OS

# Operating Systems

## Lecture 13: File System

Department of Computer Science & Technology
Tsinghua University

- Basic Concepts

- Virtual File System

- Data Block Caching

- Data Structures for Open Files

- File Allocation

- Free-Space List

- Management of Multiple Disks – RAID

◆ File System & File

◆ File Descriptor

◆ Directory

◆ File Aliasing

◆ Types of File System

# File System and File

- File system: an OS abstraction for using persistent storage
  - Π Organizing, manipulating, navigating, accessing, and retrieving data on the persistent storage

- Most computer systems have file systems
  - Π PCs, servers, laptops
  - Π iPod, Tivo/set-top-box, cellphones/PDAs
  - Π Google is made possible by a file system

- File: an OS abstraction for a unit of related data in the file system

# File System Functionality

- ## Allocate disk storage to files
  - Managing **file blocks** (which blocks belong to which file)
  - Managing **free space** (which blocks are free)
  - Allocation **algorithms** (policies)

- ## Manage the collection of files
  - **Locate** files and their contents
  - **Naming**: interface to find files by name
  - Most common: **hierarchical file system**
  - File system type (different ways to organize files)

- ## Provide convenience and features
  - **Protection**: layers to keep data secure
  - **Reliability/Durability**: Keeping of files durable despite crashes, media failures, attacks, etc

# File and Blocks

◆ **File attributes**

    π Name, type, location, size, protection, creator, creation time, last-modified-time, …

◆ **File header**

    π On-storage metadata storing information on each file

    π Storing the file attributes

    π Tracking which blocks of the storage belong at which offsets within the logical file structure

◆ File System & File

◆ File Descriptor

◆ Directory

◆ File Aliasing

◆ Types of File System

# Open File and File Descriptor

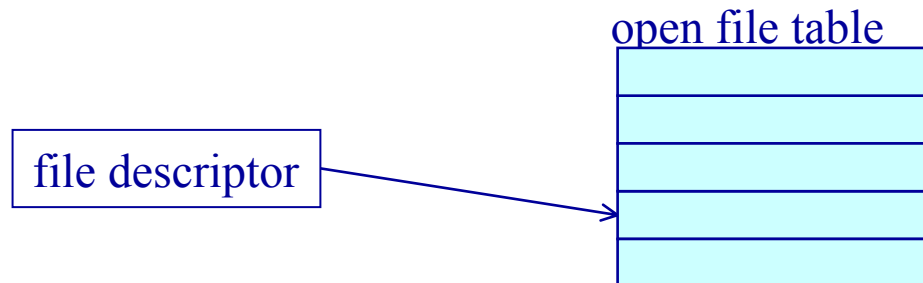- File use model
    - п User program must "open" a file before use

        f = open(name, flag);

        …

        … = read(f, …);

        …

        close(f);
- Kernel keeps track of open files for each process
    - п OS maintains an open file table per process
    - п An open file descriptor is an index into this table

    open file table

    | file descriptor |

# File Descriptor

- Several pieces of data are needed to manage open files:

  - π File pointer:  pointer to last read/write location, per process that has the file open

  - π File-open count: counter of number of times a file is open – to allow removal of data from open-file table when last processes closes it

  - π Disk location of the file: cache of data access information

  - π Access rights: per-process access mode information

# User vs System View of a File

- ## User's view:
  - π Durable data structures

- ## At system call interface
  - π Collection of bytes (UNIX)
  - π Doesn't matter to system what kind of data structures you want to store on disk!

- ## OS's internal view
  - π Collection of blocks (a block is a logical transfer unit, while a sector is the physical transfer unit)
  - π Block size 鑠 sector size; in UNIX, block size is 4KB
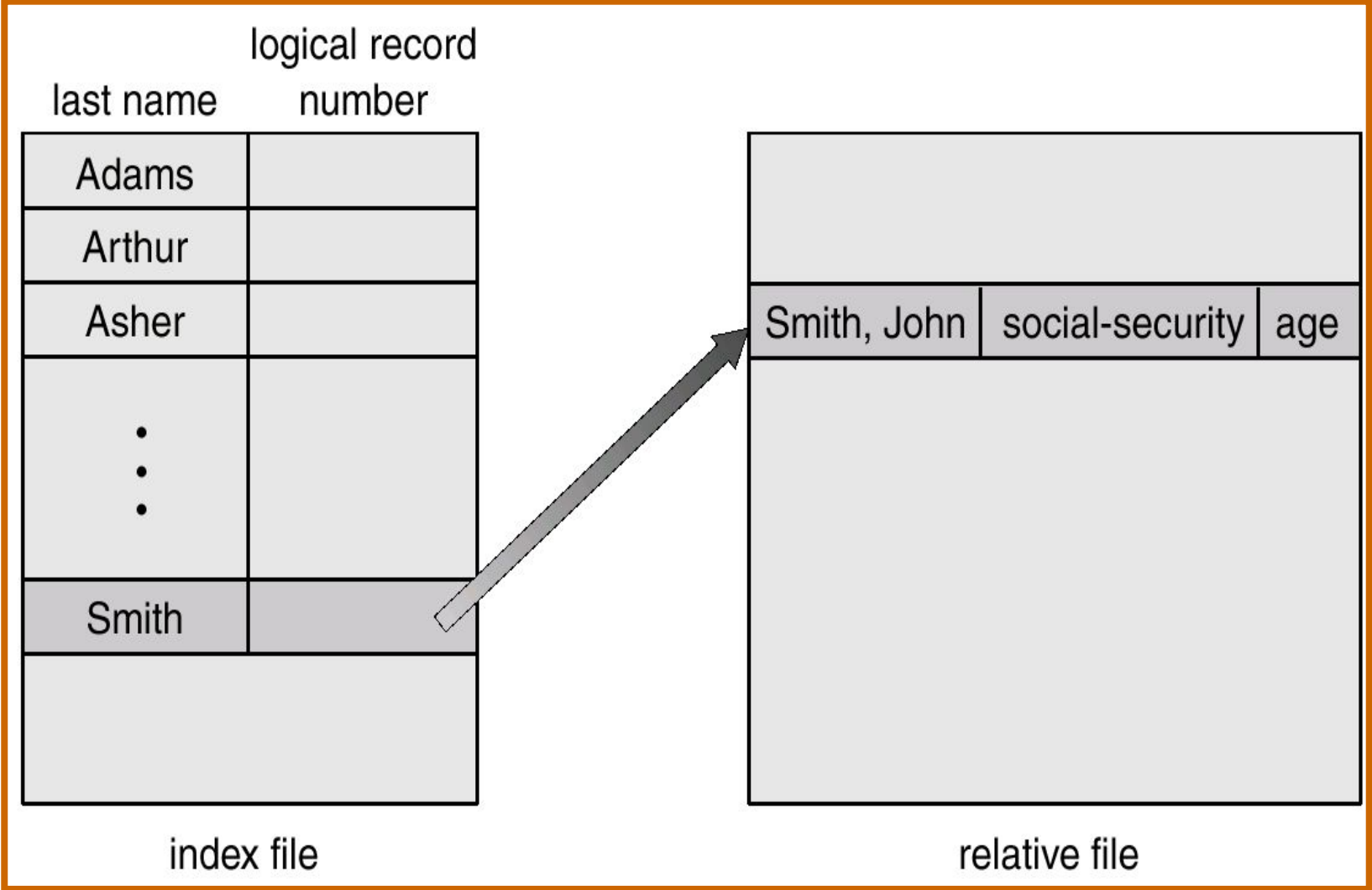
# Translating from User to System View

