

File System Implementation



Prof. Yongtae Kim

Computer Science and Engineering
Kyungpook National University

The Way to Think

■ We introduce a simple file system implementation vsfs

- Specifically, vsfs (Very Simple File System) will be implemented
- The file system is pure software; no hardware feature is required

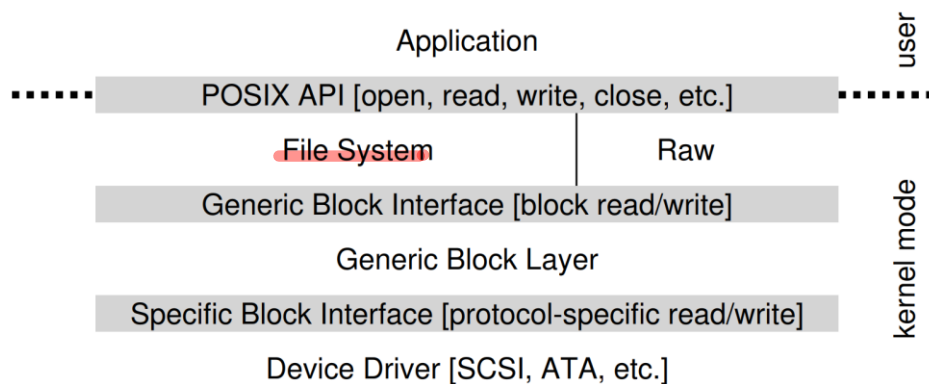
■ We need to understand how the file system works in two aspects

- 1) The first is the data structures of the file system

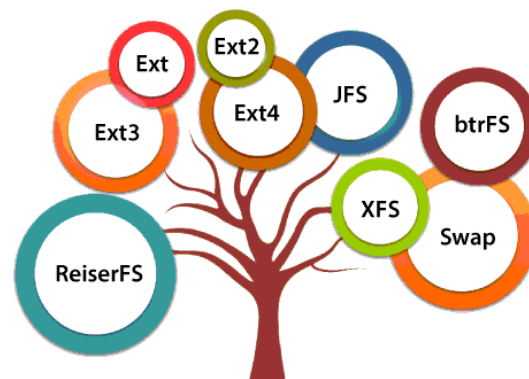
In other words, what types of on-disk structures are utilized by the file system to organize its data and metadata? → array, tree, etc

- 2) The second aspect of a file system is its access methods

How does it map the calls made by a process, such as `open()`, `read()`, `write()`, etc., onto its structures?

