



### This lecture

#### • Implementing file system abstraction

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



# User vs. system view of a file

- User's view
  - Durable data structures
- System's 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 >= sector size.

#### File structure

- None sequence of words, bytes
- Simple record structure
  - Lines
  - Fixed length
  - Variable length
- Complex structures
  - Formatted document
  - Relocatable load file
- Can simulate last two with first method by inserting appropriate control characters.
- Who decides:
  - Operating system
  - Program

### File attributes

- Name only information kept in human-readable form.
- Type needed for systems that support different types.
- Location pointer to file location on device.
- Size current file size.
- Protection controls who can do reading, writing, executing.
- Time, date, and user identification data for protection, security, and usage monitoring.
- Information about files are kept in the directory structure, which is maintained on the disk.

#### File operations

- ♦ create
- write
- read
- reposition within file file seek
- delete
- truncate
- open(F<sub>i</sub>) search the directory structure on disk for entry F<sub>i</sub>, and move the content of entry to memory.
- close (F<sub>i</sub>) move the content of entry F<sub>i</sub> in memory to directory structure on disk.

File types - name, extension					
File Type	Usual extension	Function			
Executable	exe, com, bin or none	ready-to-run machine- language program			
Object	obj, o	complied, machine language, not linked			
Source code	c, p, pas, 177, asm, a	source code in various languages			
Batch	bat, sh	commands to the command interpreter			
Text	txt, doc	textual data documents			
Word processor	wp, tex, rrf, etc.	various word-processor formats			
Library	lib, a	libraries of routines			
Print or view	ps, dvi, gif	ASCII or binary file			
Archive	arc, zip, tar	related files grouped into one file, sometimes compressed.			











# File system design

#### Data structures

- Directories: file name -> file metadata
   \* Store directories as files
- File metadata: how to find file data blocks
- Free map: list of free disk blocks
- How do we organize these data structures?
  - Device has non-uniform performance

# Design challenges

- Index structure
  - How do we locate the blocks of a file?
- Index granularity
  - What block size do we use?
- Free space
  - How do we find unused blocks on disk?
- Locality
  - How do we preserve spatial locality?
- Reliability
  - What if machine crashes in middle of a file system op?

	FAT	FFS	NTFS	ZFS
Index structure	Linked list	Tree (fixed, assym)	Tree (dynamic)	Tree (COW) dynamia
granularity	block	block	extent	block
free space allocation	FAT array	Bitmap (fixed location)	Bitmap (file)	Space m (log- structure
Locality	defragmenta tion	Block groups + reserve space	Extents Best fit defrag	Write- anywher Block-gro







# Data structures for disk management

- A file header for each file (part of the file metadata)
  - Disk sectors associated with each file
- A data structure to represent free space on disk
  - Bit map
    - \* 1 bit per block (sector)
    - \* blocks numbered in cylinder-major order, why?
  - Linked list
  - Others?
- How much space does a bitmap need for a 4G disk?

## Contiguous allocation

- Request in advance for the size of the file
- Search bit map or linked list to locate a space
- File header
  - first sector in file
  - number of sectors
- Pros
  - Fast sequential access
  - Easy random access
- Cons
  - External fragmentation
  - Hard to grow files









# FAT

- ♦ Pros:
  - Easy to find free block
  - Easy to append to a file
  - Easy to delete a file
- ♦ Cons:
  - Small file access is slow
  - Random access is very slow
  - Fragmentation
    - \* File blocks for a given file may be scattered
    - \* Files in the same directory may be scattered
    - \* Problem becomes worse as disk fills











### FFS inode

- Metadata
  - File owner, access permissions, access times, ...
- Set of 12 data pointers
  - With 4KB blocks => max size of 48KB files
- Indirect block pointer
  - pointer to disk block of data pointers
- Indirect block: 1K data blocks => 4MB (+48KB)

## FFS inode

- Metadata
  - File owner, access permissions, access times, ...
- Set of 12 data pointers
  - With 4KB blocks => max size of 48KB
- Indirect block pointer
  - pointer to disk block of data pointers
  - 4KB block size => 1K data blocks => 4MB
- Doubly indirect block pointer
  - Doubly indirect block => 1K indirect blocks
  - 4GB (+ 4MB + 48KB)

