# File Systems
# CS 111
# Operating System Principles
# Peter Reiher

# Outline

- File systems:
  - Why do we need them?
  - Why are they challenging?

- Basic elements of file system design

- Designing file systems for disks
  - Basic issues
  - Free space, allocation, and deallocation

# Introduction

- Most systems need to store data persistently
  - So it's still there after reboot, or even power down
- Typically a core piece of functionality for the system
  - Which is going to be used all the time
- Even the operating system itself needs to be stored this way
- So we must store some data persistently
  - Most commonly on a disk drive

# Our Persistent Data Options

- Use raw persistent storage to store the data
  - Hard for users to work with
  - Not much easier for OS developers

- Use a database to store the data
  - Probably more structure (and possibly overhead) than we need or can afford

- Use a file system
  - Some organized way of structuring persistent data
  - Which makes sense to users and programmers

# File Systems

- Originally the computer equivalent of a physical filing cabinet

- Put related sets of data into individual containers

- Put them all into an overall storage unit

- Organized by some simple principle
    - E.g., alphabetically by title
    - Or chronologically by date

- Goal is to provide:
    - Persistence
    - Ease of access
    - Good performance

# The Basic File System Concept

- Organize data into natural coherent units
  - Like a paper, a spreadsheet, a message, a program
- Store each unit as its own self-contained entity
  - A *file*
  - Store each file in a way allowing efficient access
- Provide some simple, powerful organizing principle for the collection of files
  - Making it easy to find them
  - And easy to organize them
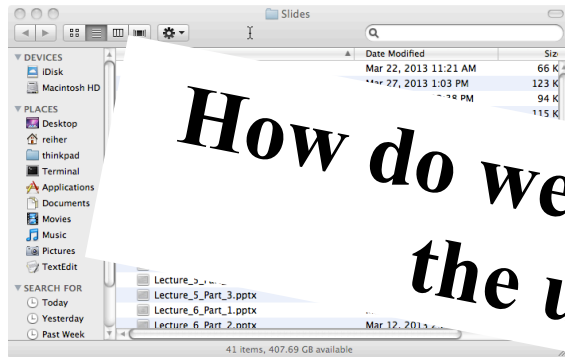
# File Systems and Hardware

- File systems are typically stored on hardware providing persistent memory

  – Disks, tapes, flash memory, etc.

- With the expectation that a file put in one "place" will be there when we look again

- Performance considerations will require us to match the implementation to the hardware

- But ideally, the same user-visible file system should work on any reasonable hardware

# Data and Metadata

- File systems deal with two kinds of information

- *Data* – the information that the file is actually supposed to store

  – E.g., the instructions of the program or the words in the letter

- *Metadata* – Information about the information the file stores

  – E.g., how many bytes are there and when was it created
  – Sometimes called *attributes*

- Ultimately, both data and metadata must be stored persistently

  – And usually on the same piece of hardware

# Bridging the Gap

We want something like . . .

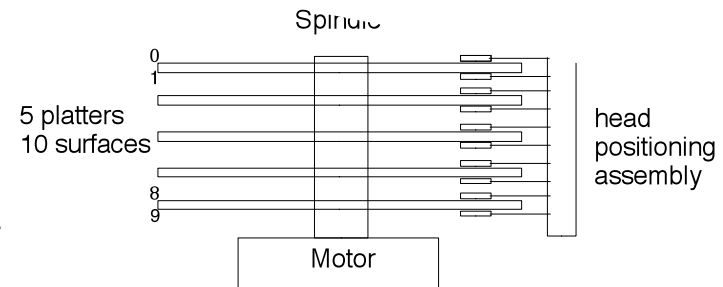But we've got something like . . .

How do we get from the hardware to the useful abstraction?

Or . . .

Or at least

```
drwxr-xr-x     8 root   wheel     272 May   4   2010 X11
lrwxr-xr-x     1 root   wheel       3 May   4   2010 X11R6 -> X11
drwxr-xr-x   913 root   wheel   31042 Apr  21 12:21 bin
drwxr-xr-x   336 root   wheel   11424 Mar  17 09:13 lib
drwxr-xr-x   103 root   wheel    3502 Apr  21 12:23 libexec
drwxr-xr-x     7 root   wheel     238 Jan  16 23:00 local
drwxr-xr-x   238 root   wheel    8092 Mar  17 09:13 sbin
drwxr-xr-x    59 root   wheel    2006 Apr  21 12:21 share
drwxr-xr-x     4 root   wheel     136 May   4   2010 standalone
```