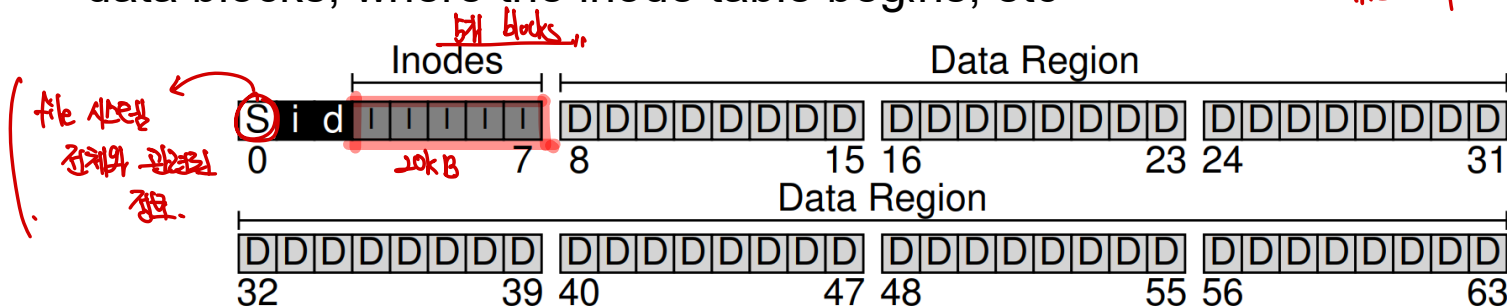


Types of Linux File System



Overall Organization

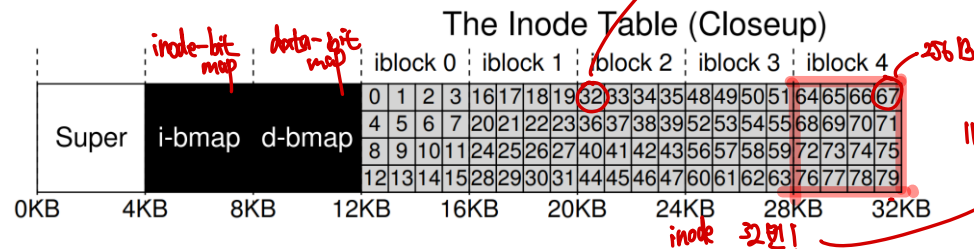
- The first thing we'll need to do is divide the disk into **block**
 - Simple file systems use just one block size, commonly-used size of **4KB**
 - Our view of the disk partition is a series of blocks (0~N-1); Assume **64 blocks**.
- Think about what we store in these blocks to build vsfs
 - To store **user data**, the **data region of 56 blocks** on the disk is reserved
 - To store **metadata**, which contains (file size, owner, etc) file systems have **inode**; **5 blocks of inode table**; 256 byte per inode, 20KB can hold 80 total inodes
 - **Allocation structures** to track if a block is free are required in file systems
 - Allocation-tracking methods: **free list** and **bitmap**, etc \rightarrow free list: inode, data를 저장할 수 있음. bitmap is simple and thus vsfs will use **bitmap for inode and data**
 - The last block is reserved for **superblock**, which contains how many inodes and data blocks, where the inode table begins, etc \rightarrow superblock



File Organization: The Inode

- One of the most important on-disk structures of file system is the **inode**, which is short for **index node**

- Each inode is implicitly referred to by a number (i-number; low-level name)



- The **inode location** can be calculated using the i-number

e.g.) i-number: 32 $\rightarrow 32 \times \text{sizeof}(\text{inode}) + 12\text{KB} = 32 \times 256\text{B} + 12\text{KB} = 20\text{KB}$

To find the corresponding sector $20\text{KB} / 512 = 40$

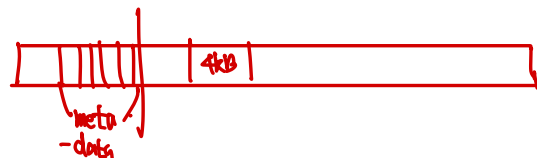
- A general equation to find the sector:

```
blk = (inumber * sizeof(inode_t)) / blockSize;
sector = ((blk * blockSize) + inodeStartAddr) / sectorSize;
```

- Each inode includes the file's metadata, such as its size, protection, time, etc (e.g. simplified ext2 inode on right figure)

Size	Name	What is this inode field for?
2	mode	can this file be read/written/executed?
2	uid	who owns this file?
4	size	how many bytes are in this file?
4	time	what time was this file last accessed?
4	ctime	what time was this file created?
4	mtime	what time was this file last modified?
4	dtime	what time was this inode deleted?
2	gid	which group does this file belong to?
2	links_count	how many hard links are there to this file?
4	blocks	how many blocks have been allocated to this file?
4	flags	how should ext2 use this inode?
4	osd1	an OS-dependent field
60	block	a set of disk pointers (15 total)
4	generation	file version (used by NFS)
4	file_acl	a new permissions model beyond mode bits
4	dir_acl	called access control lists