- **What happens if user says: give me bytes 2—12?**
  - π Fetch block corresponding to those bytes
  - π Return just the correct portion of the block
- **What about: write bytes 2—12?**
  - π Fetch block
  - π Modify portion
  - π Write out Block
- **Everything inside File System is in whole size blocks**
  - π For example, getc(), putc() 鑛 buffers something like 4096 bytes, even if interface is one byte at a time

# Access Patterns

- ## How do users access files?
  - ∏ Need to know type of access patterns user is likely to throw at system
- ## Sequential access: bytes read in order
  - ∏ Almost all file access are of this flavor
- ## Random Access: read/write element out of middle
  - ∏ Less frequent, but still important. For example, virtual memory backing file: page of memory stored in file
  - ∏ Want this to be fast – don't want to have to read all bytes to get to the middle of the file
- ## Content-based Access: by characteristics
  - ∏ Many systems don't provide this; instead, databases are built on top of disk access to index content (requires efficient random access)

# **File Internal Structure**

- ◆ No structure
  - π Sequence of words, bytes

- ◆ Simple record structure
  - π Lines
  - π Fixed length
  - π Variable length

- ◆ Complex structures
  - π Formatted document (e.g., MS Word, PDF)
  - π Executable file
  - π …

# File Sharing and Access Control

- **Sharing of files on multi-user systems is desirable**

- **Access control**
  - π Who can have what type accesses to what files
  - π Types of access: read, write, execute, delete, list, etc.

- **Per-file access control list (ACL)**
  - π <entity, permission>

- **Unix model**
  - π <user|group|world, read|write|execute>
  - π User IDs identify users, allowing permissions and protections to be per-user
  - π Group IDs allow users to be in groups, permitting group access rights

# Consistency Semantics

- Specify how multiple users/clients are to access a shared file simultaneously
  - π Similar to process synchronization algorithms
  - π Less complex due to disk I/O and network latency
- Unix file system (UFS) semantics
  - π Writes to an open file are visible immediately to other users of the same open file
  - π Sharing file pointer to allow multiple users to read and write concurrently
- Session semantics
  - π Writes only visible after the file is closed
- Locking
  - π Provided by some OS and file systems

◆ File System & File

◆ File Descriptor

◆ Directory

◆ File Aliasing

◆ Types of File System

# Hierarchical File System

- Files are organized in directories
- Directory is a kind of special files
  - π Each contains a <name, pointer to file header> table
- Tree structure for directories and files
  - π Some early file systems are flat (single-level directory)
- Hierarchical name space



/spell/mail/prt/first

/programs/p/list

# Operations Performed on Directory

- Typical operations
  - π Search for a file
  - π Create a file
  - π Delete a file
  - π List a directory
  - π Rename a file
  - π Traverse a path in the file system

- OS should only allow kernel mode to modify a directory
  - π Ensure integrity of the mapping
  - π Application programs can read directory (e.g., ls)

# **Directory Implementation**

◆ Linear list of file names with pointer to the data blocks

- Π simple to program
- Π time-consuming to execute

◆ Hash Table – linear list with hash data structure

- Π decreases directory search time
- Π collisions – situations where two file names hash to the same location
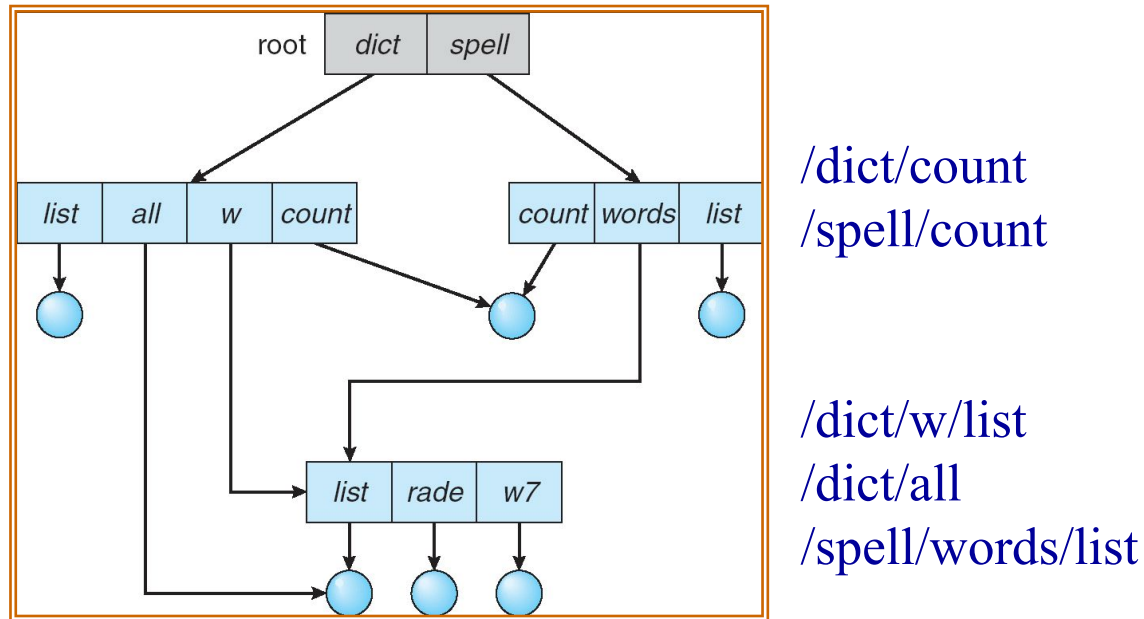- Π fixed size

◆ File System & File

◆ File Descriptor

◆ Directory

◆ File Aliasing

◆ Types of File System

# File Aliasing

♦ **Two or more different names referring same file**



/dict/count
/spell/count

/dict/w/list
/dict/all
/spell/words/list

♦ **Hard Links**: multiple directory entries point at the same file

♦ **Soft Links**: "shortcut" pointers to other files

π Implemented by storing the logical name of actual file

# The Dangling Pointer Problem in File Aliasing

- ## What if one delete the file pointed by one name
  - Π The name alias becomes "dangling pointer"

- ## Backpointers solution:
  - Π Each file has a list of backpointers, so we can delete all pointers
  - Π Backpointers using a daisy chain organization

- ## Add a level of indirection: directory entry data structure
  - Π Link – another name (pointer) to an existing file
  - Π Resolve the link – follow pointer to locate the file

# Cycles in Directory



/avi/book/avi/book/avi/book/avi/book/avi/book/avi/book/avi/book/avi/book/avi/…

◆ How do we guarantee no cycles?

