

CS 162 Operating Systems and Systems Programming

Professor: Anthony D. Joseph
Spring 2003

Lecture 17: File Systems and Disk Management

17.0 Main Points

Implementing file system abstraction
Comparison among disk management algorithms

Physical Reality	File System Abstraction
Block oriented	Byte oriented
Physical sector #'s	Named files
No protection	Users protected from each other
Data might be corrupted if machine crashes	Robust to machine failures

17.1 File System Components

Disk management: how to arrange collection of disk blocks into files

Naming: user gives file name, not track 50, platter 5, etc.

Protection: keep information secure

Reliability/durability: when system crashes, lose stuff in memory, but want files to be durable.

17.2 User vs. System View of a File

User's view:

- Durable data structures.

Systems' view (system call interface):

- Collection of bytes (UNIX)

System's view (inside OS):

- Collection of blocks (a **block** is a logical transfer unit, while a sector is the physical transfer unit. Block size \geq sector size; in UNIX, block size is 4KB.)

17.2.1 Translating from user to system view

What happens if user says: give me bytes 2 – 12?

- a. Fetch block corresponding to those bytes
- b. Return just the correct portion of the block

What about: write bytes 2 – 12?

- a. Fetch block
- b. Modify portion
- c. Write out block

Everything inside file system is in whole size blocks. For example, `getc`, `putc` \Rightarrow buffers 4096 bytes, even if interface is one byte at a time.

From now on, file is collection of blocks.

17.3 File Usage Patterns

How do users access files?

1. Sequential access – bytes read in order (give me the next X bytes, then give me next)
2. Random access – read/write element out of middle of array (give me bytes i – j)
3. Content-based access – “find me 100 bytes starting with JOSEPH”

Many file systems don't provide #3; instead, database is built on top to index content (requires efficient random access)

How are files typically used?

1. Most files are small (for example, `.login`, `.c` files)

2. Large files use up most of the disk space

3. Large files account for most of the bytes transferred to/from disk

Bad news: need everything to be efficient.

- Need small files to be efficient, since lots of them.
- Need large files to be efficient, since most of the disk space, most of the I/O due to them

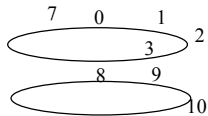
17.4 Disk Management Policies

How do we organize file on disk?

17.4.1 Common data structures

Need a “file header,” one for each file: which disk sectors are associated with each file.

Also, need bitmap to represent free space on disk, one bit per block (or sector). Blocks numbered in cylinder-major order, so that adjacent numbered blocks can be accessed without seeks or rotational delay.



Track 0, surface 0, sector 0, 1, ... | surface 1, sector 0, 1 ... | ... | track 1, surface 1, sector 0...

Caching: every OS today keeps a cache of recently used disk blocks in memory, to avoid having to go to disk. Common to all organizations. For now, assume no cache, and add it later.

17.4.2 Contiguous allocation

User says in advance how big file will be

Search bit map (using best fit/first fit) to locate space for file

File header contains:

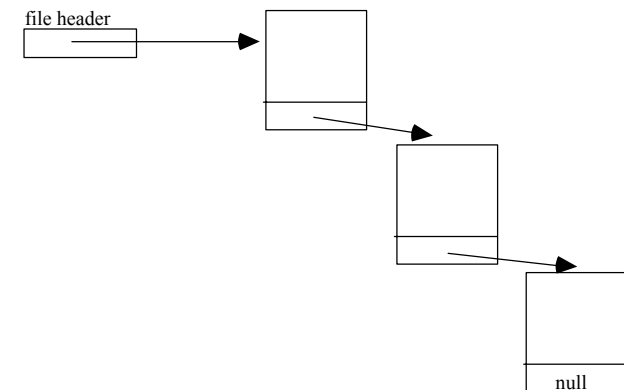
- First sector in file
- File size (# of sectors)

Pros & cons:

- + Fast sequential access
- + Easy random access
- External fragmentation
- Hard to grow files

17.4.3 Linked files

Each block, pointer to next on disk (Alto)



File header pointer to first block on disk

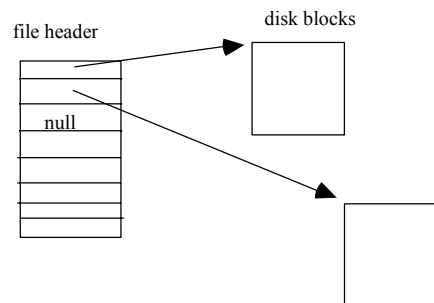
Pros & Cons:

- + Can grow files dynamically
- + Free list managed same as file
- Sequential access: seek between each block
- – Random access: horrible
- Unreliable (lose block, lose rest of file)

MS-DOS used a similar linked approach but instead of embedding links in pages, they used a separate structure called the File Allocation Table (FAT). The FAT has an entry for each block on the disk and the entries corresponding to the blocks of a particular file are linked up.

