

File Systems

- Files
- Directories
- Implementation
- Examples

Long-Term Information Storage: Goals

- Must store large amounts of data
- Information must survive process termination
- Multiple processes must have concurrent access to information

What a file is/is not

- File: a named collection of bytes stored on disk
 - contained within its own contiguous logical address space
- User view:
 - collection of records
- OS view:
 - bunch of blocks stored on a device
 - File system packs bytes into disk blocks on write
 - File system unpacks bytes from disk blocks on read

Why files?

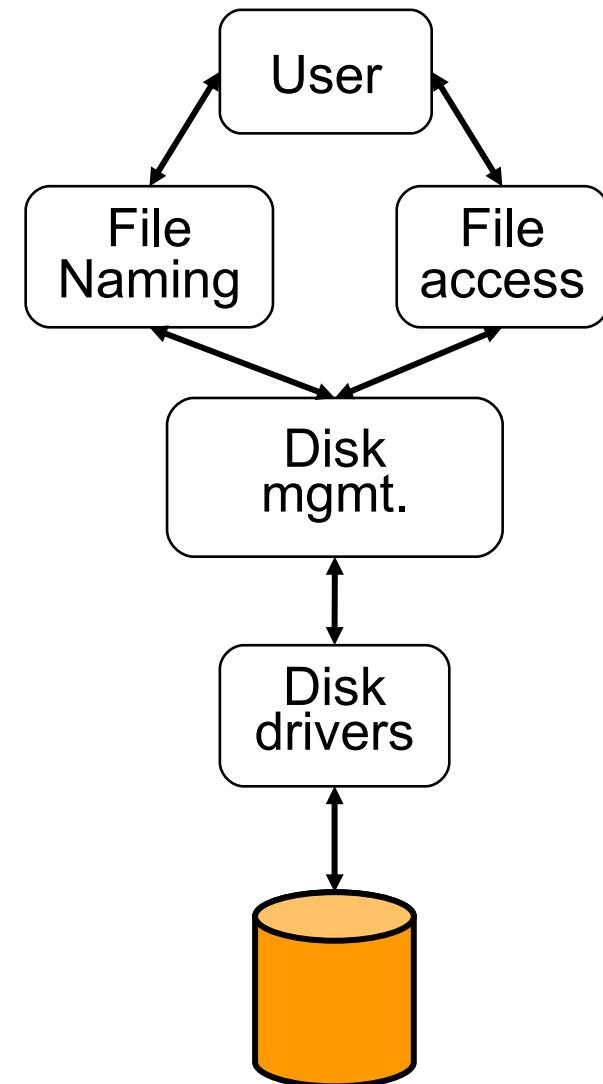
- Physical reality
 - block oriented
 - physical sector numbers
 - no protection between users
 - data corruption due to crashes, etc.
- File system abstraction
 - byte oriented
 - named files
 - users protected
 - some recovery possible after failure

Sounds simple, what's the catch?

- Files grow and shrink
 - sizes span orders of magnitude
- Little/no a-priori knowledge about contents
- Users desire efficiency
- Must overcome disk performance
 - it is awful
- Must cope with failure

High-Level Components/Abstractions

- Disk management
 - arrange collection of disk blocks into files
- Naming
 - allow user to name files, rather than give sector/track numbers, to locate data
- Security
 - keep information secure
- Reliability/durability
 - ensure crashes aren't catastrophic

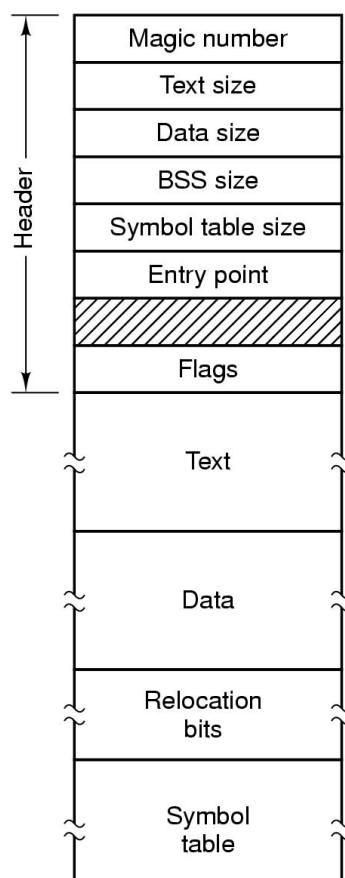


File Structure

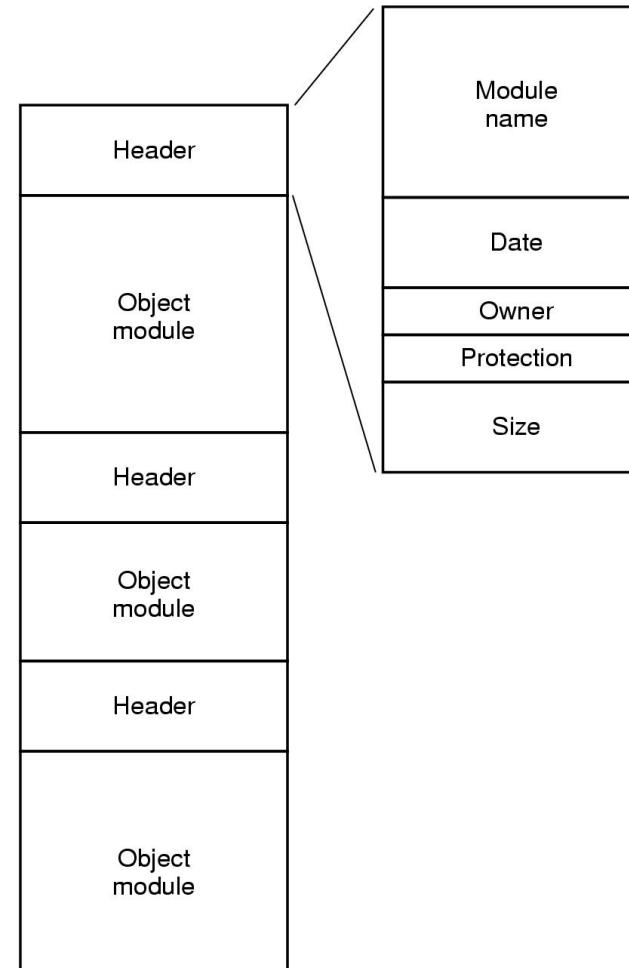
- None
 - sequence of words/bytes
- Simple
 - lines
 - fixed/variable length
- Complex
 - formatted document
 - relocatable load file, etc
- Who decides?
 - OS and programs

Examples of structured files

- Executable (a) vs archive (b)



(a)



(b)

Sampling of file types (by 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

Sampling of file attributes

Attribute	Meaning
Protection	Who can access the file and in what way
Password	Password needed to access the file
Creator	ID of the person who created the file
Owner	Current owner
Read-only flag	0 for read/write; 1 for read only
Hidden flag	0 for normal; 1 for do not display in listings
System flag	0 for normal files; 1 for system file
Archive flag	0 for has been backed up; 1 for needs to be backed up
ASCII/binary flag	0 for ASCII file; 1 for binary file
Random access flag	0 for sequential access only; 1 for random access
Temporary flag	0 for normal; 1 for delete file on process exit
Lock flags	0 for unlocked; nonzero for locked
Record length	Number of bytes in a record
Key position	Offset of the key within each record
Key length	Number of bytes in the key field
Creation time	Date and time the file was created
Time of last access	Date and time the file was last accessed
Time of last change	Date and time the file has last changed
Current size	Number of bytes in the file
Maximum size	Number of bytes the file may grow to

List of basic file operations

- Open/close
- Create/delete
- Read/write/append
- Seek
- Get/set attributes
- Rename

FS System Calls example (1)

FS System Calls example (2)

```
/* Open the input file and create the output file */
in_fd = open(argv[1], O_RDONLY); /* open the source file */
if (in_fd < 0) exit(2); /* if it cannot be opened, exit */
out_fd = creat(argv[2], OUTPUT_MODE); /* create the destination file */
if (out_fd < 0) /* if it cannot be created, exit */

/* Copy loop */
while (TRUE) {
    rd_count = read(in_fd, buffer, BUF_SIZE); /* read a block of data */
    if (rd_count <= 0) break; /* if end of file or error, exit loop */
    wt_count = write(out_fd, buffer, rd_count); /* write data */
    if (wt_count <= 0) exit(4); /* wt_count <= 0 is an error */
}

/* Close the files */
close(in_fd);
close(out_fd);
if (rd_count == 0) /* no error on last read */
    exit(0);
else
    exit(5); /* error on last read */

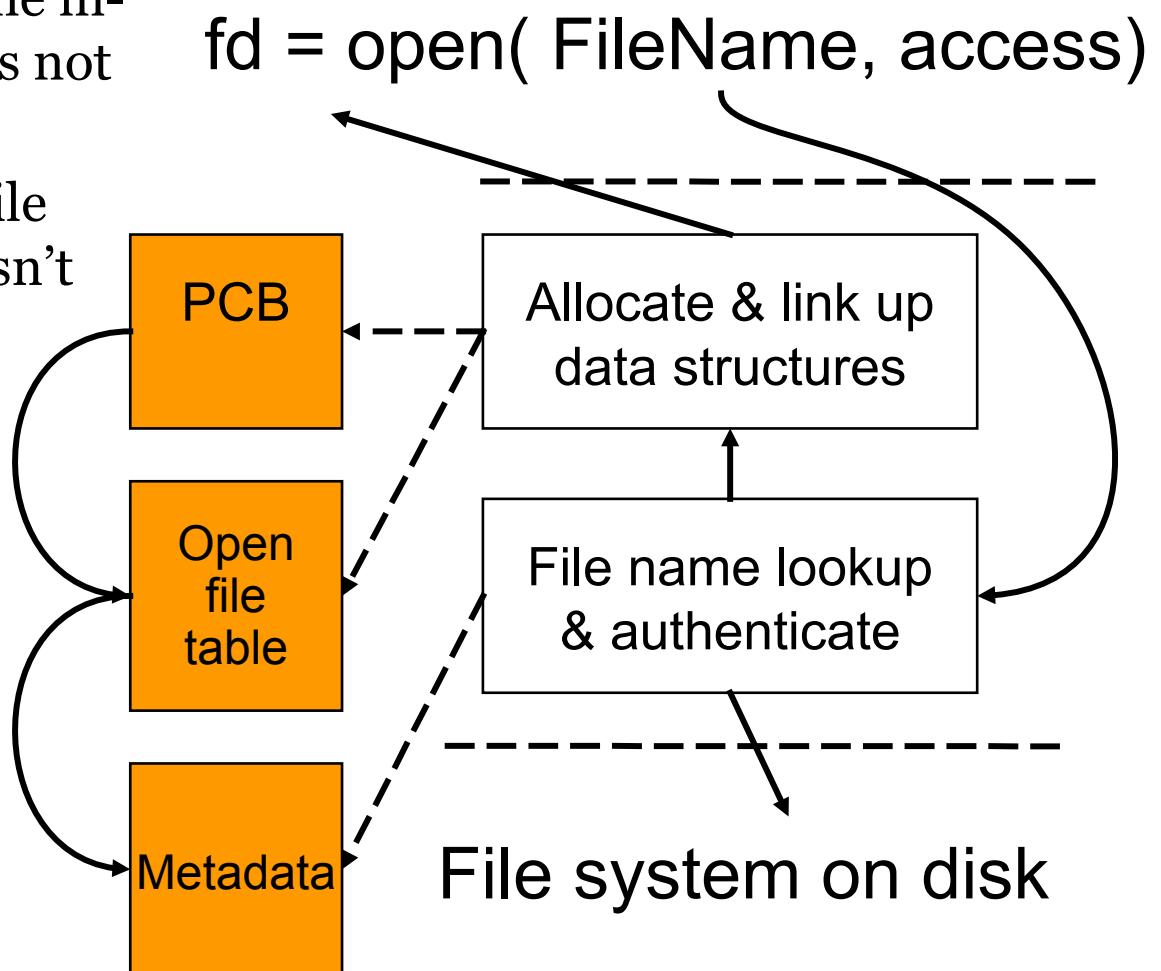
}
```

Open Files

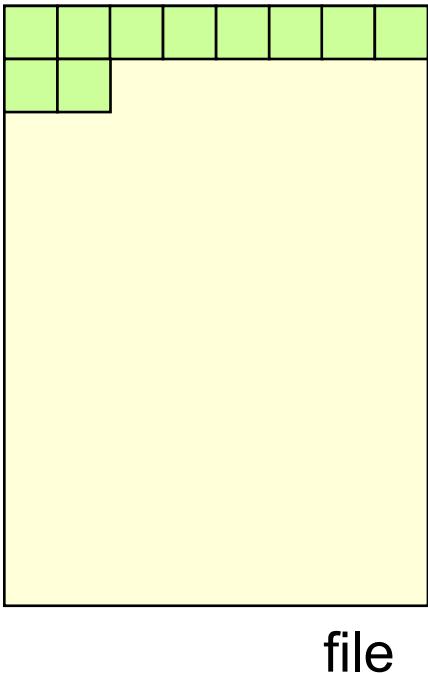
- Need several pieces of information to manage open files
 - file pointer: pointer to last read/write location, for each process which has opened the file
 - file-open count: how many times has this file been opened
 - needed to remove file from open-file table when last process closes it
 - disk location: where is the file?
 - access rights: per-process access mode information
- Some systems allow open file locking
 - mediated access to files ^ limit concurrent access

Opening A File

- File name lookup and authentication
- Copy the file metadata into the in-memory data structure, if it is not in yet
- Create an entry in the open file table (system wide) if there isn't one
- Create an entry in PCB
- Link up the data structures
- Return a pointer to user



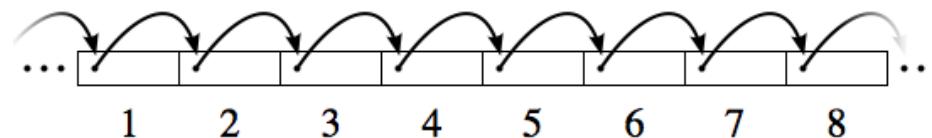
Open / Create / Write Example



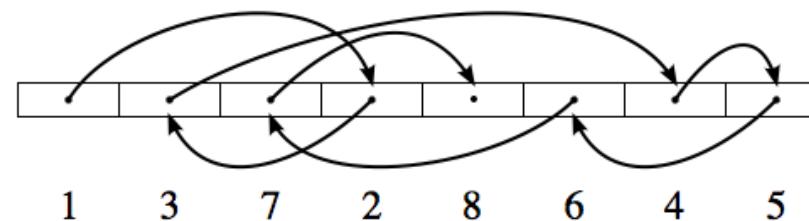
```
ThreadA () {  
    FILE* f;  
    int i=0;  
    CREAT ("file");  
    f=FOPEN ("file");  
    for (;i<10;i++) {  
        FWRITE (f,i);  
    }  
    FCLOSE (f);  
}
```

File access types

- Sequential
 - read all bytes/records from beginning
 - no jumping around



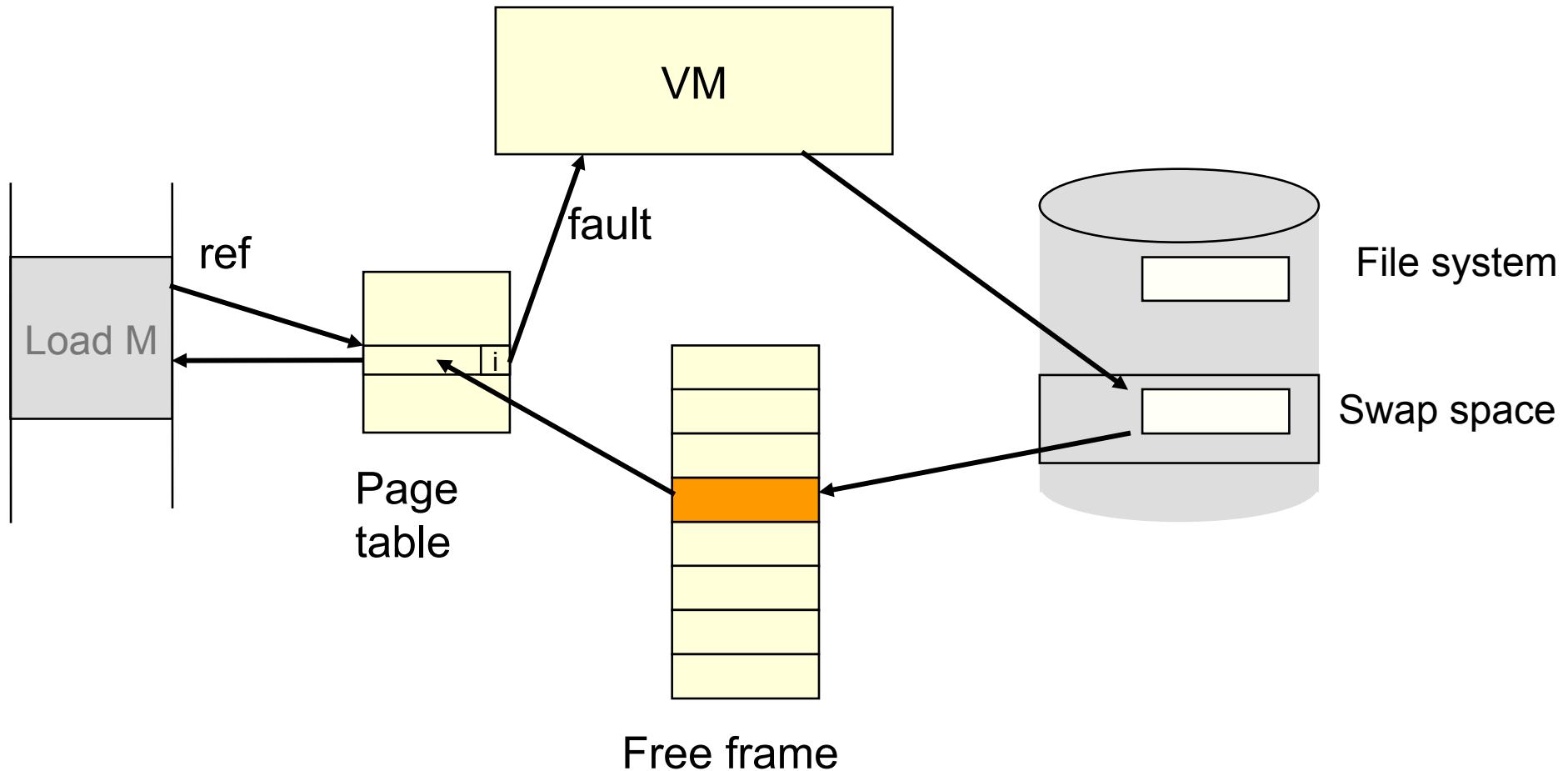
- Random
 - read bytes/records in any order



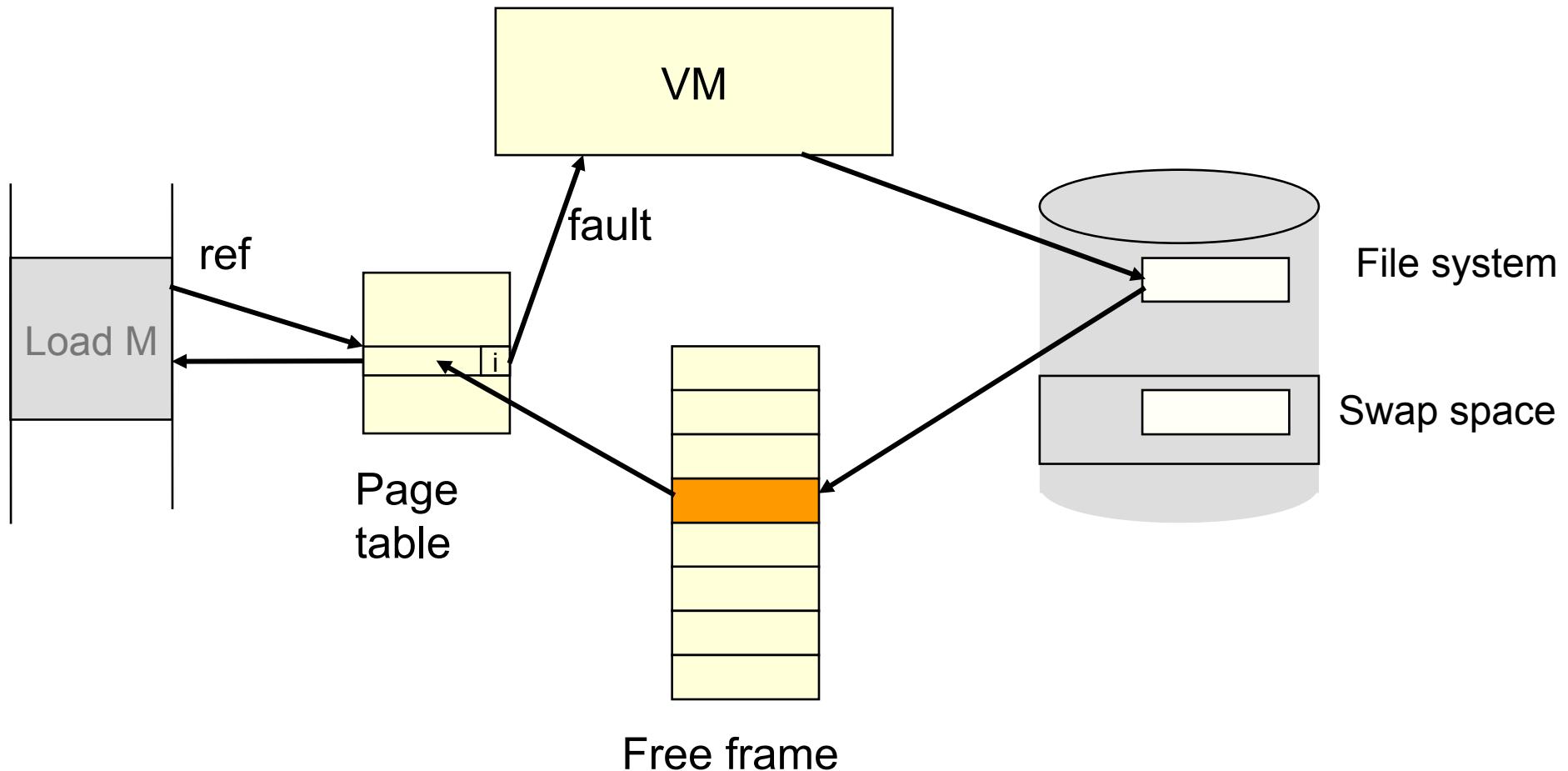
File system interfaces

- Typically, explicit read/write operations
 - file data explicitly copied between disk file and process mem.
 - programs do not directly access file data
- Memory-mapped files
 - map files into process virtual memory
 - directly access file data through regular memory access ops.

