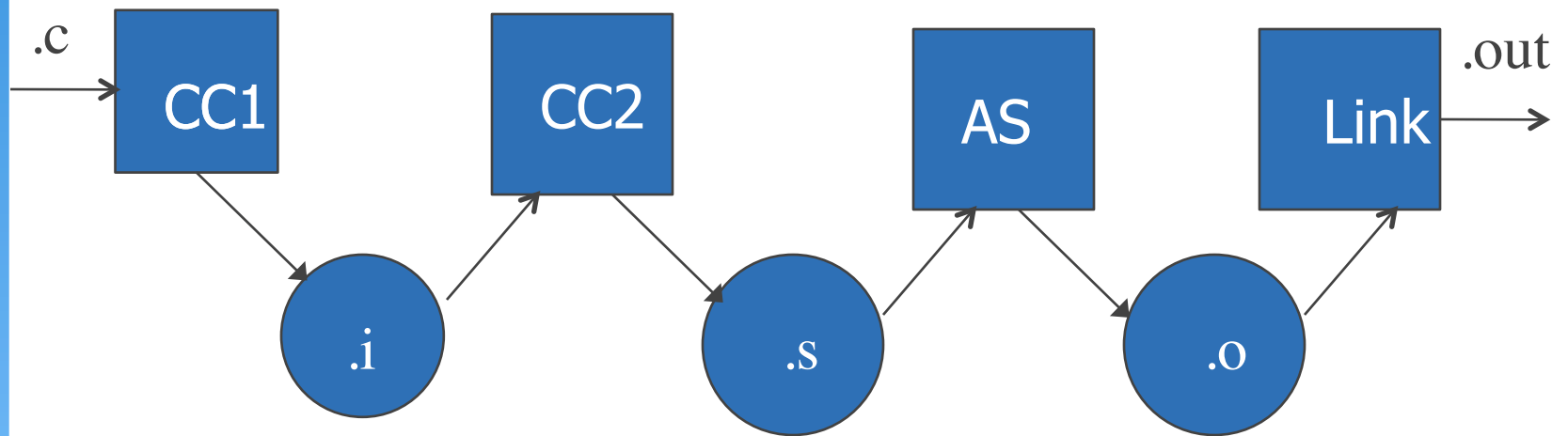


# Tmpfs/memory file system

- ▶ Use virtual memory to build a file system
  - Will not survive reboots
  - Contents might be written back to disk as part of VM
  - Temporary files, that need to survive reboots can be fast because nothing ever goes to disk



