

Operating System Design and Implementation

Filesystem and Block I/O

Shiao-Li Tsao

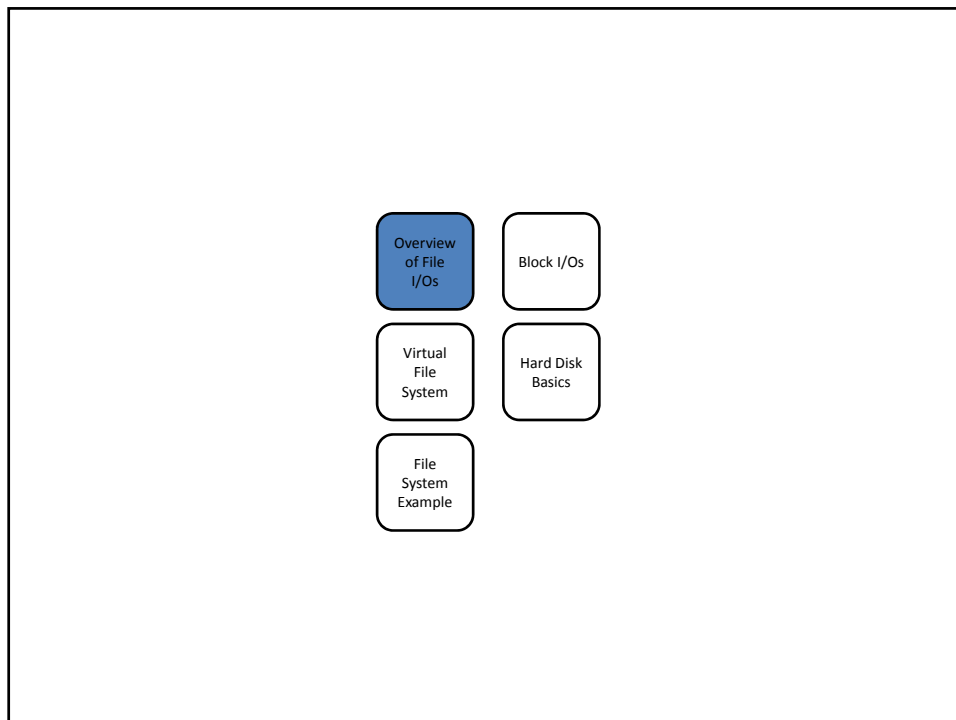
Overview
of File
I/Os

Block I/Os

Virtual
File
System

Hard Disk
Basics

File
System
Example



```
#include <stdio.h>
#include <stdlib.h>

int main()
{
    char ch, file_name[25];
    FILE *fp;

    printf("Enter the name of file you wish to see\n");
    gets(file_name);

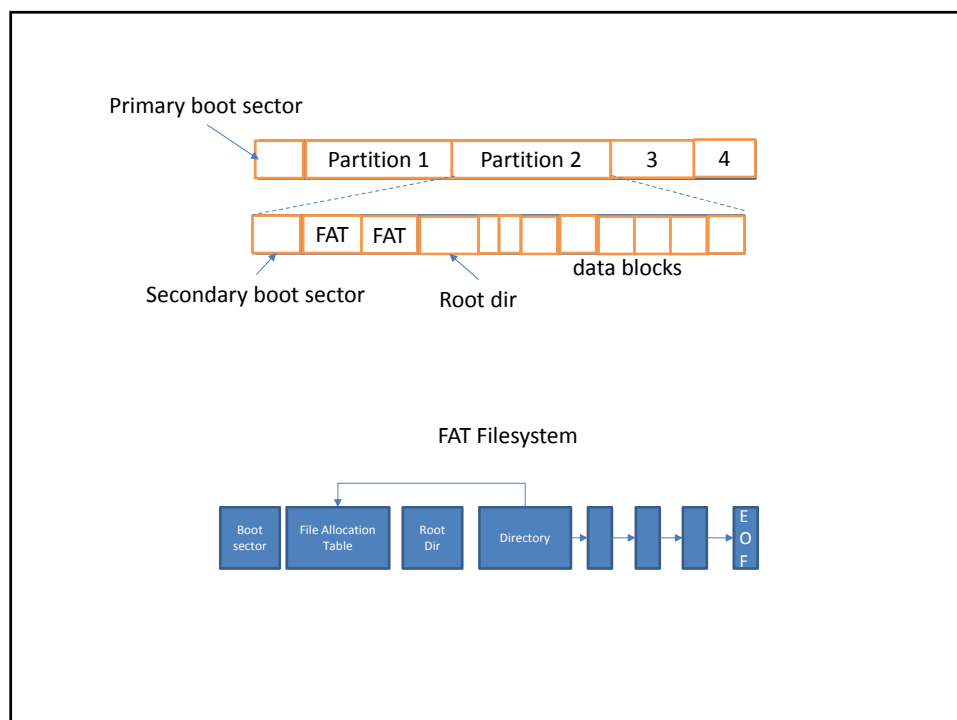