Remember Demand Paging in VM?



Memory-Mapped Files



Memory-Mapped Files: properties

- File mapped by initializing virtual memory so that the file directly serves as the backing store for a region in the process' address space
 - Hence no explicit copying
 - Hence direct access
- Processes mapping the same file share physical memory that caches file data

Memory-Mapped Files: example

- use mmap()/munmap()

```
fd = open(argv[1], O_RDWR);
fstat(fd, &sb);
addr = mmap(NULL, sb.st_size, PROT_READ | PROT_WRITE,
            MAP_SHARED, fd, 0);
int i;
for (i = 0; i < 10; i++)
    printf("%c", addr[i]);

munmap(addr, sb.st_size);
fclose(fd);
```

Memory-Mapped Files: overview

- Benefits
 - explicit file I/O instructions not needed
 - simple to program
 - with possibly different semantics
 - e.g. `read()` copies data ^ local modifications are private
 - `mmap()` does not ^ must explicitly copy to make modifications private
- Caveat
 - `mmap()` only works on page-size granularity
 - can't really append to files this way (OS doesn't know where the file ends, only that it uses a given page)

File Systems

- Files
- Directories
- Implementation
- Examples

Directory Organization

- Directory = system file
 - keeps track of other files
- Possible approaches to directories
 - each place in directory holds info about a file and pointer to address where file data is stored
 - each place in directory holds a file name and pointer to a data structure, which then holds info about the file and pointers to data location

Directory Operations

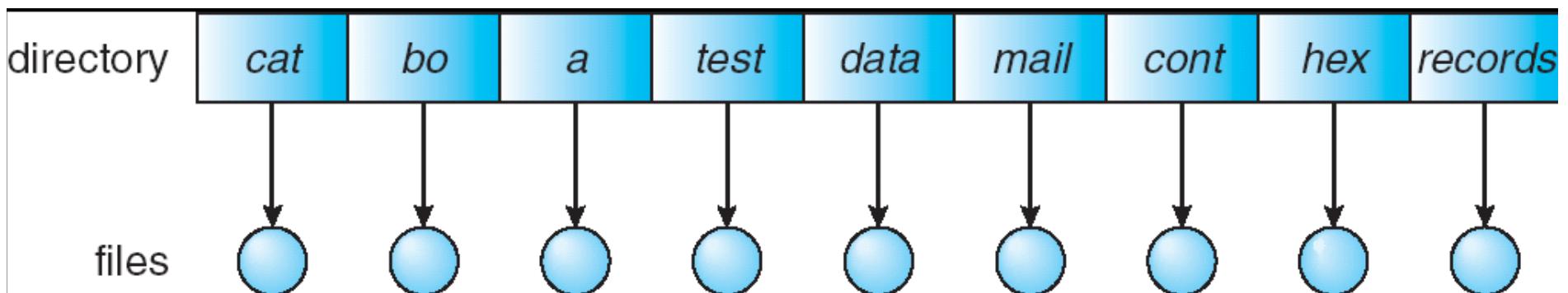
- Create/delete
- Open/close
- Read
- Rename
- Link/unlink

Directory Organization: goals

- Efficiency
 - can you locate a file quickly?
- Naming
 - user convenience
 - users can have different names for same file/same names for different file
- Grouping
 - allows for grouping of files by some properties (code, utilities, games, questionable pictures, etc.)

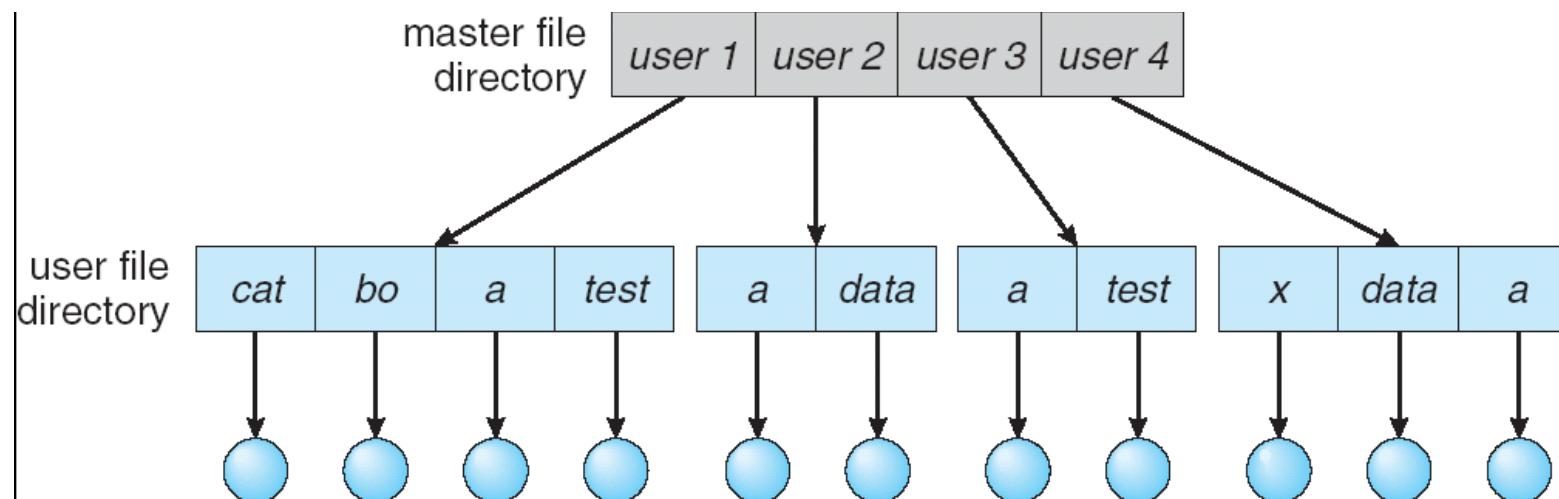
Single-Level Organization

- All users see single directory
 - naming is inconvenient
 - grouping is a problem
 - but, if there are no users, who cares?



Two-Level Organization

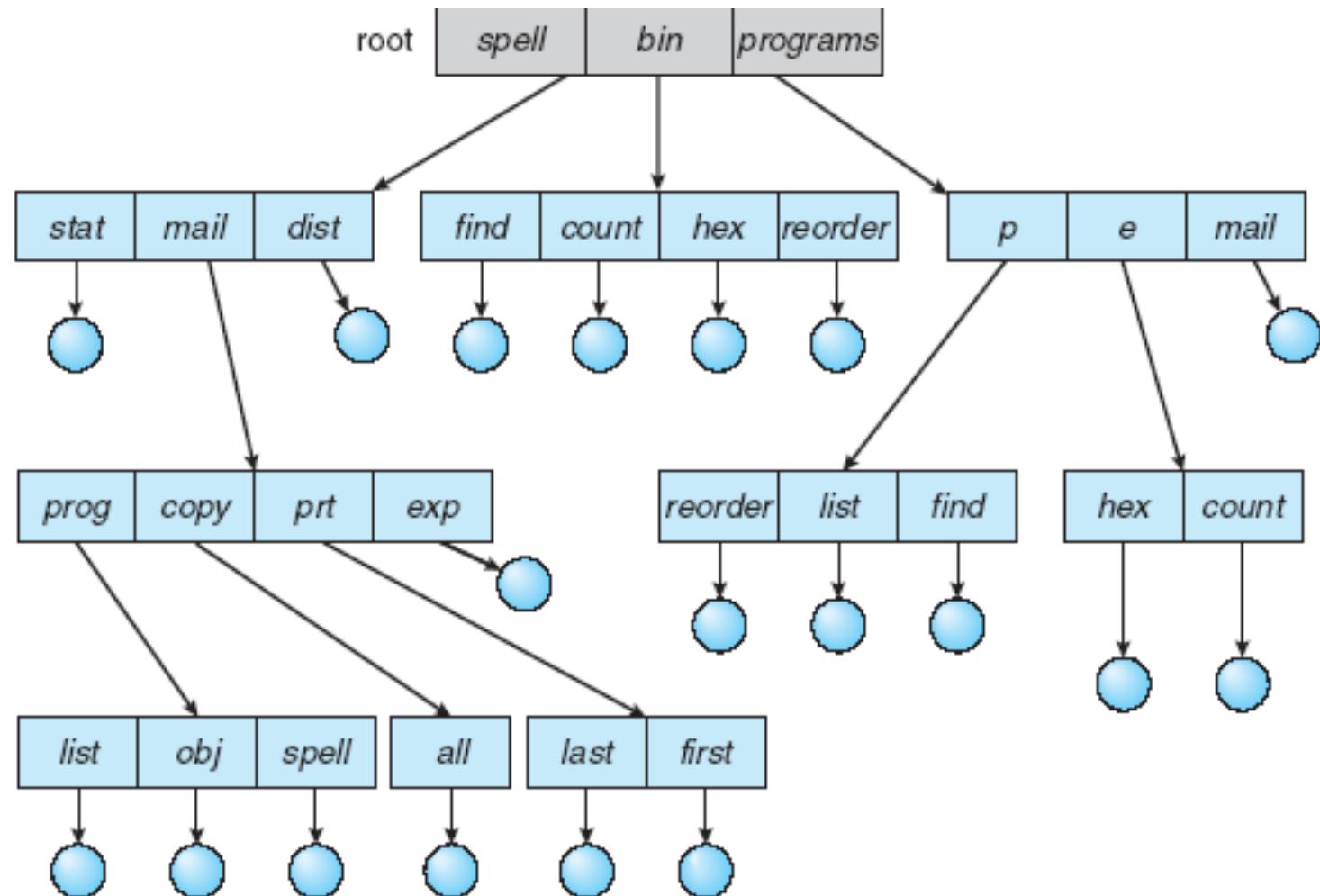
- Separate directory for each user
 - can have same/different filenames for different/same files
 - searching is (more) efficient
 - single user's files are still in one directory



Tree-Structured Organization

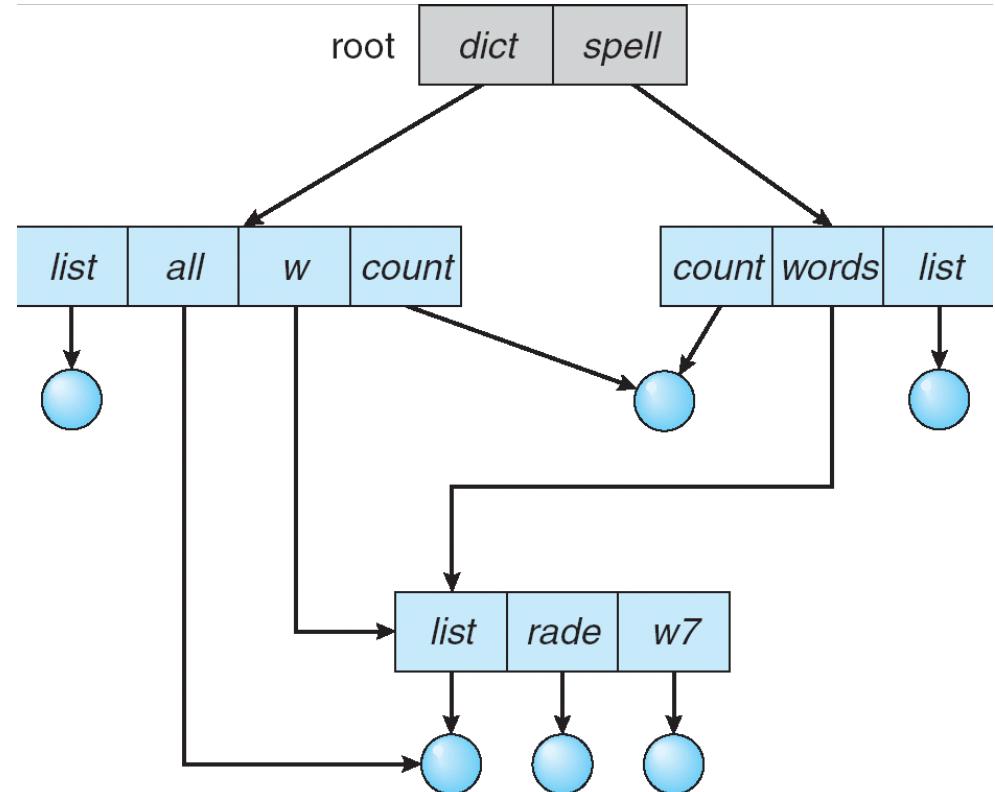
- Directories
 - contain files, directories
 - arbitrary grouping ^ flexibility
 - easily expand to hold large number of files
 - drawback: more work to find stuff
- File naming looks like: /dir1/dir2/dir3/filename
 - that's an absolute pathname... describes exact location starting at root
 - relative pathname:
 - if you are “in” dir2, then dir3/filename is location of file
 - note no leading ‘/’

Tree-Structured Organization: example



Linking

- Link: another name/pointer to existing file
- Resolving a link: following pointer to locate the file
- Linking
 - saves space
 - makes updates easier
- Hard links
 - links file into a directory structure
- Soft links
 - just a pointer to a file (e.g. shortcut in windows)



Implementing directories

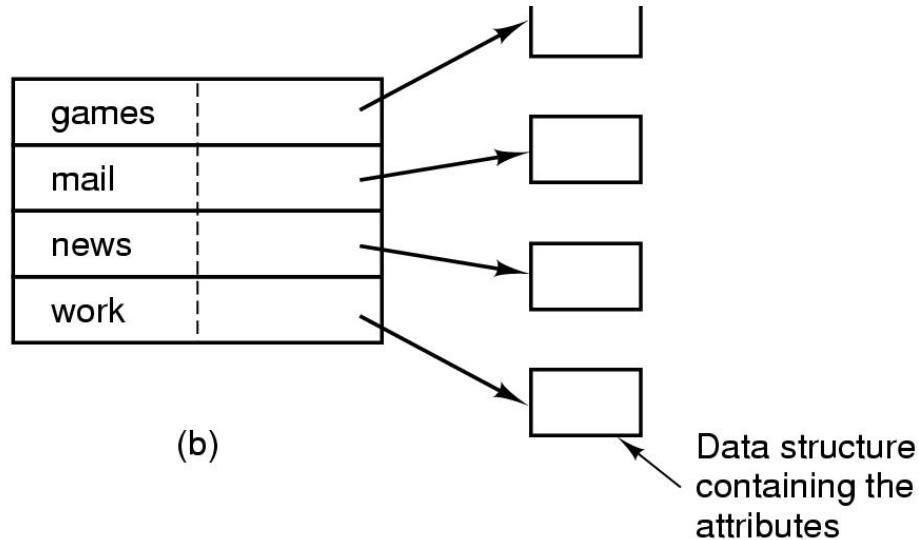
- Linux ^ specially formatted file
 - can read it just like a file
 - it's a directory only because we treat it that way
- Contains names of files
 - <name, fd_index> pairs in no particular order
 - fd_index may be another directory
 - root has no name

Possible approaches redux

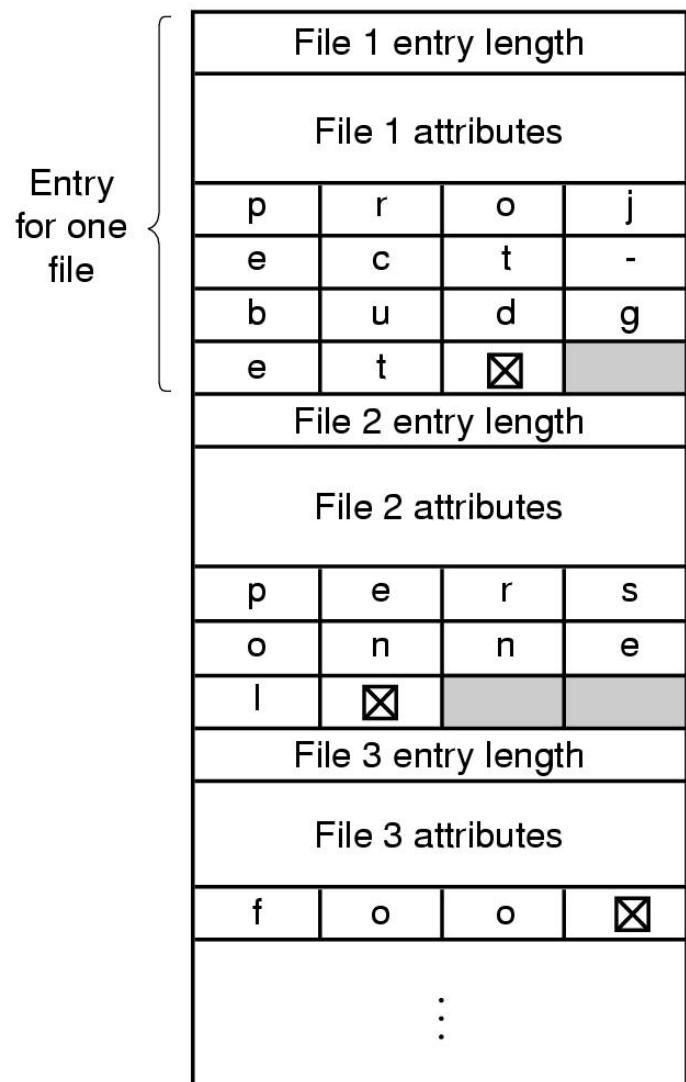
- Possible approaches to directories
 - each place in directory holds info about a file and pointer to address where file data is stored
 - each place in directory holds a file name and pointer to a data structure, which then holds info about the file and pointers to data location

games	attributes
mail	attributes
news	attributes
work	attributes

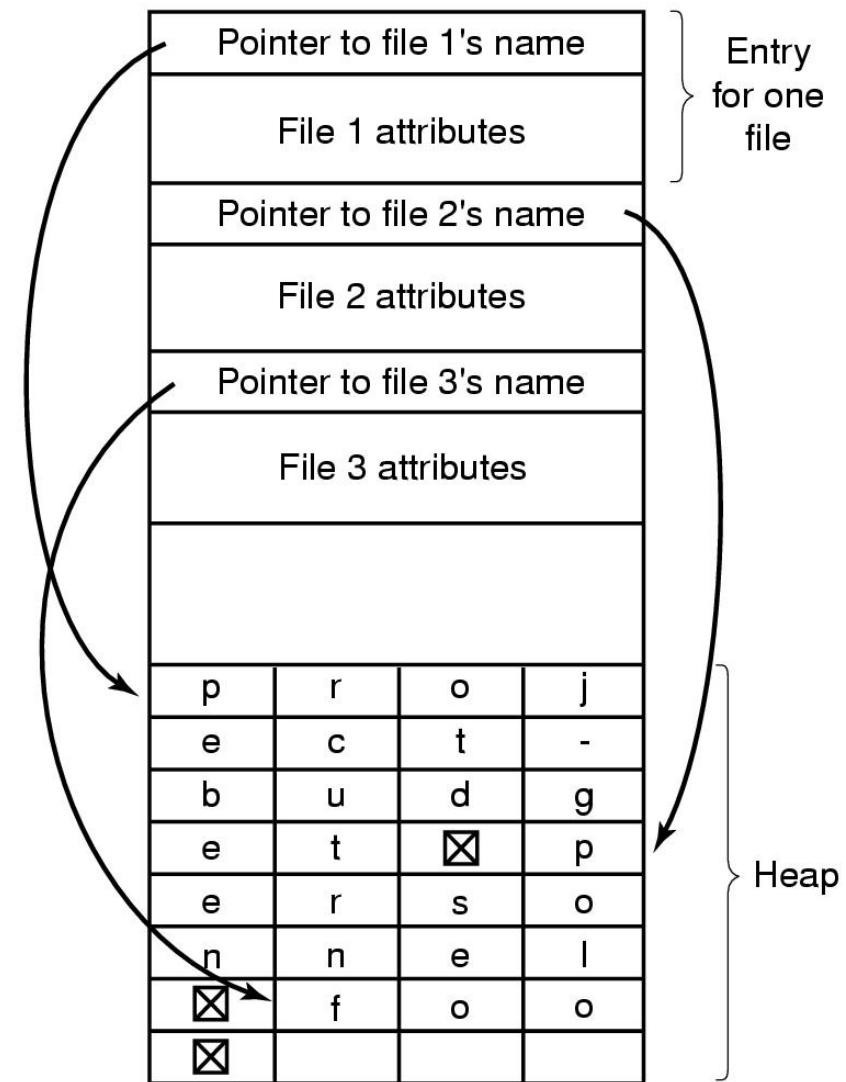
(a)



Filename handling



(a)



(b)

So, how might this work...