`gcc -O hello.c`

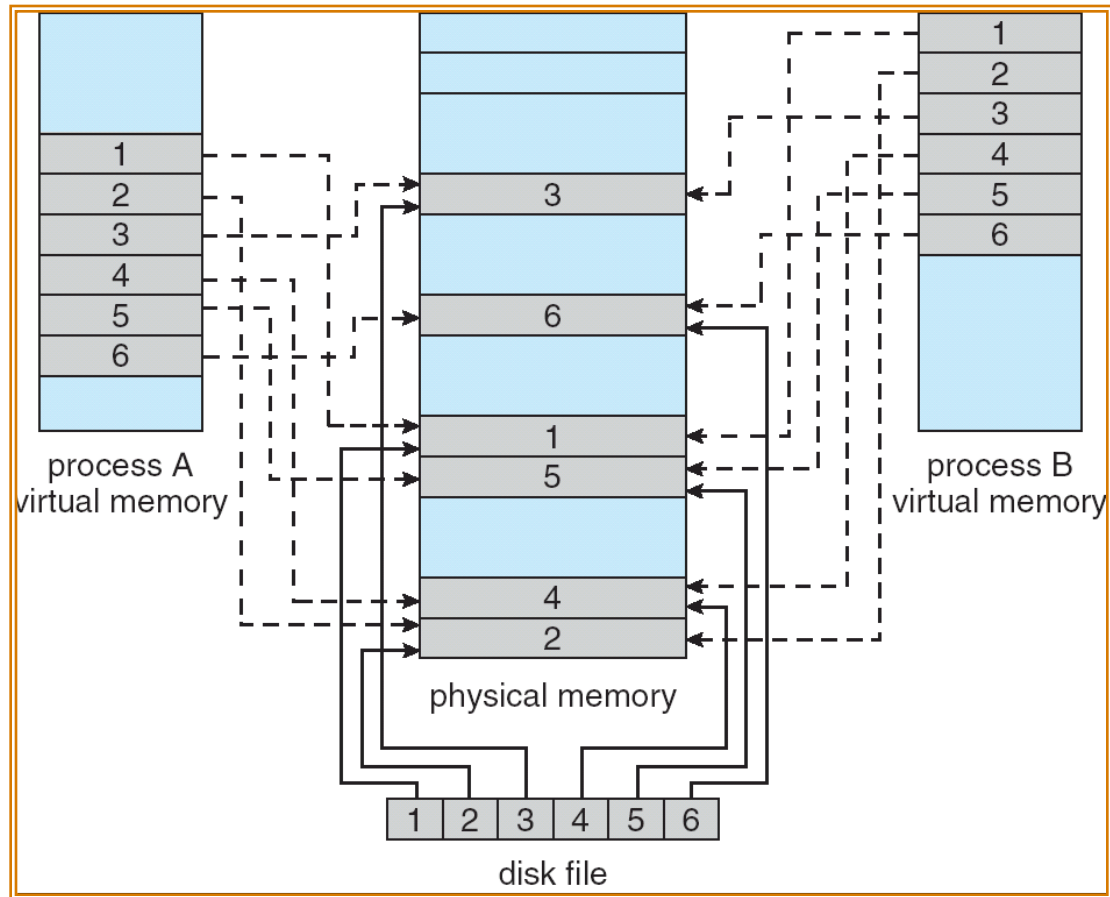


# Memory-Mapped Files

- ▶ Memory-mapped file I/O allows file I/O to be treated as routine memory access by mapping a disk block to a page in memory
- ▶ A file is initially read using demand paging. A page-sized portion of the file is read from the file system into a physical page. Subsequent reads/writes to/from the file are treated as ordinary memory accesses.
- ▶ Simplifies file access by treating file I/O through memory rather than `read()` `write()` system calls
- ▶ Also allows several processes to map the same file allowing the pages in memory to be shared



# Memory Mapped Files



# Sample code using mmap

```
#include <sys/mman.h>
#include <sys/stat.h>
#include <fcntl.h>
#include <unistd.h>
```

```
main(int argc, char *argv[], char *envp[]) {
    int fd;
    char *ptr, *path = (argc == 2) ? argv[1] : "file";

    /* Open a file and write some contents. If file already exists,
       delete old contents */
    fd = open(path, O_WRONLY | O_CREAT | O_TRUNC, 0660);
    write(fd, "hello", strlen("hello"));
    write(fd, " world", strlen(" world"));
    close(fd);
}
```



## (continued)

```
fd = open(path, O_RDWR);

// mmap(addr, len, prot, flags, fildes, off);
ptr = mmap(0, 4, PROT_READ|PROT_WRITE,
  MAP_SHARED, fd, 0);
ptr+=2;
memcpy(ptr, "lp ", 3);
munmap(ptr, 4);
close(fd);
}
```

- ▶ Transform “hello world” into “help world”

