

操作系统原理与设计

第 11 章文件系统的实现

陈香兰

中国科学技术大学计算机学院

May 21, 2014

提纲

- 1 File-System Structure
- 2 FS Implementation
- 3 Directory Implementation
- 4 Allocation Methods (分配方法)
- 5 Free-Space Management
- 6 Efficiency and Performance
- 7 Recovery
- 8 Log Structured File Systems
- 9 小结和作业

Outline

- 1 File-System Structure
- 2 FS Implementation
- 3 Directory Implementation
- 4 Allocation Methods (分配方法)
- 5 Free-Space Management
- 6 Efficiency and Performance
- 7 Recovery
- 8 Log Structured File Systems
- 9 小结和作业

File-System Structure

- File structure
 - Logical storage unit
 - Collection of related information
- FS resides on secondary storage (disks)
- FS organization
 - How FS should look to the user
 - How to map the logical FS onto the physical secondary-storage devices
- FS organized into layers

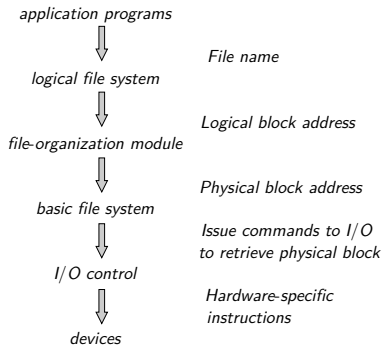


Figure: Layered File System

Outline

- 1 File-System Structure
- 2 FS Implementation**
- 3 Directory Implementation
- 4 Allocation Methods (分配方法)
- 5 Free-Space Management
- 6 Efficiency and Performance
- 7 Recovery
- 8 Log Structured File Systems
- 9 小结和作业

FS Implementation

- Structures and operations used to implement file system operation, OS- & FS-dependment
 - ① On-disk structures
 - ② In-memory structures

① On-disk structures

① Boot control block

- To boot an OS from the partition (volume)
- If empty, no OS is contained on the partition

② Volume control block

③ Directory structure

④ Per-file FCB

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

Figure: A typical file control block

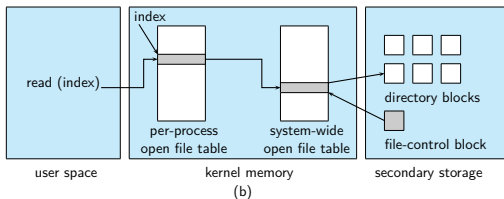
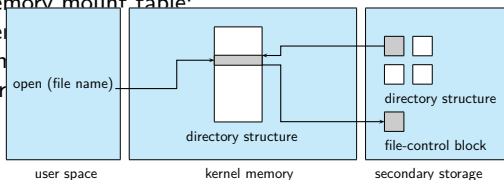
- ② **In-memory information:** For both FS management and performance improvement via caching
 - Data are loaded at mount time and discarded at dismount
 - Structures include:
 - in-memory mount table;
 - in-memory directory-structure cache
 - system-wide open-file table;
 - per-process open-file table

FS Implementation

② In-memory information: For both FS management and performance improvement via caching

- Data are loaded at mount time and discarded at dismount
- Structures include:

- in-memory mount table
- in-memory directory structure
- system-wide open file table
- per-process open file table

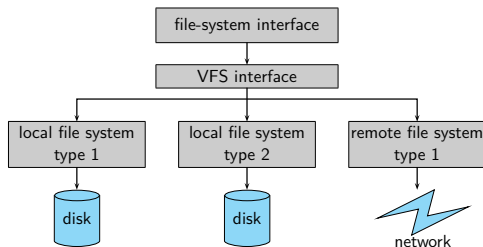


Partitions and mounting

- Partition (分区)
 - Raw (E.g. UNIX swap space & some database) VS. cooked
 - **Boot** information, with its own format
 - Boot image
 - Boot loader unstanding multiple FSes & OSes
Dual-boot
- **Root partition** is mounted at boot time
- **Others** can be automatically mounted at boot or manually mounted later

Virtual File Systems (虚拟文件系统)

- **Virtual File Systems (VFS, 虚拟文件系统)** 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.