Spindle

0
1

5 platters
10 surfaces

head
positioning
assembly

8
9

Motor

# A Further Wrinkle

- We want our file system to be agnostic to the storage medium
- Same program should access the file system the same way, regardless of actual storage medium
  - Otherwise hard to write portable programs
- Should work the same for disks of different types
- Or if we use a RAID instead of one disk
- Or if we use flash instead of disks
- Or if even we don't use persistent memory at all
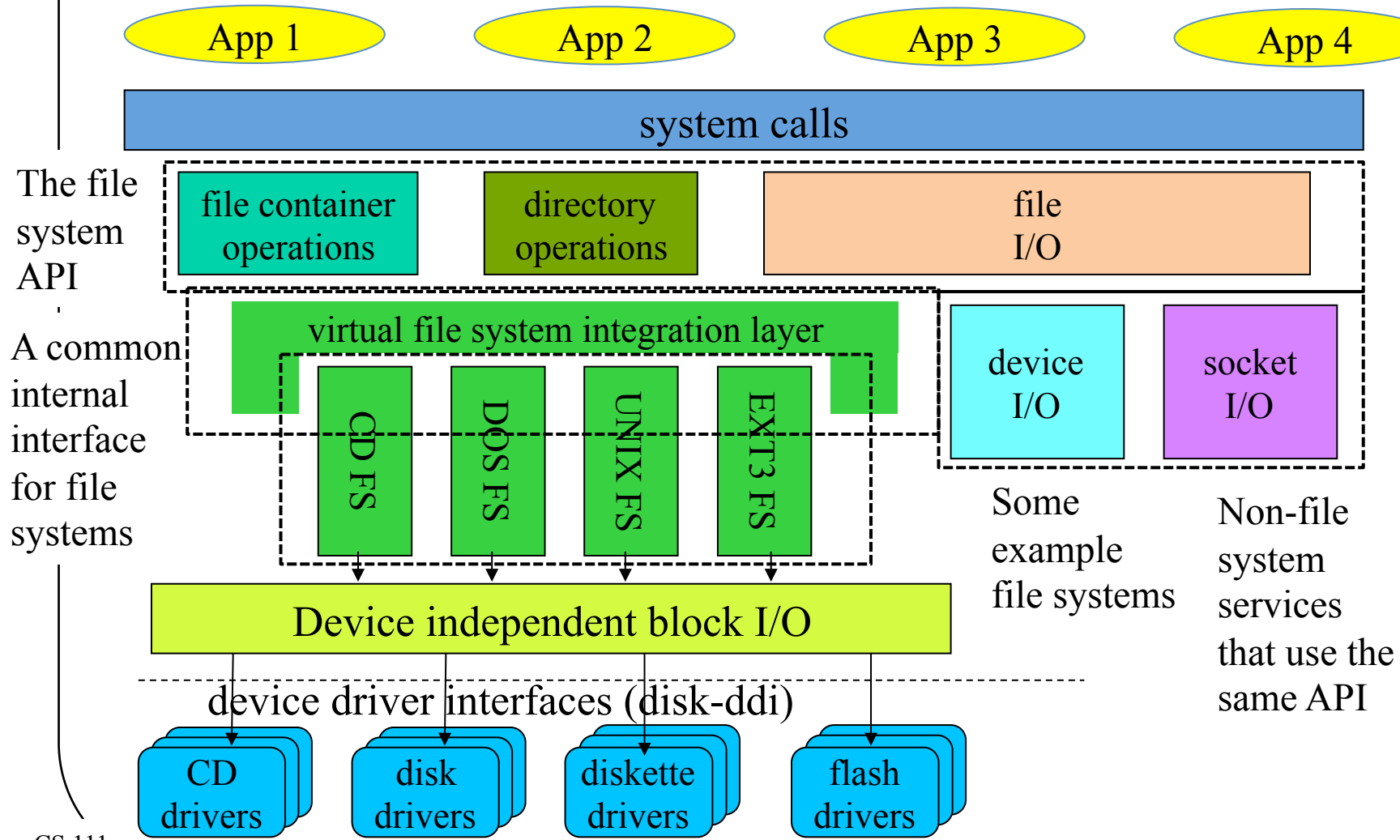  - E.g., RAM file systems