- Finding '/dir1/dir2/dir3/readme.txt':
 - Fetch root inode
 - Start loading root directory data blocks
 - Walk directory data until you find dir1
 - Get inode # for dir1 from directory file
 - Fetch dir1's inode
 - Start loading dir1's directory data blocks
 - Walk directory data until you find dir2
 - Get inode # for dir2 from directory file
 - ...
 - Repeat process until you have inode # for file

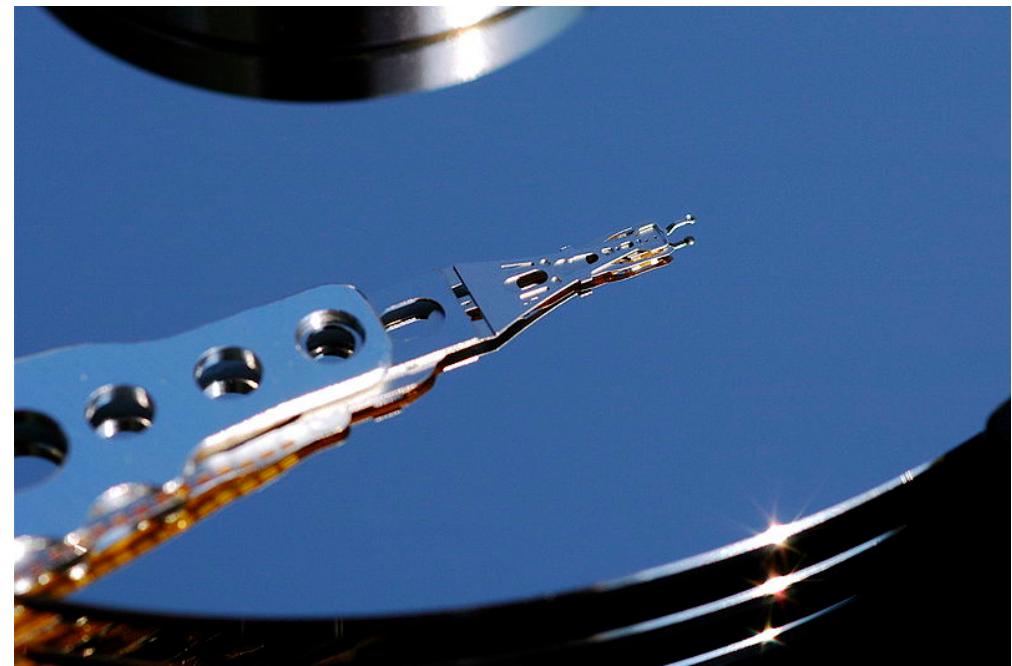
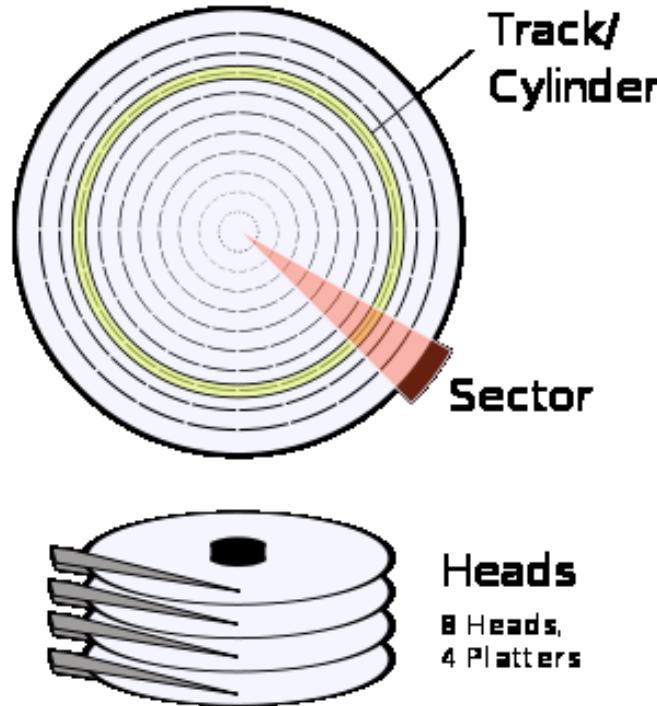
Observations

- Name lookup is slow
 - lots of CPU time
 - lots of fetches (inodes + data blocks)
 - but... also lots of repetition
- So
 - cache name lookups
 - cache recently used inodes and directory blocks
 - reserve some kernel memory for this

File Systems

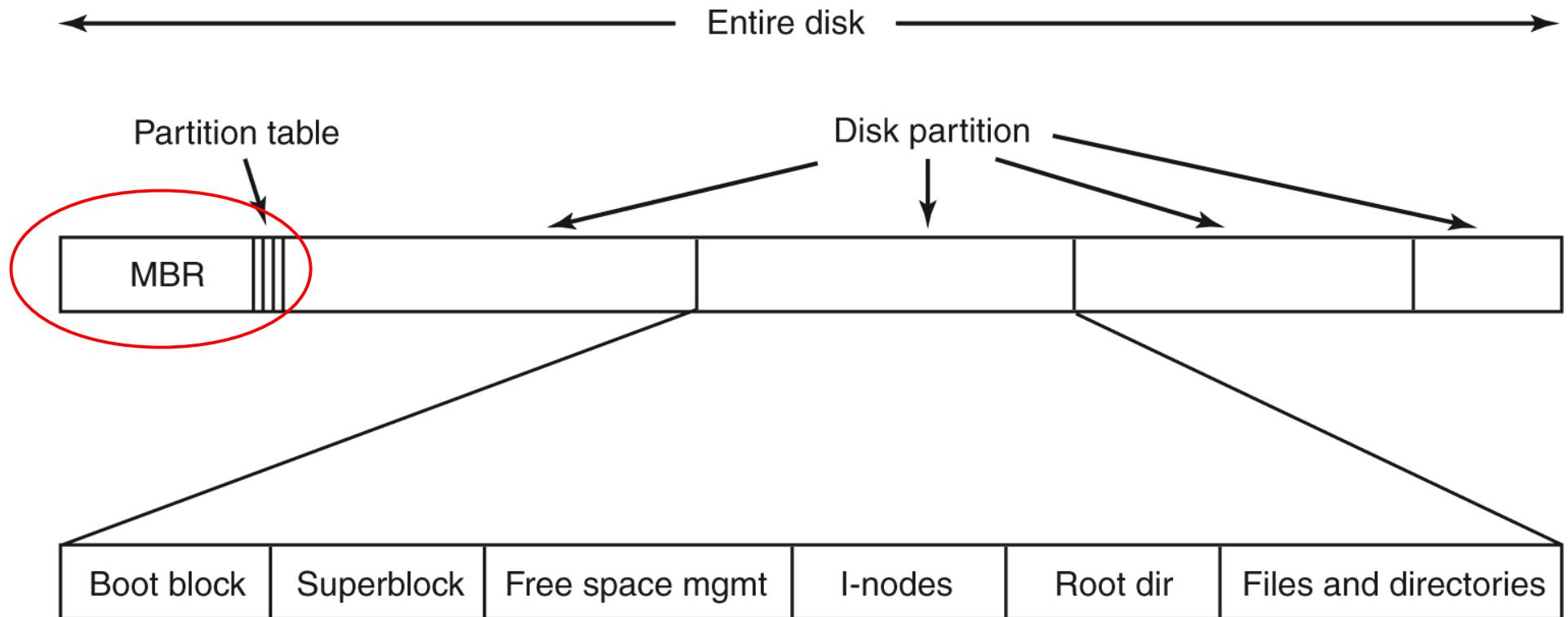
- Files
- Directories
- Implementation
- Examples

A Bit About Hard Drives



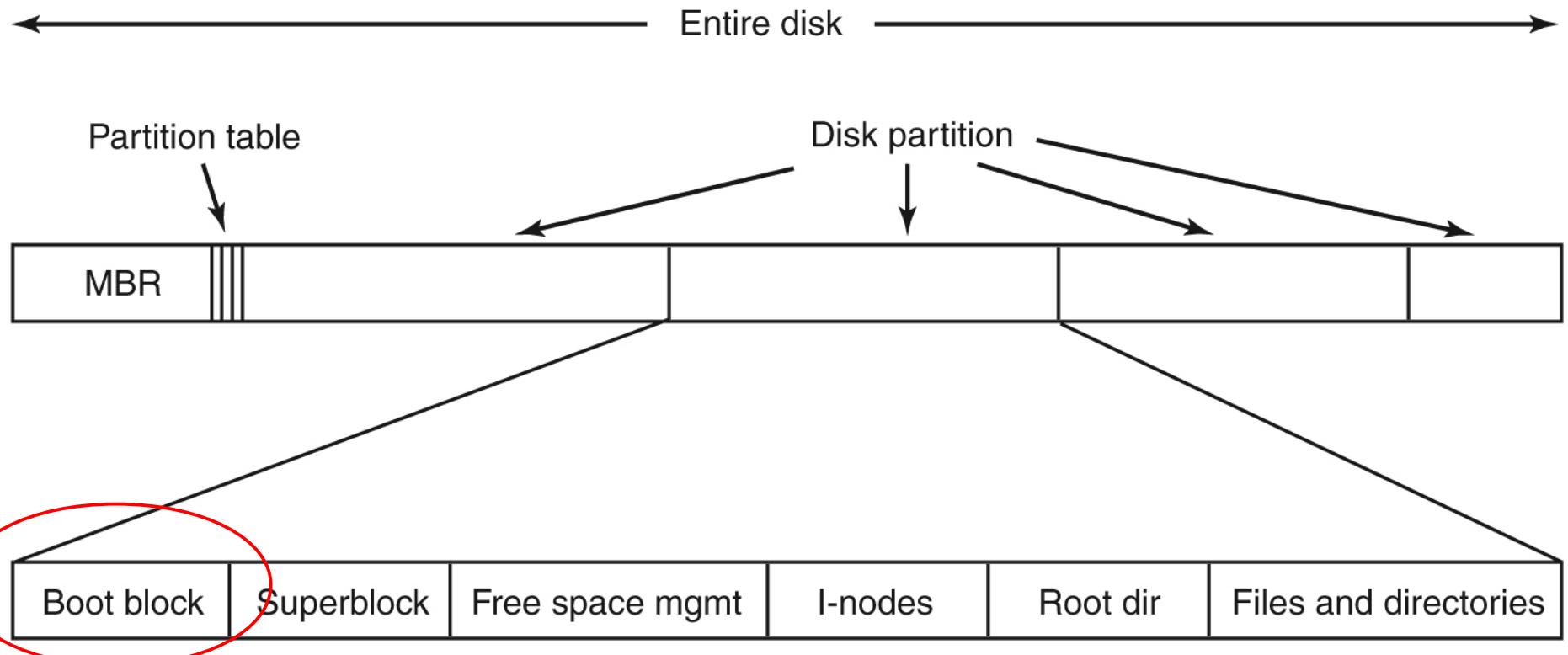
- Physical performance determined by seek time and rotational latency
- To be efficient, file systems must play to hard disk strengths and weaknesses

File System Layout



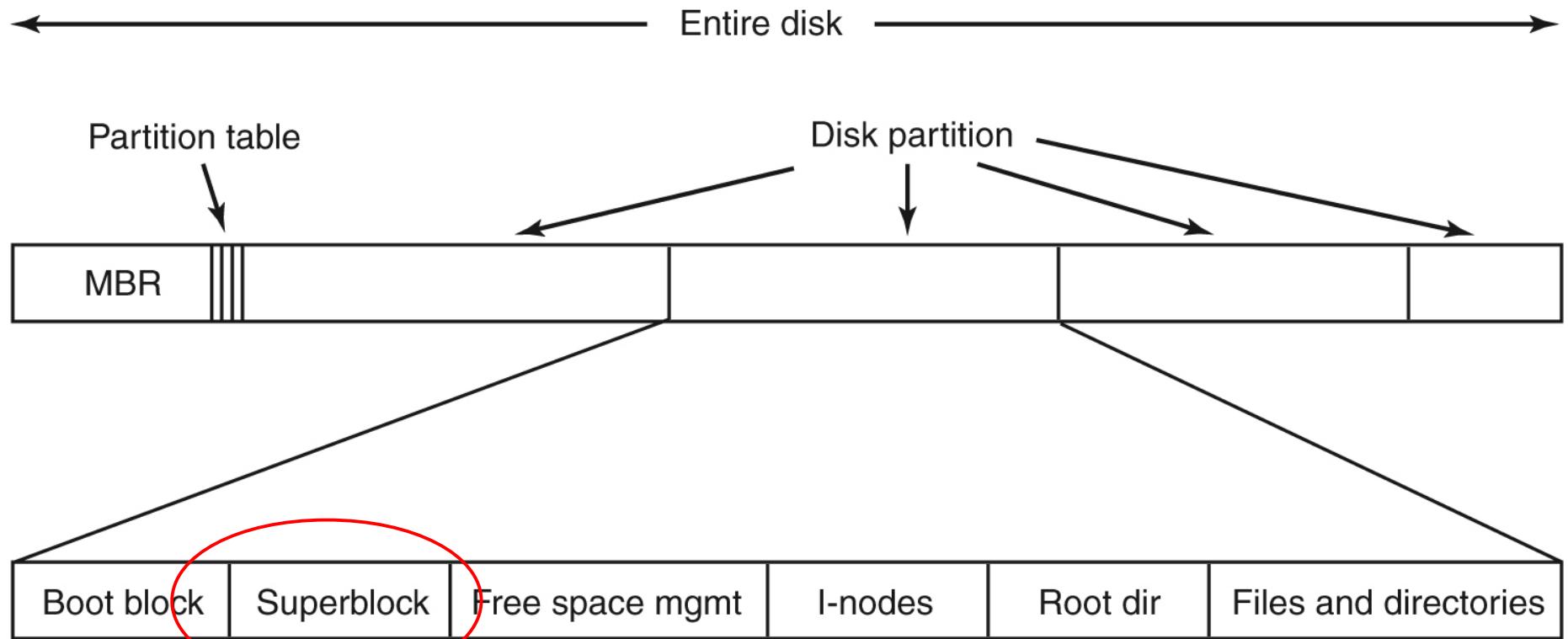
- **MBR**
 - Sector 0 of disk
 - Contains partition table (start, end, active)
 - Boot from active partition

File System Layout



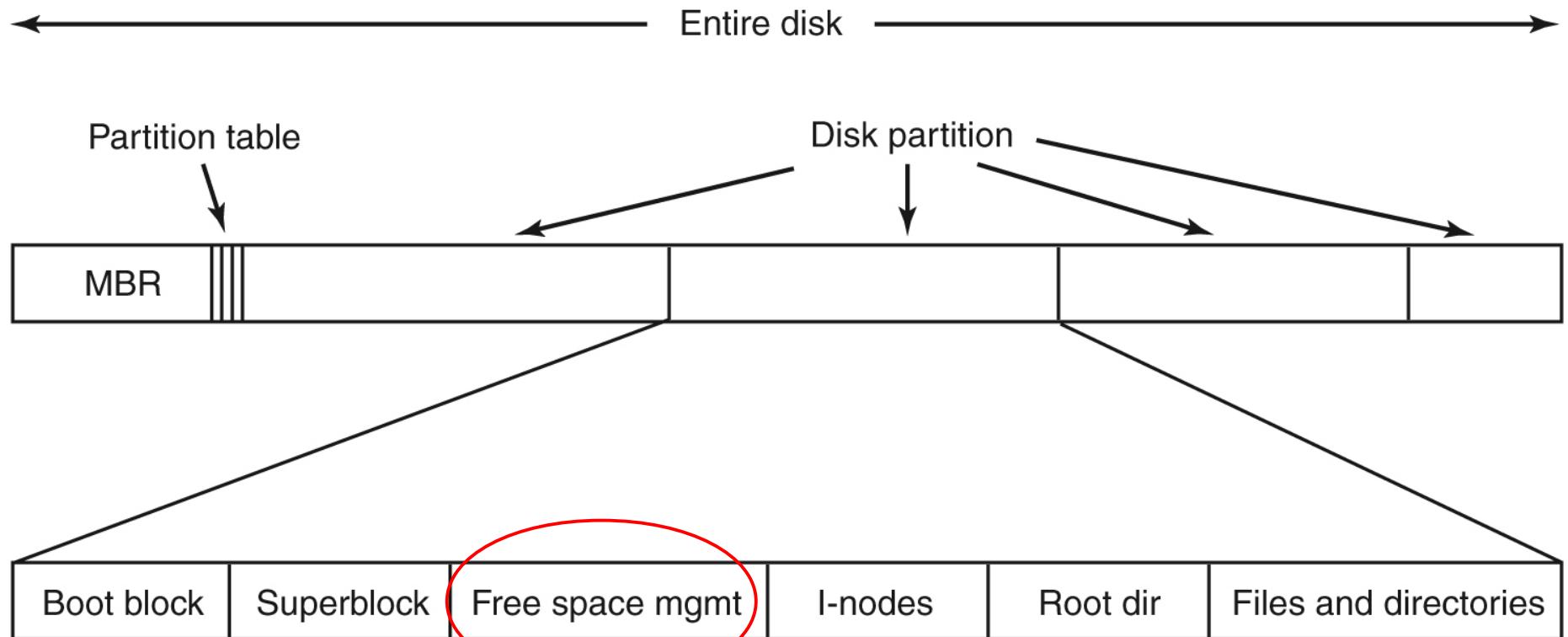
- Boot block
 - Loader for OS in that partition