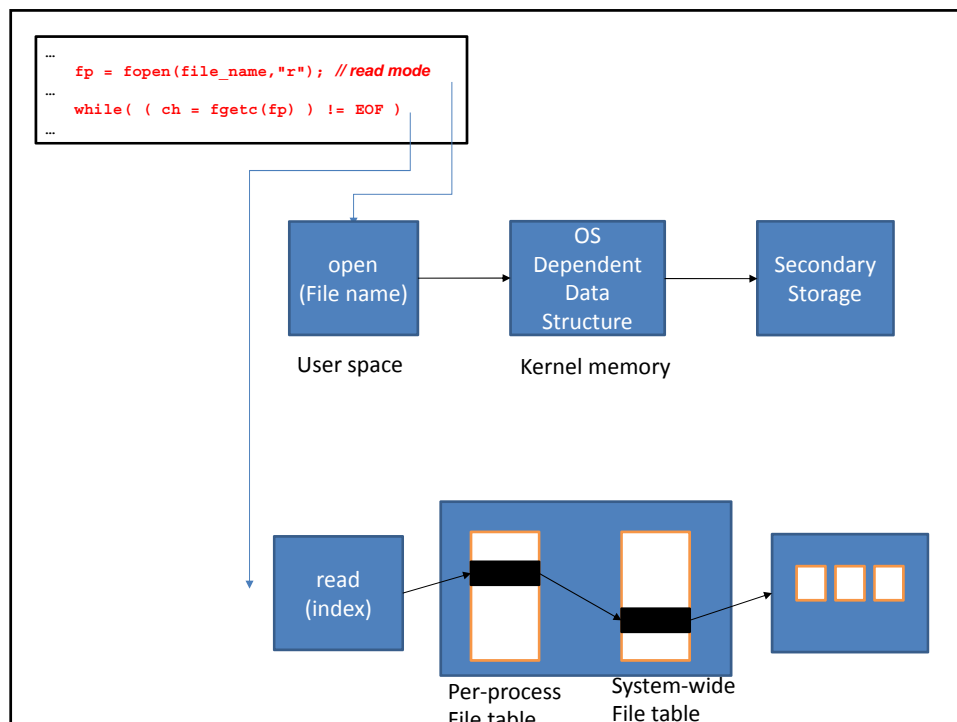
    fp = fopen(file_name, "r"); //read mode

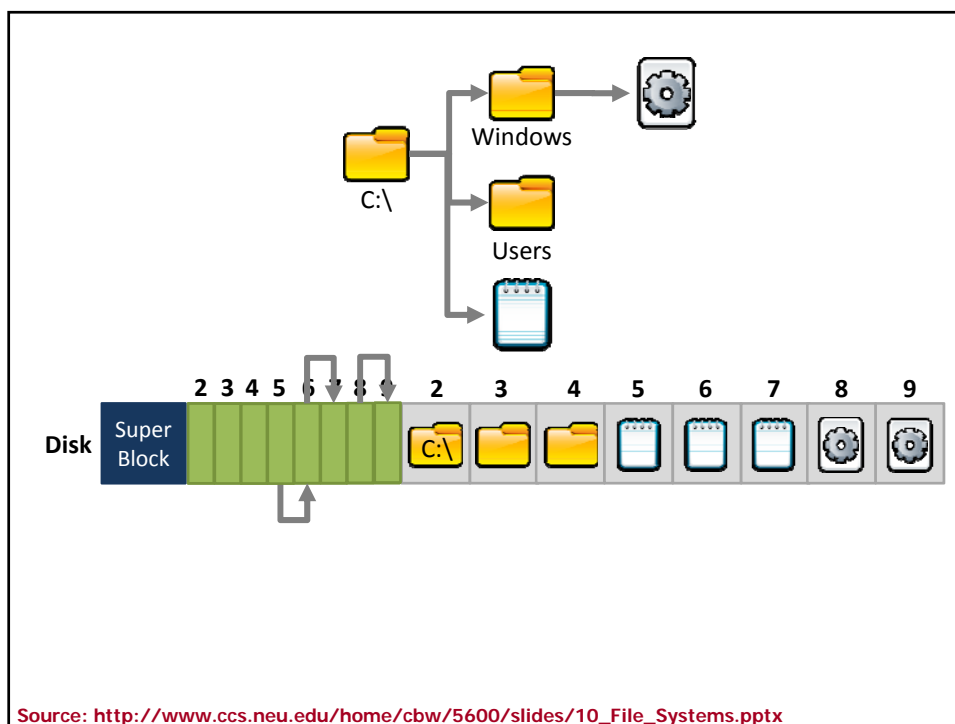
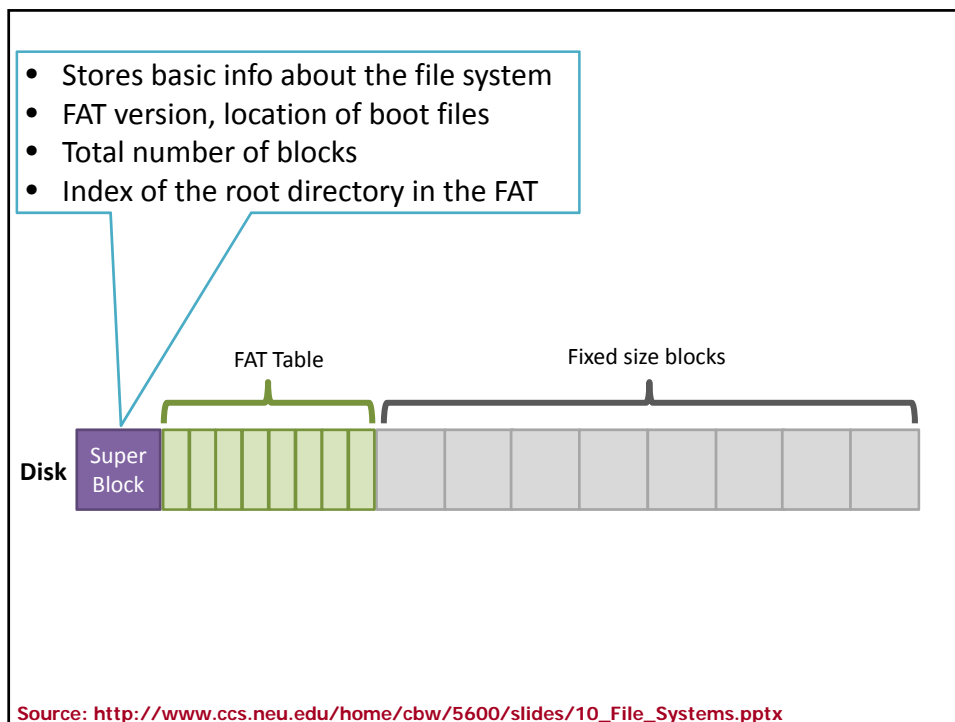
    if( fp == NULL )
    {
        perror("Error while opening the file.\n");
        exit(EXIT_FAILURE);
    }

    printf("The contents of %s file are :\n", file_name);

    while( ( ch = fgetc(fp) ) != EOF )
        printf("%c",ch);

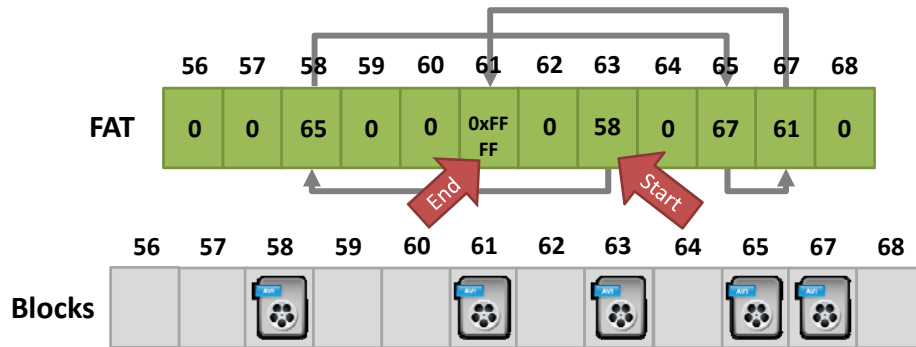
    fclose(fp);
    return 0;
}
```