# Desirable File System Properties

- What are we looking for from our file system?
  - Persistence
  - Easy use model
    - For accessing one file
    - For organizing collections of files
  - Flexibility
    - No limit on number of files
    - No limit on file size, type, contents
  - Portability across hardware device types
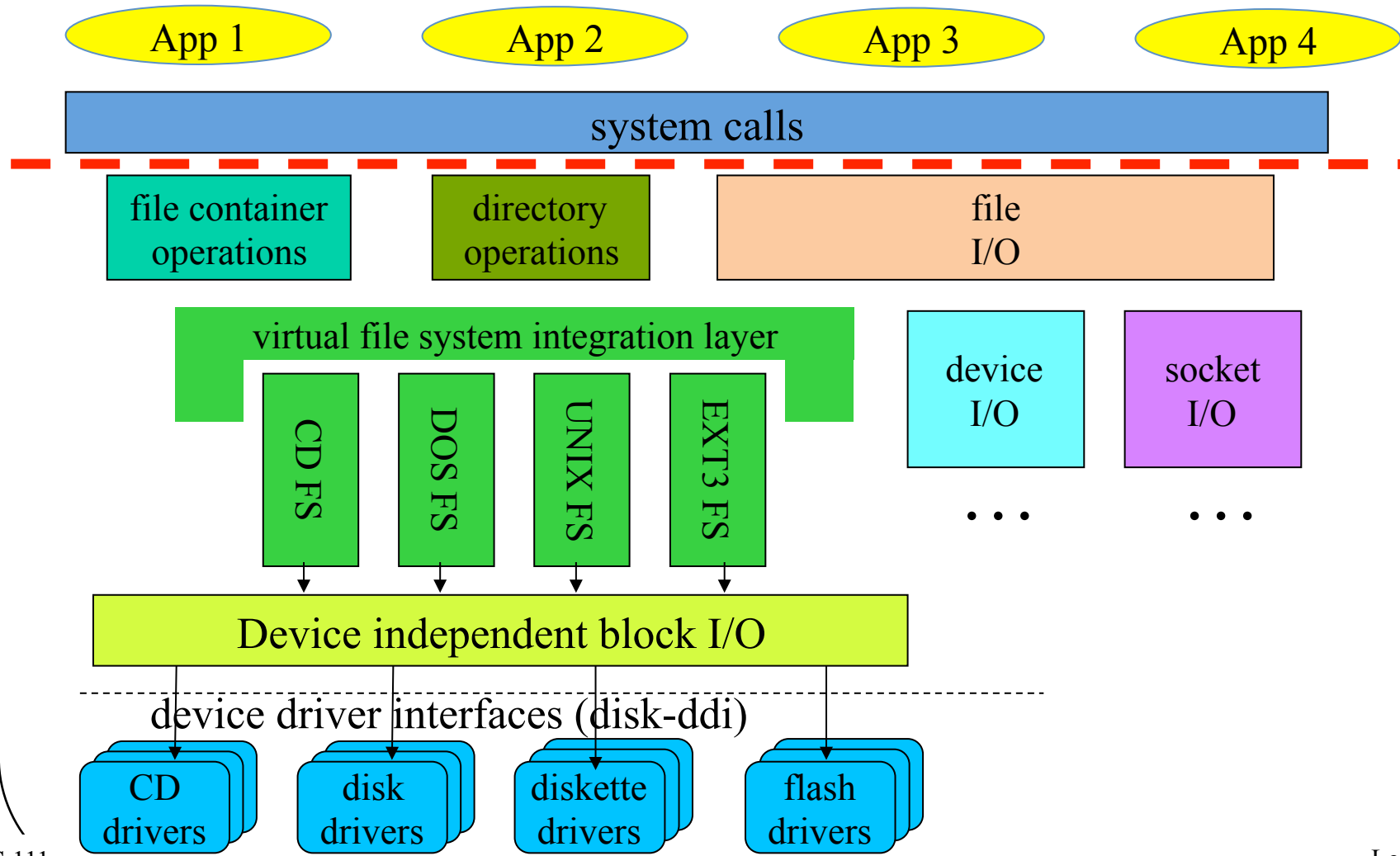  - Performance
  - Reliability
  - Suitable security

# Basics of File System Design

- Where do file systems fit in the OS?

- File control data structures

# File Systems and the OS

App 1     App 2     App 3     App 4

system calls

The file system API

file container operations

directory operations

file I/O

A common internal interface for file systems

virtual file system integration layer

CD FS

DOS FS

UNIX FS

EXT3 FS

device I/O

socket I/O

Some example file systems

Non-file system services that use the same API

Device independent block I/O

device driver interfaces (disk-ddi)

CD drivers

disk drivers

diskette drivers

flash drivers

# The File System API

App 1    App 2    App 3    App 4

system calls

file container operations     directory operations     file I/O

virtual file system integration layer

CD FS    DOS FS    UNIX FS    EXT3 FS

device I/O    socket I/O

. . .    . . .

