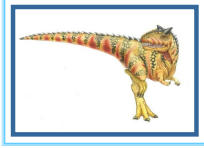


Chapter 14: File-System Implementation



Chapter 14: Implementing File Systems

- File-System Structure
- File-System Implementation
- Directory Implementation
- Allocation Methods
- Free-Space Management



Objectives

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

文件的分配方法有三种主要类型：连续分配、链接分配和索引分配。连续分配提供最佳性能，但可能导致外部碎片；链接分配避免了外部碎片，但不支持随机访问；索引分配通过索引块集中管理指针，支持随机访问，但可能浪费空间。



File-System Structure

- Disks provide most of the secondary storage on which file systems are maintained.
- **Two characteristics of disks** make them convenient for this usage:
 - ① A disk can be rewritten in place; it is possible to read a block from the disk, modify the block, and write it back onto the same place on the disk
 - ② A disk can access directly any block of information it contains. Thus it is simple to access any file either sequentially or randomly, and switching from one file to another requires only moving the read-write heads and waiting for disk to rotate – discussed in Chapter 11
- To improve I/O efficiency, I/O transfers between memory and disk **are performed in units of blocks**. Each **block** contains one or more sectors. A sector size varies from 32 bytes to 4,096 bytes (4KB), usually 512 bytes (0.5 KB)

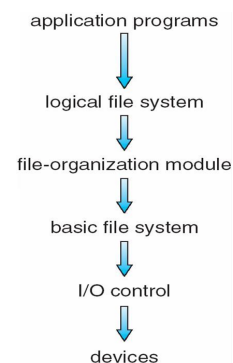


File-System Structure (Cont.)

- File structure
 - Logical storage unit
 - Collection of related information
- **File system** resides on secondary storage – hard drive or disks
 - It provides efficient and convenient **access** to disk by allowing data to be stored, located, and retrieved easily
 - It provides **user interface**: file and file attributes, operations on files, directory for organizing files **界面**
 - It provides **data structures** and algorithms for mapping logical file system onto physical secondary storage devices
- File systems are organized into different layers



Layered File System





File System Layers

- **I/O control** and **device drivers** manage I/O devices at the I/O control layer
 - It consists of **device drivers** and **interrupt handlers** to transfer information between memory and disks
 - Given commands like "read drive 1, cylinder 72, track 2, sector 10" (disk physical address), into memory location 1060" outputs **low-level hardware specific commands** to hardware controller
- **Basic file system** issues generic commands to the appropriate device driver to read and write physical blocks on the disk
 - given commands like "retrieve block 123", translates it to a specific device driver
 - It manages memory buffers and caches that hold various file-system, directory, and data block. (allocation, freeing, replacement)
 - Buffers hold data in transit. A block in memory buffer is allocated before the transfer of a disk block occurs.
 - Caches hold frequently used file-system metadata to improve performance
- **File organization module** knows files, and their logical blocks, as well as physical blocks
 - Translates logical block # (address) to physical block #, pass this to basic file system for transfer
 - Manages free disk space, disk block allocation



File System Layers (Cont.)

- **Logical file system** manages metadata information
 - **Metadata** includes all of the file-system structures except the actual data, i.e., the contents of files
 - It manages directory structure to provide the information needed by file-organization module.
 - Translates file name into file number or file handle, location by maintaining **file control blocks**
 - A **file control block (FCB)** (called **inode** in Unix file systems) contains all information about a file including ownership, permissions, and location of the file contents (on the disk)
 - It is also responsible for file protection
- **分层的优点:** Layering useful for **reducing complexity and redundancy**, but adds overhead and can decrease performance **缺点:**
- Many file systems are in use today, and most operating systems support more than one file system
 - Each with its own format - CD-ROM is ISO 9660; Unix has **UFS** (Unix File System) based on FFS; Windows has FAT, FAT32, NTFS (or Window NT File System) as well as floppy, CD, DVD Blu-ray, Linux has more than 40 types of file systems, with **extended file system** ext2 and ext3; plus distributed file systems, etc.
 - New ones still arriving – ZFS, GoogleFS, Oracle ASM, FUSE



File-System Implementation

Several **on-disk and in-memory structures** are used to implement a file system.