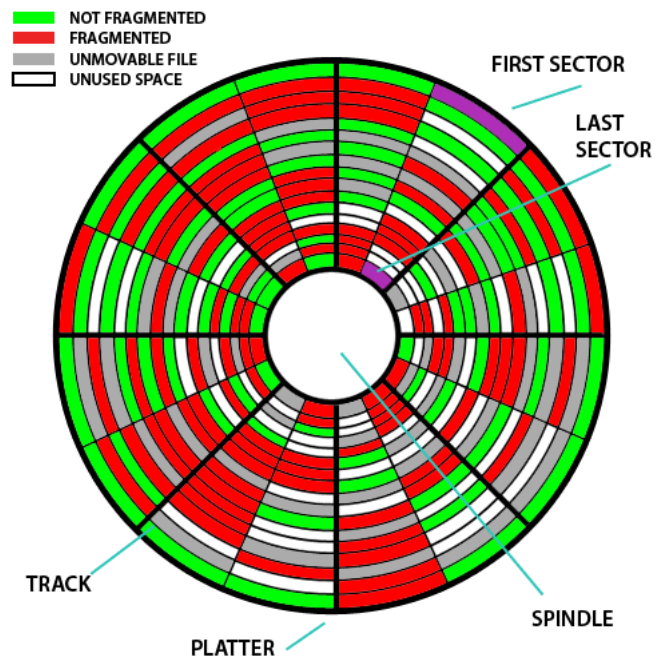


Fragmentation

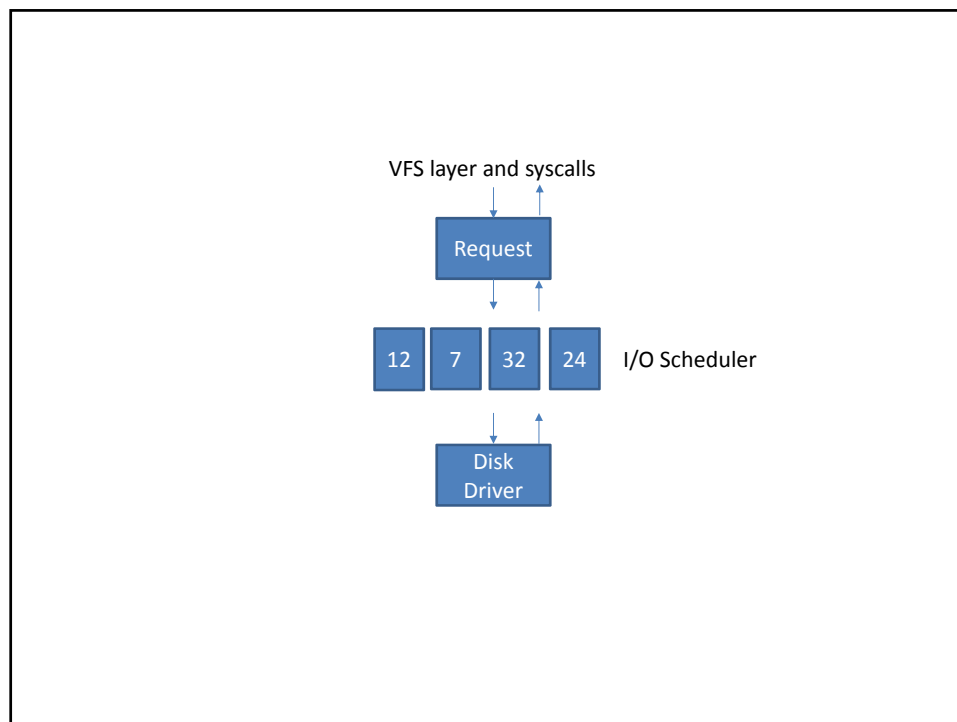
- Blocks for a file need not be contiguous



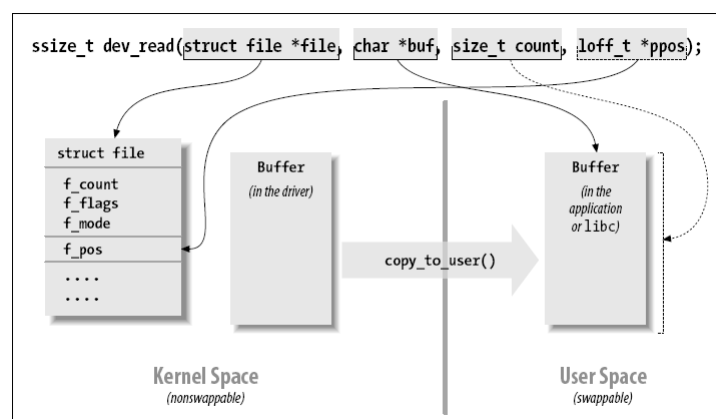
Source: http://www.ccs.neu.edu/home/cbw/5600/slides/10_File_Systems.pptx



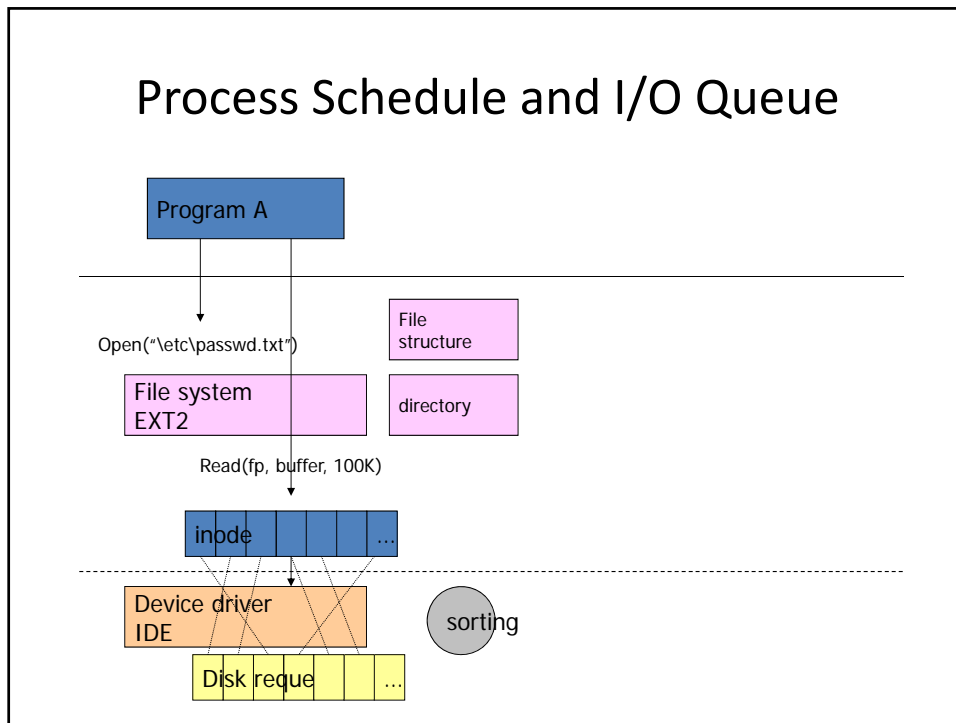
Source: http://www.atomicpages.net/blog/category/tech_tips/software/