Device independent block I/O

device driver interfaces (disk-ddi)

CD drivers    disk drivers    diskette drivers    flash drivers

# The File System API

- Highly desirable to provide a single API to programmers and users for all files

- Regardless of how the file system underneath is actually implemented

- A requirement if one wants program portability
  - Very bad if a program won't work because there's a different file system underneath

- Three categories of system calls here

  1. File container operations

  2. Directory operations

  3. File I/O operations

# File Container Operations

- Standard file management system calls
  - Manipulate files as objects
  - These operations ignore the contents of the file
- Implemented with standard file system methods
  - Get/set attributes, ownership, protection ...
  - Create/destroy files and directories
  - Create/destroy links
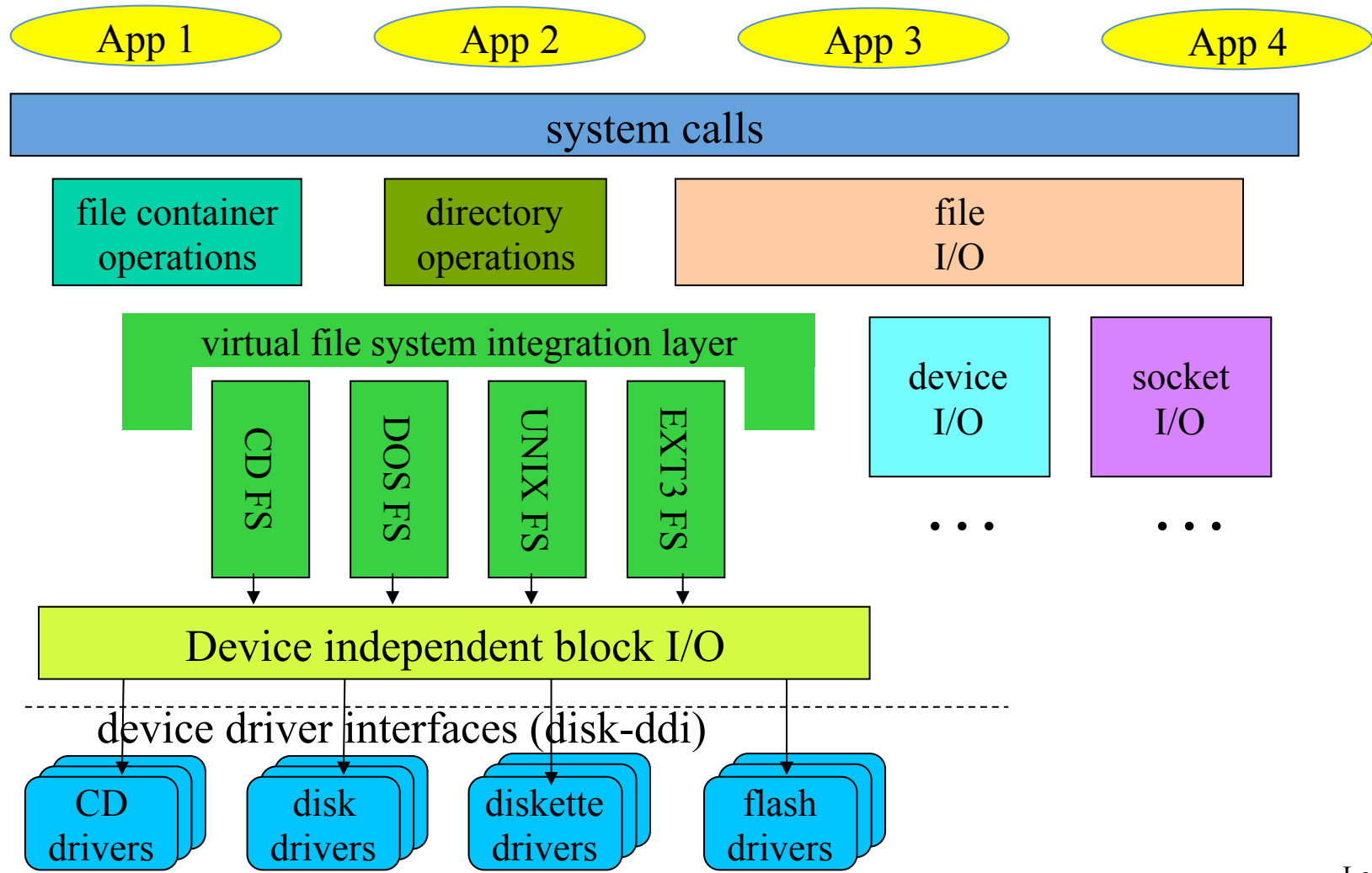- Real work happens in file system implementation