17.4.4 Indexed files (Nachos, VMS)

User declares max file size; system allocates a file header to hold an array of pointers big enough to point to file size number of blocks.



Pros & Cons:

- + Can easily grow up to space allocated for descriptor
- + Random access is fast
- Clumsy to grow file bigger than table size
- Still lots of seeks: blocks can be spread all over the disk, so sequential access is slow.

17.4.5 Multilevel indexed (UNIX 4.1)

Key idea: efficient for small files, but still allow big files

File header contains 13 pointers (fixed size table, pointers not all equivalent) (the header is called an “inode” in UNIX)

First ten are pointers to data blocks. (If file is small enough, some pointers will be NULL.)

What if you allocate 11th block?

Pointer to an indirect block – a block of pointers to data blocks. Gives us 256 blocks, + 10 (from file header) = 1/4 MB

What if you allocate a 267th block?

Pointer to a doubly indirect block – a block of pointers to indirect blocks (in turn, block of pointers to data blocks). Gives us about 64K blocks => 64MB

What if want a file bigger than this? One last pointer – what should it point to?

Instead, pointer to triply indirect block – block of pointers to doubly indirect blocks (which are ...)

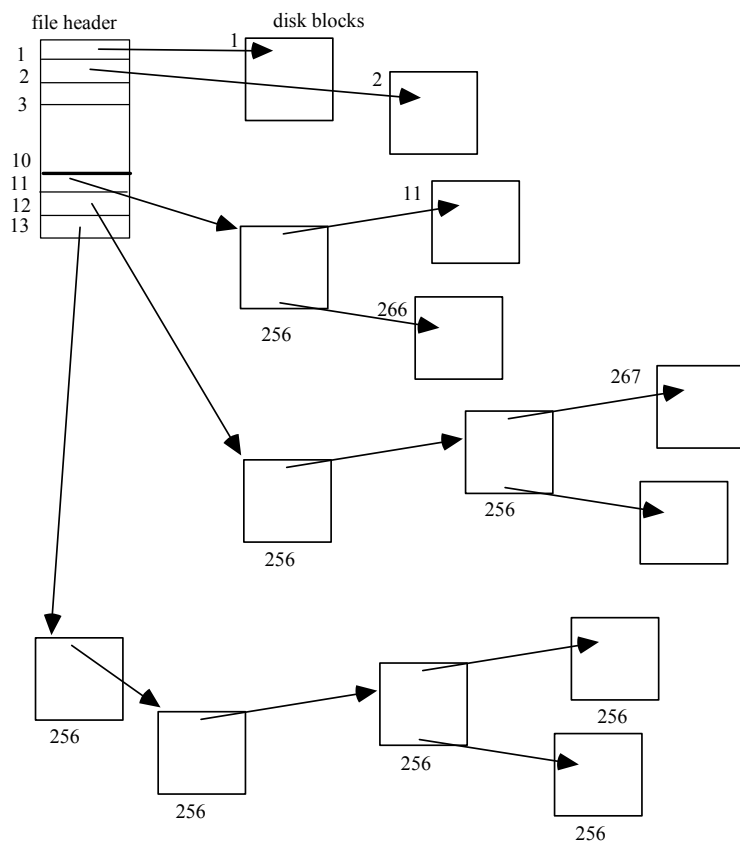
Thus, file header is:

First 10 data block pointers – point to one block each, so 10 blocks

11 indirect block pointer – points to 256 blocks

12 doubly indirect block pointer – points to 64K blocks

13 triply indirect block pointer – points to 16M blocks



1. Bad news: Still an upper limit on file size ~ 16 GB.
2. Pointers get filled in dynamically: need to allocate indirect block only when file grows > 10 blocks. If small file, no indirection needed.
3. How many disk accesses to reach block #23? Which are they?

How about block # 5?

How about block # 340?

UNIX Pros & Cons:

- + Simple (more or less)
- + Files can easily expand (up to a point)
- + Small files particularly cheap and easy
- Very large files, spend lots of time reading indirect blocks
- Lots of seeks

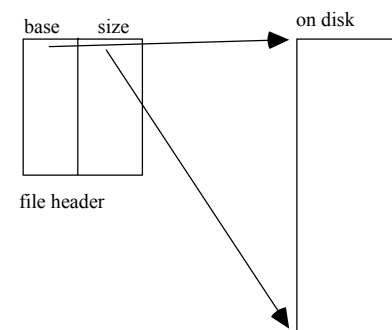
17.4.6 DEMOS

OS for Cray-1, mid to late 70's. File system approach corresponds to segmentation.