Schematic View of Virtual File System

Outline

- 1 File-System Structure
- 2 FS Implementation
- 3 Directory Implementation**
- 4 Allocation Methods (分配方法)
- 5 Free-Space Management
- 6 Efficiency and Performance
- 7 Recovery
- 8 Log Structured File Systems
- 9 小结和作业

Directory Implementation

- ① **Linear list** of file names with pointer to the data blocks.
 - Simple to program
 - Time-consuming to execute
- ② **Hash Table** – linear list with hash data structure.
 - Decreases directory search time
 - **Collisions** – situations where two file names hash to the same location
 - Fixed & variable size or chained-overflow hash table

Outline

- 1 File-System Structure
- 2 FS Implementation
- 3 Directory Implementation
- 4 Allocation Methods (分配方法)**
- 5 Free-Space Management
- 6 Efficiency and Performance
- 7 Recovery
- 8 Log Structured File Systems
- 9 小结和作业

Allocation Methods (分配方法)

- An allocation method refers to **how disk blocks are allocated** for files so that disk space is **utilized effectively** & files can be **accessed quickly**
 - ① Contiguous allocation (连续分配)
 - ② Linked allocation (链接分配)
 - ③ Indexed allocation (索引分配)
 - ④ Combined (组合方式)

1. Contiguous Allocation (连续分配) I

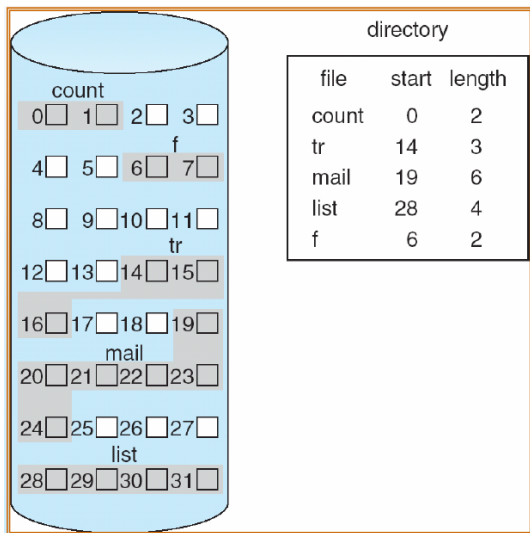
- Each file occupies a set of **contiguous** blocks on the disk
- Simple - directory entry only need
 - **starting location (block #)**
 - **& length (number of blocks)**
- Mapping from logical to physical

$$\text{LogicalAddress}/512 = Q \dots R$$

$$\text{Block to be accessed} = Q + \text{starting address}$$

$$\text{Displacement into block} = R$$

1. Contiguous Allocation (连续分配) II



1. Contiguous Allocation (连续分配) III

- Advantages:
 - Support **both random & sequential** access
 - Start block: b ;
Logical block number: i
 \Rightarrow physical block number: $b + i$
 - **Fast** access speed, because of short head movement
- Disadvantages:
 - **External fragmentation**
 - **Wasteful of space** (dynamic storage-allocation problem).
 - **Files cannot grow**,
or File size must be known in advance.
 \Rightarrow **Internal fragmentation**

- 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.