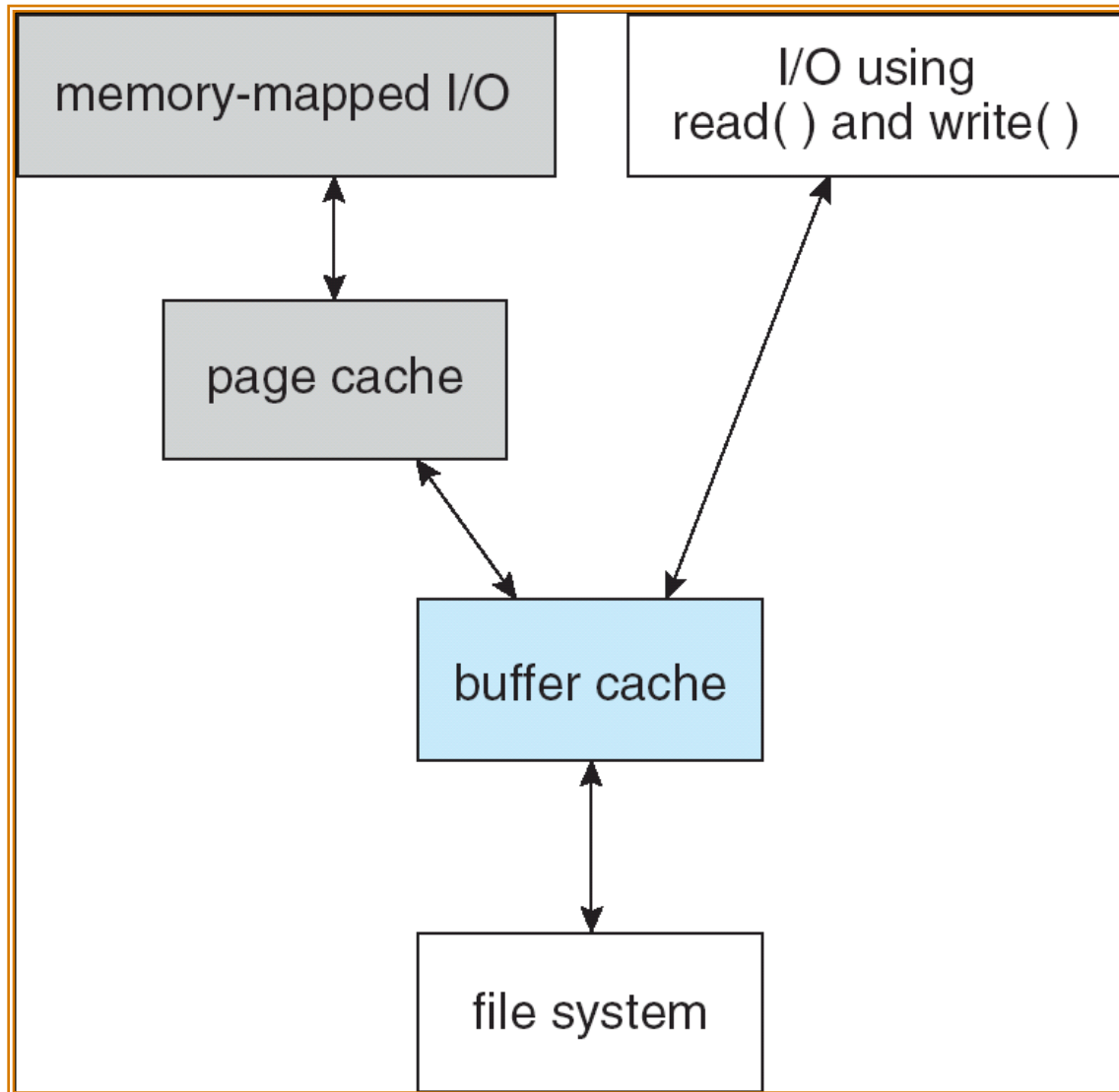


# Page Cache

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



# I/O Without a Unified Buffer Cache



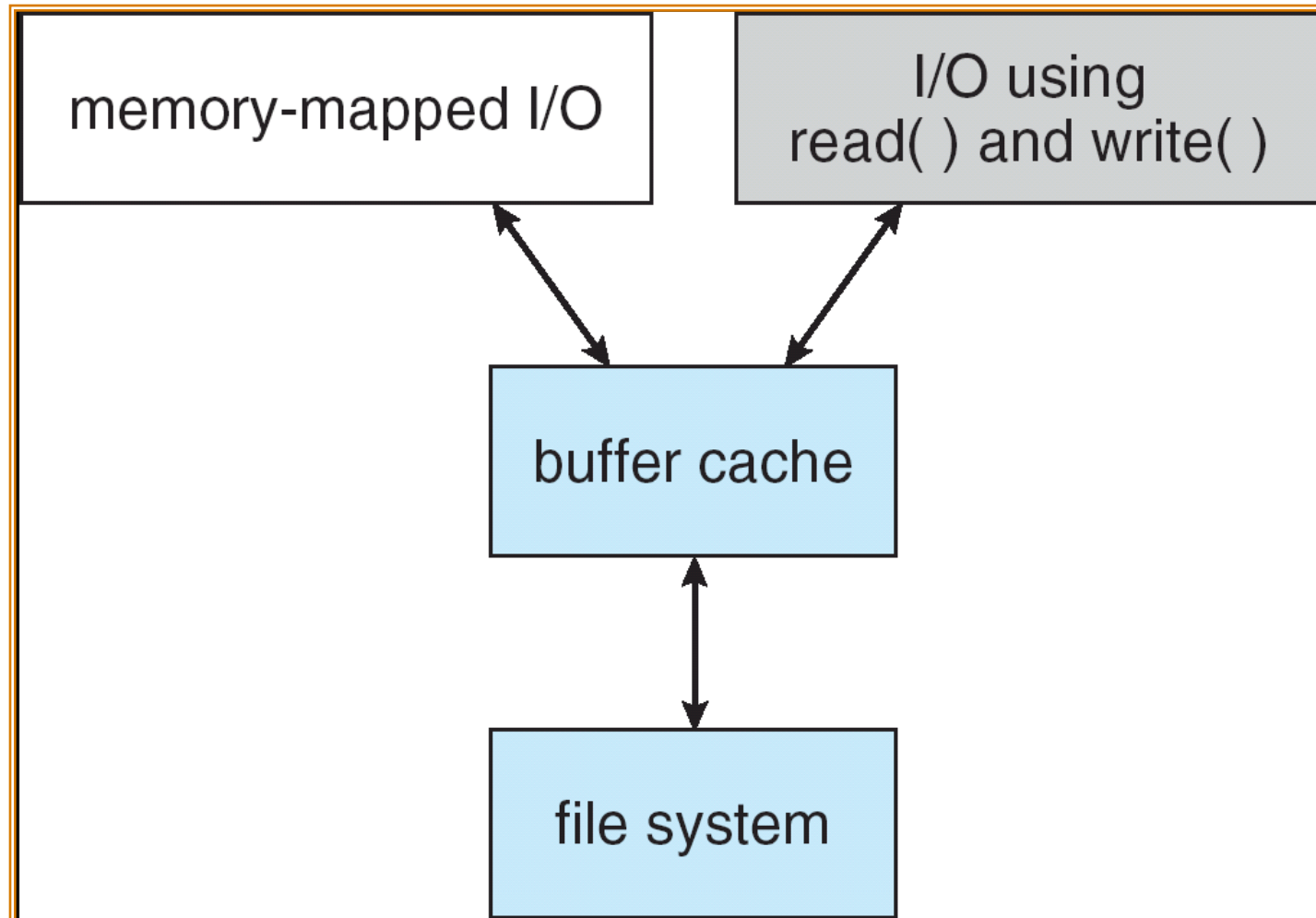
# Unified Buffer Cache

- ▶ 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



# Recovery

- ▶ Consistency checking – compares data in directory structure with data blocks on disk, and tries to fix inconsistencies
  - scandisk in DOS, fsck in unix
- ▶ 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