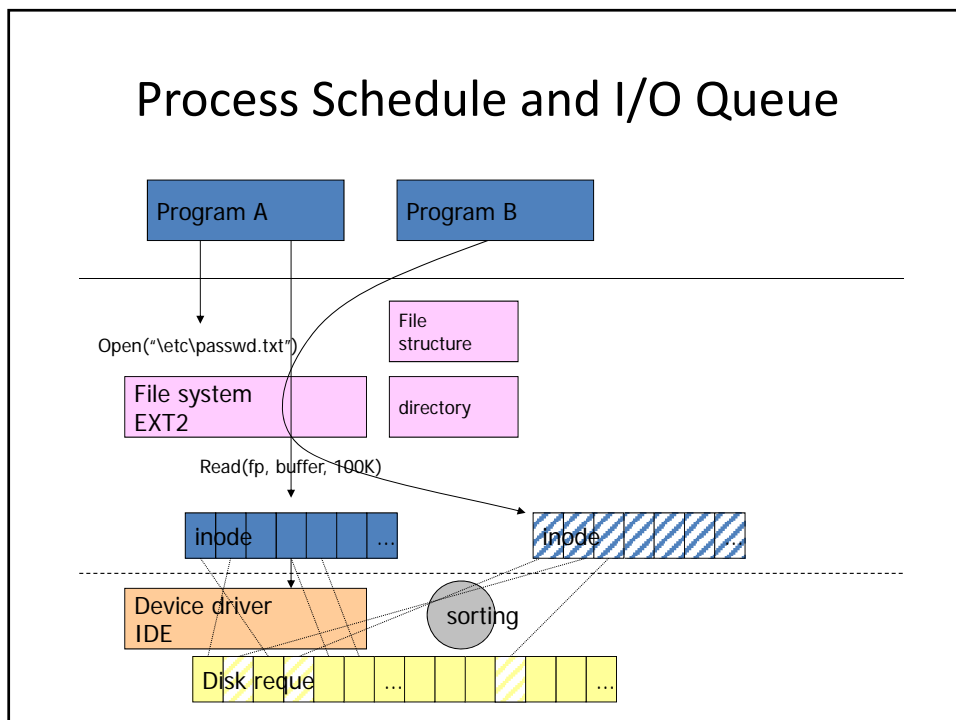
Device read

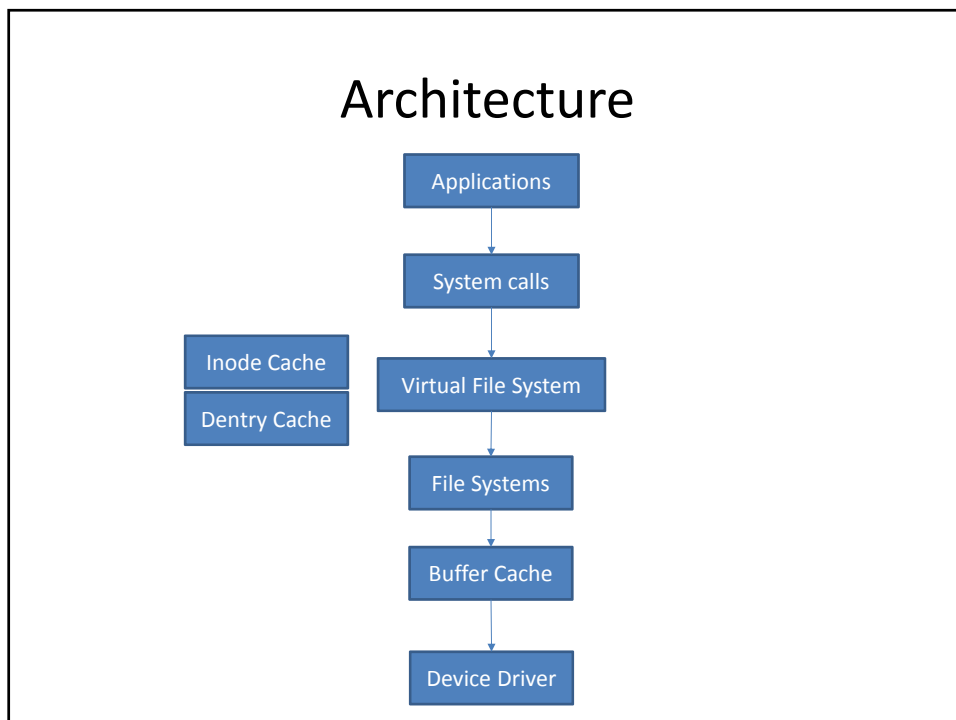
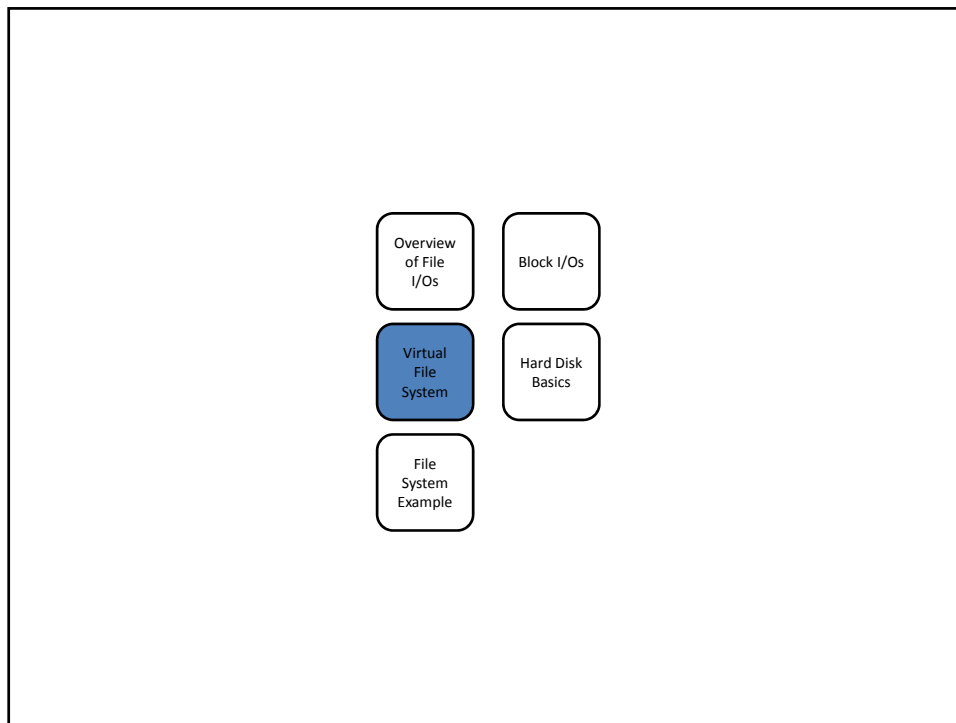