The Multi-Level Index

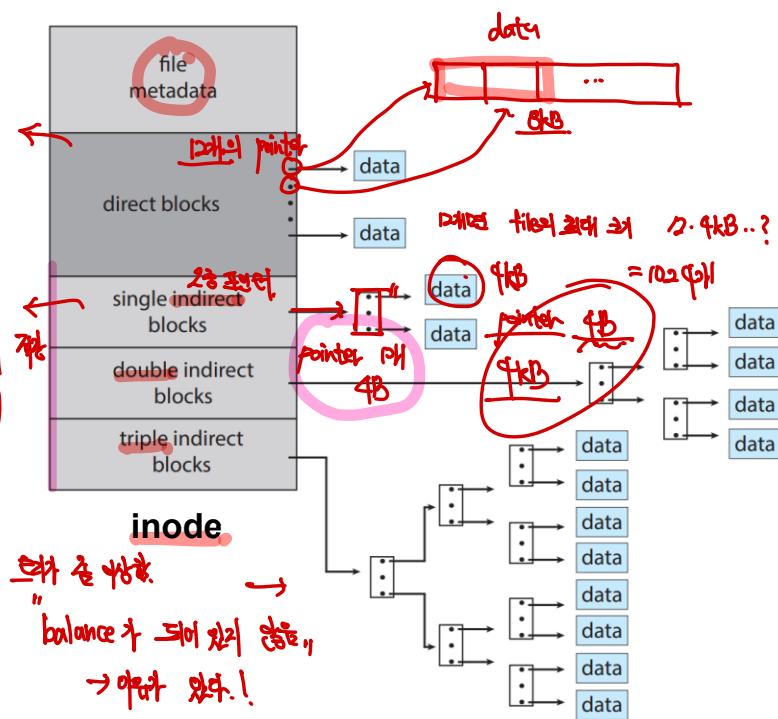


■ To support bigger files, file systems uses an indirect pointer

- Instead of directly pointing to a block that contains user data, it points to a block that contains more pointers, each of which point to user data
- Each inode of vsfs contains 12 direct pointers and 1 indirect pointers; under 4-byte disk address, the file can grow to be $(12 + 1024) \times 4KB = 4144KB$ // 최대 크기 ~ 4MB
- To have even larger files, we can add another pointer to inode: double indirect pointer $\rightarrow (12 + 1024 + 1024^2) \times 4KB \approx 4GB$ Double indirect 배
- Even more? triple indirect pointer (배치)
- This imbalanced tree is referred to as the multi-level index approach
- Why use an imbalanced tree like this? 1024개의 들것에 들것 (pointer)

One finding is that most files are small

Most files are small	파일이 대부분 그렇게 크기 작아. ~2K is the most common size
Average file size is growing	Almost 200K is the average
Most bytes are stored in large files	A few big files use most of space
File systems contains lots of files	Almost 100K on average
File systems are roughly half full	Even as disks grow, file systems remain ~50% full
Directories are typically small	Many have few entries; most have 20 or fewer



Directory Organization & Free Space Management

여러 inode와 레코드 블록이 서로 연결되어 있는 bitmap 저장

Directories basically contain a list of (entry name, i-number) pair

- e.g.) a directory (i-number 5) has files (foo, bar, foobar_is_a_pretty_longname) with inode number 12, 13, and 24, respectively

i-node.

inum	reclen	strlen	name
5	12	2	.
2	12	3	..
12	12	4	foo
13	12	4	bar
24	36	28	foobar_is_a_pretty_longname

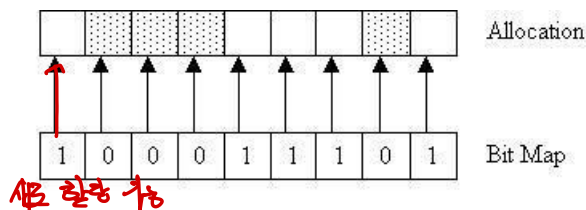
meta.

하나라도 있.

- Each entry has an inode number, record length (total bytes for name + any left over space), string length (name's actual length), and the file name

Free space management is important for all file systems

- vsfs uses two bitmaps to track which inodes and data blocks are free
- When creating a file, the file system allocates the an inode and search through bitmap for a free inode, and allocate it to file; similar activity for data blocks



inode 번호, data 번호
둘 다 찾아야 한다.