# 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

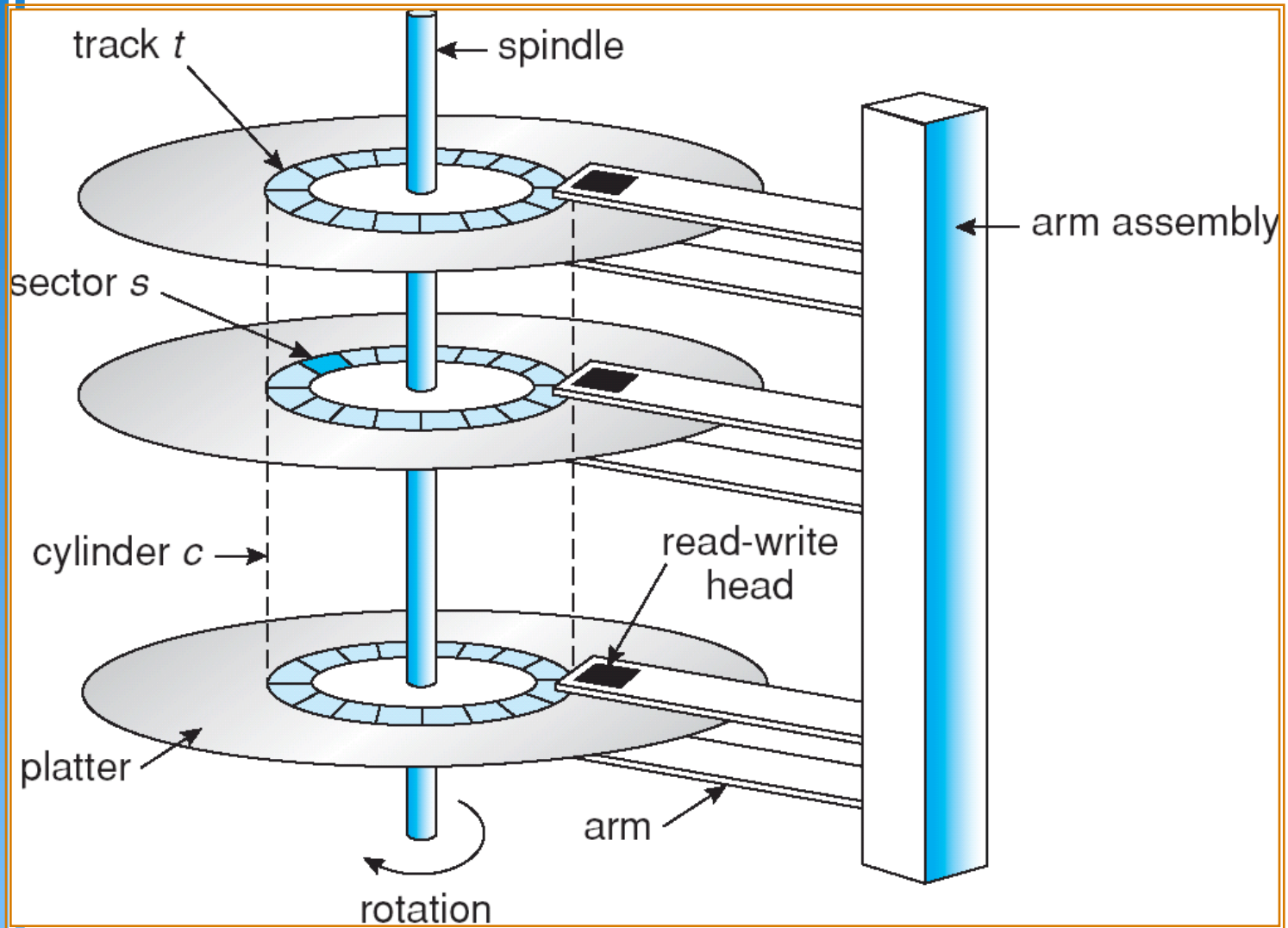


# Overview of Mass Storage Structure

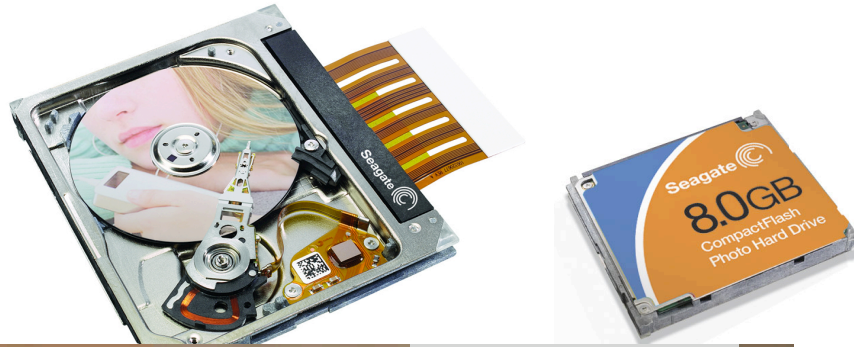
- ▶ Magnetic disks provide bulk of secondary storage
  - Drives rotate at 70 to 250 times per second
    - Ipod disks: 4200 rpm
    - Laptop disks: 4200, 5400 rpm or 7200 rpm
    - Desktop disks: 7200 rpm
    - Server disks: 10000 rpm or 15000 rpm
  - **Transfer rate** is rate at which data flow between drive and computer
  - **Positioning time (random-access time)** is time to move disk arm to desired cylinder (**seek time**) and time for desired sector to rotate under the disk head (**rotational latency**)
  - **Head crash** results from disk head contacting disk surface
    - That's bad
- ▶ Disks can be removable
- ▶ Drive attached to computer via **I/O bus**
  - Busses vary, including **EIDE, ATA, SATA, Firewire, USB, Fibre Channel, SCSI**
  - **Host controller** in computer uses bus to talk to **disk controller** built into drive or storage array