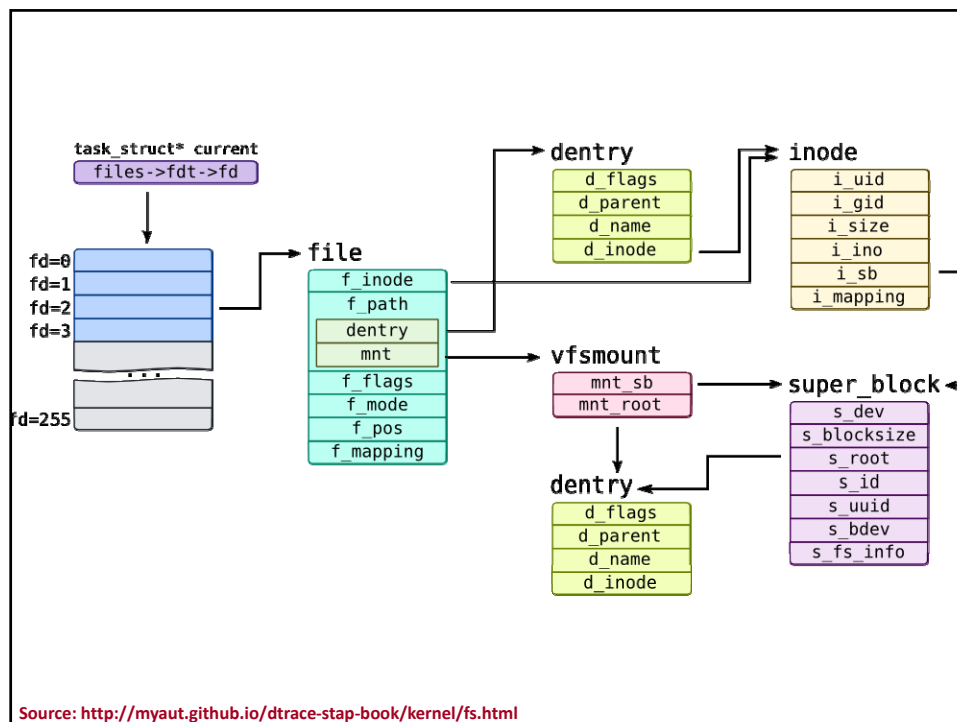
Process Schedule and I/O Queue



Process Schedule and I/O Queue





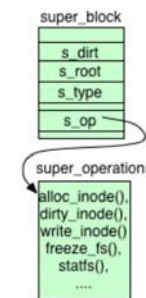


VFS Objects and Data Structures

- superblock object
 - represents a specific mounted filesystem
- inode object
 - represents a specific file
- dentry object
 - represents a directory entry, a single component of a path
- file object
 - represents an open file as associated with a process

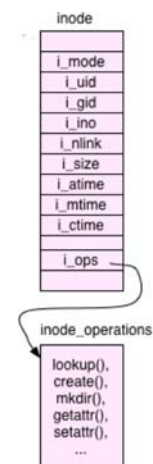
The Superblock Object

- The superblock object is implemented by each filesystem, and is used to store information describing that specific filesystem



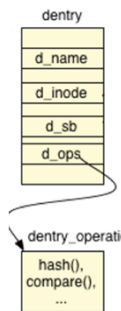
The Inode Object

- The inode object represents all the information needed by the kernel to manipulate a file or directory



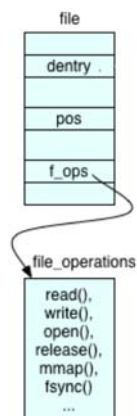
The Dentry Object

- VFS treats directories as files
 - e.g. In the path /bin/vi, both bin and vi are files - bin being the special directory file and vi being a regular file
- An inode object represents these files
- Path name lookup involves translating each component of a path
- To facilitate this, the VFS employs the concept of a directory entry (dentry)
- A dentry is a specific component in a path
 - e.g. /, bin, and vi are all dentry objects

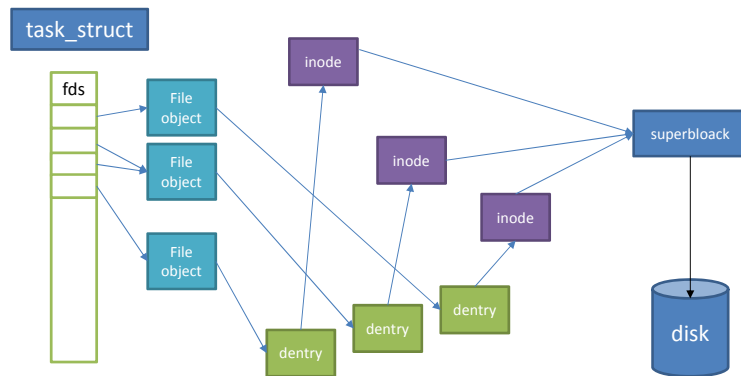


The File Object

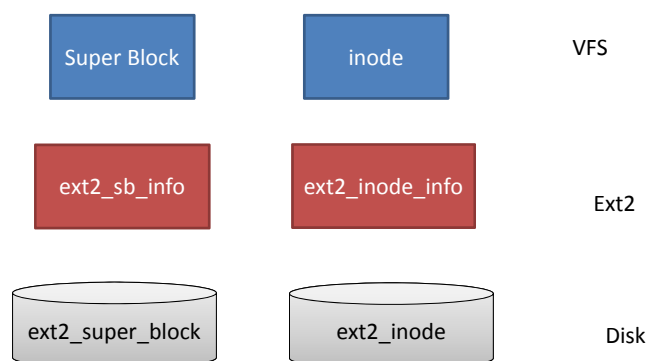
- The file object is used to represent a file opened by a process



