)
- **Virtual File System (VFS) layer** – distinguishes local files from remote ones, and local files are further distinguished according to their file-system types

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

95

## Three Major Layers of NFS Architecture

- The VFS activates file-system-specific operations to handle local requests according to their file-system types
- Calls the NFS protocol procedures for remote requests
- **NFS service layer** – bottom layer of the architecture
  - Implements the NFS protocol

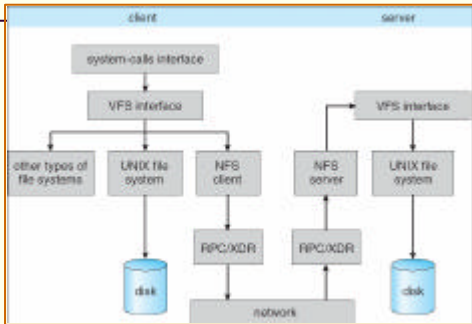
April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

96



## Schematic View of NFS Architecture



April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

97

## NFS Path-Name Translation

- Performed by breaking the path into component names and performing a separate NFS lookup call for every pair of component name and directory vnode
- To make lookup faster, a directory name lookup cache on the client's side holds the vnodes for remote directory names

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

98

## NFS Remote Operations

- Nearly one-to-one correspondence between regular UNIX system calls and the NFS protocol RPCs (except opening and closing files)
- NFS adheres to the remote-service paradigm, but employs buffering and caching techniques for the sake of performance
- **File-blocks cache** – when a file is opened, the kernel checks with the remote server whether to fetch or revalidate the cached attributes

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

99

## NFS Remote Operations

- Cached file blocks are used only if the corresponding cached attributes are up to date
- **File-attribute cache** – the attribute cache is updated whenever new attributes arrive from the server
- Clients do not free delayed-write blocks until the server confirms that the data have been written to disk

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

100

## Example: WAFL File System

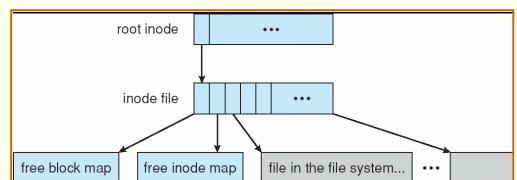
- Used on Network Appliance “Filers” – distributed file system appliances
- “Write-anywhere file layout”
- Serves up NFS, CIFS, http, ftp
- Random I/O optimized, write optimized
  - NVRAM for write caching
- Similar to **Berkeley Fast File System**, with extensive modifications

April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

101

## The WAFL File Layout

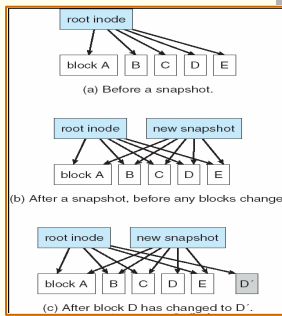


April 05

Prof. Ismael H. F. Santos - ismael@tegraf.puc-rio.br

102

## Snapshots in WAFL



April 05

103

## Permissions

- file permissions
- file dates (create, access, write)
- file owner, group, ACL
- file size
- file data blocks or pointers to file data blocks

April 05

Prof. Ismael H. F. Santos - ismael@teograf.puc-rio.br

104