2. Linked Allocation (链接分配)

- Each file is a linked list of disk blocks: blocks may be scattered anywhere on the disk.
- Two types
 - ① **Implicit (隐式链接)**
 - ② **Explicit (显式链接)**

2. Linked Allocation (链接分配)

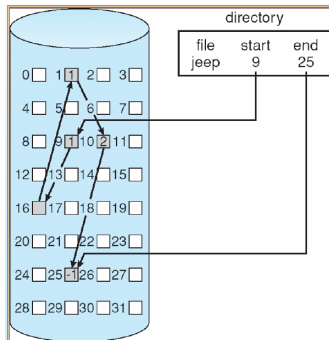
① Implicit (隐式链接)

- **Directory** contains a pointer to the first block & last block of the file.
- **Each block** contains a pointer to to the next block.

a block =

pointer

- **Allocate as needed, link together**
 - Simple – need only starting address
 - Free-space management system – **no waste of space**



2. Linked Allocation (链接分配)

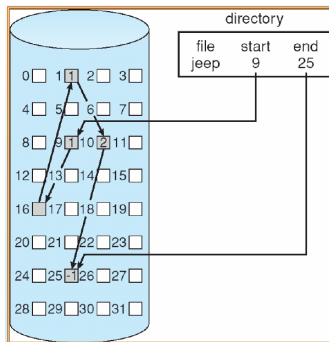
① Implicit (隐式链接)

• Disadvantage:

- No random access
- Link pointers need disk sapce
E.g.: 512 per block, 4 per pointer \Rightarrow 0.78%

Solution: clusters

\Rightarrow disk throughput \uparrow
But internal fragmentation \uparrow



2. Linked Allocation (链接分配)

① Implicit (隐式链接)

• Mapping:

Suppose

- ① block size=512bytes,
- ② block pointer size=1byte, using the first byte of a block
- ③ Logical address in the file to be accessed = A

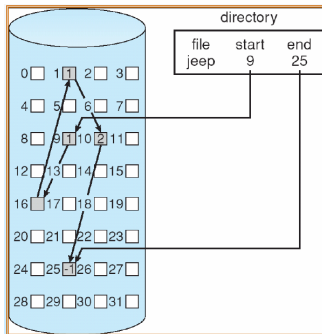
we have

- ① Data size for each block = $512 - 1 = 511$
- ② $A/511 = Q \dots R$

then

- ① Block to be accessed is the Q^{th} block in the linked chain of blocks representing the file.
- ② Displacement into block = $R + 1$

• How to reduce searching time?

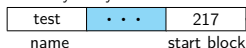


2. Linked Allocation (链接分配)

② **Explicit linked allocation: File Allocation table, FAT**

Disk-space allocation used by MS-DOS and OS/2

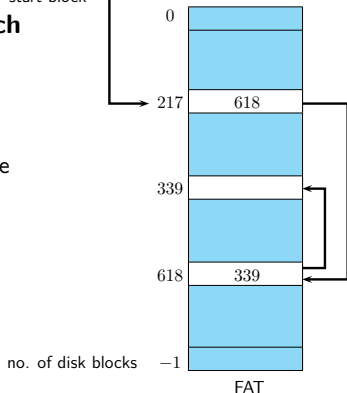
directory entry



- **A section of disk** at the beginning of each partition is set aside to contain the **FAT**

- Each disk block one entry
- The entry contains
 - (1) **the index of the next block** in the file
 - (2) **end-of-file**, for the last block entry
 - (3) **0**, for unused block

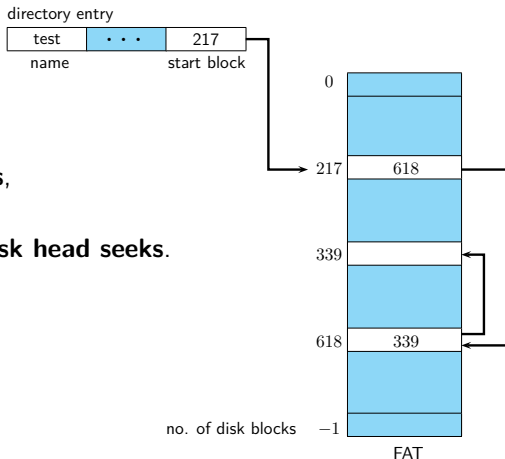
- **Directory entry** contains the first block number



2. Linked Allocation (链接分配)

② **Explicit linked allocation: File Allocation table, FAT**

Disk-space allocation used by MS-DOS and OS/2



- Now **support random access**, but still not very efficient
- May result in a **significant disk head seeks**.