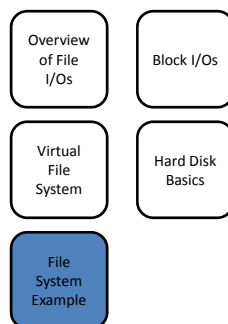
Data Structure Overview of VFS



Data Structure Overview of VFS



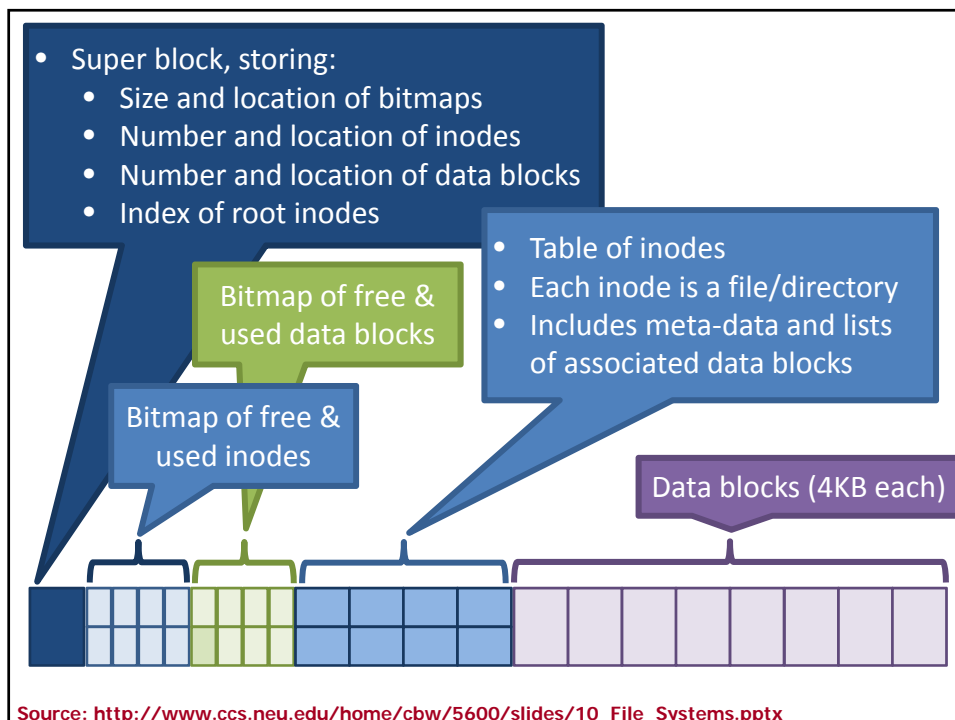
- Ref:
https://kaiwantech.files.wordpress.com/2009/08/vfs_msdos_31.png



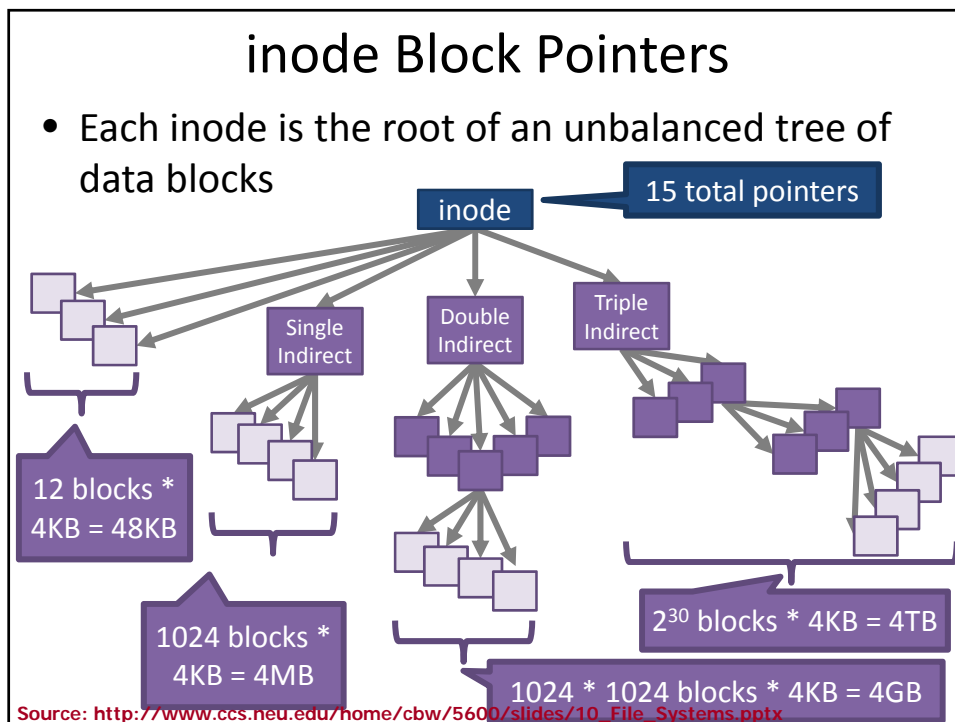
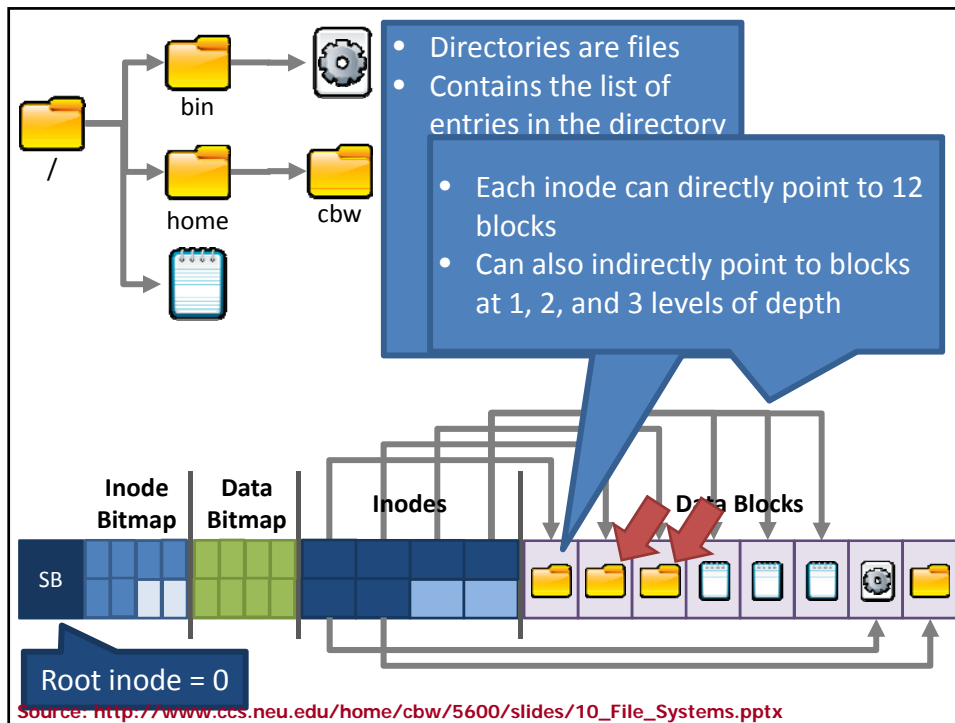
Status Check