Access Paths: Reading and Writing

채널 다스 블록

Opening, reading, closing a file /foo/bar whose size is 12KB

- The file system finds the inode for **bar** from **root** (its i-number is already known)
- open()**: **root inode** (pointer) → root directory data (**foo i-number**) → **foo inode** (pointer) → **foo directory data** (**bar i-number**) → **bar inode**
- read()**: read **bar** data → **update the inode** (last access time, etc)
- close()**: close the file and deallocate the file descriptor

① open: root directory의 inode도 채워 풀러가 어떤 정보를 올릴
 ↓
 ② root directory의 data 블록에 foo라는 이름, inode 번호를
 넣고, 해당 블록을 가리키는 포인터를 inode에 올릴 것

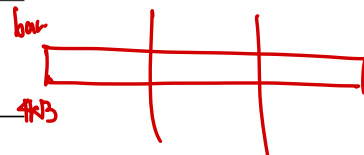
↓
 ③ foo directory의 inode를 가리키는 foo directory의 data 블록으로 갈
 ↓
 ④ foo directory의 data 블록을 가리키는 bar라는 포인터와 바이트 정보를
 올리고 bar inode를 올릴 것

↓
 ⑤ bar inode로 가리키는 블록을 읽어
 바이트 정보를 올릴 수 있다
 ↓
 바이트 정보를 올린 후 바이트 정보를 올릴
 bar inode의 바이트 정보를 올릴 것

	data bitmap	inode bitmap	root inode	foo inode	bar inode	root data	foo data	bar data	bar data	bar data
								[0]	[1]	[2]
			① read			② read				
				③ read			④ read			
					⑤ read					
read()					⑥ read			⑦ read		
					⑧ write					
					⑨ read					
					⑩ write					
					⑪ read					
					⑫ write					
					⑬ read					
					⑭ write					
					⑮ read					
					⑯ write					
					⑰ read					
					⑱ write					
					⑲ read					
					⑳ write					
					㉑ read					
					㉒ write					
					㉓ read					
					㉔ write					
					㉕ read					
					㉖ write					
					㉗ read					
					㉘ write					
					㉙ read					
					㉚ write					
					㉛ read					
					㉜ write					
					㉝ read					
					㉞ write					
					㉟ read					
					㊱ write					
					㊲ read					
					㊳ write					
					㊴ read					
					㊵ write					
					㊶ read					
					㊷ write					
					㊸ read					
					㊹ write					
					㊺ read					
					㊻ write					
					㊼ read					
					㊽ write					
					㊾ read					
					㊿ write					

블록 2개
 ↑

root directory
 ↓ (inode)
 foo



- Unlike reading, writing allocates a block unless the block is being overwritten
- **create()**: inode bitmap (free inode) → inode (initialization) → directory data (link filename to i-number) → directory inode (update)
- **write()**: data bitmap (free block) → inode (update) → actual write
- **close()**: close the file and deallocate the file descriptor
- 12KB write to **/foo/bar**: 10 I/O for file create and 5 I/O for each block write



Caching and Buffering

- Reading and writing files can be expensive due to many disk I/O**

↑
inode 정보를 다 disk에서 가져와야 함.

 - With a long pathname (e.g., /1/2/3/ ... /100/file.txt), the file system would literally perform hundreds of reads just to open the file
 - Most file systems aggressively use memory (DRAM) to cache important blocks

→ 어떤 블록에 대한 inode 정보를 cache.
→ 비휘발성으로 쓰지 않으면 휘발성 (I/O 성능 향상)
- Early file systems introduced fixed-size cache for popular blocks**

static partition

 - This static partitioning of memory can be wasteful (e.g. 10% of memory)

→ 이런 양함
 - Modern systems, in contrast, employ dynamic partitioning approach; The OS integrates virtual memory pages and file system pages into unified page cache.

pageing
- Imagine the file open example with caching**
 - Opening a file generate many disk I/O (very slow) to read inode, directory data
 - What if opening again the same file or files in the same directory? → cache hit

(no slow disk I/O) → performance ↑

disk 가져 오지 않음.
- Consider the effect of caching on writes**

→ 데이터 불변성 문제나
고려되어야 함.

 - Write buffering certainly has a number of performance benefits:
 - batch update, 2) write scheduling, 3) avoid write (e.g. temporary files)

→ 컴파일 단계 필요함?