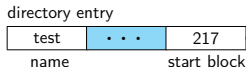
Solution: **Cached FAT**

2. Linked Allocation (链接分配)

② Explicit linked allocation: File Allocation table, FAT

Disk-space allocation used by MS-DOS and OS/2

• How to compute FAT size?



Suppose

- ① Disk space = 80 GB
- ② Block size = 4 KB

Then

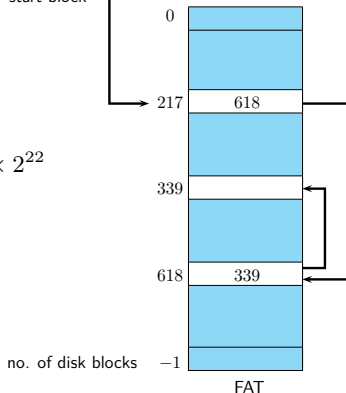
- ① Total block number = $80 \times 2^{30} / 2^{12} = 5 \times 2^{22}$
- ② $4 \times 2^{22} = 2^{24} < 5 \times 2^{22} < 8 \times 2^{22} = 2^{25}$

• Length of each FAT entry?

(25bits? 28bits? 32bits?)

• Length of FAT?

$(5 \times 2^{22} \times 4B = 80MB = 80GB/2^{10})$



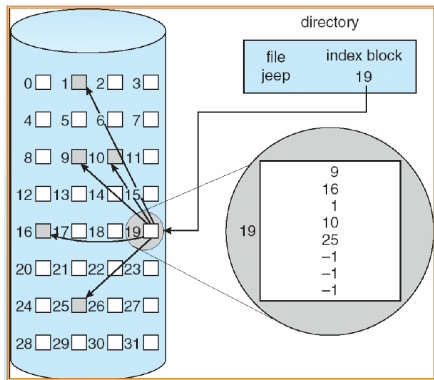
3. Indexed Allocation (索引分配)

- **Indexed Allocation (索引分配):**

Brings all pointers together into one location – **the index block**.

- **Each file** has its own index block
- **Directory entry** contains the index block address
- **Each index block:** An array of pointers (an index table)

*Logical block number i
= the i^{th} pointer*



3. Indexed Allocation (索引分配)

- **Indexed Allocation (索引分配):**

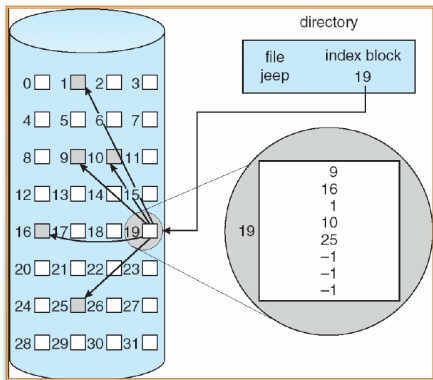
Brings all pointers together into one location – **the index block**.

- **Advantage:**

- **Random** access
- Dynamic access **without external fragmentation**

- **Disadvantage:**

- have **overhead** of index block.
- File **size limitation**, since one index block can contains limited pointers



3. Indexed Allocation (索引分配)

- Indexed Allocation (索引分配):**

Brings all pointers together into one location – **the index block**.