# Directory Operations

- Directories provide the organization of a file system

  - Typically hierarchical

  - Sometimes with some extra wrinkles

- At the core, directories translate a name to a lower-level file pointer

- Operations tend to be related to that

  - Find a file by name

  - Create new name/file mapping

  - List a set of known names

# File I/O Operations

- Open – map name into an open instance
- Read data from file and write data to file
  - Implemented using logical block fetches
  - Copy data between user space and file buffer
  - Request file system to write back block when done
- Seek
  - Change logical offset associated with open instance
- Map file into address space
  - File block buffers are just pages of physical memory
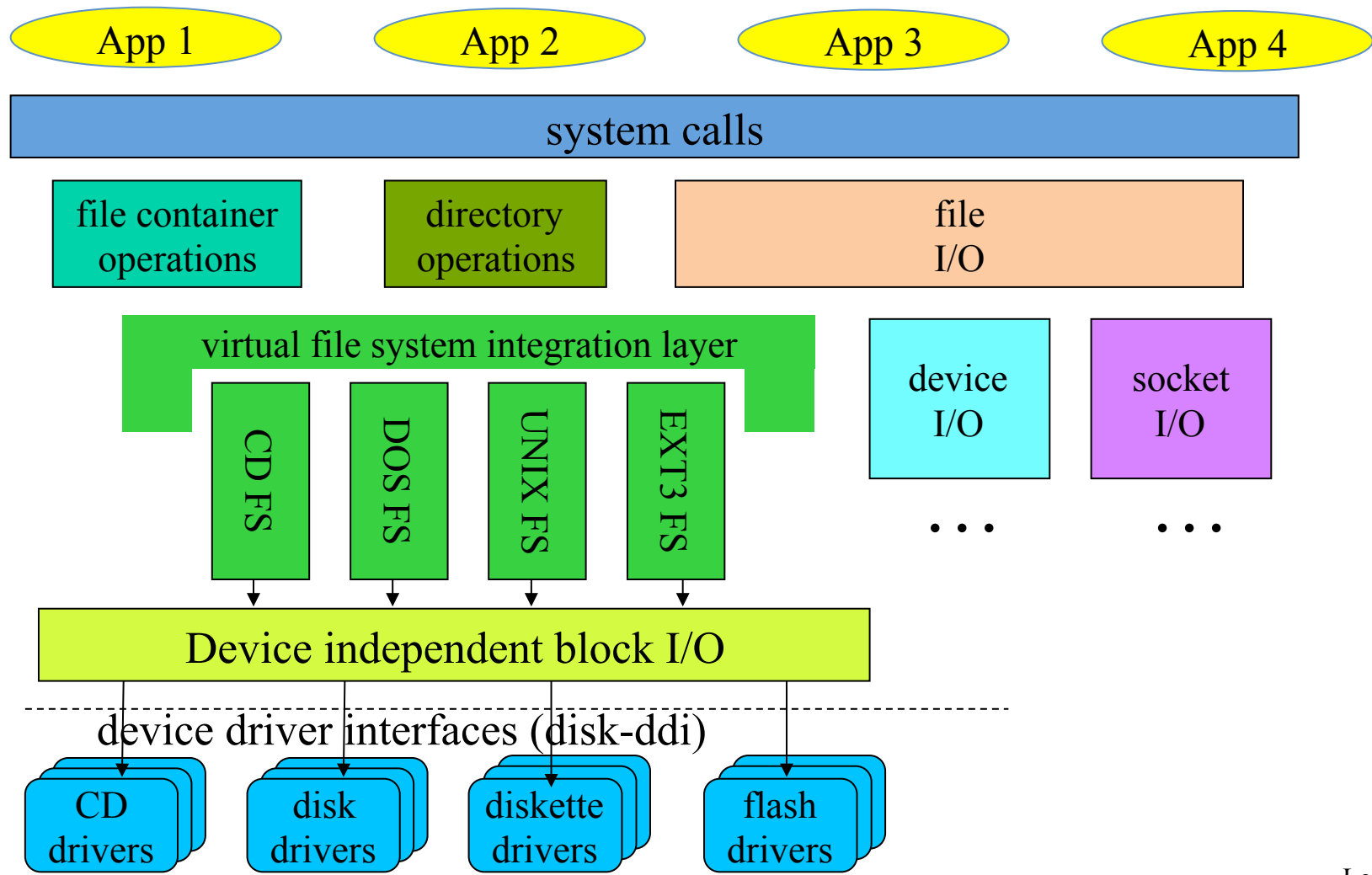  - Map into address space, page it to and from file system

# The Virtual File System Layer

App 1  App 2  App 3  App 4

system calls

| file container operations | directory operations | file I/O |
|---|---|---|

virtual file system integration layer

CD FS  DOS FS  UNIX FS  EXT3 FS

device I/O  socket I/O

. . .  . . .