　π Allow only links to file not subdirectories

　π Every time a new link is added use a cycle detection
　　algorithm to determine whether it is OK

◆ More practical

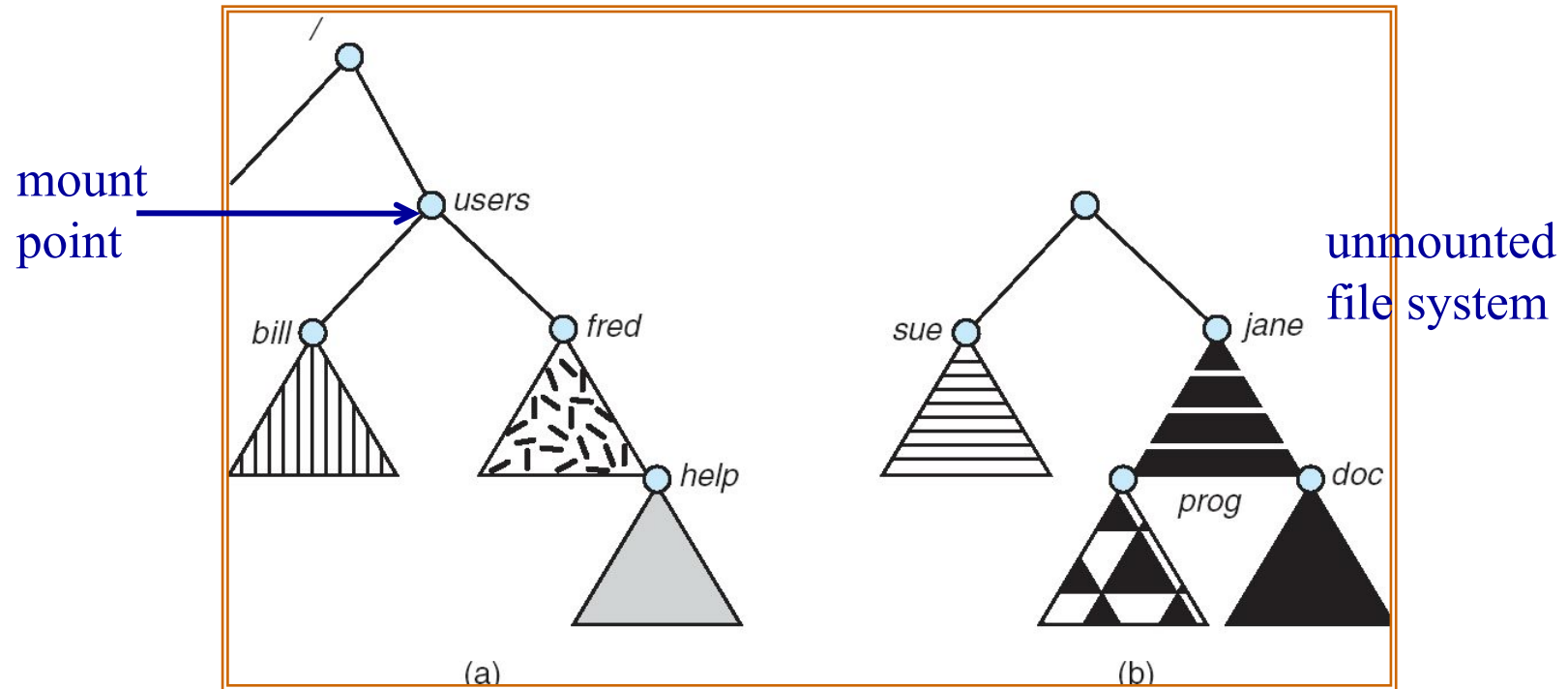　π Limit the number of directories that a path can traverse

# Name Resolution: Path Traversal

- Name resolution: the process of converting a logical name into a physical resource (like a file)
  - π In file system: file name (path) to actual file
  - π Traverse succession of directories until reach target file
- Example: resolving "/bin/ls"
  - π Read in file header for root (fixed spot on disk)
  - π Read in data block for root; search for "bin" entry
  - π Read in file header for "bin"
  - π Read in data block for "bin"; search for "ls"
  - π Read in file header for "ls"
- Present working directory (PWD)
  - π Per-process pointer to a directory for resolving file name
  - π Allows user to specify relative path instead of absolute path (say PWD="/bin" can resolve "ls")

# File System Mounting

- ◆ A file system must be mounted before it can be accessed
- ◆ A unmounted file system is mounted at a mount point

mount
point →

unmounted
file system



(a)

(b)

◆ **File System & File**

◆ **File Descriptor**

◆ **Directory**

◆ **File Aliasing**

◆ **Types of File System**

# Types of File Systems

- **Disk file systems**
  - π Files on a data storage device, like disk.
  - π Example: FAT, NTFS, ext2/3, ISO9660, etc.

- **Database file systems**
  - π Files are addressable (resolution) by characteristics
  - π Example: WinFS

- **Transactional file systems**
  - π Changes/events to file systems are logged
  - π Example: journaling file system

- **Network/distributed file systems**
  - π Example: NFS, SMB, AFS, GFS

- **Special/virtual file systems**

http://en.wikipedia.org/wiki/Comparison_of_file_systems

# Network/Distributed File Systems

- **Files may be shared across a network**
  - Π Files located at remote servers
  - Π Clients to mount remote file systems from servers
  - Π Standard OS file calls are translated into remote calls
  - Π Standard file sharing protocols: NFS for Unix, CIFS for Windows
- **Distributed system problems**
  - Π Client and user-on-client identification complicated
  - Π For example, NFS is insecure
  - Π Consistency problem
  - Π Dealing with failure mode
- **Truly distributed file systems is still a research**
  - Π Examples: Andrew File System (AFS)

# **Outline**

- Basic Concepts
- Virtual File System
- Data Block Caching
- Data Structures for Open Files
- File Allocation
- Free-Space List
- Management of Multiple Disks – RAID

# File System Implementation in an OS

◆ **Layering structure**

  π Upper layer: virtual (logical) file system

  π Lower layer: specific file system modules

| File/File System API | | | | | | |
|---|---|---|---|---|---|---|
| Virtual File System layer | | | | | | |
| ext2 | fat | iso9660 | nfs | smb | | |
| Device I/O | | | Network I/O | | | |

# Virtual File System (VFS) Layer

- ## Purpose
  - Π Abstraction for all different file system implementations

- ## Functions
  - Π Provide the same file and file system interface
  - Π Manage all file and file system related data structures
  - Π Routines for efficient lookup, traverse the file system
  - Π Interact with specific file system modules

# File System Basic Data Structures

- **Volume Control Block (Unix: "superblock")**
  - π One per file system
  - π Detail information about the file system
  - π # of blocks, block size, free-block count/pointer, etc.

- **File Control Block (Unix: "vnode" or "inode")**
  - π One per file
  - π Detail information about the file
  - π Permission, owner, size, data block locations, etc.