- Mapping** from logical to physical

Suppose

(1) Block size = 1KB

(2) Index size = 4B

Then for logical address LA , we have

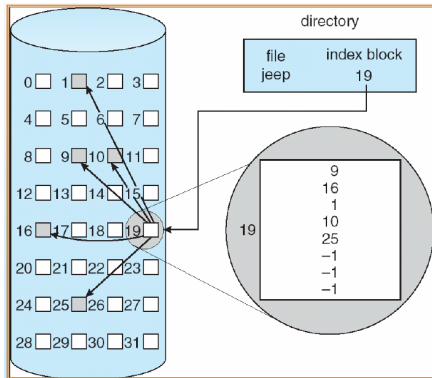
$$LA/512 = Q...R$$

(3) Q = the index of the pointer

(4) R = displacement into block

We also have **Max file size**

$$= 2^{10}/4 \times 1KB = 256KB$$



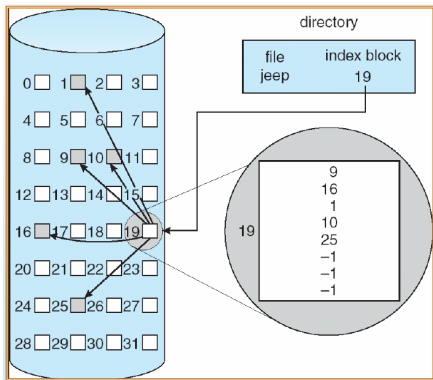
3. Indexed Allocation (索引分配)

- **Indexed Allocation (索引分配):**

Brings all pointers together into one location – **the index block**.

- **How to support a file of unbounded length?**

- ① **linked scheme**
- ② **multi-level index scheme**



3. Indexed Allocation (索引分配)

① **Linked scheme**

- **Link blocks of index table (no limit on size).**
- **Mapping**

Suppose

(1) Block size=1KB

(2) Index or link pointer size = 4B

Then

$$LA / (1KB \times (1K/4 - 1)) = Q_1 \dots R_1$$

(3) Q_1 = block of index table

(4) R_1 is used as follows:

$$R_1 / 1K = Q_2 \dots R_2$$

(5) Q_2 = index into block of index table

(6) R_2 = displacement into block of file:

3. Indexed Allocation (索引分配)

② multi-level index scheme

Example: **Two-level index** (maximum file size is ?)

- We have

$$LA / (1K \times 1K / 4) = Q_1 \dots R_1$$

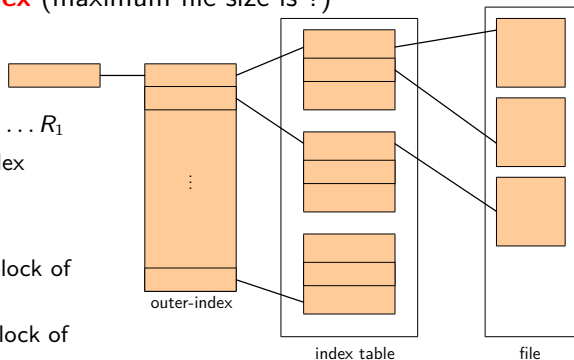
(1) Q_1 = index into outer-index

(2) R_1 is used as follows:

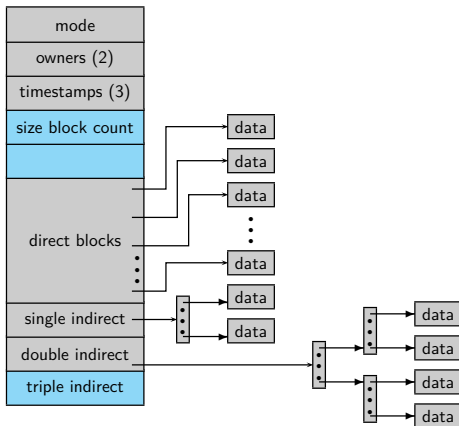
$$R_1 / 1KB = Q_2 \dots R_2$$

(3) Q_2 = displacement into block of index table

(4) R_2 = displacement into block of file



4. Combined Scheme (组合方式): UNIX (4K bytes per block) I



4. Combined Scheme (组合方式): UNIX (4K bytes per block) II

- if 4KB per block, and 4B per entry

$$\text{Direct blocks} = 10 \times 4KB = 40KB$$

$$\text{Number of entries per block} = 4KB/4B = 1K$$

$$\text{Single indirect} = 1K \times 4KB = 4MB$$

$$\text{Double indirect} = 1K \times 4MB = 4GB$$

$$\text{Triple indirect} = 1K \times 4GB = 4TB$$

Maximum file size = ?