Device independent block I/O

device driver interfaces (disk-ddi)

CD drivers  disk drivers  diskette drivers  flash drivers

# The Virtual File System (VFS) Layer

- Federation layer to generalize file systems
  - Permits rest of OS to treat all file systems as the same
  - Support dynamic addition of new file systems
- Plug-in interface for file system implementations
  - DOS FAT, Unix, EXT3, ISO 9660, network, etc.
  - Each file system implemented by a plug-in module
  - All implement same basic methods
    - Create, delete, open, close, link, unlink,
    - Get/put block, get/set attributes, read directory, etc.
- Implementation is hidden from higher level clients
  - All clients see are the standard methods and properties
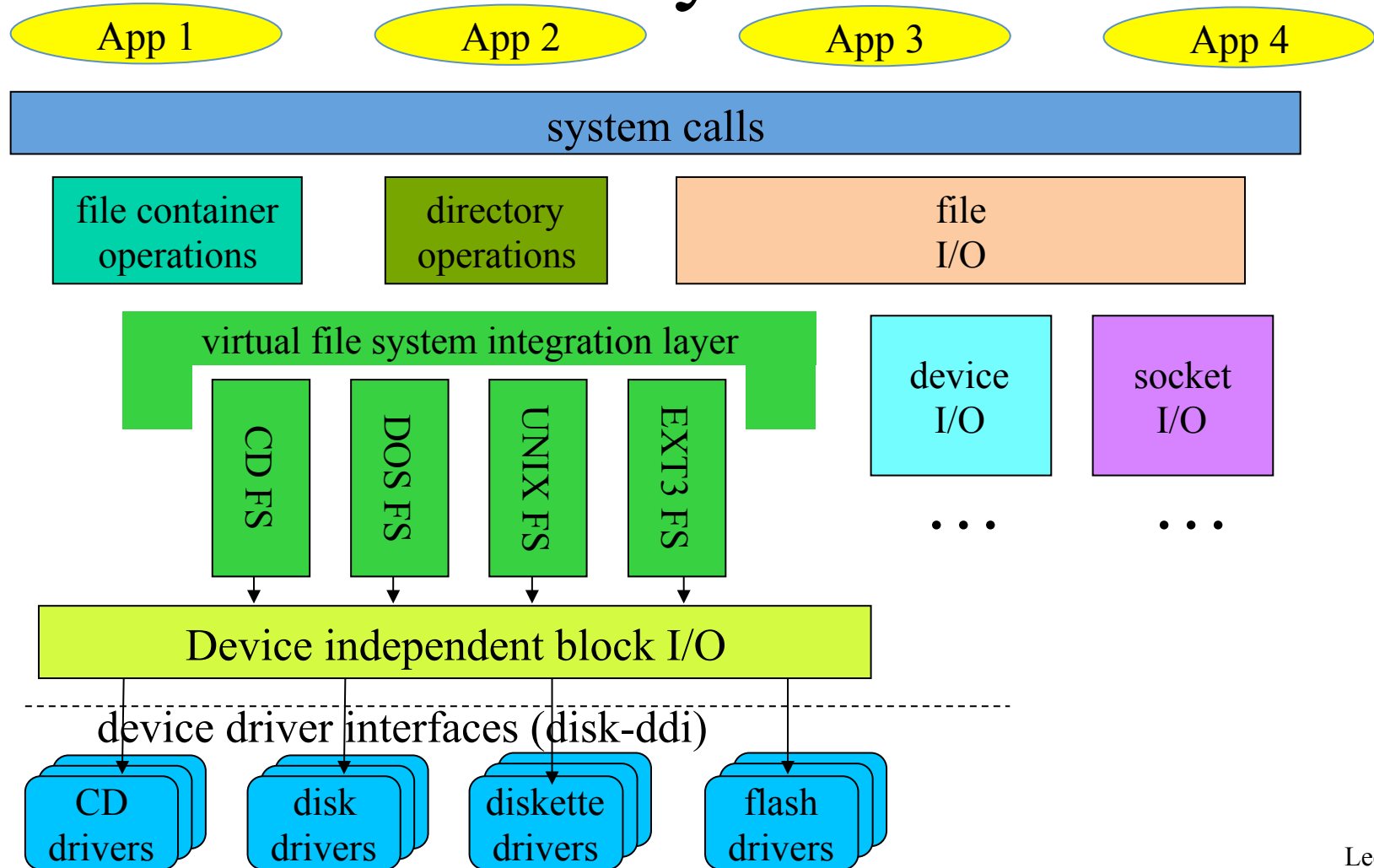
# The File System Layer

App 1    App 2    App 3    App 4

system calls

| file container operations | directory operations | file I/O |

virtual file system integration layer

CD FS    DOS FS    UNIX FS    EXT3 FS

device I/O    socket I/O

. . .    . . .

Device independent block I/O

device driver interfaces (disk-ddi)

CD drivers    disk drivers    diskette drivers    flash drivers

# The File Systems Layer

- Desirable to support multiple different file systems

- All implemented on top of block I/O
  - <u>Should</u> be independent of underlying devices

- All file systems perform same basic functions
  - Map names to files
  - Map <file, offset> into <device, block>
  - Manage free space and allocate it to files
  - Create and destroy files
  - Get and set file attributes
  - Manipulate the file name space

# Why Multiple File Systems?

- Why not instead choose one "good" one?