- **Directory Node (Linux: "dentry")**
  - π One per directory entry (directory or file)
  - π A tree data structure to encode the directory structure and tree layout
  - π Pointer to file control block, parent, list of entries, etc.

# Where are the Data Structure Stored

- **File system data structures**
  - π Volume control block (one per file system)
  - π File control block (one per file)
  - π Directory node (one per directory entry)
- **Persistently stored on the secondary storage**
  - π In data block(s) allocated in the storage
- **Loaded to memory when needed**
  - π Volume control block: in memory if file system is mounted
  - π File control block: if the file is accessed
  - π Directory node: during traversal of a file path

# Storage View

# **Outline**

- Basic Concepts

- Virtual File System

- Data Block Caching

- Data Structures for Open Files

- File Allocation

- Free-Space List

- Management of Multiple Disks – RAID

# Data Block Caching

- ◆ Data blocks are read into memory on-demand
    - Π To serve a read() operation
    - Π Read-ahead: prefetch subsequent data blocks
- ◆ Data blocks are cached after used
    - Π Under assumption that they may be used again
    - Π Writes may be buffered and delayed
- ◆ Two methods of caching data block
    - Π Normal buffer cache
    - Π Page cache: unified caching for data blocks and memory pages

# Remember Demanded Paging Memory Model?

- ## Demand paging
  - π Bring a page into memory only when it is needed

- ## Backing store
  - π A page (in virtual address space) can be mapped to a location in a file (in secondary storage)
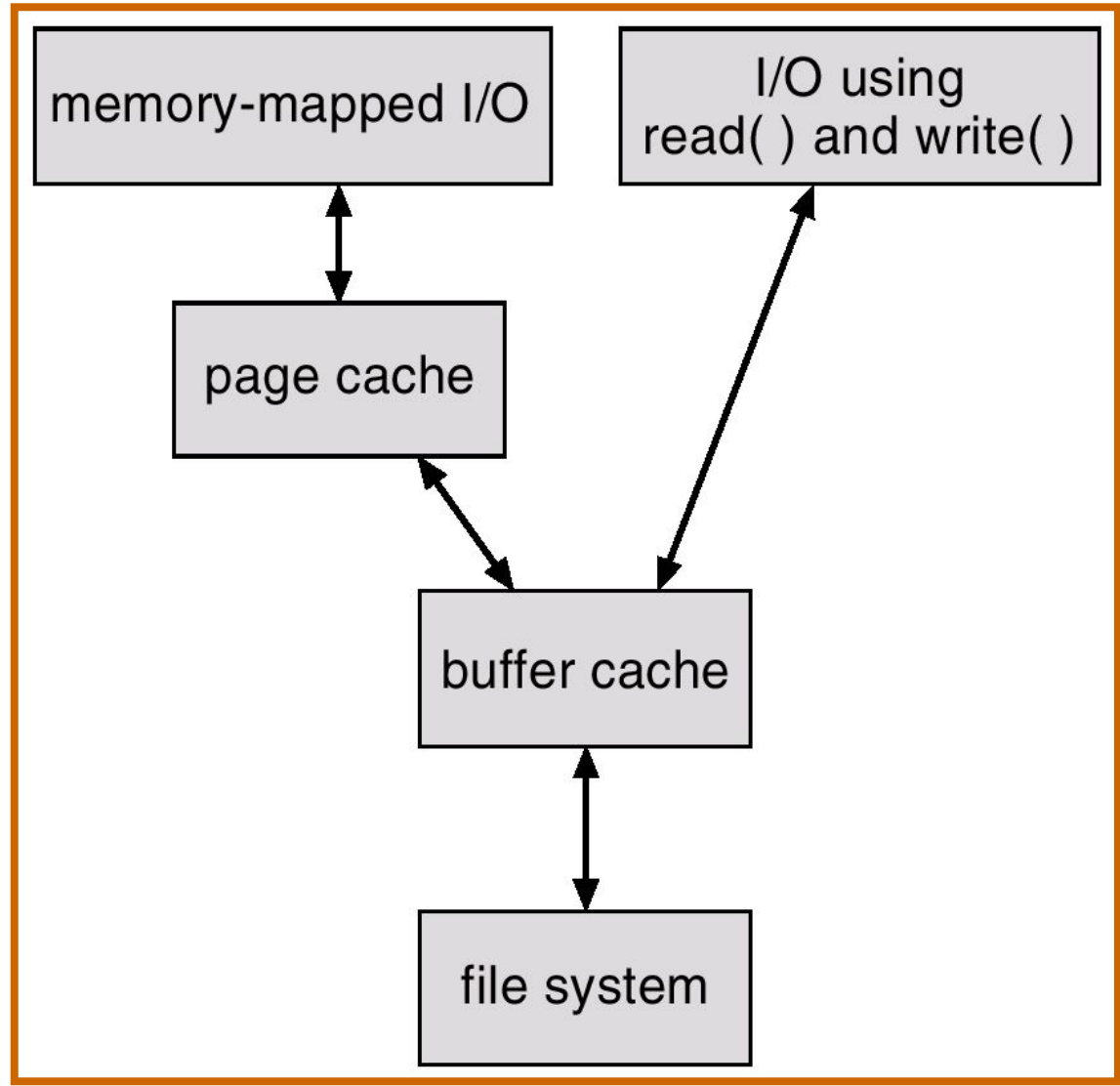


virtual address

OS kernel memory management

main memory

swap file | file | secondary storage

backing store

# Page Cache

◆ Page cache for file data blocks

- Π A file data block is mapped to a page in virtual memory
- Π File read/write op is translated to memory access
- Π May cause page-fault and/or set the page dirty
- Π Issue: page replacement – taken from processes or file page cache?

virtual address

OS kernel memory management

Page cache →    main memory

files    secondary storage

# 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.

- Basic Concepts
- Virtual File System
- Data Block Caching
- Data Structures for Open Files
- File Allocation
- Free-Space List
- Management of Multiple Disks – RAID

# File System Data Structures for Open Files

- ## Open file descriptor
  - Π One per open file
  - Π Information about the file status
  - Π Directory entry, current file pointer, set of file ops, etc.

- ## Open file tables
  - Π One per process
  - Π One system-wide
  - Π Each volume control block should keep a list too
  - Π So that it wouldn't dismount if still open file(s)

# Open-File Tables

# Open File Locking

- **Provided by some operating systems and file systems**

- **Mediates access to a file**

- **Mandatory or advisory:**

  - Π **Mandatory** – access is denied depending on locks held and requested

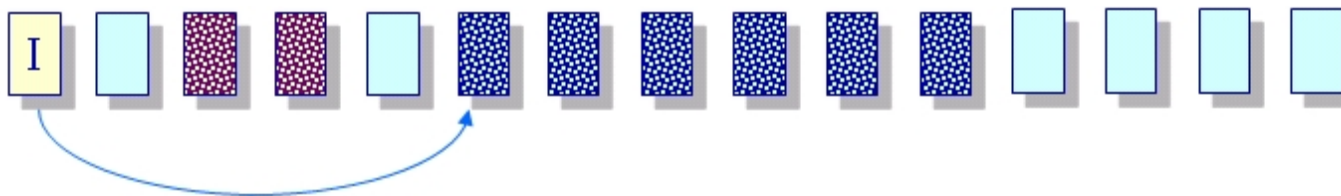  - Π **Advisory** – processes can find status of locks and decide what to do

# Outline

- Basic Concepts
- Virtual File System
- Data Block Caching
- Data Structures for Open Files
- File Allocation
- Free-Space List
- Management of Multiple Disks – RAID

# Use Pattern

- ## Most files are small.