- **On-disk structure**, it may contain information about how to boot an operating system stored on the disk, the total number of blocks, number and location of free blocks, directory structure, and individual files **引导**
- **In-memory information** used for both file-system management and performance improvement via **caching**. The data are loaded at mount time, updated during file-system operations, and discarded at dismount time. **文件系统管理** **挂载**



On-Disk File-System Structure

- **启动** **Boot control block** (per volume **卷**) contains info needed by system to boot OS from that volume
 - If the disk does not contain an OS, this block can be empty
 - Usually the first block of a volume. In UFS, it is called the **boot block**. In NTFS, it is the **partition boot sector** **分区引导扇区**
- **Volume control block** (per volume) contains volume (or partition) details
 - Total # of blocks, # of free blocks, block size, free block count and pointers, a free FCB count and pointer
 - In UFS, this is called **superblock**. In NTFS, it is stored in the **master file table**
- **目录结构** **A directory structure** (per file system) is used to organize files
 - In UFS, this includes file names and associate inode numbers (FCB in Unix). In NTFS, it is stored in the master file table
- **Per-file File Control Block (FCB)** contains many details about a file
 - It has a unique identifier number to associate with a directory entry.
 - In UFS, **inode** number, permissions, size, dates
 - NTFS stores into in master file table using relational DB structure, with a row per file



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



In-Memory File System Structures

The in-memory information is used for both file-system management and performance improvement **via caching**. The data are loaded at mount time, updated during file-system operations, and discarded at dismount

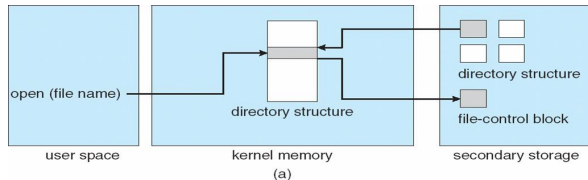
- An in-memory **mount table** contains information about each mounted volume
- An in-memory **directory-structure cache** holds the directory information of recently accessed directories.
- The **system-wide open-file table** contains a copy of the FCB of each open file to locate the files, as well as other information
- The **per-process open-file table** contains a pointer to the appropriate entry in the system-wide open-file table, as well as other information such as per-process file protection and access rights.
- Buffers hold file-system blocks when they are being read from disk or written to disk





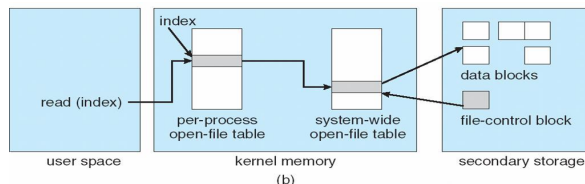
In-Memory File System Structures

- Figure below refers to opening a file
 - The **open()** operation passes a file name to the logical file system.
 - If the file is already in use by another process, only the **per-process open-file table entry** is created pointing to the corresponding entry of this file in the **system-wide open-file table**.
 - If not, it searches the **directory structure** (part of it may be cached in memory). Once the file is found, an entry in system-wide open-file table and in per-process open-file table are created, respectively. The **FCB** is copied into the system-wide open-file table in memory for subsequent access.
 - The system-wide open-file table **not only stores the FCB but also tracks the number of processes that have open**, and thus are using this file – **file count**.
 - The other fields in the per-process open-file table may include a pointer to the current location in the file (for next read() or write() operation) and access mode in which the file is open.