Cray-1 had 12 ns cycle time, so CPU:disk speed ratio about the same as today (a few million instructions = 1 seek).

Idea: reduce disk seeks by using contiguous allocation in normal case, but allow flexibility to have non-contiguous allocation.

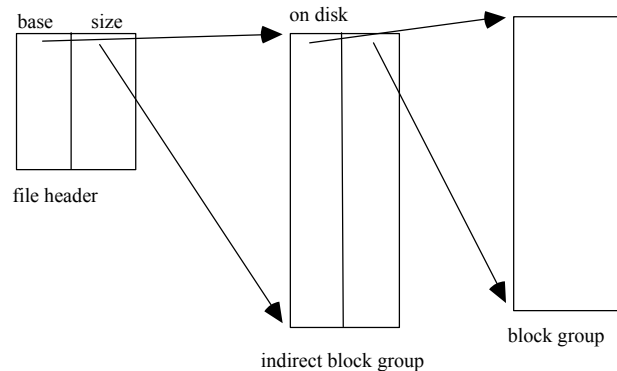
File header: table of base & size (10 "block group" pointers)



Each "block group" – a contiguous region of blocks

Are 10 block group pointers enough? No. If need more than 10 block groups, set flag in file header: BIGFILE.

Each table entry now points to an indirect block group – a block group of pointers to block groups.



Can get huge files this way: Suppose 1000 blocks in a block group (can be bigger or smaller) => 80 GB max file size

How do you find an available block group? Use bit map to find block of 0's

Pros & cons:

- + Easy to find free block groups
- + Free areas merge automatically
- When disk fills up:
 - a. No long runs of blocks (fragmentation)
 - b. High CPU overhead to find free block

In practice, disks are always full.

Solution:

Don't let disk get full – keep portion in reserve

Free count = # blocks free in bitmap.

Normally, don't even try allocate if free count = 0.

Change this to: don't allocate if free count < reserve

Why do this?

Tradeoff: pay for more disk space, get contiguous allocation

How much of a reserve do you need?

In practice, 10% seems like enough.

17.4.7 UNIX BSD 4.2

Exactly the same as BSD 4.1 (same file header and triply indirect blocks), except incorporated some ideas from DEMOS:

- Uses bitmap allocation in place of free list
- Attempt to allocate files contiguously
- 10% reserved disk space
- Skip sector positioning

Problem: when you create a file, don't know how big it will become (in UNIX, most writes are by appending to the file). So how much contiguous space do you allocate for a file, when it's created?

In Demos, power of 2 growth: once it grows past 1 MB, allocate 2MB, etc.

In BSD 4.2, just find some range of free blocks, put each new file at the front of a different range. When need to expand a file, you first try successive blocks in bitmap.

Also, rotational delay can cause a problem, even with sequential files.

Read one block, do processing, and read next block. In the meantime, disk has continued turning. If have to wait for entire rotation, problem! Go from reading at disk bandwidth, to reading one sector every rotation.

Two solutions:

- Skip sector positioning (BSD 4.2)
- Read ahead/disk track buffers – read next block right after first, even if application hasn't asked for it yet. This could be done either by OS (read ahead) or by disk itself (track buffers).

17.5 Disk scheduling

Disk can do only one request at a time; what order do you choose to do requests?

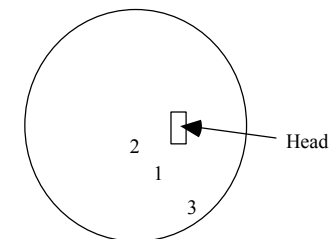
If 0 or 1 request is queued, the choice is easy. But what if more than 1? Try to arrange requests in some order that reduces seek time.

17.5.1 FIFO order

Fair among requesters, but order of arrival may be to random spots on the disk => long seeks

17.5.2 SSTF

SSTF: shortest seek time first. Pick the request that's closest on the disk. (Although called SSTF, today, include rotational delay in calculation, since rotation can be as long as seek.)



Order requests will be serviced using SSTF.

SSTF is good at reducing seeks, but may get starvation

17.5.3 SCAN

SCAN implements an elevator algorithm: take the closest request in the direction of travel. No starvation, but retains flavor of SSTF.

Circular-Scan (C-SCAN): only goes in one direction --- it skips any requests on the way back. This is fairer than SCAN, which is biased towards pages in the middle of the disk.