  - Π Need strong support for small files.

  - Π Block size can't be too big.

- ## Some files are very large.

  - Π Must allow large files (64-bit file offsets).
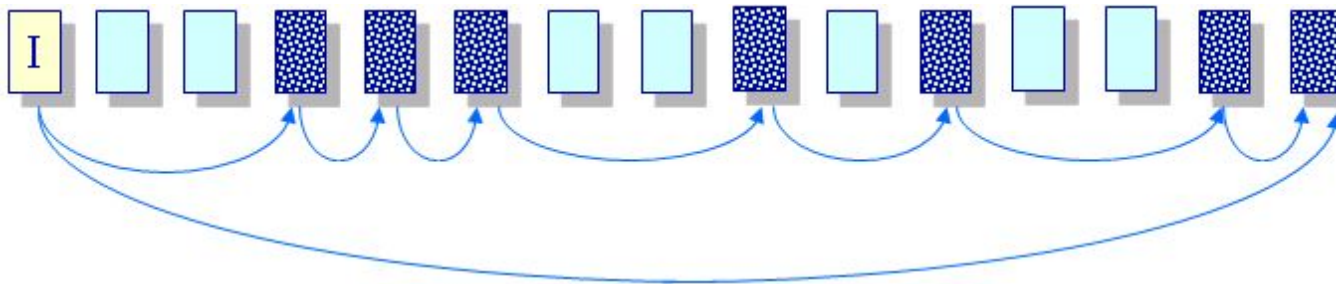
  - Π Large file access should be reasonably efficient.

# File Allocation

◆ How to allocate data blocks to each file

◆ Allocation methods

   π Contiguous allocation

   π Linked allocation

   π Indexed allocation

◆ Metrics

   π Efficiency: e.g., storage utilization (external fragmentation)

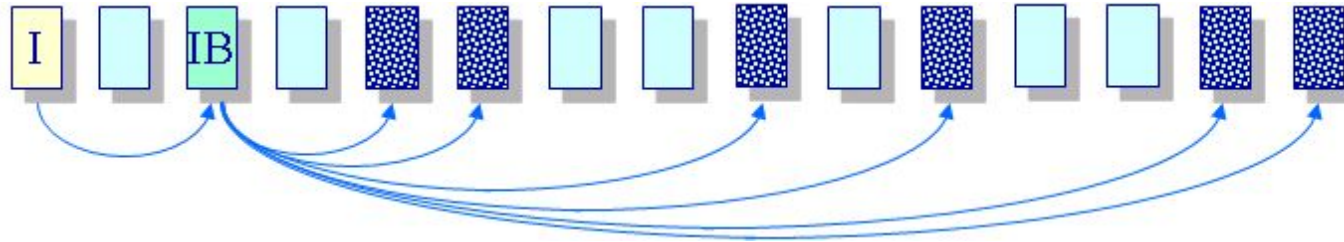   π Performance: e.g., access speed

# Contiguous Allocation



- ◆ File header specifies starting block & length
- ◆ Placement/Allocation policies
  - π First-fit, best-fit, ...

u Pluses
  - π Best file read performance
  - π Efficient sequential & random access

u Minuses
  - π Fragmentation!
  - π Problems with file growth
    - �147 Pre-allocation?
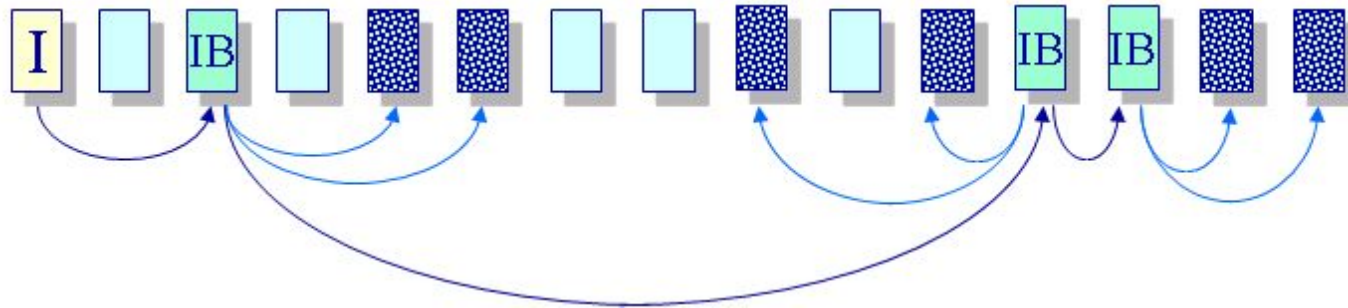    - �147 On-demand allocation?

# Linked Allocation



◆ Files stored as a linked list of blocks

◆ File header contains a pointer to the first and last file blocks

u Pluses
- Π Easy to create, grow & shrink files
- Π No fragmentation

鐙 Minuses
- Π Impossible to do true random access
- Π Reliability
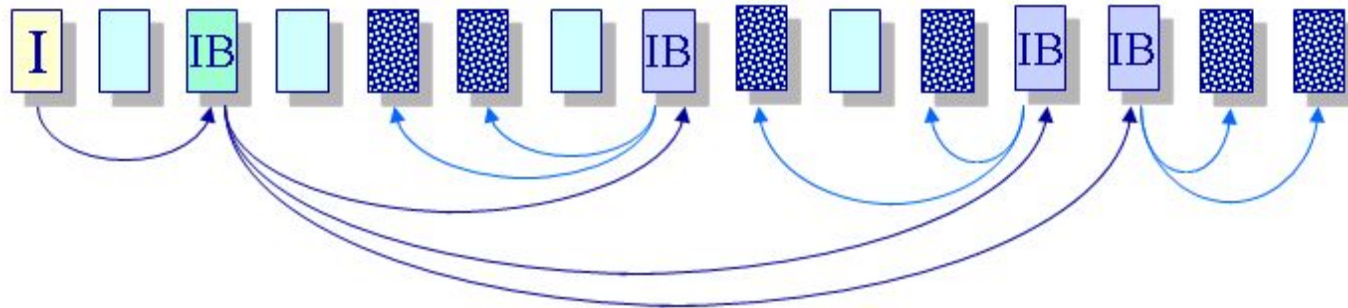  - 鐚 Break one link in the chain and...