In-Memory File System Structures

- Figure below refers to reading a file
 - The **open()** call returns a pointer to the appropriate entry in the per-process open-file table. In another word, **open()** operation creates the corresponding entries in both per-process and system-wide open-file tables, and subsequent operations use this pointer to locate file content without the need to access directory structure – this speeds up the subsequent operations on the file.
 - All file operations are then performed via this pointer. UNIX systems refer to it as a **file descriptor**; Windows refers to it as a **file handle**.
 - Data from read() eventually copied to specified user process memory address (part of process address space)



Directory Implementation 目录实现

- The selection of directory-allocation and directory management algorithms significantly affects the efficiency, performance, and reliability of the file system.
- Linear list** of file names with pointer to the data blocks
 - Simple to program
 - Time-consuming to execute
 - The major disadvantage of a linear list is that finding a file requires a linear search time.
 - Cache in memory the frequently used directory information
 - Could keep ordered alphabetically via linked list or use B+ tree
- Hash Table** – linear list with hash data structure
 - Decreases directory search time
 - Collisions** – situations where two file names hash to the same location
 - Only good if entries are fixed size, or use **chained-overflow** method
 - The major difficulties with a hash table are its generally fixed size and the dependence of the hash function on that size.

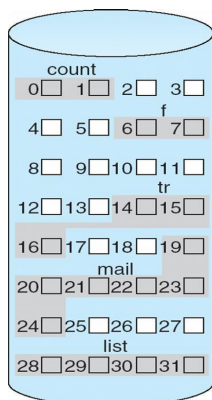


Allocation Methods - Contiguous

- An **allocation method** refers to how disk blocks are allocated for files, such that the disk space is utilized effectively, and files can be accessed quickly.
- There are three major methods of allocating disk space that are widely in use, **contiguous**, **linked** and **indexed**.
- Contiguous allocation** – each file occupies a set of contiguous blocks of the disk
 - Best performance in most cases – support sequential and direct access easily
 - Simple – only starting location (block #) and length (number of blocks) are required
 - Problems with finding space for a new file, and when file size grows
 - This is also a **dynamic storage-allocation problem** discussed earlier, which involves how to satisfy a request of size n (variable) from a list of free variable-sized holes – external fragmentation exists
 - Best-fit and first-fit are common strategies and shown to be more efficient than worst-fit.
 - The cost of compaction is particularly high for large disk, which may take hours. Some system require that compaction be done only when off-line, with the file system unmounted
 - This is preferred for files that the file size must be known at the time of file creation – overestimation leads to large amount of internal fragmentation



Contiguous Allocation of Disk Space

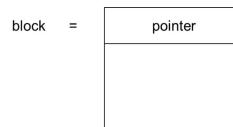


directory		
file	start	length
count	0	2
tr	14	3
mail	19	6
list	28	4
f	6	2



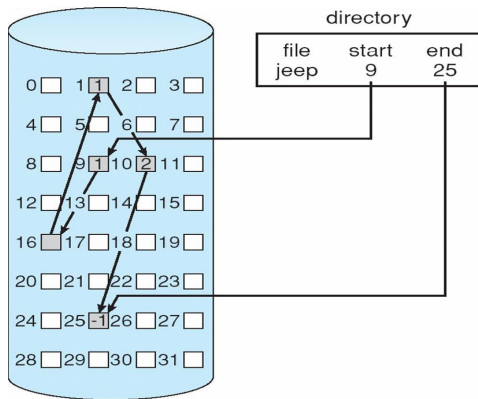
Allocation Methods - Linked

- Linked allocation** – each file consists of a linked-list of blocks
 - Each file is a linked list of disk blocks, which may be scattered anywhere on the disk
 - The directory contains a pointer to the first and last blocks of a file
 - File ends at null pointer (the end-of-list pointer value)
 - Each block contains pointer to next block
 - No compaction needed, and no external fragmentation
 - A file can continue to grow as long as free blocks are available
- It is inefficient to support direct access of the file, only good for sequential access
- Extra disk space required for the pointers. If a pointer requires 4 bytes out of a 512-byte block, then 0.78% of the disk space is being used for pointers.
- Reliability can be a problem; for instance, what happen if a pointer is lost or damaged.