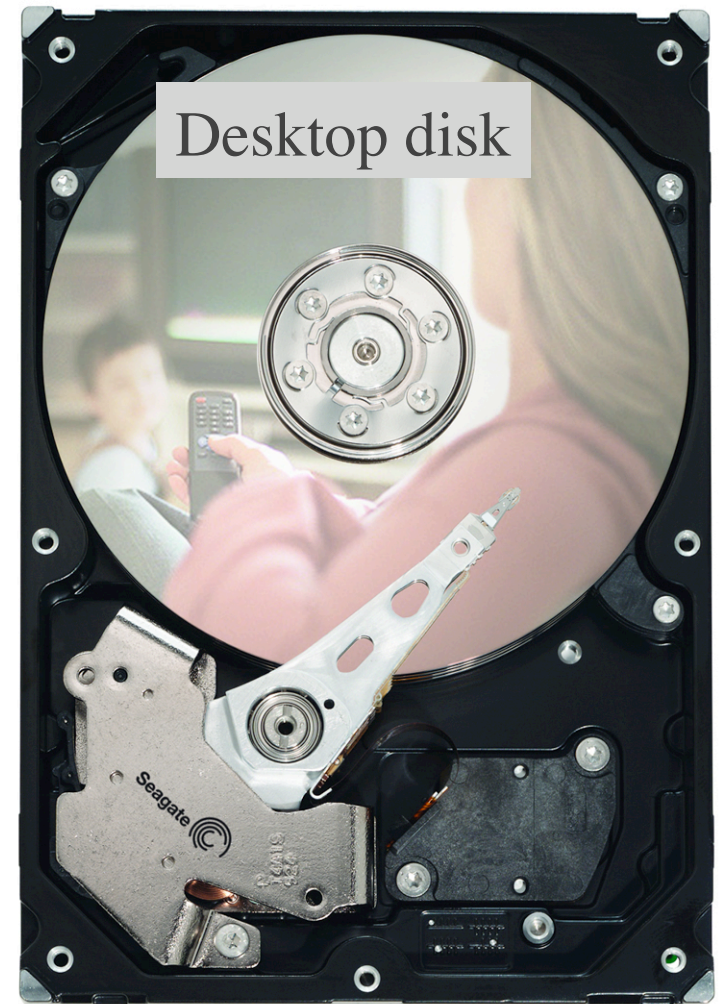
# Moving-head Disk Mechanism



# Disk drives



Server disk



Desktop disk



# Hard disk head, platter and disk crash



# Disk Structure

- ▶ Disk drives are addressed as large 1-dimensional arrays of *logical blocks*, where the logical block is the smallest unit of transfer.
- ▶ The 1-dimensional array of logical blocks is mapped into the sectors of the disk sequentially.
  - Sector 0 is the first sector of the first track on the outermost cylinder.
  - Mapping proceeds in order through that track, then the rest of the tracks in that cylinder, and then through the rest of the cylinders from outermost to innermost.