Outline

- 1 File-System Structure
- 2 FS Implementation
- 3 Directory Implementation
- 4 Allocation Methods (分配方法)
- 5 Free-Space Management**
- 6 Efficiency and Performance
- 7 Recovery
- 8 Log Structured File Systems
- 9 小结和作业

- **Disk Space**: limited
 - Free space management: **To keep track of free disk space**
 - **How?** Free-space list?
 - **Algorithms**
 - ① **Bit vector**
 - ② **Linked list**
 - ③ **Grouping (成组链接法)**
 - ④ **Counting**

① Bit vector

- **Free-space list** is implemented as a **bit map** or **bit vector**
 - **1 bit for each block**
1=free;
0=allocated
 - Example:
a disk where blocks 2,3,4,5,8,9,10,11,12,13,17,18,25,26,27 are free and the rest blocks are allocated. The bitmap would be

0011 1100 1111 1100 0110 0000 0111 0000 0...

- **Bit map length.**
For n blocks, if the base unit is word, and the size of word is 16 bits, then

$$\text{bit map length} = (n + 15)/16$$

U16 bitMap[bitMapLength];

① Bit vector

- How to find the first free block or n consecutive free blocks on the disk?
 - Many computers supply **bit-manipulation instructions**
 - To find the first free block:
Suppose: base unit = word (**16** bits) or other
 - (1) find **the first non-0 word**
 - (2) find **the first 1 bit** in the first non-0 word
 - If first K words is 0, & $(K + 1)^{th}$ word > 0 ,
the first $(K + 1)^{th}$ word's first 1 bit has offset L ,
then

$$\text{first free block number } N = K \times 16 + L$$

① Bit vector

- Simple
- Must be kept on disk

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)

- **Solution: Clustering**

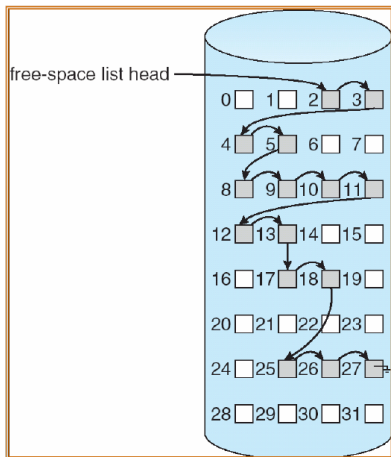
① Bit vector

- **Efficient** to get the first free block or n consecutive free blocks, **if we can always store the vector in memory.**
 - But **copy in memory and disk may differ.**
E.g. $\text{bit}[i] = 1$ in memory & $\text{bit}[i] = 0$ on disk
 - **Solution:**
 - Set $\text{bit}[i] = 1$ in memory.
 - Allocate $\text{block}[i]$
 - Set $\text{bit}[i] = 1$ in disk
- Need to protect:
 - Pointer to free list
 - Bit map

Free-Space Management

② Linked Free Space List on Disk

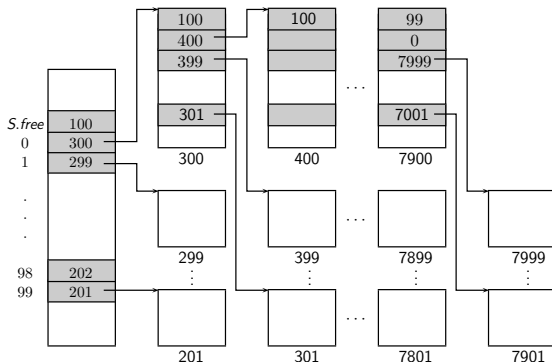
- **Link together all the free disk blocks**
 - **First** free block
 - **Next** pointer
- Not efficient
- **Cannot get contiguous space easily**
- No waste of space