# Indexed Allocation



- Create a non-data block for each file called the index block
  - π A list of pointers to file blocks
- File header contains the index block

u Pluses
  - π Easy to create, grow & shrink files
  - π No fragmentation
  - π Supports direct access

u Minuses
  - π Overhead of storing index when files are small
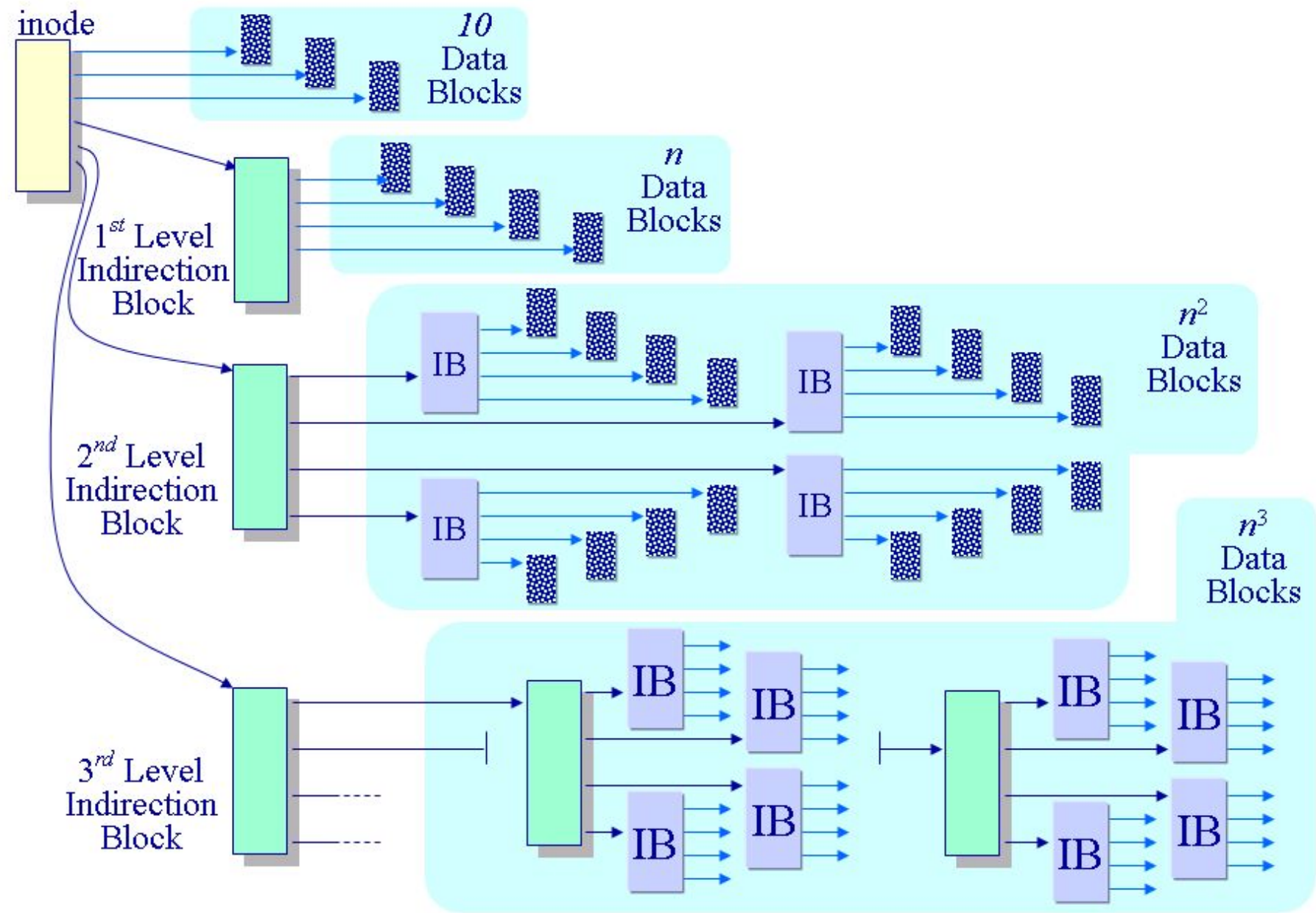  - π How to handle large files?

# Indexed Allocation for Large Files

- Linked index blocks (IB+IB+…)



- Multilevel index blocks (IB*IB*…)

# Multi-level Indexed Allocation in UNIX

# Multi-level Indexed Allocation in UNIX

- File header contains 13 pointers
  - π 10 pointes to data blocks;
  - π 11th pointer → indirect block;
  - π 12th pointer → doubly-indirect block;
  - π 13th pointer → triply-indirect block

- Implications
  - π Upper limit on file size
  - π Blocks are allocated dynamically, files can easily expand
  - π Small files are cheap
  - π Allocate indirect blocks only for large files, and large files require a lot of seek to access indirect blocks

# OS
## **Outline**

- Basic Concepts
- Virtual File System
- Data Block Caching
- Data Structures for Open Files
- File Allocation
- Free-Space List
- Management of Multiple Disks – RAID

# Free-Space List

- Keep track of all unallocated blocks in the storage

- Where is free-space list stored?

- What is a good data structure for free-space list?

# Free-Space List: Bit Map

- Represent the list of free blocks as a bit map:
  - π 111111111111111001110101011101111...
  - π If bit i = 0 then block i is free, otherwise it is allocated

- Simple to use but this can be a big vector:
  - π 160GB disk -> 40M blocks -> 5MB worth of bits
  - π However, if free sectors are uniformly distributed across the disk then the expected number of bits that must be scanned before finding a "0" is n/r, where
    - 鐚 n = total number of blocks on the disk
    - 鐚 r = number of free blocks
  - π If a disk is 90% full, then the average number of bits to be scanned is 10, independent of the size of the disk

◆ Need to protect:

  π Pointer to free list

  π Bit map

   ▩ Must be kept on disk

   ▩ Copy in memory and disk may differ.

   ▩ Cannot allow for block[$i$] to have a situation where bit[$i$] = 1 in memory and bit[$i$] = 0 on disk.
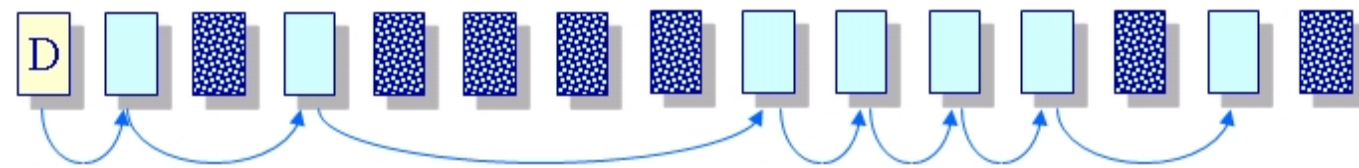
  π Solution:

   ▩ Set bit[$i$] = 1 in disk.

   ▩ Allocate block[$i$]

   ▩ Set bit[$i$] = 1 in memory          ?