- The efficiency of FAT is very low
 - Lots of seeking over file chains in FAT
 - Only way to identify free space is to scan over the entire FAT
- Linux file system uses more efficient structures
 - Extended File System (ext) uses **index nodes (inodes)** to track files and directories

Source: http://www.ccs.neu.edu/home/cbw/5600/slides/10_File_Systems.pptx



Source: http://www.ccs.neu.edu/home/cbw/5600/slides/10_File_Systems.pptx



Advantages of inodes

- Optimized for file systems with many small files
 - Each inode can directly point to 48KB of data
 - Only one layer of indirection needed for 4MB files
- Faster file access
 - Greater meta-data locality → less random seeking
 - No need to traverse long, chained FAT entries
- Easier free space management
 - Bitmaps can be cached in memory for fast access
 - inode and data space handled independently

Source: http://www.ccs.neu.edu/home/cbw/5600/slides/10_File_Systems.pptx

File Reading Example

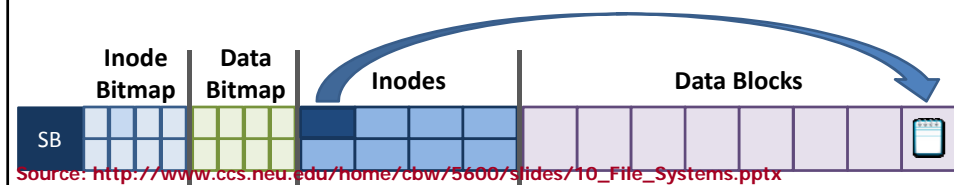
	Bitmaps		inodes			Data Blocks				
	data	inode	root	tmp	file	root	tmp	file[0]	file[1]	file[3]
Time ↓	open("/tmp/file")		read							
						read				
				read						
					read		read			
	read()				read			read		
					write					
	read()				read				read	
					write					
	read()				read					read
					write					