File System Layout



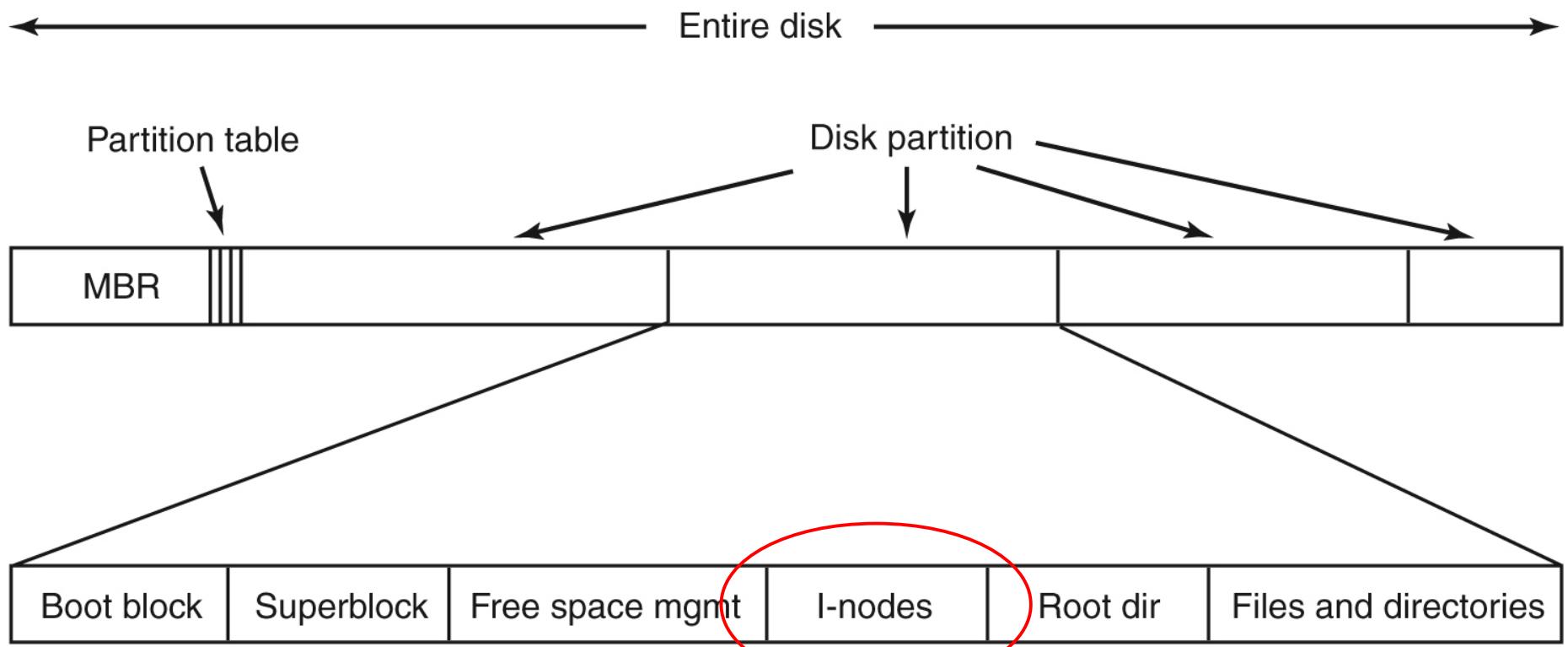
- Superblock
 - File system parameters (magic number, number of blocks in file system, etc)

File System Layout



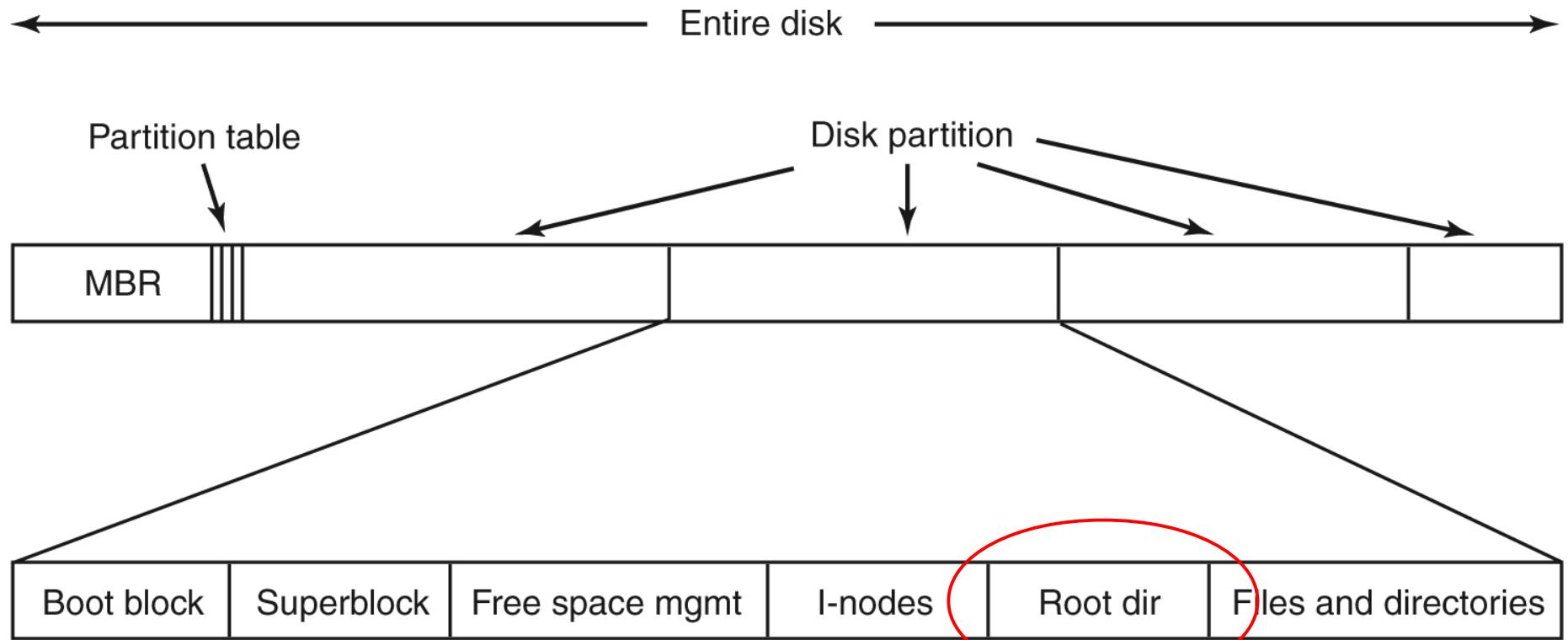
- Information about free blocks in system
 - Bitmap
 - List of pointers

File System Layout



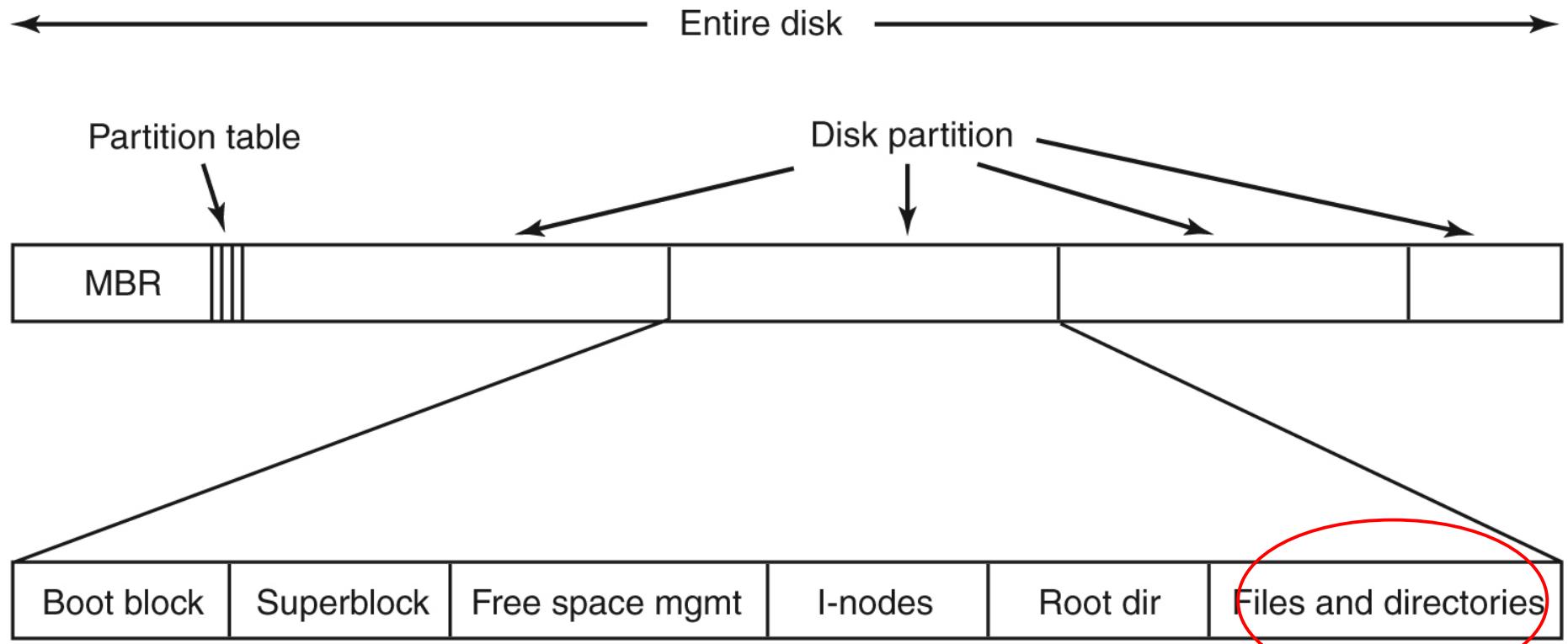
- inodes
 - Array of structures describing the file

File System Layout



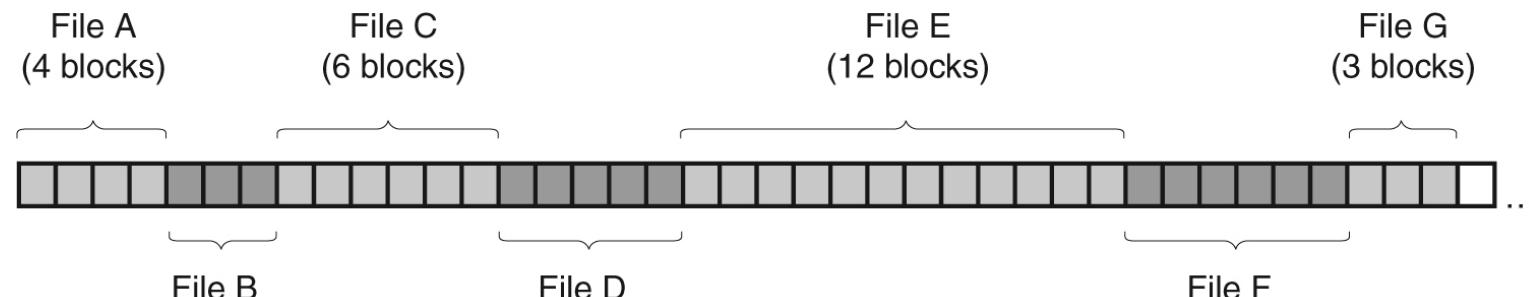
- Top of the file system

File System Layout

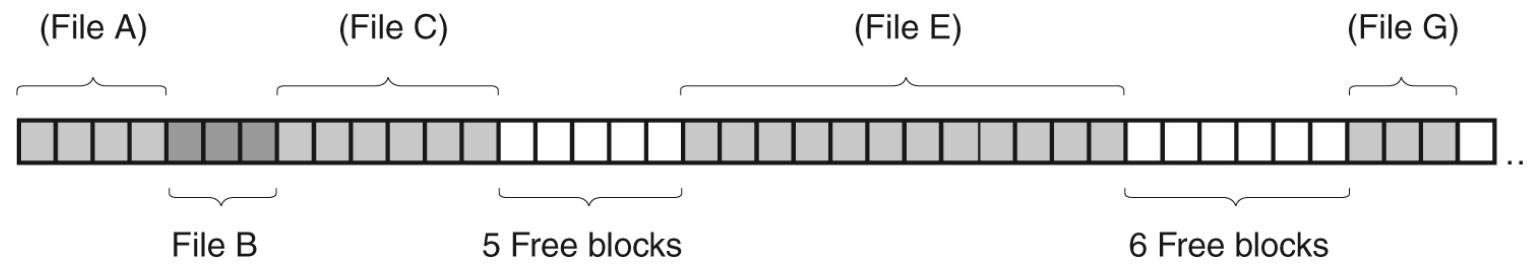


- All other directories and files

Files (Contiguous Allocation)



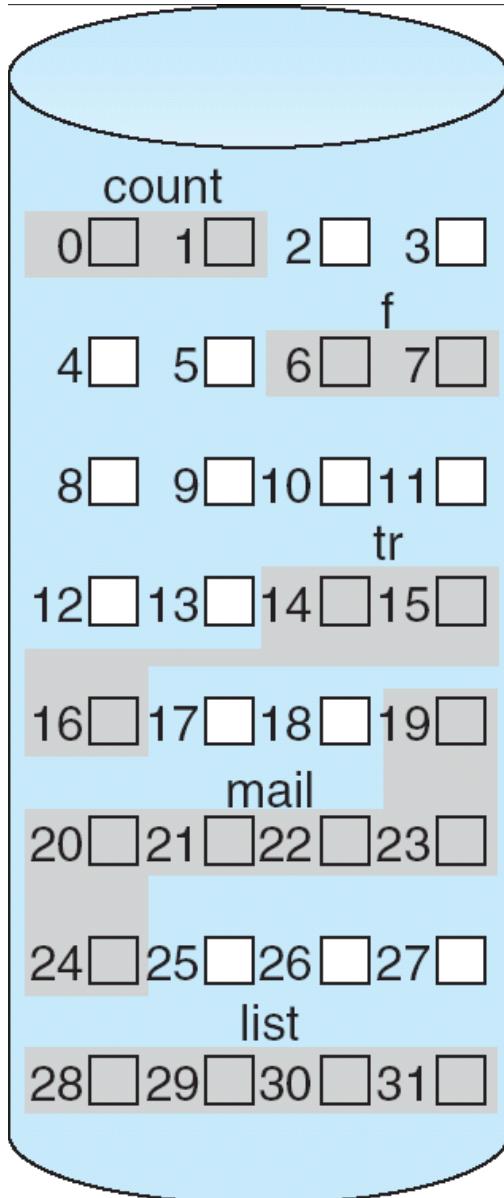
(a)



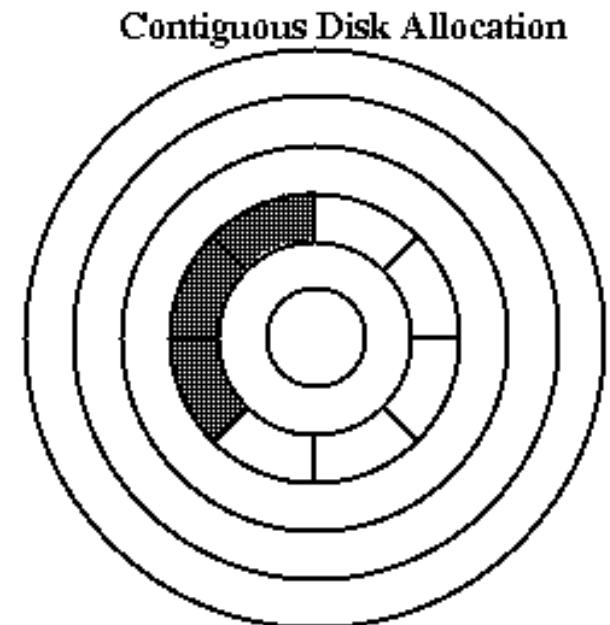
(b)

- Files begin at start of new block
- Simple: keep starting block and # of blocks
- Read performance excellent (few seeks)
- Fragmentation, have to know length of files

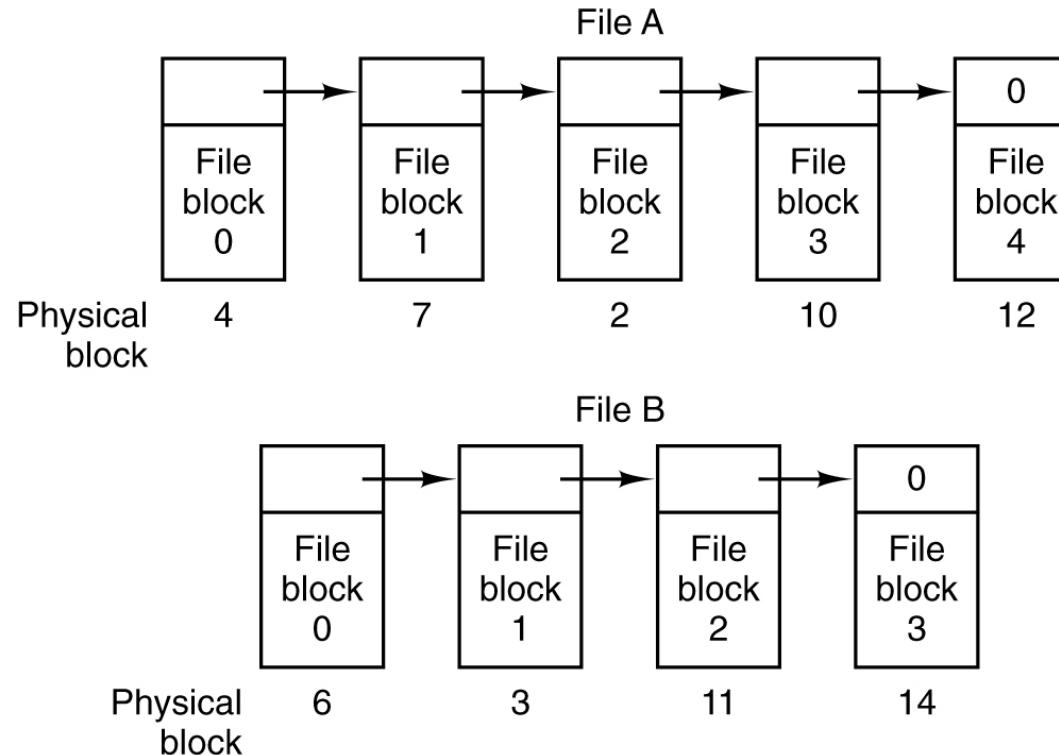
Contiguous Allocation of Disk Space



directory		
file	start	length
count	0	2
tr	14	3
mail	19	6
list	28	4
f	6	2

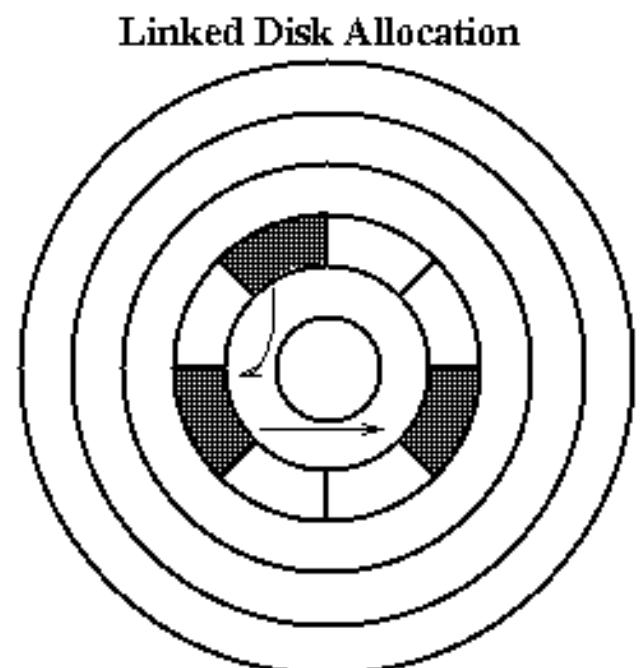
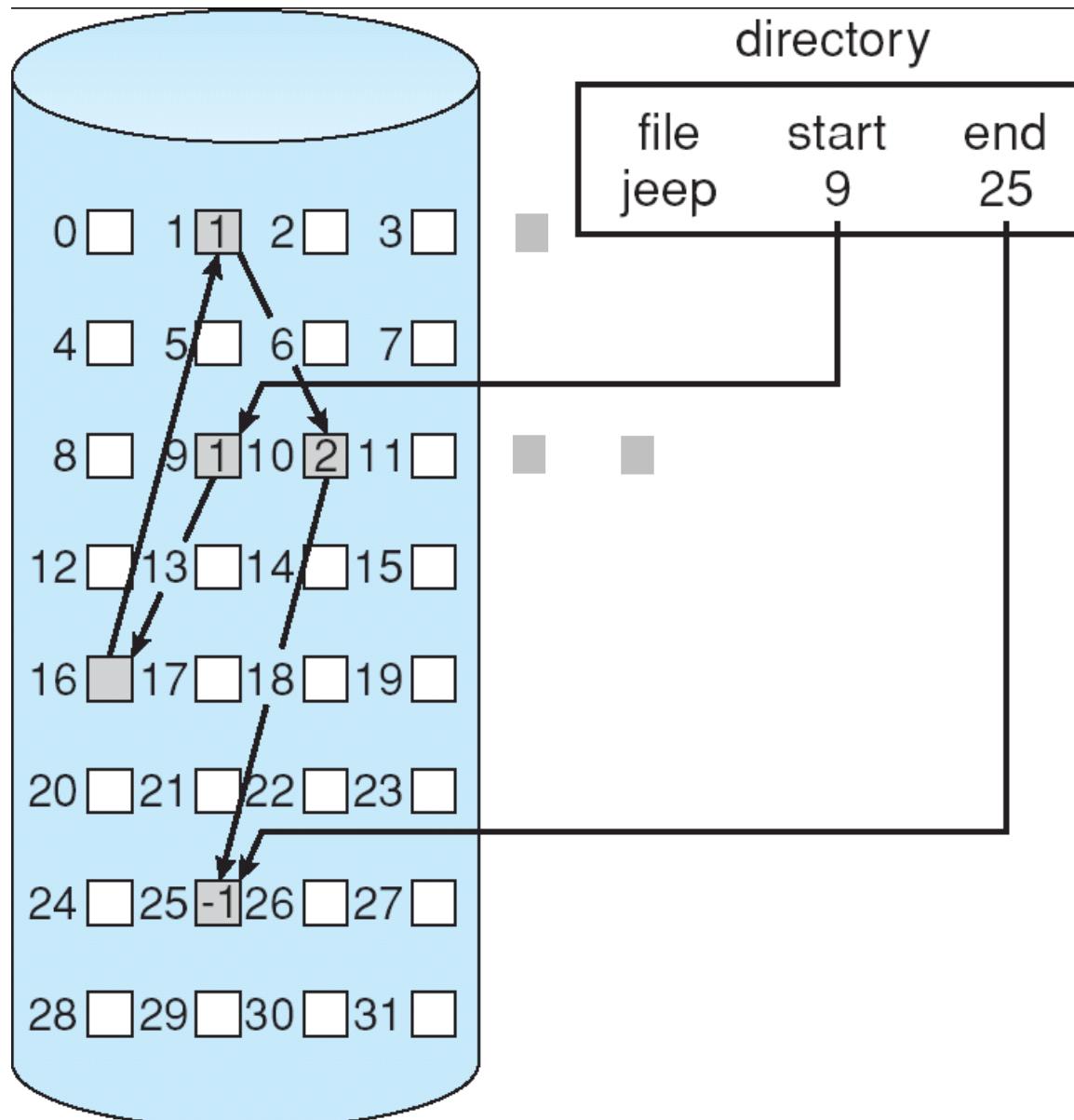