Free-Space Management

③ Grouping(成组链接法)

- To store the addresses of n free blocks (a group) in the first free block
 - First n-1 group members are actually free
 - Last one contain the next group
 - And so on
- E.g.: UNIX



④ Counting

- Assume:
 - Several contiguous blocks may be allocated or freed simultaneously
- Each = first free block number & a counter (number of free blocks)
- Shorter than linked list at most time, counter > 1

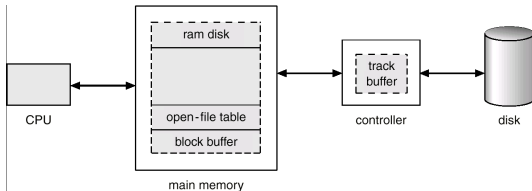
Outline

- 1 File-System Structure
- 2 FS Implementation
- 3 Directory Implementation
- 4 Allocation Methods (分配方法)
- 5 Free-Space Management
- 6 Efficiency and Performance**
- 7 Recovery
- 8 Log Structured File Systems
- 9 小结和作业

- **Efficiency** dependent on:
 - Disk allocation and directory algorithms
 - **Various approaches**
 - Inodes distribution
 - Variable cluster size
 - Types of data kept in file's directory entry
 - Large pointers provides larger file length, but cost more disk space

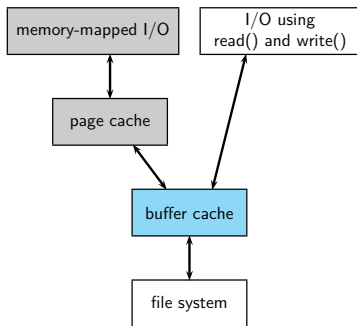
● Performance

- **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**



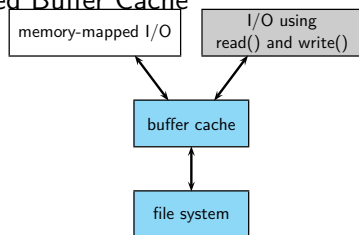
Page Cache I

- 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
- This leads to the following figure
 - I/O Without a Unified Buffer Cache



Unified Buffer Cache

- A unified buffer cache uses the same page cache to cache both memory-mapped pages and ordinary file system I/O
- I/O Using a Unified Buffer Cache



Outline

- 1 File-System Structure
- 2 FS Implementation
- 3 Directory Implementation
- 4 Allocation Methods (分配方法)
- 5 Free-Space Management
- 6 Efficiency and Performance
- 7 Recovery
- 8 Log Structured File Systems
- 9 小结和作业

- Consistency checking (一致性检查)
 - compares data in directory structure with data blocks on disk, and tries to fix inconsistencies
- Backup & restore
 - 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

Outline

- 1 File-System Structure
- 2 FS Implementation
- 3 Directory Implementation
- 4 Allocation Methods (分配方法)
- 5 Free-Space Management
- 6 Efficiency and Performance
- 7 Recovery
- 8 Log Structured File Systems**
- 9 小结和作业

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

Outline

- 1 File-System Structure
- 2 FS Implementation
- 3 Directory Implementation
- 4 Allocation Methods (分配方法)
- 5 Free-Space Management
- 6 Efficiency and Performance
- 7 Recovery
- 8 Log Structured File Systems
- 9 小结和作业

小结

- 1 File-System Structure
- 2 FS Implementation
- 3 Directory Implementation
- 4 Allocation Methods (分配方法)
- 5 Free-Space Management
- 6 Efficiency and Performance
- 7 Recovery
- 8 Log Structured File Systems
- 9 小结和作业

谢谢!