- linked lists

- Grouped lists



Next group block

Allocated block    Empty block
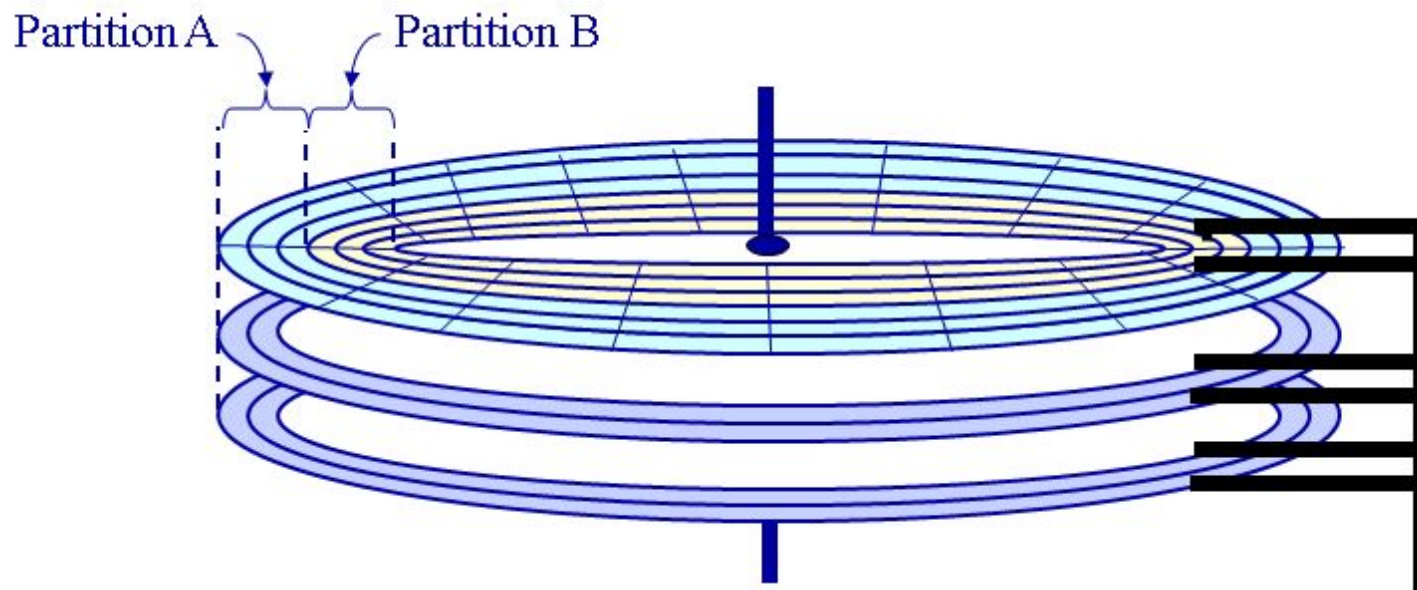
- Basic Concepts

- Virtual File System

- Data Block Caching

- Data Structures for Open Files

- File Allocation

- Free-Space List

- Management of Multiple Disks – RAID

# Disk Partitioning for Performance

◆ **Disks are typically partitioned to minimize the largest possible seek time**

- π A partition is a collection of cylinders
- π Each partition is a logically separate disk

# A Typical Disk File-System Organization

◆ Partition: a division of hard disk to apply OS-specific formatting

◆ Volume: a single accessible storage area with a single instance of a filesystem

  Π Typically resident on a single partition of a hard disk

# Management of Multiple Disks

- Use multiple parallel disks to increase
  - π Throughput (through parallelism)
  - π Reliability and availability (through redundancy)
- RAID - Redundant Array of Inexpensive Disks
  - π A variety of disk-organization techniques
  - π RAID levels: different RAID scheme (e.g., RAID-0, RAID-1, RAID-5)
- Implementation
  - π In OS kernel: storage/volume management
  - π In hardware RAID controller (I/O)

# RAID-0: Disk Striping for Throughput

- **Blocks broken into sub-blocks that are stored on separate disks**

  - π similar to memory interleaving

- **Provides for higher disk bandwidth through a larger effective block size**

OS disk block

| 8 | 9 | 10 | 11 |
| 12 | 13 | 14 | 15 |
| 0 | 1 | 2 | 3 |

8  9  10 11       12 13 14 15       0  1  2  3

## Physical disk blocks

# Raid-1: Disk Mirroring for Reliability

◆ Reliability is increased exponentially

◆ Read performance goes up linearly

  π Write to both disks, read from either.