Files (Linked-List Allocation)

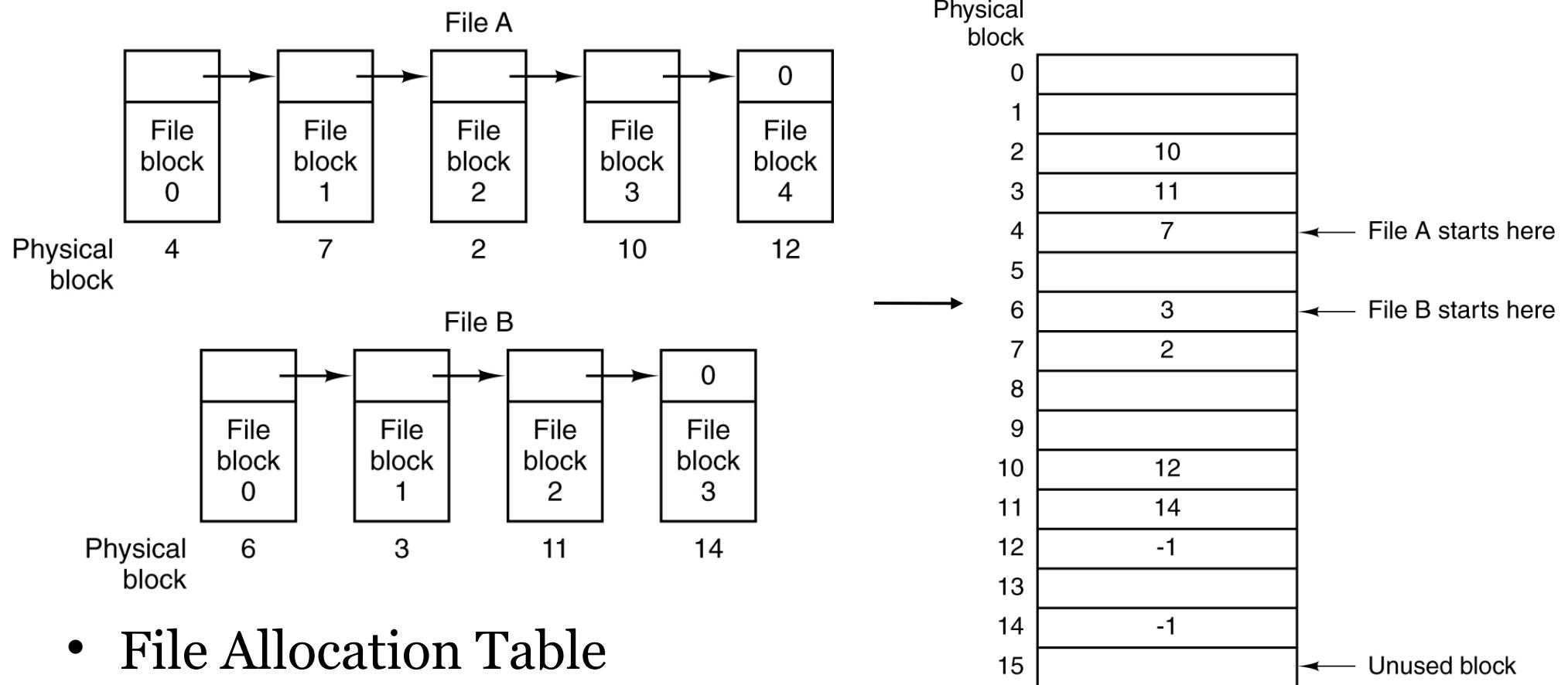


- First word in block points to next block
 - No external fragmentation
 - Random access is slow, lots of seeks
 - Reading block size of data requires 2 blocks

Linked Allocation of Disk Space

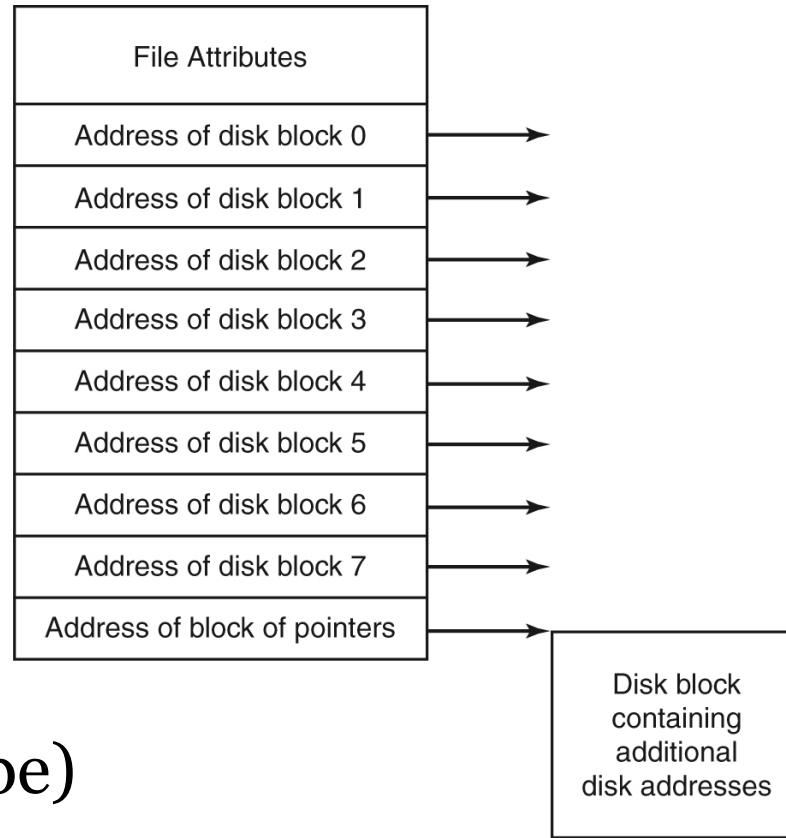


Files (FAT)



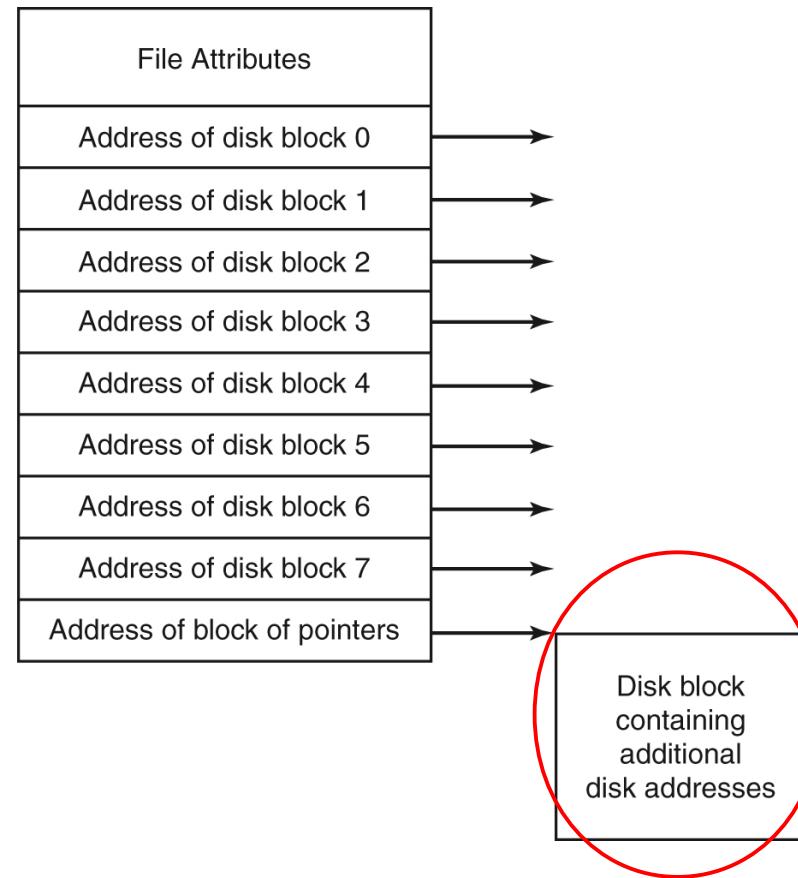
- **File Allocation Table**
 - Solves problem of wasting block space for pointer
 - Can still locate all blocks from starting block
 - Entire table in memory (e.g. 800MB for 200GB disk)

Files (inodes)



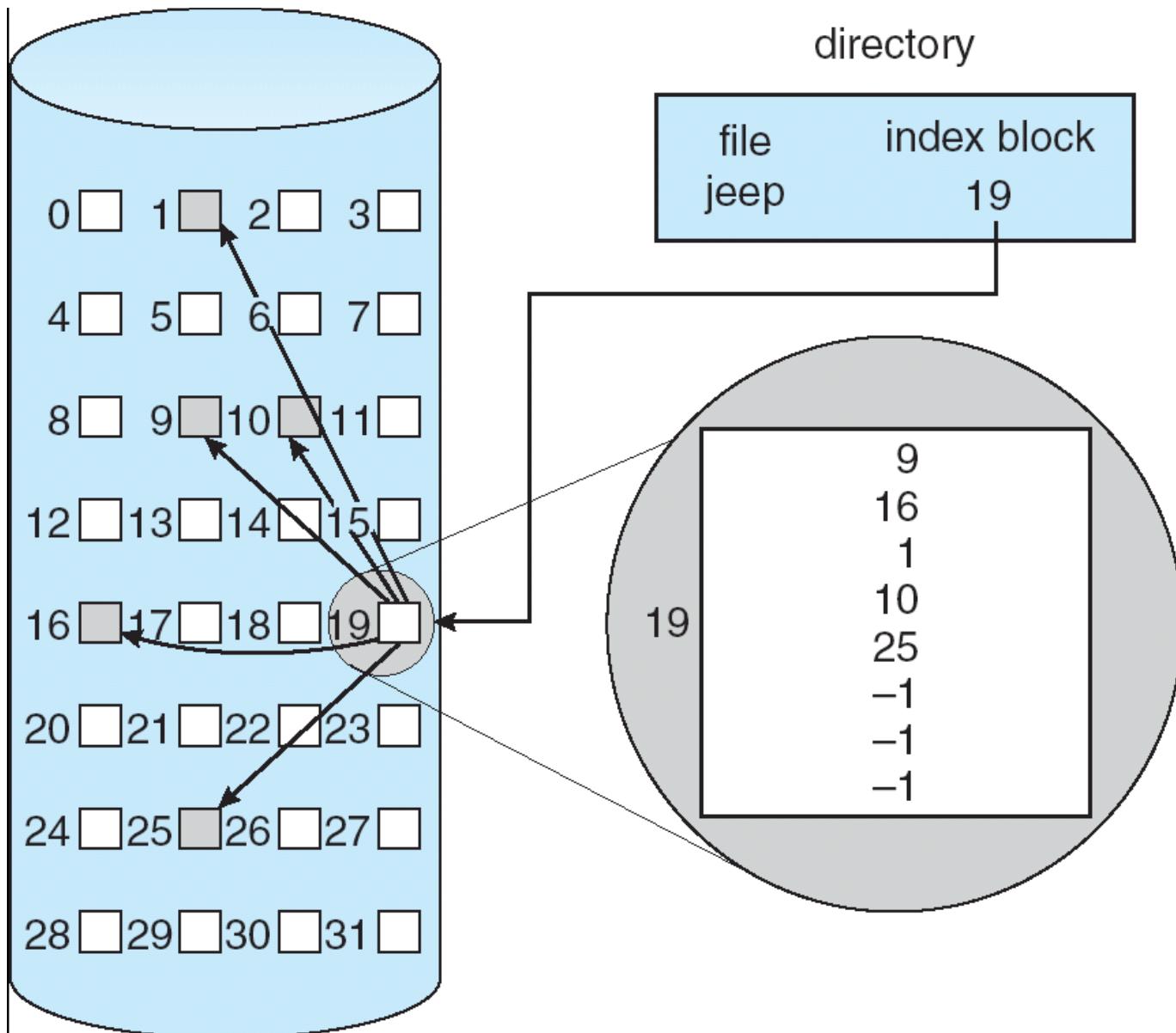
- Index-node? (maybe)
 - Lists attributes and disk addresses of all blocks
 - File's inode only in memory when file is used

Files (inodes)

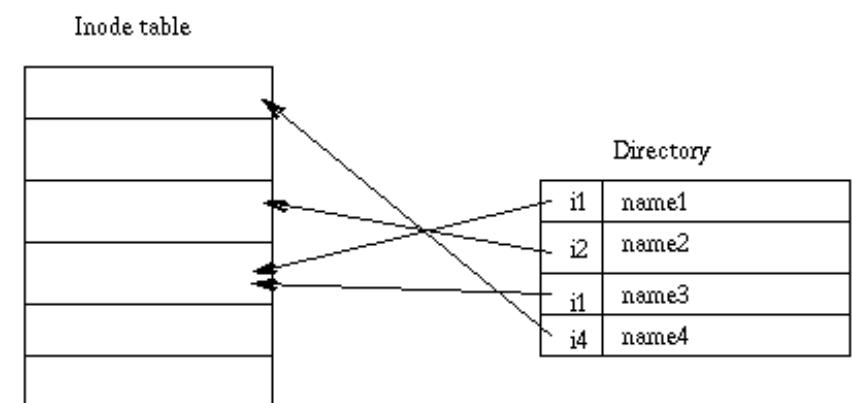
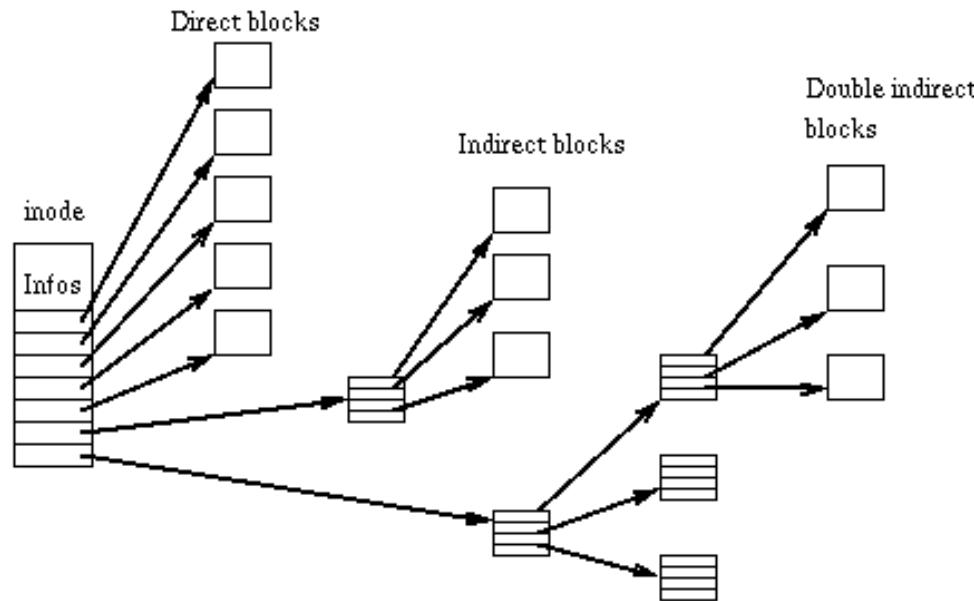


- Why this?

Example of Indexed Allocation



Summary so far...



FS Management and Optimization

- Files usually stored as collections of fixed-size blocks
 - Need not be stored adjacently on disk
- Ok, but what's a good block size?
 - Depends...

Choosing a Block Size

- Whatever is best fit

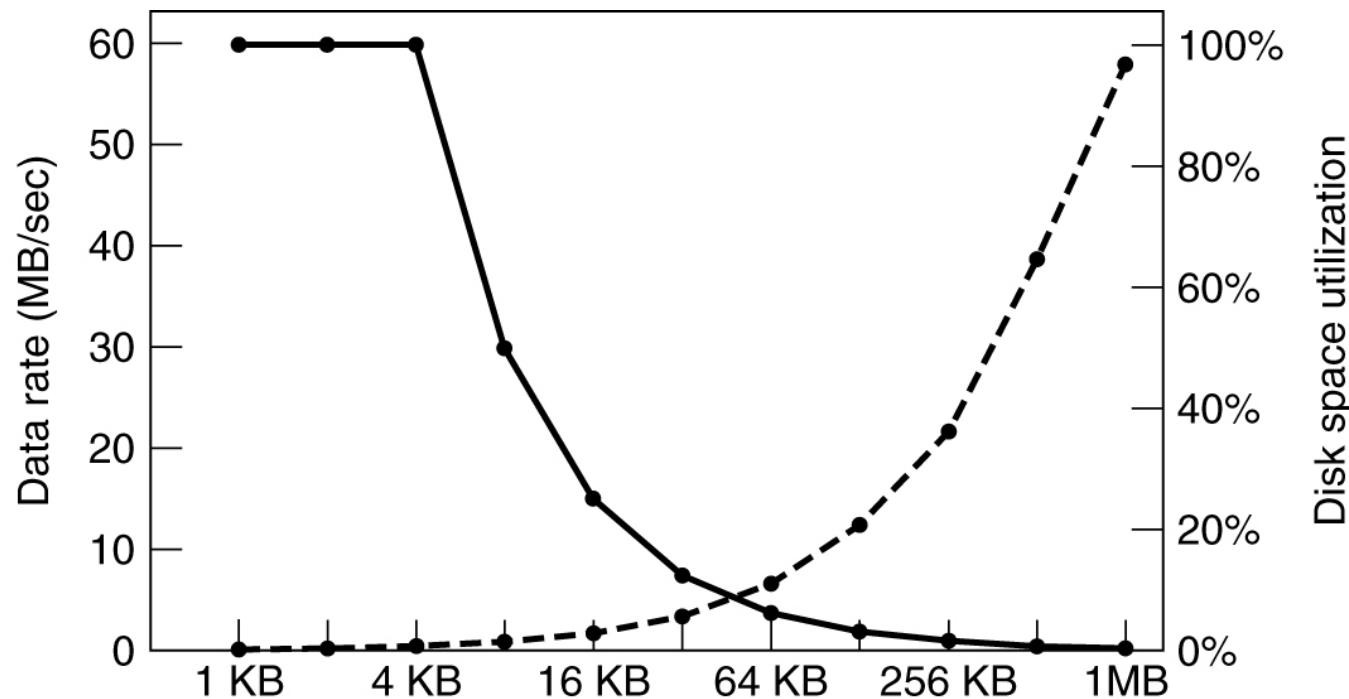
Length	VU 1984	VU 2005	Web
1	1.79	1.38	6.67
2	1.88	1.53	7.67
4	2.01	1.65	8.33
8	2.31	1.80	11.30
16	3.32	2.15	11.46
32	5.13	3.15	12.33
64	8.71	4.98	26.10
128	14.73	8.03	28.49
256	23.09	13.29	32.10
512	34.44	20.62	39.94
1 KB	48.05	30.91	47.82
2 KB	60.87	46.09	59.44
4 KB	75.31	59.13	70.64
8 KB	84.97	69.96	79.69

Length	VU 1984	VU 2005	Web
16 KB	92.53	78.92	86.79
32 KB	97.21	85.87	91.65
64 KB	99.18	90.84	94.80
128 KB	99.84	93.73	96.93
256 KB	99.96	96.12	98.48
512 KB	100.00	97.73	98.99
1 MB	100.00	98.87	99.62
2 MB	100.00	99.44	99.80
4 MB	100.00	99.71	99.87
8 MB	100.00	99.86	99.94
16 MB	100.00	99.94	99.97
32 MB	100.00	99.97	99.99
64 MB	100.00	99.99	99.99
128 MB	100.00	99.99	100.00