Linked Allocation



Allocation Methods - FAT (文件分配表)

- FAT (File Allocation Table)** – an important variation on **linked allocation**
 - This simple but efficient method of disk space allocation was used by the MS-DOS.
 - A section of disk at the beginning of volume is set aside to contain a table called **FAT**.
 - The table has one entry for each disk block and is indexed by block number.
 - The FAT is used in much the same way as a linked list.
 - The directory entry contains the block number of the first block of the file. The table entry indexed by that block number contains the block number of the next block in the file. This continues until it reaches the last block, which has a special end-of-file value as the table entry.
 - An unused block is indicated by a table entry value 0. Allocating a new block to a file is a simple matter of finding the first 0-value table entry.
 - FAT can be cached. Random (direct) access time is improved, because the disk head can find the location of any block by reading the information in the FAT, instead of moving through blocks stored on the disk in the linked-allocation scheme.

4.2 (3 points) What are the advantages of using FAT over a linked allocation?

Answer:

Linked allocation performs the worst when accessing a block that is stored at the middle or near the end of a file as it needs to access all of the individual blocks of the file in a sequential manner to find the pointer to the target block starting from the first block. This may involve considerable overhead in moving disk arm and waiting for rotation. This is much easier in FAT by chasing the pointers stored within the FAT, with no or significantly less disk arm movement. In addition, most of the FAT can be cached in memory and therefore the pointers can be determined with just memory accesses instead of having to access the disk blocks.

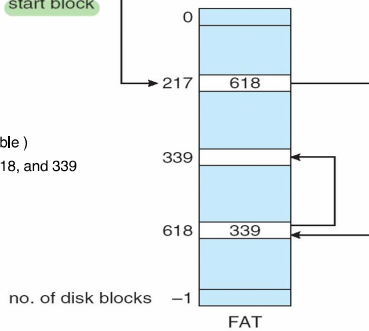


File-Allocation Table

directory entry

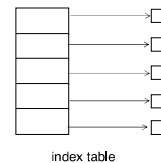
test	...	217
name		start block

- An example of FAT (File Allocation Table)
- The file consists of disk blocks 217, 618, and 339



Allocation Methods - Indexed

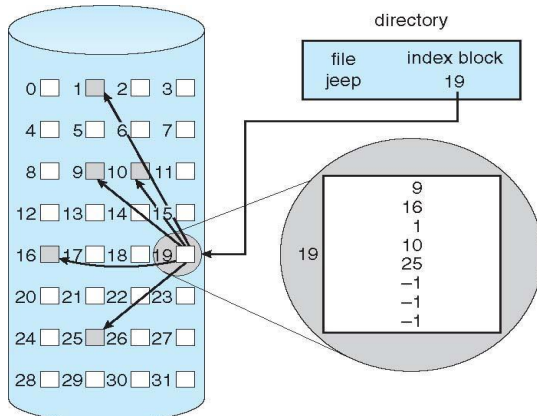
- Indexed allocation** – brings all the pointers together into one location, the **index block**
 - Each file has its own **index block**, which contains an array of pointers to its data blocks, or disk-block addresses.
- Logical view