- There may be multiple storage devices
  - E.g., hard disk and flash drive
  - They might benefit from very different file systems

- Different file systems provide different services, despite the same interface
  - Differing reliability guarantees
  - Differing performance
  - Read-only vs. read/write

- Different file systems used for different purposes
  - E.g., a temporary file system

# Device Independent Block I/O Layer

App 1  App 2  App 3  App 4

system calls

| file container operations | directory operations | file I/O |
|---|---|---|

virtual file system integration layer

CD FS  DOS FS  UNIX FS  EXT3 FS

device I/O   socket I/O

. . .   . . .

Device independent block I/O

device driver interfaces (disk-ddi)

CD drivers  disk drivers  diskette drivers  flash drivers

# File Systems and Block I/O Devices

- File systems typically sit on a general block I/O layer

- A generalizing abstraction – make all disks look same

- Implements standard operations on each block device
  - Asynchronous read (physical block #, buffer, bytecount)
  - Asynchronous write (physical block #, buffer, bytecount)

- Map logical block numbers to device addresses
  - E.g., logical block number to <cylinder, head, sector>

- Encapsulate all the particulars of device support
  - I/O scheduling, initiation, completion, error handlings
  - Size and alignment limitations

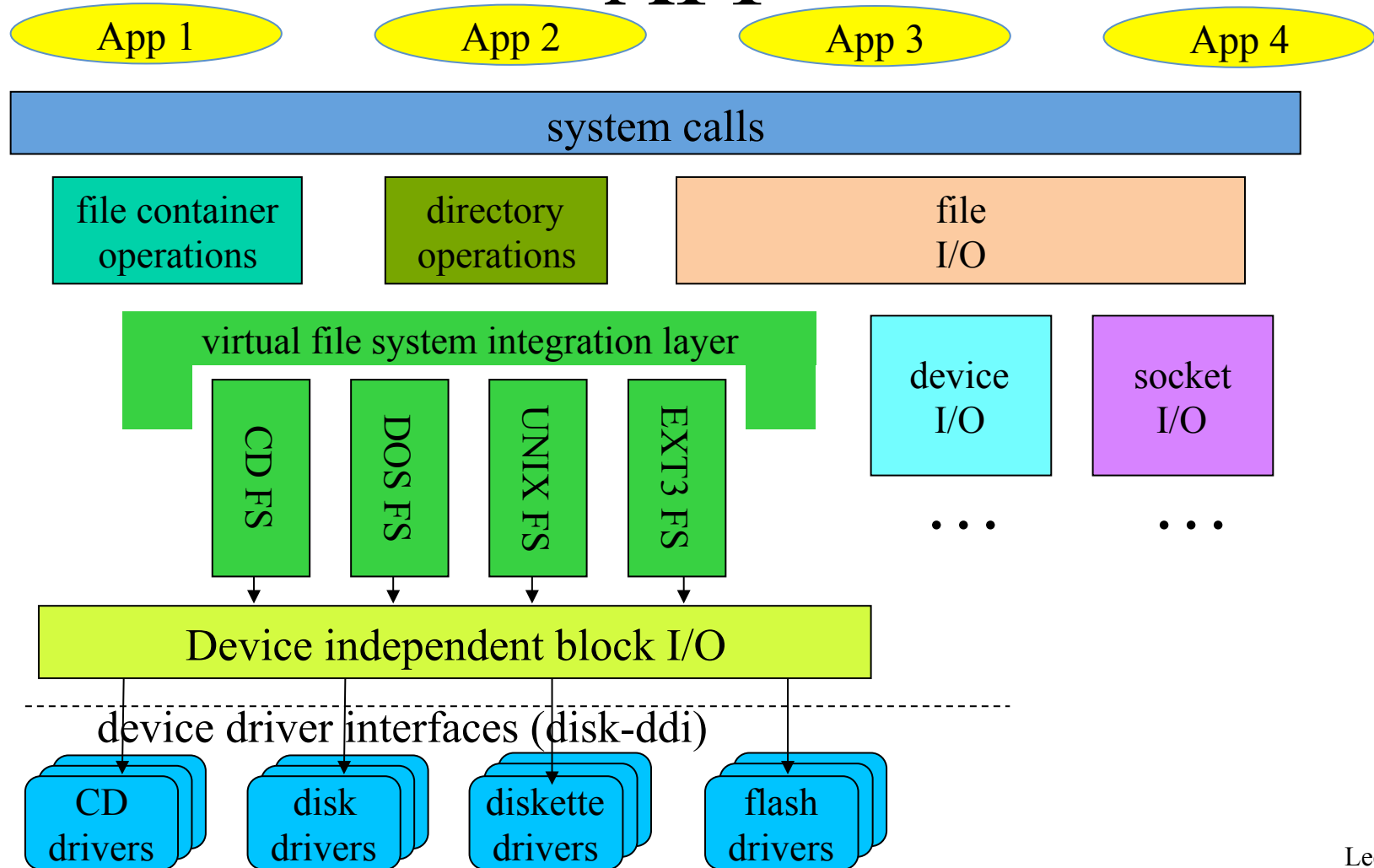# Why Device Independent Block I/O?