Update the last
accessed time
of the file

Source: http://www.ccs.neu.edu/home/cbw/5600/slides/10_File_Systems.pptx

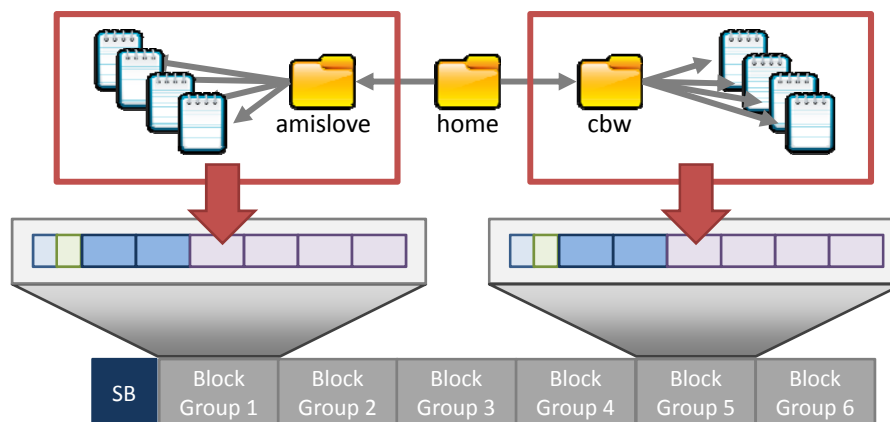
ext: The Good and the Bad

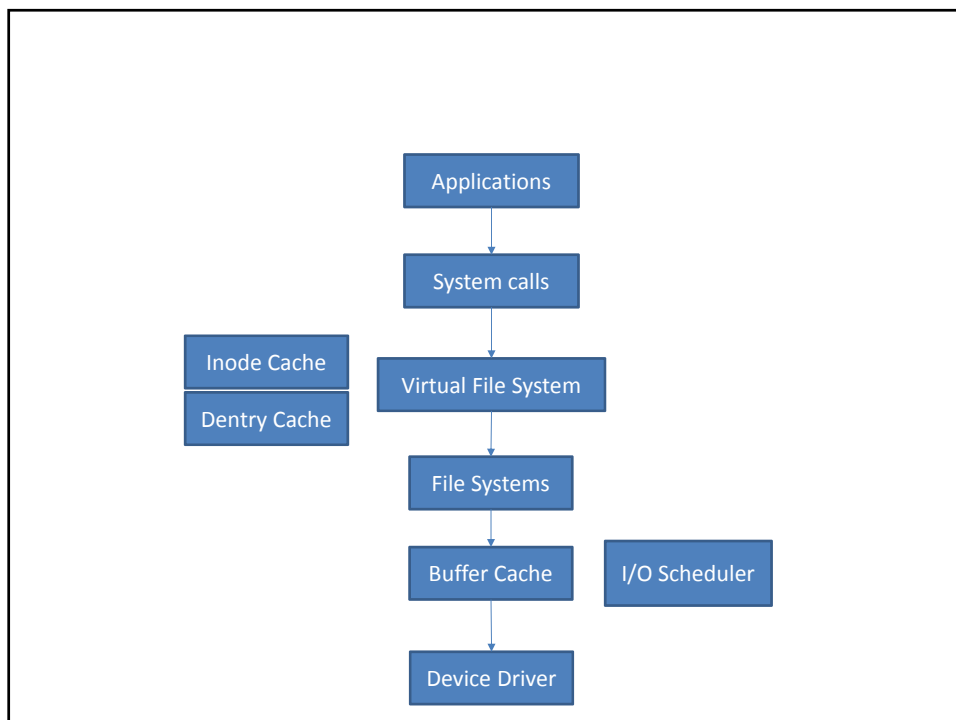
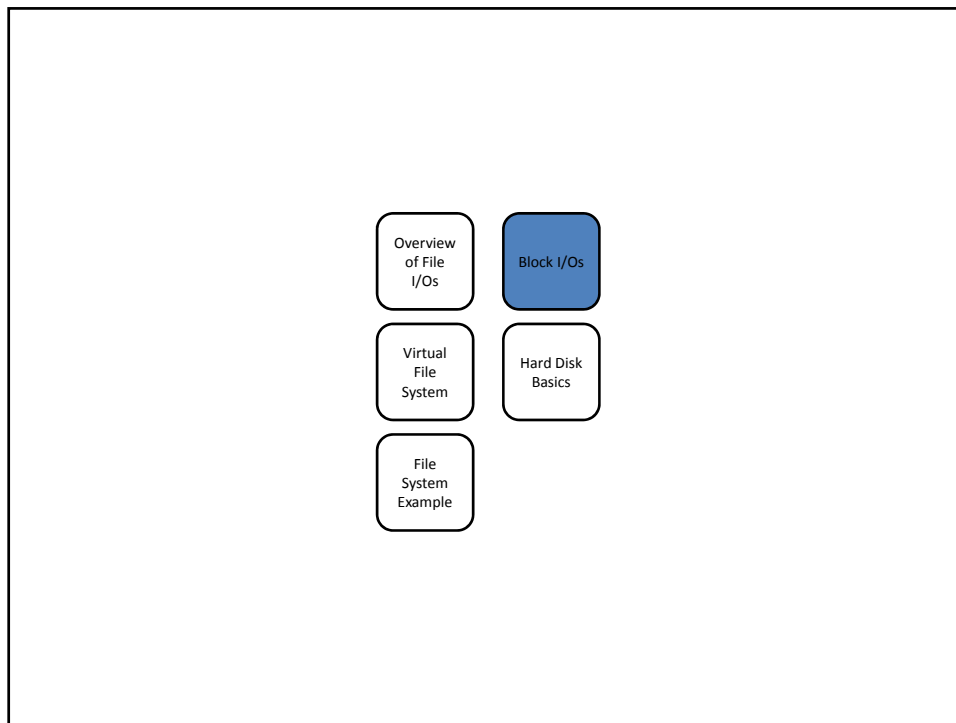
- The Good – ext file system (inodes) support:
 - All the typical file/directory features
 - Hard and soft links
 - More performant (less seeking) than FAT
- The Bad: poor locality
 - ext is optimized for a particular file size distribution
 - However, it is not optimized for spinning disks
 - inodes and associated data are far apart on the disk!



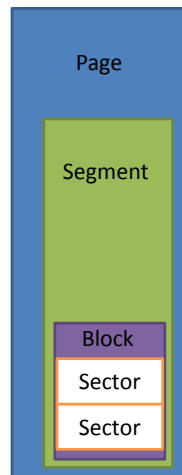
Allocation Policy

- ext2 attempts to keep related files and directories within the same block group





Typical Layout of Pages/Blocks



Block device driver

- I/O Scheduler
 - Sending out requests to the block devices, as soon as it issues them, results in awful performance
 - Minimizing seeks is absolutely crucial to the system's performance
 - Performs operations called merging and sorting to greatly improve the performance

38

Block device driver

- I/O Scheduler
 - Merge
 - Merge multiple requests together if they are adjacent sectors
 - Sort
 - Insert requests to the list based on the location (to avoid extra seek)

39

Block device driver

- The Linus Elevator
 - Used in 2.4
 - Algorithm
 - First, if a request to an adjacent on-disk sector is in the queue, the existing request and the new request are merged into a single request.
 - Second, if a request in the queue is sufficiently old, the new request is inserted at the tail of the queue to prevent starvation of the other, older, requests.
 - Next, if there is a suitable location sector-wise in the queue, the new request is inserted there. This keeps the queue sorted by physical location on disk.
 - Finally, if no such suitable insertion point exists, the request is inserted at the tail of the queue

40

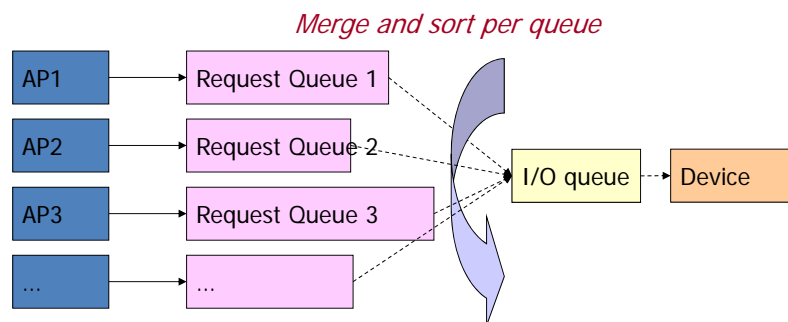
Block device driver

- Deadline I/O scheduler
 - Used in 2.6
 - Consider deadline
 - read: 500ms, write: 5s
 - Two queues
 - Sector queue: Sorted by sector (read/write queues)
 - FIFO queue: Sorted by expire time (read/write queues)
 - Sector queue first, if FIFO queue expire, schedule FIFO queue
 - Sort and merge request without starvation
 - Prefer read requests (deadline)

41

Block device driver

- Complete fair queuing I/O scheduler
 - Used in 2.6

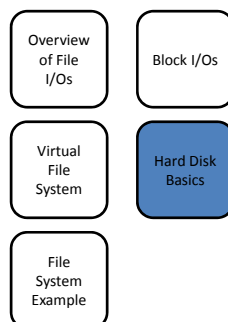


42

Block device driver

- Noop I/O scheduler
 - Used in 2.6
 - No sorting, and merging only
 - Merge with adjacent request only
 - Near-FIFO order
 - For truly random-access device

43



Source: <http://www.theadatarecovery.co.uk/how-do-hard-disk-drives-work/>

Source:
<https://www.microsoft.com/resources/documentation/windowsnt/4/workstation/reskit/en-us/diskdesc.msp?mfr=true>

Source: <http://how-computers-work.blogspot.tw>

Source: <http://www.marvell.com/storage/hdd-soc/>

Source: <http://www.mipi.org/>