Choosing a Block Size

- Keep in mind:
 - Smaller blocks ^ more seeks
 - Bigger blocks ^ more wasted space
 - Time vs space tradeoff

The Tradeoff



- Data rate increases with block size
 - Access time for a block dominated by seek and rotational latency
- Space wastage increases with block size

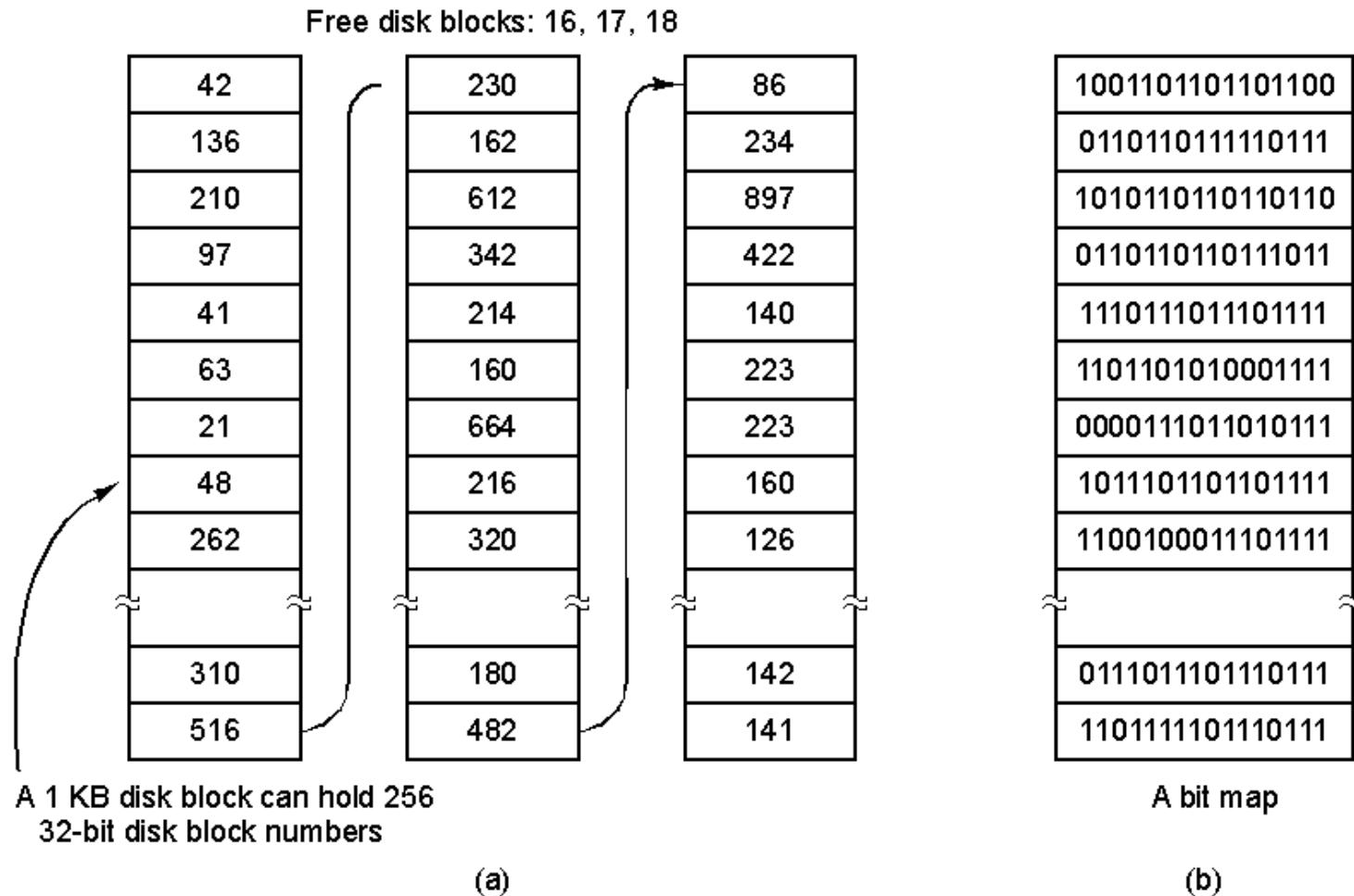
I've Chosen a Block Size, Now What

- Need to keep track of free blocks

Free disk space management

- Linked list
 - A linked-list list of unused blocks is used to keep track of free space
 - 32 bits per list entry
 - Many entries per list in empty blocks (linked)
- Bitmap
 - A bitmap of all blocks keeps track of which blocks are in use and which are free
 - Better than linked list if the entire bitmap can be kept in main memory at all times

List vs. Bit Map



- (a) Storing the free list on a linked list
(b) A bit map

Performance

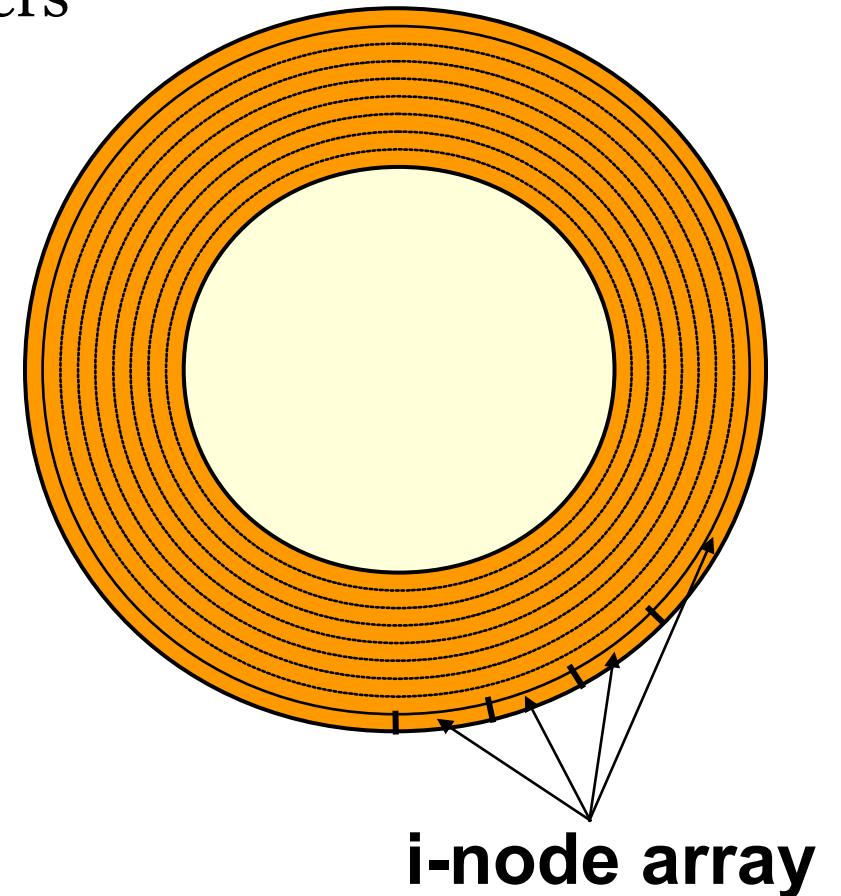
- How do users access files?
 - Sequential: bytes read/written in order
 - Random: read/write blocks in middle of file
 - May access whole or partial file
- How are files used?
 - Most are small
 - Large files take up most of disk space
 - Large files account for most bytes transferred
- Everything needs to be efficient

Performance

- Hard drive
 - Disk allocation and directory algorithms
 - Types of data kept in file's directory entry
 - Data/metadata layout
- Operating system
 - Disk cache – separate section of main memory for frequently used blocks
 - Delayed-writes – aggregation / higher priority to reads
 - Read-ahead – Prefetching of sequential blocks
 - Section of memory as virtual disk (RAM disk)

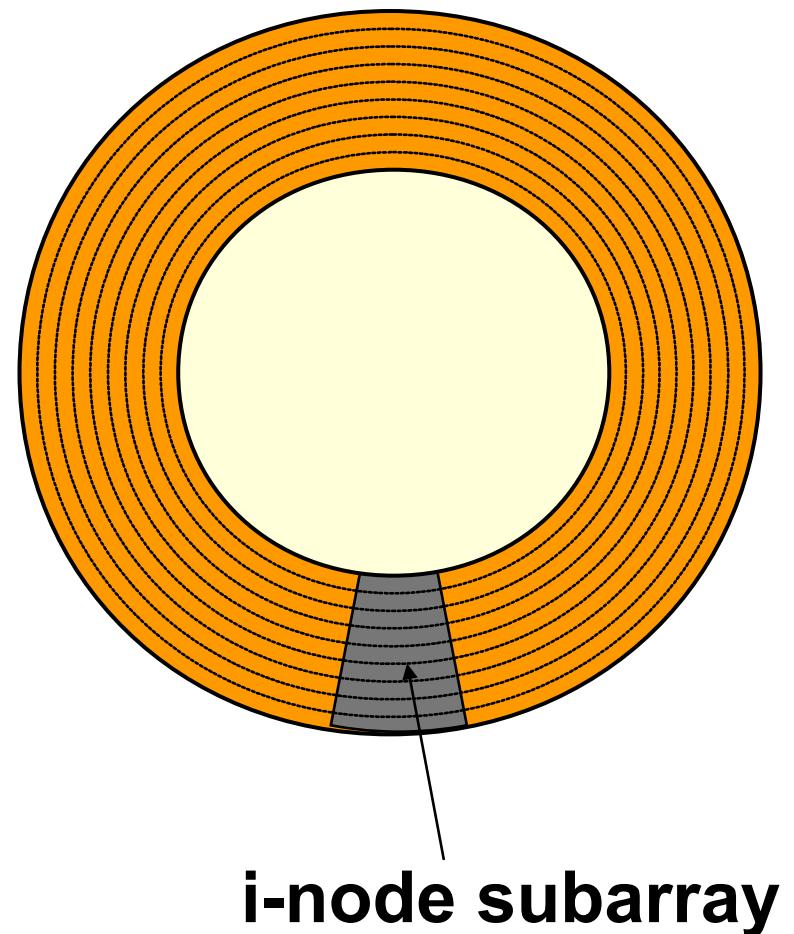
Early Unix Disk Layout

- Array of inodes in outermost cylinders
- inodes are far away from data
- at least 2 seeks per I/O



Modern Disk Layout (very broadly)

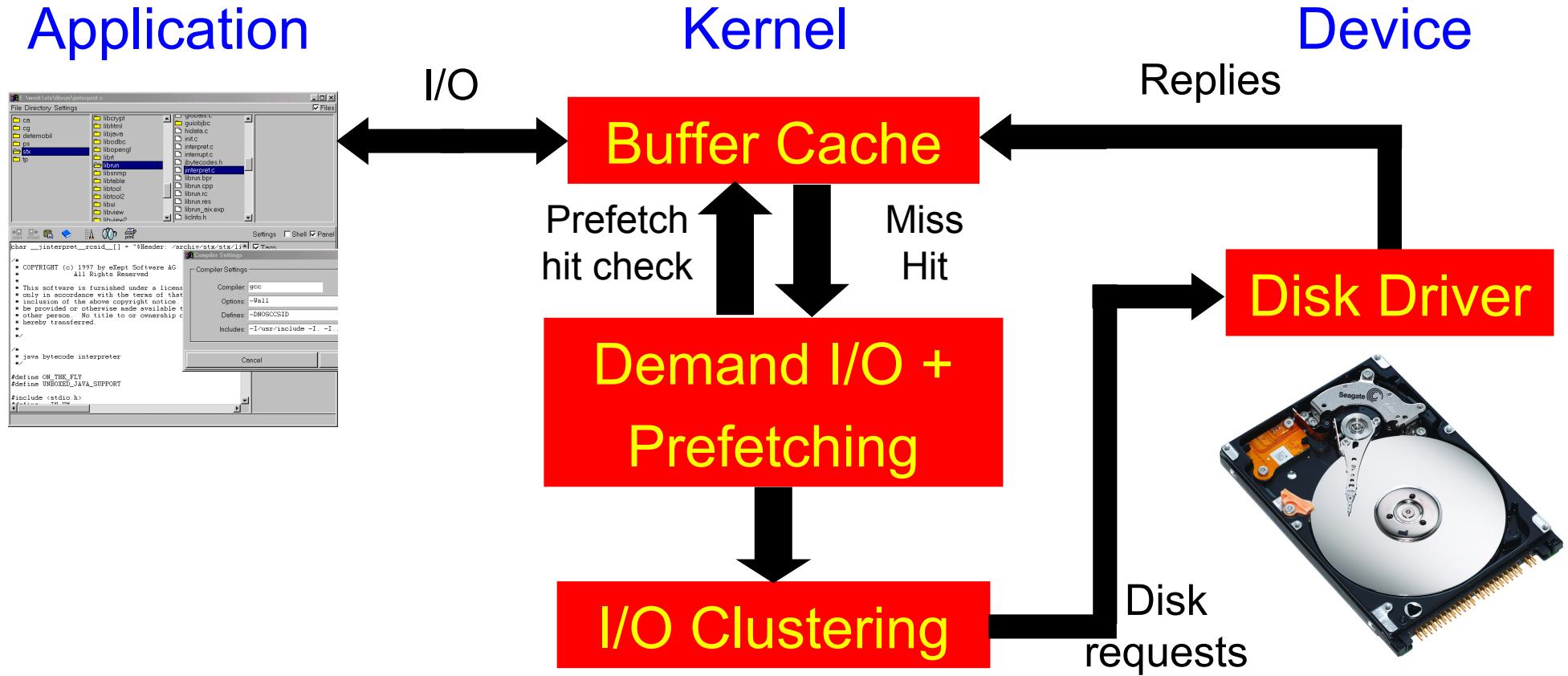
- Portion of inode array on each cylinder
- Multiple inodes on same block
- One seek to access inode AND data



File Caching

- Locality of reference in file accesses
 - Yet another application of the principle of locality
 - What were the earlier instances?
- Keep a number of disk blocks in the much faster memory
 - when accessing disk, check the cache first!

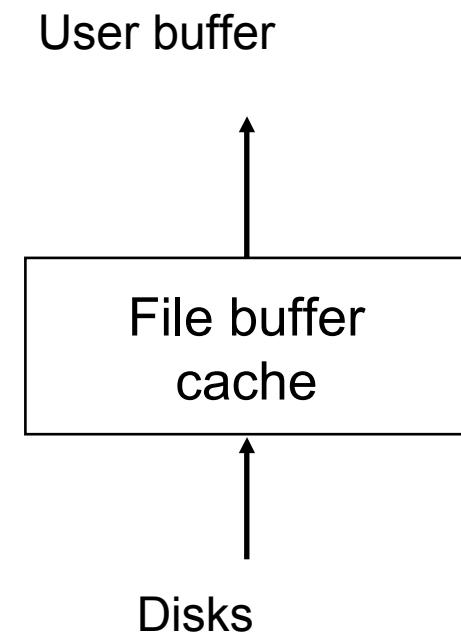
I/O path in modern OS



- Variable portion of main memory (in Linux); processes and file system compete for physical memory

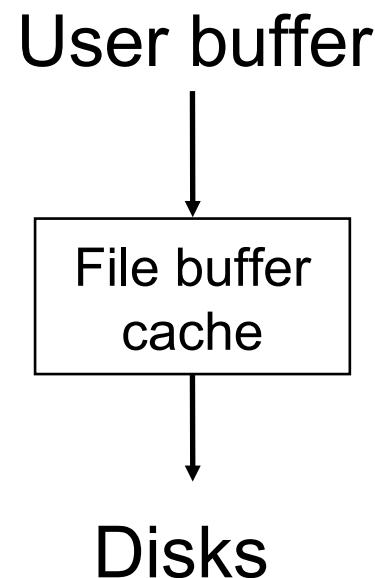
Handling Reads

- `read(fd, buf, n)`
- On a hit
 - copy from the buffer cache to a user buffer
- On a miss
 - replacement when full
 - read a file into the buffer cache



Handling Writes

- `write(fd, buffer, n)`
- On a hit
 - write to buffer cache
- On a miss
 - read the file into buffer cache
(possible replacement)
 - write to buffer cache

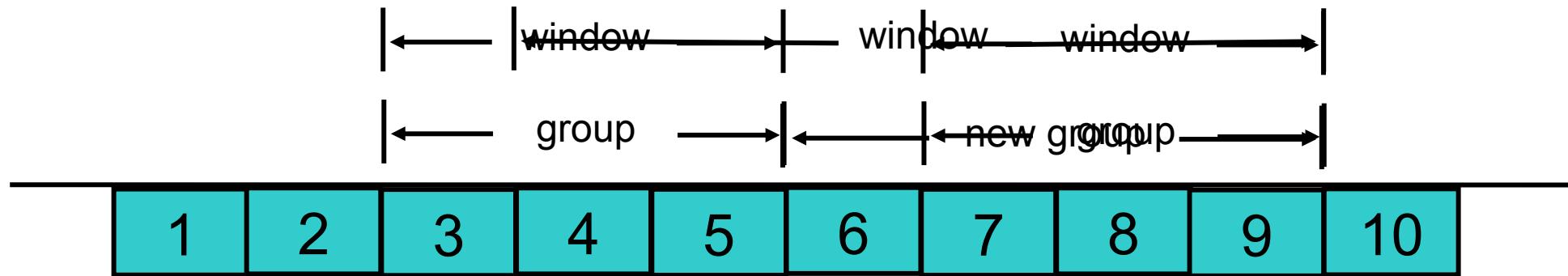


File prefetching

- Read-ahead:
 - For sequential access, read the requested block and the following N blocks together (why is this a good idea?)
- Drawbacks:
- Potential for cache pollution
 - Data prefetched that will not be accessed
 - Useful data replaced
- Wrong prefetches can saturate disk I/O bandwidth

Prefetching in Linux Kernel

- Targets consecutive blocks
- Read-ahead group
 - Consecutive blocks fetched (demand + prefetched)
- Read-ahead window
 - Consecutive blocks in the previous and current read-ahead groups



Observations

- Kernel prefetching can:
 - Narrow the performance gap of replacement algorithms
 - Change the relative performance benefits of algorithms
- Implication: Execution time is the ultimate metric

File System Consistency

- System crashes
 - Are all blocks written out?
 - Can lead to inconsistent state
- Utilities exist that check for consistency
 - UNIX: fsck
 - Win: scandisk, chkdsk
- Block consistency
 - Count blocks used
 - Count blocks in free list
 - Consistency means a block is either used or it is not

Consistency Checks

- System crashes can leave the filesystem in a number of inconsistent states

Block number																
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	
1	1	0	1	0	1	1	1	1	0	0	1	1	1	0	0	
Blocks in use																
0	0	1	0	1	0	0	0	0	1	1	0	0	0	1	1	
Free blocks																

(a)

Block number																
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	
1	1	0	1	0	1	1	1	1	0	0	1	1	1	1	0	0
Blocks in use																
0	0	0	0	1	0	0	0	0	1	1	0	0	0	1	1	
Free blocks																

(b)

Block number																
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	
1	1	0	1	0	1	1	1	1	0	0	1	1	1	1	0	0
Blocks in use																
0	0	1	0	2	0	0	0	0	1	1	0	0	0	1	1	
Free blocks																

(c)

Block number																
0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	
1	1	0	1	0	2	1	1	1	0	0	1	1	1	1	0	0
Blocks in use																
0	0	1	0	1	0	0	0	0	1	1	0	0	0	1	1	
Free blocks																

(d)

Consistency Checks

- Also need to check directories for inode references
 - Walk the directory tree, keep counts of references
 - Compare reference counts to those stored in inodes

Some File Systems

- Media based
 - Ext3/4 - Linux native
 - ufs – BSD
 - fat - DOS FS
 - vfat - win 95
 - hpfs – OS/2
 - minix – early Linux
 - Isofs – CDROM
 - sysv - Sysv Unix
 - hfs – Macintosh
 - affs - Amiga Fast FS
 - NTFS - NT's FS
 - adfs - Acorn-strongarm
- Network
 - nfs
 - Coda
 - AFS - Andrew FS
 - smbfs – LanManager
 - ncpfs - Novell
- Special ones
 - procfs -/proc
 - umsdos - Unix in DOS
 - userfs - redirector to user