| 0 1 1 0 0 |
| 1 1 1 0 1 |
| 0 1 0 1 1 |

Primary disk

| 0 1 1 0 0 |
| 1 1 1 0 1 |
| 0 1 0 1 1 |

Mirror disk

OS

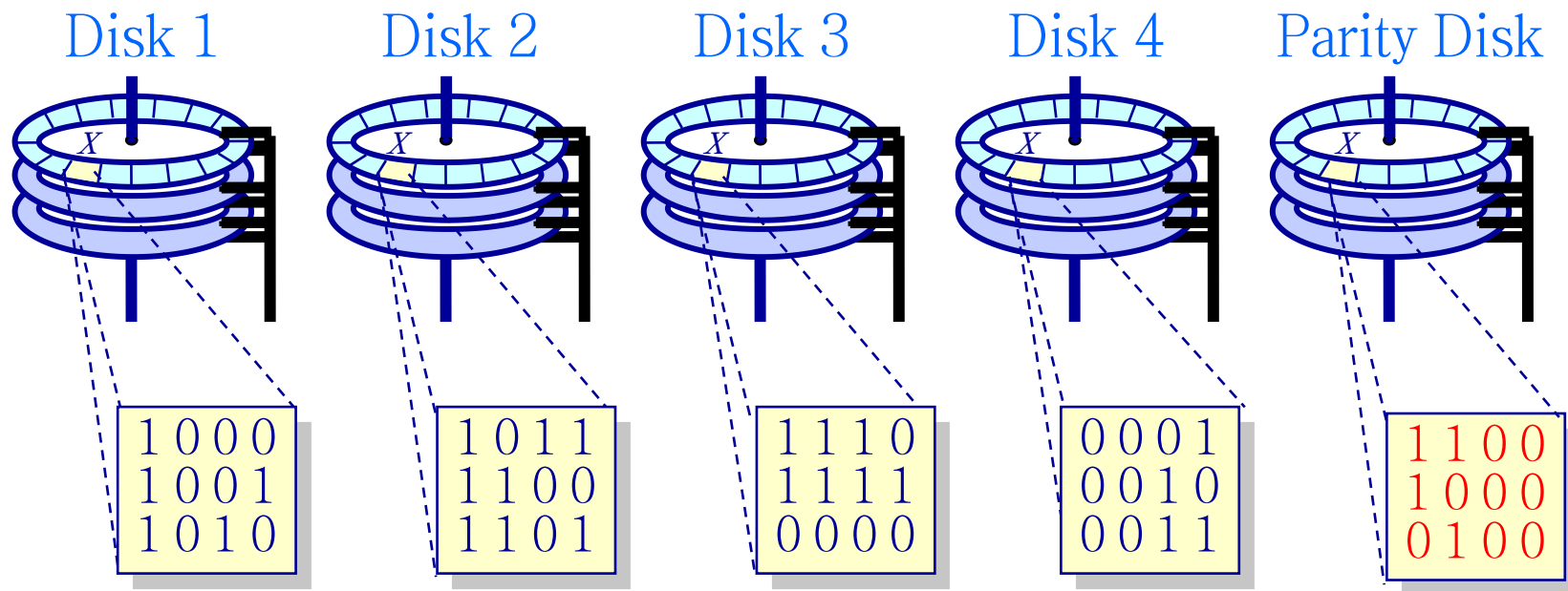# RAID-4: Parity Disk for Reliability

◆ **Block-level striping with a dedicated parity disk**

  π Allows one to recover from the crash of any one disk

  π Example: storing 8, 9, 10, 11, 12, 13, 14, 15, 0, 1, 2, 3
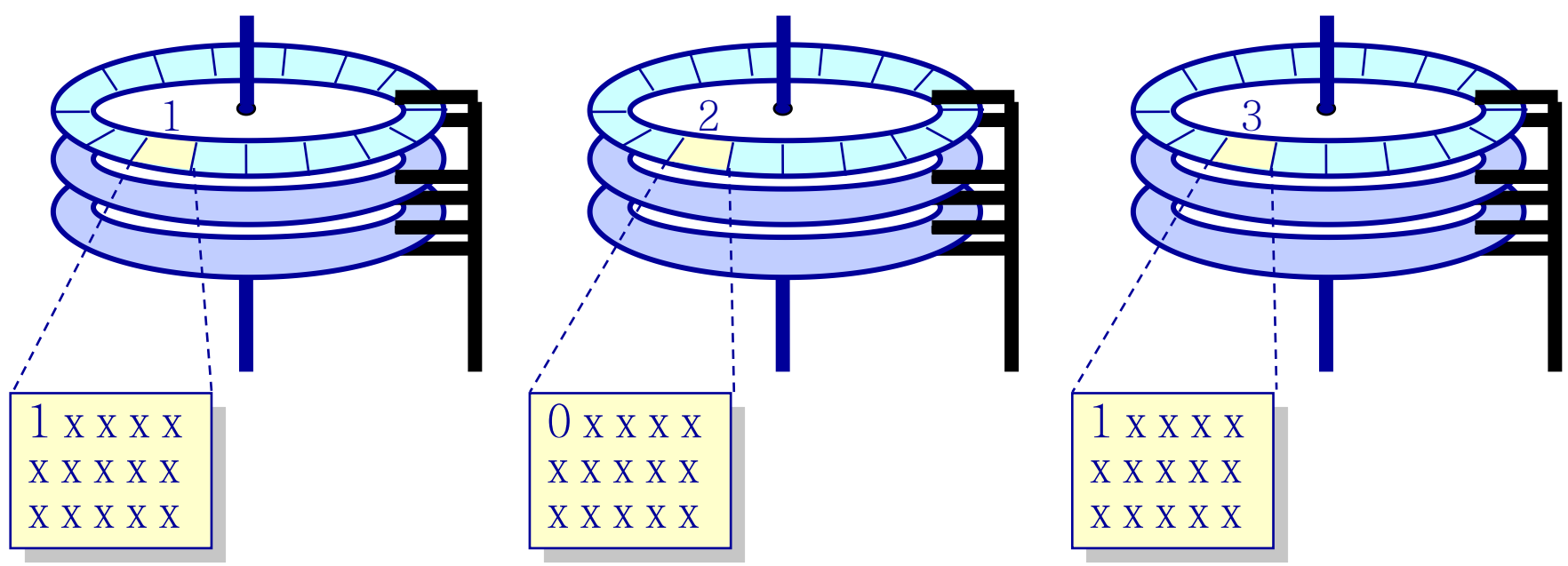
| Disk 1 | Disk 2 | Disk 3 | Disk 4 | Parity Disk |
|--------|--------|--------|--------|-------------|

```
1000    1011    1110    0001    1100
1001    1100    1111    0010    1000
1010    1101    0000    0011    0100
```

# RAID-5: Block-interleaved Distributed Parity



Disk 1  Disk 2  Disk 3  Disk 4  Disk 5

Block X

| Disk 1 | Disk 2 | Disk 3 | Disk 4 | Disk 5 |
|--------|--------|--------|--------|--------|
| 8<br>9<br>10 | 11<br>12<br>13 | 14<br>15<br>0 | 1<br>2<br>3 | Block<br>X<br>Parity |

Block x+1

| Block<br>x+1<br>Parity | a<br>b<br>c | d<br>e<br>f | g<br>h<br>i | j<br>k<br>l |

Block x+2

| m<br>n<br>o | Block<br>x+2<br>Parity | p<br>q<br>r | s<br>t<br>u | v<br>w<br>x |

Block x+3

| y<br>z<br>aa | bb<br>cc<br>dd | Block<br>x+3<br>Parity | ee<br>ff<br>gg | hh<br>ii<br>jj |

# Bit-wise vs Block-wise Disk Striping

◆ Striping and parity can be done byte-by-byte or bit-by-bit

  π RAID-0/4/5: block-wise

  π RAID-3: bit-wise

◆ Example: storing bit-string 101 in RAID-3 system
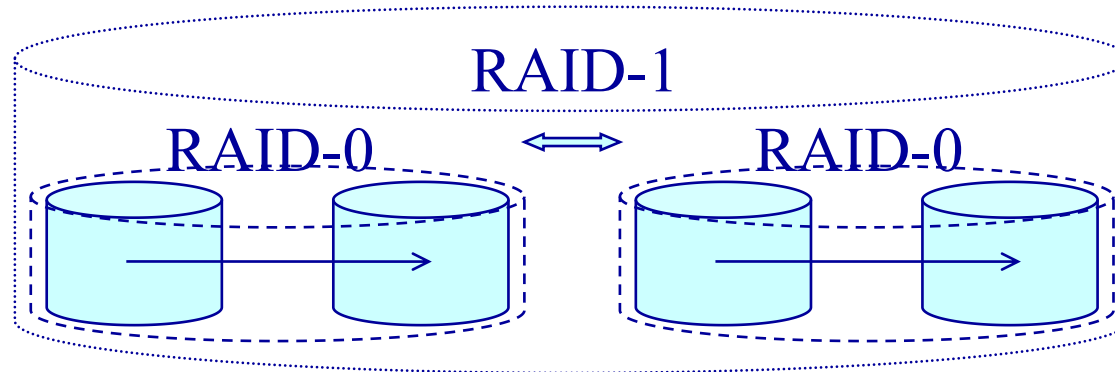


| 1 x x x x |
| x x x x x |
| x x x x x |

| 0 x x x x |
| x x x x x |
| x x x x x |

| 1 x x x x |
| x x x x x |
| x x x x x |

# Tolerating Two Disk Failures

◆ RAID-5: single parity block per striping data block

  π  tolerating one disk failure

◆ RAID-6: two redundancy blocks

  π  With a special coding scheme

  π  tolerating two disk failures

# Nested RAID Levels

◆ RAID 0+1

RAID-1

RAID-0 ⟺ RAID-0

◆ RAID 1+0

RAID-0

RAID-1 ⟶ RAID-1