index table



Example of Indexed Allocation

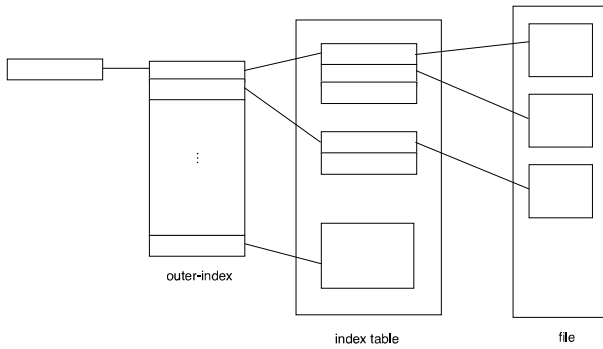


Indexed Allocation 索引分配

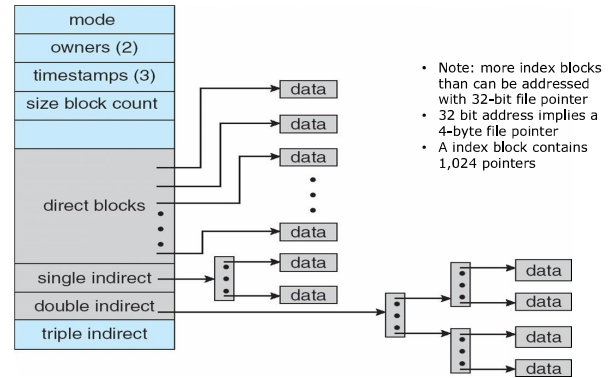
- Indexed allocation supports **direct access**, without suffering from **external fragmentation**, because any free block on the disk can satisfy a request for more space.
- Indexed allocation suffers from some of the same performance problem as linked allocation. Specifically, **index block(s)** can be cached in memory, but data blocks may still be spread all over a volume (disk).
- Indexed allocation **does suffer from wasted space**. The pointer overhead of the index block is generally greater than the pointer overhead in the linked allocation.
 - Suppose a file only has one or two blocks. Indexed allocation loses an entire index block, while linked allocation loses the space of only one pointer per block.
- So far, an index block occupies one disk block. What happens if the file is too large such that one index block is too small to hold enough pointers?
 - Linked scheme** – to link several together index blocks
 - 多级 Multilevel scheme** – a first-level index block points to a set of second-level index blocks, which in turn point to the file blocks. This could be continued to a third or fourth level, depending on the desired maximum file size. With a 4,096-byte block, we could store 1,024 **four-byte** pointers in one index block. Two levels of index allow 1,048,576 data blocks, and a file size up to 4GB.
 - Combined scheme** – **direct blocks** for **small files**, and **indirect blocks** (single indirect, double indirect, and triple indirect blocks) for **larger files**, used in UNIX-based file systems.



Indexed Allocation – Multilevel Scheme



Combined Scheme: UNIX UFS (4K bytes per block, 32-bit addresses)



自由空间管理 Free-Space Management

- File system maintains **free-space list** to track free disk space
 - (Using term "block" for simplicity)
- Bit vector** or **bit map** (n blocks)



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

- The main advantage is its simplicity and efficiency in finding the first free blocks or n consecutive free blocks on the disk
- Bit vector is **inefficient** unless the entire vector can be kept in memory, but requires extra space

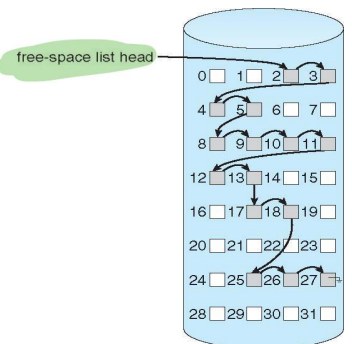
特殊:

block size = 4KB = 2^{12} bytes
disk size = 2^{40} bytes (1 terabyte)
 $n = 2^{40}/2^{12} = 2^{28}$ bits (or 32MB)

Linked Free Space List on Disk

Linked-List - link together all the free disk blocks

- Keeping a pointer to the first free block, in a special location on the disk, can be cached in memory
- The first block contains a pointer to the next free blocks, and so on
- Easily locate one free block, not easy to obtain contiguous blocks
- It is not efficient to traverse the entire list, since it must read each block, which requires substantial I/O time. Fortunately, this is not a frequent action



Free-Space Management (Cont.)

- Grouping**
 - Modify linked-list to store addresses of n free blocks in first free block. The first n-1 of these blocks are free. The last block contains addresses of another n free blocks, and so on.
 - The addresses of a large number of free blocks can be found more quickly than linked-list
- Counting** - Because several contiguous blocks may be allocated and freed simultaneously, particularly when contiguous-allocation algorithm or extents is used
 - By taking advantage of this, rather than keeping a list of n free disk addresses, we can keep address of first free block and count of following contiguous free blocks
 - Each entry in the free-space list consists of a disk address and a count
 - Although each entry requires more space than would a simple disk address, the overall list is shorter, as long as the count is generally greater than one
 - Note this method of tracking free space is similar to the extent method of allocating blocks.

End of Chapter 14