File Systems

Files

Directories

Implementation

Examples

- Virtual File System

Linux File Systems

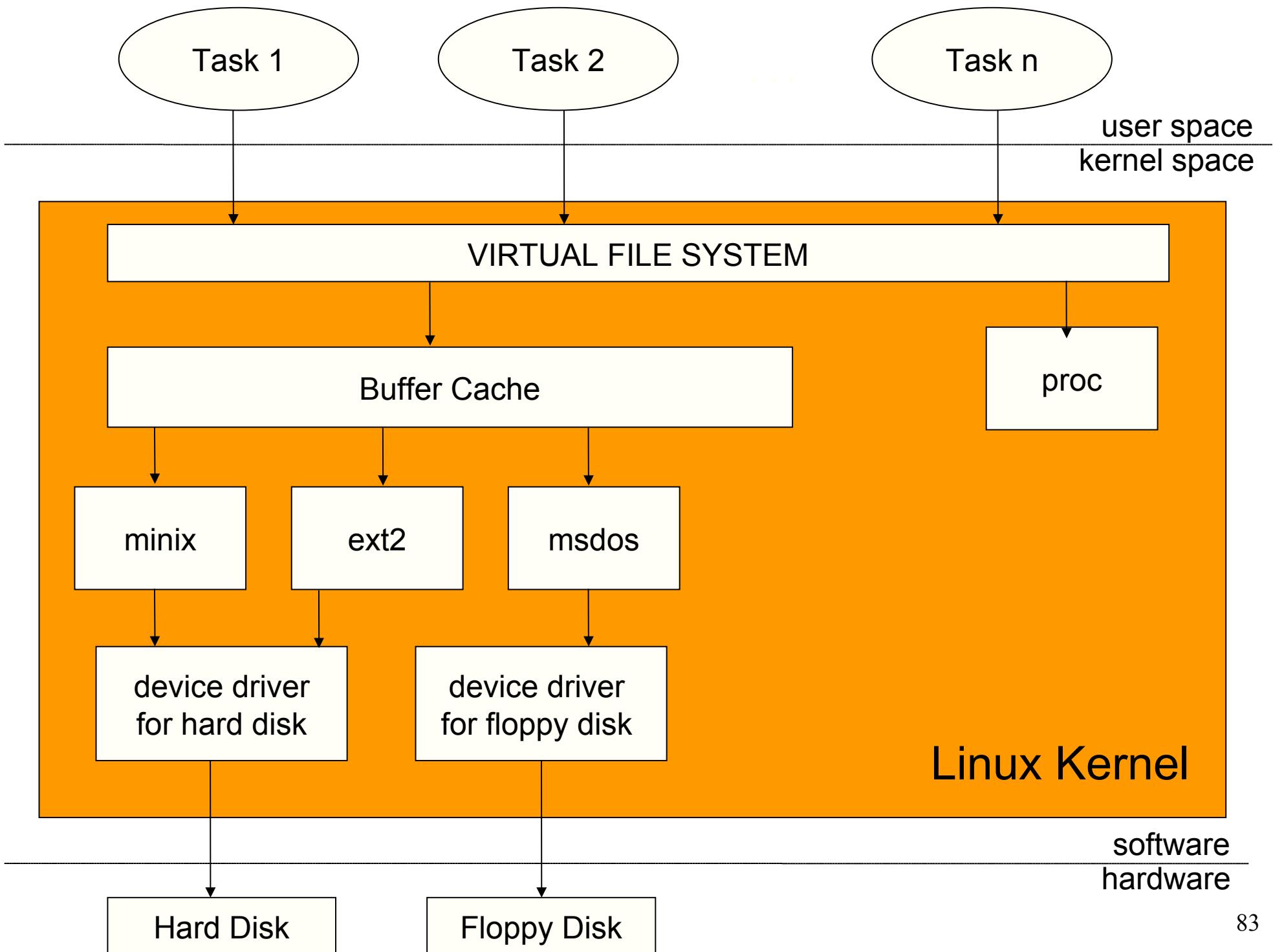
- To the user, Linux's file system appears as a hierarchical directory tree obeying UNIX semantics
- Internally, the kernel hides implementation details and manages the multiple different file systems via an abstraction layer, that is, the *virtual file system (VFS)*

Linux File Systems

- The Linux VFS is designed around object-oriented principles and is composed of two components:
 - A set of definitions that define what a file object is allowed to look like
 - The *inode-object* and the *file-object* structures represent individual files
 - the *file system object* represents an entire file system
 - A layer of software to manipulate those objects

Virtual File Systems

- Goal
 - Allow one machine to use multiple file system types
 - Unix FFS
 - MS-DOS FAT
 - CD-ROM ISO9660
 - Remote/distributed: NFS/AFS
 - Standard system calls should work transparently
 - Share common concepts and high-level operations
- Solution
 - Insert a level of indirection (VFS)!
 - Applications interact with this VFS
 - Kernel translates abstract-to-actual



File Systems

Files

Directories

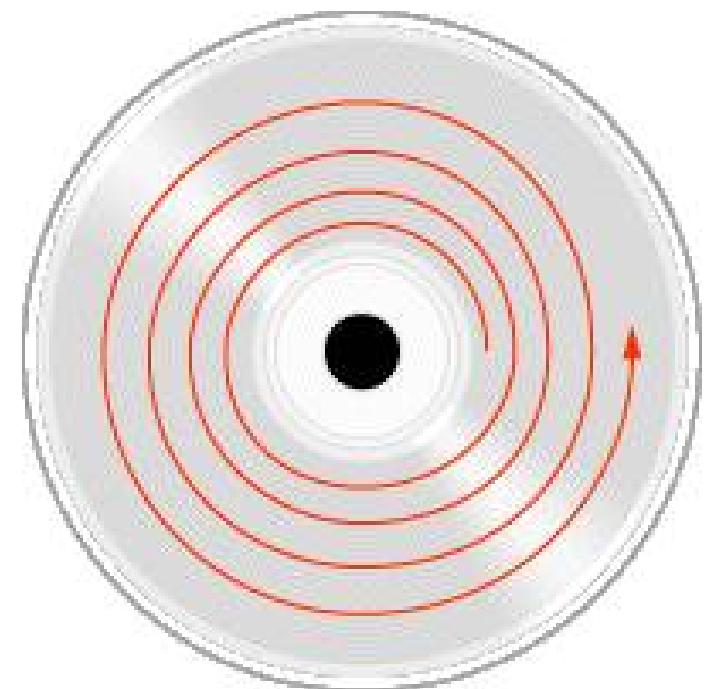
Implementation

Examples

- Virtual File System
- CD-Rom FS (ISO 9660)

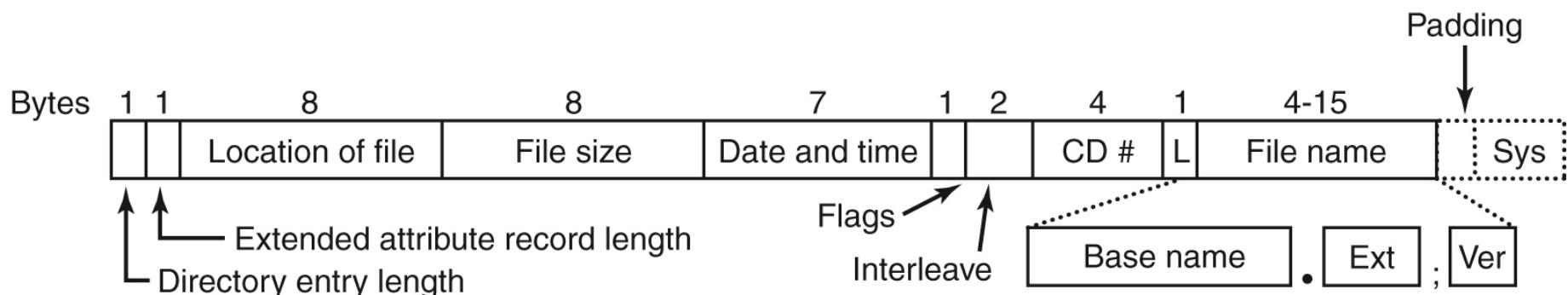
CD-ROM

- Write-once, append-only
- No tracking of free blocks
- Free space at end of CD
- Spiral contains bits in linear sequence
- 2352 byte blocks
- 2048 byte payload



CD-ROM

- 16 blocks undefined
- (block 17) Primary volume descriptor
 - Ids: system, volume, publisher, data preparer, etc
 - Root directory entry (where is it on disc?)
- Can locate rest of file system from root dir
- Directory entries:



File Systems

Files

Directories

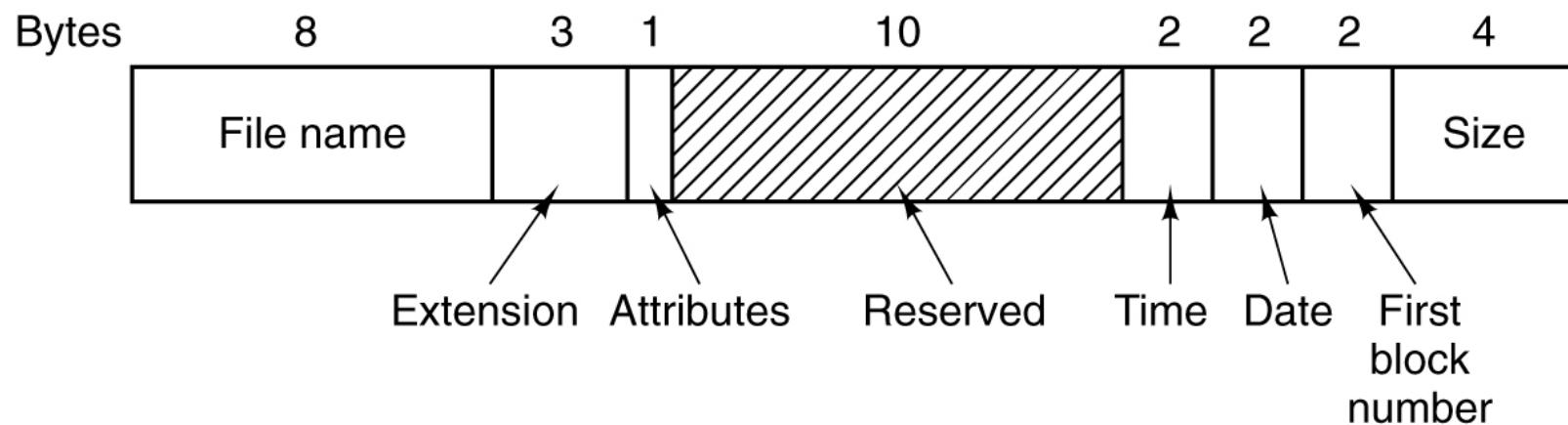
Implementation

Examples

- Virtual File System
- CD-Rom FS (ISO 9660)
- MS-DOS

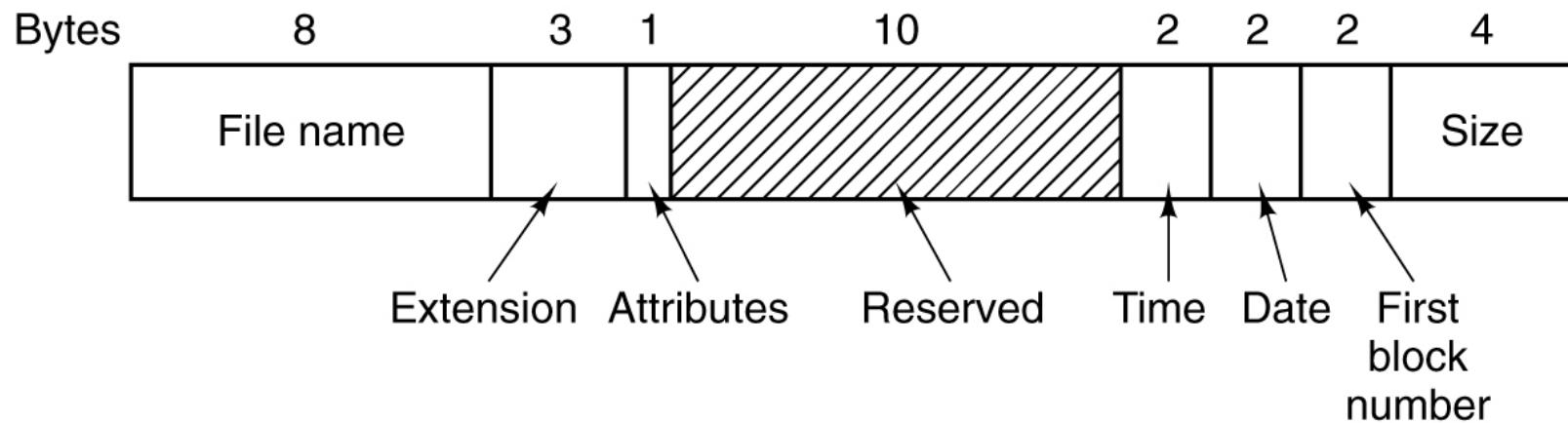
MS-DOS FS

- Used in: Windows 98-ME, floppy disks, embedded systems (FAT32), iPod
- Supported by: Windows 2000, XP, Vista, Linux
- Fixed-size dir entries:



MS-DOS FS

- Time field limited to 2 bytes
 - 65,536 unique values, but day has 86,400 seconds
- Year maxes out at 2107
 - Y2108
- File size ^ 32 bits (theoretical max 4GB)



MS-DOS FS

- Flavors: FAT-12, FAT-16, FAT-32
 - Number refers to number of bits used in disk address
 - Not all bits used for address: e.g. FAT-32 might be better called FAT-28
- Disk blocks: 512 bytes
- So...

FAT-12 w/ 512 byte blocks

Maximum partition size?

$$2^{12} * 512 = 2\text{MB}$$

MS-DOS FS

- Maximum partition sizes w/ different block sizes

Block size	FAT-12	FAT-16	FAT-32
0.5 KB	2 MB		
1 KB	4 MB		
2 KB	8 MB	128 MB	
4 KB	16 MB	256 MB	1 TB
8 KB		512 MB	2 TB
16 KB		1024 MB	2 TB
32 KB		2048 MB	2 TB

- So, if I have a 8GB disk and FAT-32, at most how many partitions do I need?
- How about FAT-16?

File Systems

Files

Directories

Implementation

Examples

- Virtual File System
- CD-Rom FS (ISO 9660)
- MS-DOS
- Ext2/Ext3

Design Goals

- No performance loss
 - should have a performance gain!
- Backward compatibility
- Reliability through
 - preservation: stable data not affected by crashes
 - predictability: known failure modes
 - atomicity: each operation fully completes or is undone after recovery

JFS

- Atomic updates
- Old and new versions of data held on disk until commit
- Undo logging
 - copy old data to log
 - write new data to disk
 - crash during update? copy old data from log
- Redo logging
 - write new data to log
 - old data remains on disk
 - crash during update? copy new data from log

JFS (2)

- Speedy recovery after a crash
 - fsck is slow on large disks
 - JFS rereads journal after a crash, from last checkpoint
- Journal:
 - contains three types of blocks
 - metadata: contents of single fs block metadata as updated by transaction
 - descriptor: describe other journal metadata blocks
 - header: head and tail of journal, sequence number
 - circular structure

JFS vs LSFS

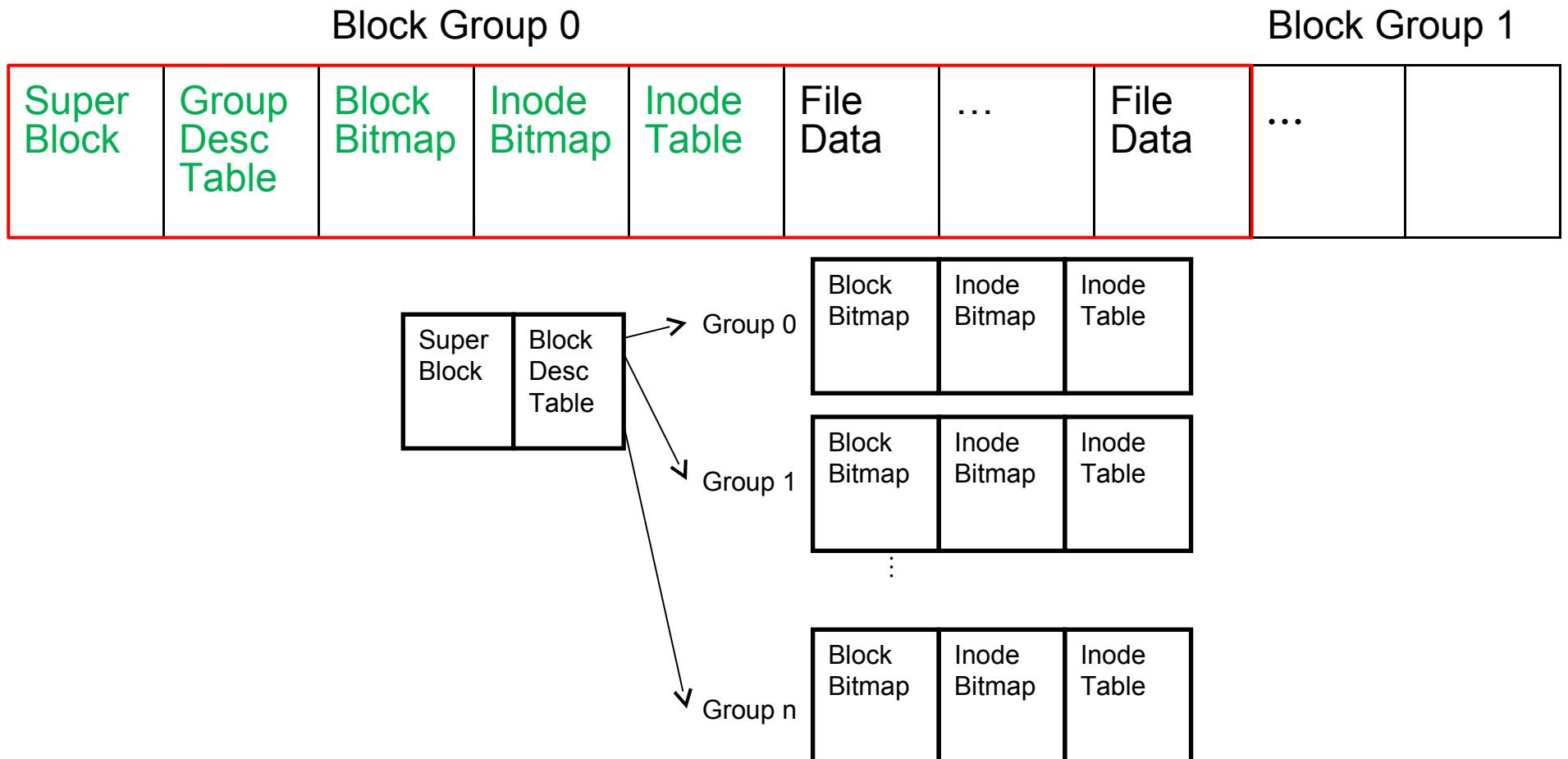
- Log structured FS ONLY contains a log
 - everything written to end
- LSFS dictates how data is stored on disk
- LFS does not

ext3 and JFS

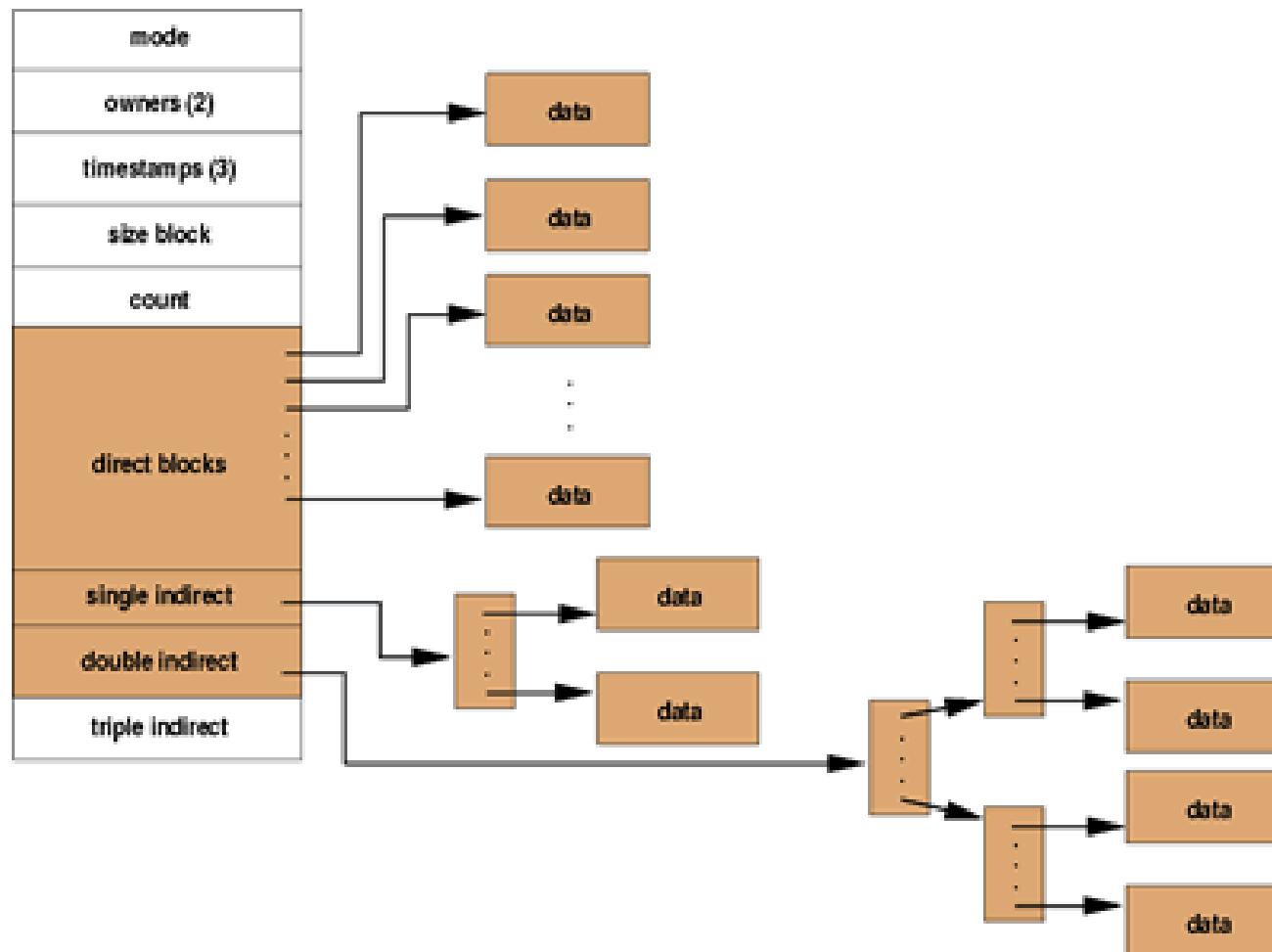
- ext3 grew out of ext2 (same code base)
 - backward compatible
- Two separate layers
 - /fs/ext3 ^ just the FS with transactions
 - /fs/jdb ^ just the journalling stuff
- ext3 calls JFS as needed
 - start/stop transaction
 - journal recovery after unclean reboot