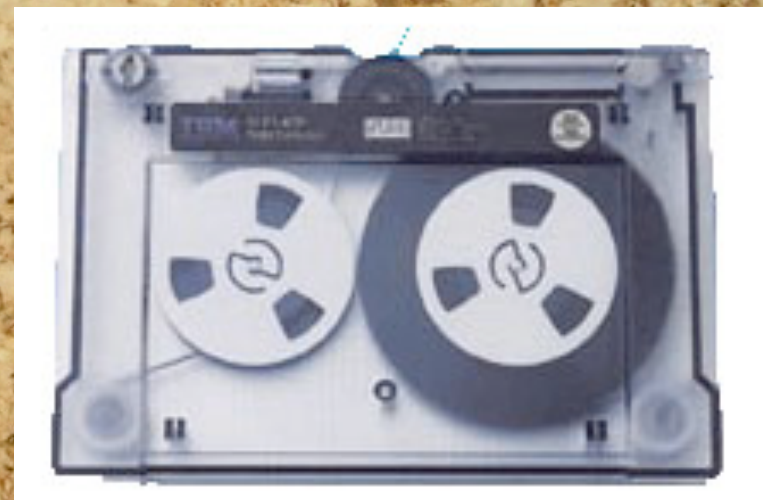


# Magnetic tape

- ▶ Was early secondary-storage medium
- ▶ Relatively permanent and holds large quantities of data
- ▶ Access time slow
- ▶ Random access ~1000 times slower than disk
- ▶ Mainly used for backup, storage of infrequently-used data, transfer medium between systems
- ▶ Kept in spool and wound or rewound past read-write head
- ▶ Once data under head, transfer rates comparable to disk
- ▶ 20-200GB typical storage
- ▶ Common technologies are 4mm, 8mm, 19mm, LTO-2 and SDLT



# Tape pictures



# Tape Drives

- ▶ The basic operations for a tape drive differ from those of a disk drive.
- ▶ **locate** positions the tape to a specific logical block, not an entire track (corresponds to **seek**).
- ▶ The **read position** operation returns the logical block number where the tape head is.
- ▶ The **space** operation enables relative motion.
- ▶ Tape drives are “append-only” devices; updating a block in the middle of the tape also effectively erases everything beyond that block.
- ▶ An EOT mark is placed after a block that is written.



# Application Interface

- ▶ Most OSs handle removable disks almost exactly like fixed disks — a new cartridge is formatted and an empty file system is generated on the disk.
- ▶ Tapes are presented as a raw storage medium, i.e., and application does not not open a file on the tape, it opens the whole tape drive as a raw device.
- ▶ Usually the tape drive is reserved for the exclusive use of that application.
- ▶ Since the OS does not provide file system services, the application must decide how to use the array of blocks.
- ▶ Since every application makes up its own rules for how to organize a tape, a tape full of data can generally only be used by the program that created it.



# Tertiary Storage Devices

- ▶ Low cost is the defining characteristic of tertiary storage.
- ▶ Generally, tertiary storage is built using *removable media*
- ▶ Common examples of removable media are floppy disks and CD-ROMs; other types are available.



# Removable Disks

- ▶ Floppy disk — thin flexible disk coated with magnetic material, enclosed in a protective plastic case.
  - Most floppies hold about 1 MB; similar technology is used for removable disks that hold more than 1 GB.
  - Removable magnetic disks can be nearly as fast as hard disks, but they are at a greater risk of damage from exposure.



# Removable Disks (Cont.)

- ▶ A magneto-optic disk records data on a rigid platter coated with magnetic material.
  - Laser heat is used to amplify a large, weak magnetic field to record a bit.
  - Laser light is also used to read data (Kerr effect).
  - The magneto-optic head flies much farther from the disk surface than a magnetic disk head, and the magnetic material is covered with a protective layer of plastic or glass; resistant to head crashes.
- ▶ Optical disks do not use magnetism; they employ special materials that are altered by laser light.



# WORM Disks

- ▶ The data on read-write disks can be modified over and over.
- ▶ WORM (“Write Once, Read Many Times”) disks can be written only once.
- ▶ Thin aluminum film sandwiched between two glass or plastic platters.
- ▶ To write a bit, the drive uses a laser light to burn a small hole through the aluminum; information can be destroyed by not altered.
- ▶ Very durable and reliable.
- ▶ *Read Only* disks, such as CD-ROM and DVD, come from the factory with the data pre-recorded.