## FFS inode

#### Metadata

- File owner, access permissions, access times, ...
- Set of 12 data pointers
   With 4KB blocks => max size of 48KB
- Indirect block pointer
  - pointer to disk block of data pointers
  - 4KB block size => 1K data blocks => 4MB
- Doubly indirect block pointer
  - Doubly indirect block => 1K indirect blocks
  - 4GB (+ 4MB + 48KB)
- Triply indirect block pointer
  - Triply indirect block => 1K doubly indirect blocks
  - 4TB (+ 4GB + 4MB + 48KB)



# FFS asymmetric tree

- Small files: shallow tree
  Efficient storage for small files
- Large files: deep tree
  Efficient lookup for random access in large files
- Sparse files: only fill pointers if needed







# FFS locality

- Block group allocation
  - Block group is a set of nearby cylinders
  - Files in same directory located in same group
  - Subdirectories located in different block groups
- inode table spread throughout disk
   inodes, bitmap near file blocks
- First fit allocation
  - Small files fragmented, large files contiguous



FFS first fit block allocation	
In-Use Free Start of <sup>Block</sup> Block Block Group	] •••







#### File header storage

- Where is file header stored on disk?
  - In (early) Unix & DOS FAT file sys, special array in outermost cylinders
- Unix refers to file by index into array --- tells it where to find the file header
  - "i-node" --- file header; "i-number" --- index into the array
- Unix file header organization (seems strange):
  - header not anywhere near the data blocks. To read a small file, seek to get header, seek back to data.
  - fixed size, set when disk is formatted.



#### Naming and directories

#### Options

- Use index (ask users specify inode number). Easier for system, not as easy for users.
- Text name (need to map to index)
- Icon (need to map to index; or map to name then to index)
- Directories
  - Directory map name to file index (where to find file header)
  - Directory is just a table of file name, file index pairs.
  - Each directory is stored as a file, containing a (name, index) pair.
  - Only OS permitted to modify directory

#### **Directory structure**

- Approach 1: have a single directory for entire system.
  - \* put directory at known location on disk
  - \* directory contains <name, index> pairs
  - \* if one user uses a name, no one else can
  - $^{\star}\,$  many older personal computers work this way.

#### Approach 2: have a single directory for each user

- \* still clumsy. And Is on 10,000 files is a real pain
- \* many older mathematicians work this way.

#### • Approach 3: hierarchical name spaces

- \* allow directory to map names to files or other dirs
- \* file system forms a tree (or graph, if links allowed)
- \* large name spaces tend to be hierarchical (ip addresses, domain names, scoping in programming languages, etc.)

















## Example: open-read-close (cont'd)

4. The kernel:

- From "fd" find the file pointer
- Based on the file system block size (let's say 1 KB), find the blocks where the bytes (file\_pointer, file\_pointer+length) lies;
- Read the inode
- For (each block) {
  - \* If the block # <11, find the disk address of the block in the entries in the inode</li>
     \* If the block # >= 11, but < 11 + (1024/4): read the "single indirect" block to find the address of the block</li>
  - \* If the block # >= 11+(1024/4) but < 11 + 256 + 256 \* 256: read the "double indirect" block and find the block's address
  - \* Otherwise, read the "triple indirect" block and find the block's address }
- Read the block from the disk
- Copy the bytes in the block to the appropriate location in the buffer
- 5. The process calls close(fd);
- 6. The kernel: deallocate the fd entry, mark it as empty.



### Example: create-write-close (cont'd)

- 3. The process calls write (fd, buffer, length);
- 4. The kernel:
  - From "fd" find the file pointer;
  - Based on the file system block size (let's say 1 KB), find the blocks where the bytes (file\_pointer, file\_pointer+length) lies;
  - Read the inode
  - For (each block) {
    - \* If the block is new, allocate a new disk block;
      - Based on the block no, enter the block's address to the appropriate places in the inode or the indirect blocks; (the indirect blocks are allocated as needed)
    - \* Copy the bytes in buffer to the appropriate location in the block }
  - Change the file size field in inode if necessary
- 5. The process calls close(fd);
- 6. The kernel deallocate the fd entry --- mark it as empty.

#### NTFS

- Master File Table
  - Flexible 1KB storage for metadata and data
- Extents
  - Block pointers cover runs of blocks
  - Similar approach in linux (ext4)
  - File create can provide hint as to size of file
- Journaling for reliability