Internals

- Physically organized into block groups
 - each group contains redundant copies of system control blocks



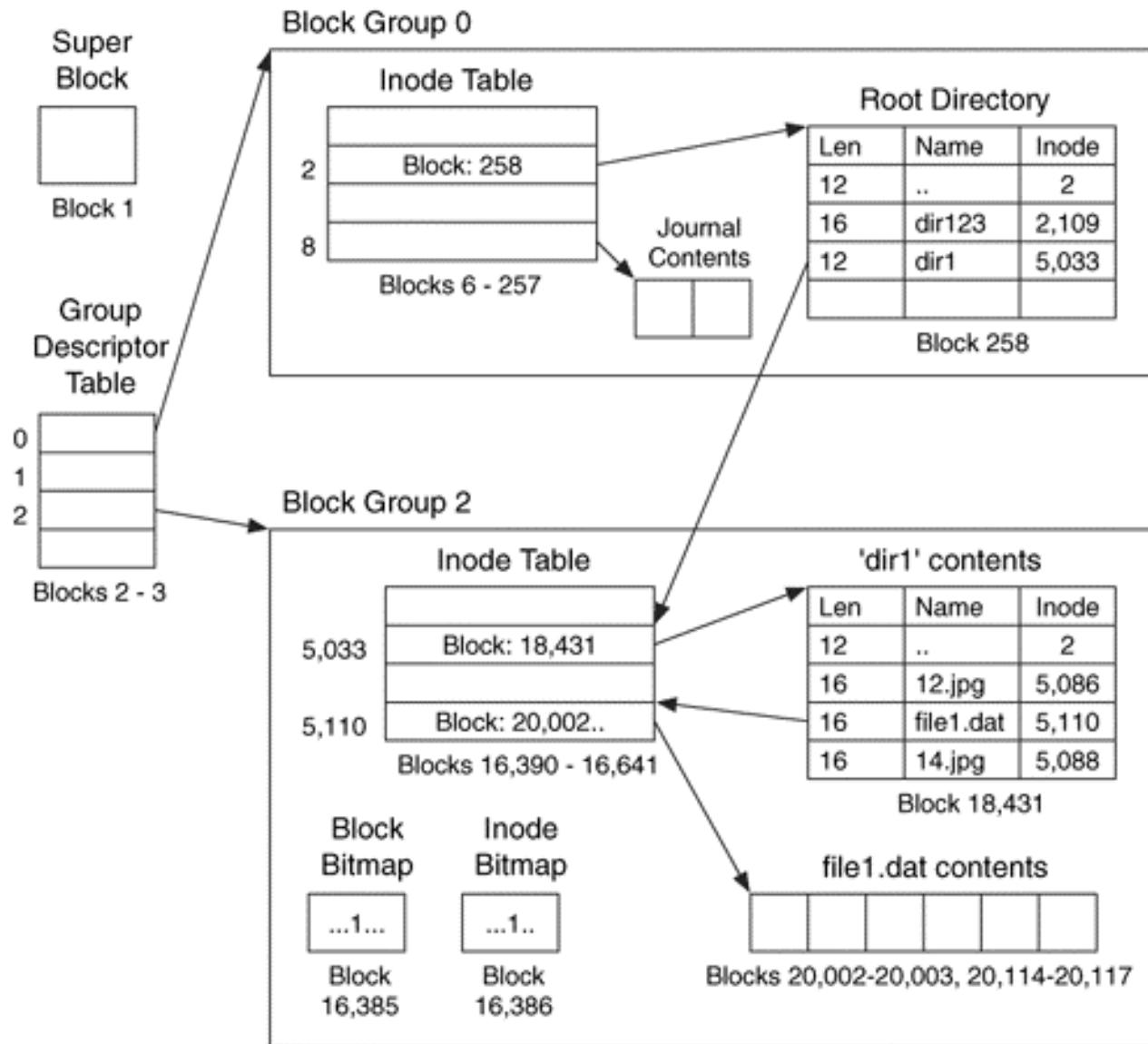
inodes



inodes

- Allocation
 - if a new inode is for a non-directory file, allocate inode in the same block group as parent directory
 - unless there is no free inode or block
 - search groups to find one whose free inodes or blocks are less than the average of free inodes and blocks per group
 - if that search fails, find group with smallest number of directories
 - why is all this necessary?

Example



File Systems

Files

Directories

Implementation

Examples

- Virtual File System
- CD-Rom FS (ISO 9660)
- MS-DOS
- NTFS
- Ext2
- Ext3
- Distributed File Systems

What Distributed File Systems Provide

- Access to data stored at servers using file system interfaces
- What are the file system interfaces?
 - Open a file, check status of a file, close a file
 - Read data from a file
 - Write data to a file
 - Lock a file or part of a file
 - List files in a directory, create/delete a directory
 - Delete a file, rename a file, add a symlink to a file
 - etc

Why DFSs are Useful

- Data sharing among multiple users
- User mobility
- Location transparency
- Backups and centralized management

Distributed systems

- Key differences from centralized systems
 - No shared memory
 - Communication: delays, unreliable
 - No common clock
 - Independent node failure modes
 - Hardware/software heterogeneity

DFS Structure

- **Service** – software entity running on one or more machines and providing a particular type of function to a priori unknown clients
- **Server** – service software running on a single machine
- **Client** – process that can invoke a service using a set of operations that forms its client interface
 - A client interface for a file service is formed by a set of primitive file operations (create, delete, read, write)
 - Client interface of a DFS should be transparent, i.e., not distinguish between local and remote files

Naming Properties

- Naming – mapping between logical and physical objects
- Location transparency:
 - Name of the file does not reveal the file's physical storage location
- Location independence:
 - Name of file does not need to be changed when file's physical location changes
- A file may be replicated in several sites
 - Mapping hides both the existence of multiple copies and their location

Two naming schemes

- Files named by combination of their host name and local name; guarantees a unique system-wide name
 - Neither location transparent
 - Nor location independent
- “Attach” remote directories to local directories, giving the appearance of a coherent directory tree
 - NFS

What is NFS?

- First commercially successful network file system:
 - Developed by Sun Microsystems for their diskless workstations
 - Designed for robustness and “adequate performance”
 - Sun published all protocol specifications
 - Many many implementations
 - Widely used today

➤ <http://nfs.sourceforge.net/>

NFS Design Objectives

- Machine and Operating System Independence
 - Could be implemented on low-end machines of the mid-80's
- Transparent Access
 - Remote files should be accessed in exactly the same way as local files
- Fast Crash Recovery
 - No global synchronization
- “Reasonable” performance
 - Robustness and preservation of UNIX semantics were much more important

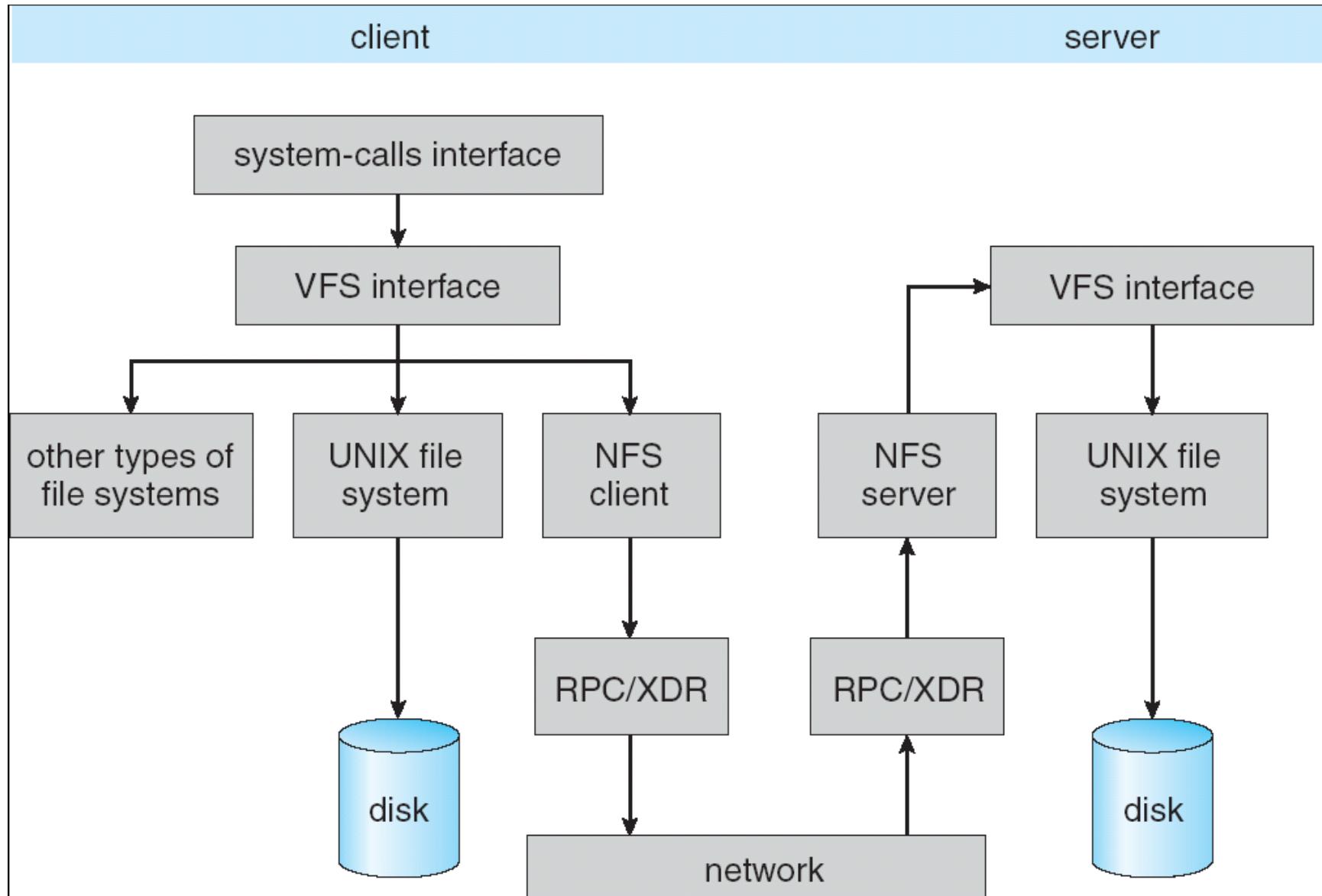
NFS Protocol

- 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
- Nearly one-to-one correspondence between regular UNIX system calls and the NFS protocol RPCs (except opening and closing files)

Three Major Layers in NFS

- Linux 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
 - 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)

NFS Architecture

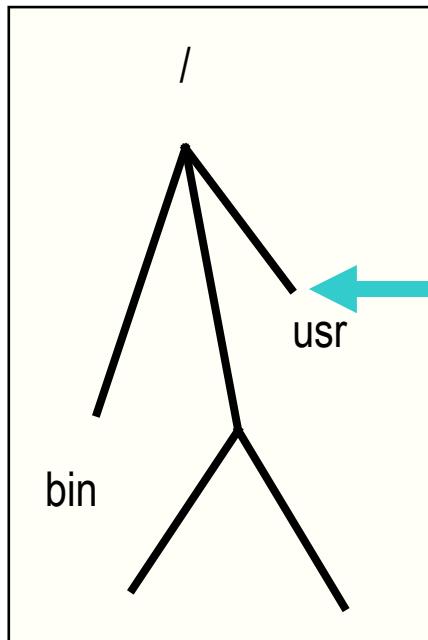


Client Side

- Transparent interface to NFS achieved by
 - Mapping between remote file names and remote file addresses via remote mount
 - Extension of UNIX mount
 - Specified in a mount table
 - Makes a remote subtree appear part of a local subtree

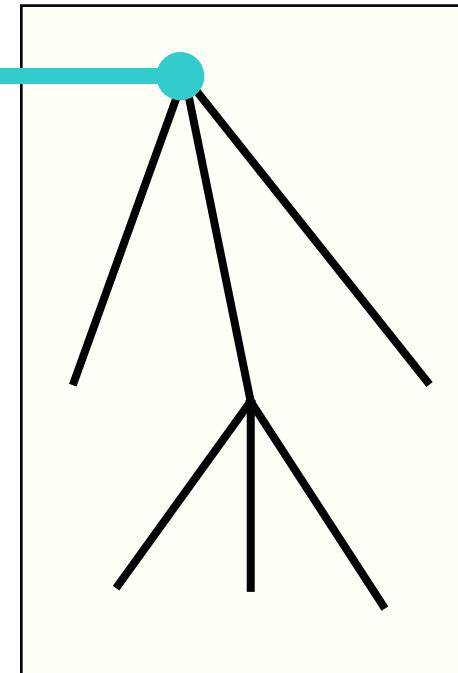
Client side: Remote Mount

Client tree



rmount

Server subtree



**After rmount, root of server subtree
can be accessed as /usr**

Mount Table on os

```
[icrk@os /]$ cat /etc/mtab
rootfs / rootfs rw 0 0
proc /proc proc rw,relatime 0 0
sysfs /sys sysfs rw,relatime 0 0
devtmpfs /dev devtmpfs rw,nosuid,size=499808k,nr_inodes=124952,mode=755 0 0
devpts /dev/pts devpts rw,relatime,gid=5,mode=620,ptmxmode=000 0 0
tmpfs /dev/shm tmpfs rw,relatime 0 0
tmpfs /run tmpfs rw,nosuid,nodev,mode=755 0 0
/dev/vda1 / ext4 rw,relatime,data=ordered 0 0
tmpfs /sys/fs/cgroup tmpfs rw,nosuid,nodev,noexec,mode=755 0 0
cgroup /sys/fs/cgroup/systemd cgroup rw,nosuid,nodev,noexec,relatime,release_agent=/lib/systemd/systemd-cgroups-agent,name=systemd 0 0
cgroup /sys/fs/cgroup/cpuset cgroup rw,nosuid,nodev,noexec,relatime,cpuset 0 0
cgroup /sys/fs/cgroup/cpu,cpuacct cgroup rw,nosuid,nodev,noexec,relatime,cpuacct,cpu 0 0
cgroup /sys/fs/cgroup/memory cgroup rw,nosuid,nodev,noexec,relatime,memory 0 0
cgroup /sys/fs/cgroup/devices cgroup rw,nosuid,nodev,noexec,relatime,devices 0 0
cgroup /sys/fs/cgroup/freezer cgroup rw,nosuid,nodev,noexec,relatime,freezer 0 0
cgroup /sys/fs/cgroup/net_cls cgroup rw,nosuid,nodev,noexec,relatime,net_cls 0 0
cgroup /sys/fs/cgroup/blkio cgroup rw,nosuid,nodev,noexec,relatime,blkio 0 0
cgroup /sys/fs/cgroup/perf_event cgroup rw,nosuid,nodev,noexec,relatime,perf_event 0 0
systemd-1 /proc/sys/fs/binfmt_misc autofs rw,relatime,fd=25,pgrp=1,timeout=300,minproto=5,maxproto=5,direct 0 0
tmpfs /var/run tmpfs rw,nosuid,nodev,mode=755 0 0
tmpfs /var/lock tmpfs rw,nosuid,nodev,mode=755 0 0
securityfs /sys/kernel/security securityfs rw,relatime 0 0
configfs /sys/kernel/config configfs rw,relatime 0 0
mqueue /dev/mqueue mqueue rw,relatime 0 0
tmpfs /media tmpfs rw,nosuid,nodev,noexec,relatime,mode=755 0 0
debugfs /sys/kernel/debug debugfs rw,relatime 0 0
hugetlbfs /dev/hugepages hugetlbfs rw,relatime 0 0
sunrpc /var/lib/nfs/rpc_pipefs rpc_pipefs rw,relatime 0 0
nfs.cs.siu.edu:/home/ /home nfs4 rw,relatime,vers=4.0,rsize=1048576,wsize=1048576,namlen=255,hard,proto=tcp,port=0,timeo=600,retrans=2,sec=sys,clientaddr=146.163.150.17,local_lock=none,addr=146.163.150.227 0 0
binfmt_misc /proc/sys/fs/binfmt_misc binfmt_misc rw,relatime 0 0
```

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, forwarded to mount server running on server
 - Export list – specifies local file systems that server exports for mounting, along with names of machines that are permitted to mount them
- Following a mount request that conforms to its export list, the server returns a file handle—a key for further accesses

NFS Mount Protocol (2)

- 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

Stateful File Service

- Mechanism
 - Client opens a file
 - Server fetches information about the file from its disk, stores it in its memory, and gives the client a connection id unique to the client and the open file
 - Identifier is used for subsequent accesses until the session ends
 - Server must reclaim the main-memory space used by clients who are no longer active

Stateless File Server

- Avoids state information by making each request self-contained
- Each request identifies the file and position in the file
- No need to establish and terminate a connection by open and close operations

Stateful vs stateless protocols

- Stateful:
 - `fd = open("/aaa/bbb/ccc", ...);`
 - `read(fd, ...)`
- Stateless:
 - `read("/aaa/bbb/ccc", offset, ...);`
 - Each procedure call contains all the information necessary to complete the call
 - Server maintains no “between call” information

Failure Recovery in Stateful

- A stateful server loses all its volatile state in a crash
 - Restore state by recovery protocol based on a dialog with clients, or abort operations that were underway when the crash occurred
 - Server needs to be aware of client failures in order to reclaim space allocated to record the state of crashed client processes (orphan detection and elimination)

Advantages of stateless

- Crash recovery is very easy:
 - When a server crashes, client just resends request until it gets an answer from the rebooted server
 - Client cannot tell difference between a server that has crashed and recovered and a slow server
- Simplifies the protocol
 - Client can always repeat any request

Consequences of stateless

- read and write RPCs must specify offset
 - Server does not keep track of current position in the file
 - Longer request messages
 - Slower request processing
- But user still uses conventional Linux APIs

Performance Optimization: Cache

- Reduce network traffic by retaining recently accessed disk blocks in a cache, so that repeated accesses to the same information can be handled locally
 - Remote-service mechanism is one transfer approach
 - If needed data not already cached, a copy of data is brought from the server to the user
 - Accesses are performed on the cached copy
 - Cache-consistency problem – keeping the cached copies consistent with the master file

Cache Update Policy

- Write-through – write data through to master's copy as soon as they are placed on any cache
 - Reliable
 - Only reads benefit
 - Poor performance
 - Network traffic increase
 - Reads stuck behind writes
 - Higher server load

Cache Update Policy

- Write-back (delayed write)– modifications written to the cache and written through to the server later
 - Data may be overwritten before written back, and so aggregated
 - Poor reliability
 - Unwritten data will be lost whenever client crashes
 - Standard OS mechanisms limit data loss
 - Scan cache at regular intervals and flush dirty blocks since the last scan
 - Write-on-close, writes data back to the server when the file is closed.
- Big benefit for files that are open for long periods and frequently modified

Cache in NFS