- A better abstraction than generic disks

- Allows unified LRU buffer cache for disk data
  - Hold frequently used data until it is needed again
  - Hold pre-fetched read-ahead data until it is requested

- Provides buffers for data re-blocking
  - Adapting file system block size to device block size
  - Adapting file system block size to user request sizes

- Handles automatic buffer management
  - Allocation, deallocation
  - Automatic write-back of changed buffers

# Why Do We Need That Cache?

- File access exhibits a high degree of reference locality at multiple levels:

  – Users often read and write a single block in small operations, reusing that block

  – Users read and write the same files over and over

  – Users often open files from the same directory

  – OS regularly consults the same meta-data blocks

- Having common cache eliminates many disk accesses, which are slow
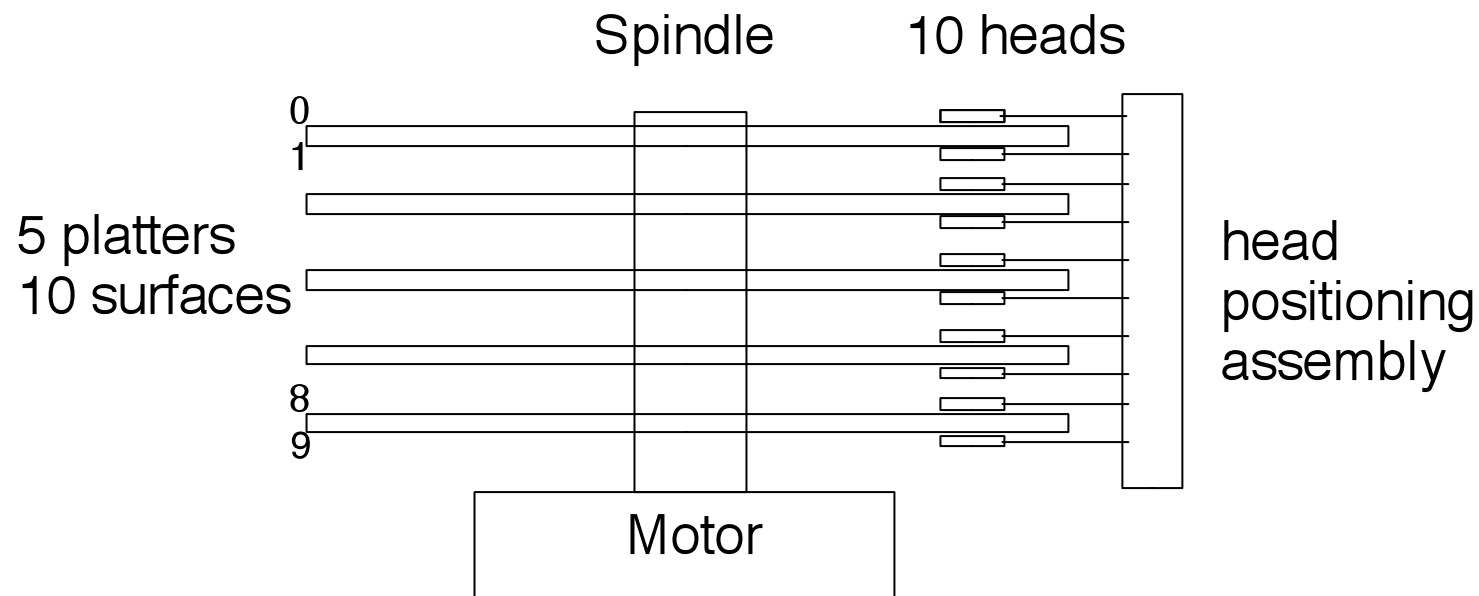
# Devices, Sockets and File System API

App 1   App 2   App 3   App 4

system calls

| file container operations | directory operations | file I/O |
|---|---|---|

virtual file system integration layer

CD FS | DOS FS | UNIX FS | EXT3 FS

device I/O

socket I/O

. . .   . . .

Device independent block I/O

device driver interfaces (disk-ddi)

CD drivers   disk drivers   diskette drivers   flash drivers

# Disk Drives

- Still the primary method of providing stable storage

  - Storage meant to last beyond a single power cycle of the computer

  - Particularly for file systems

- Getting good performance from disk drives is critical for file system performance

- A place where physics meets computer science
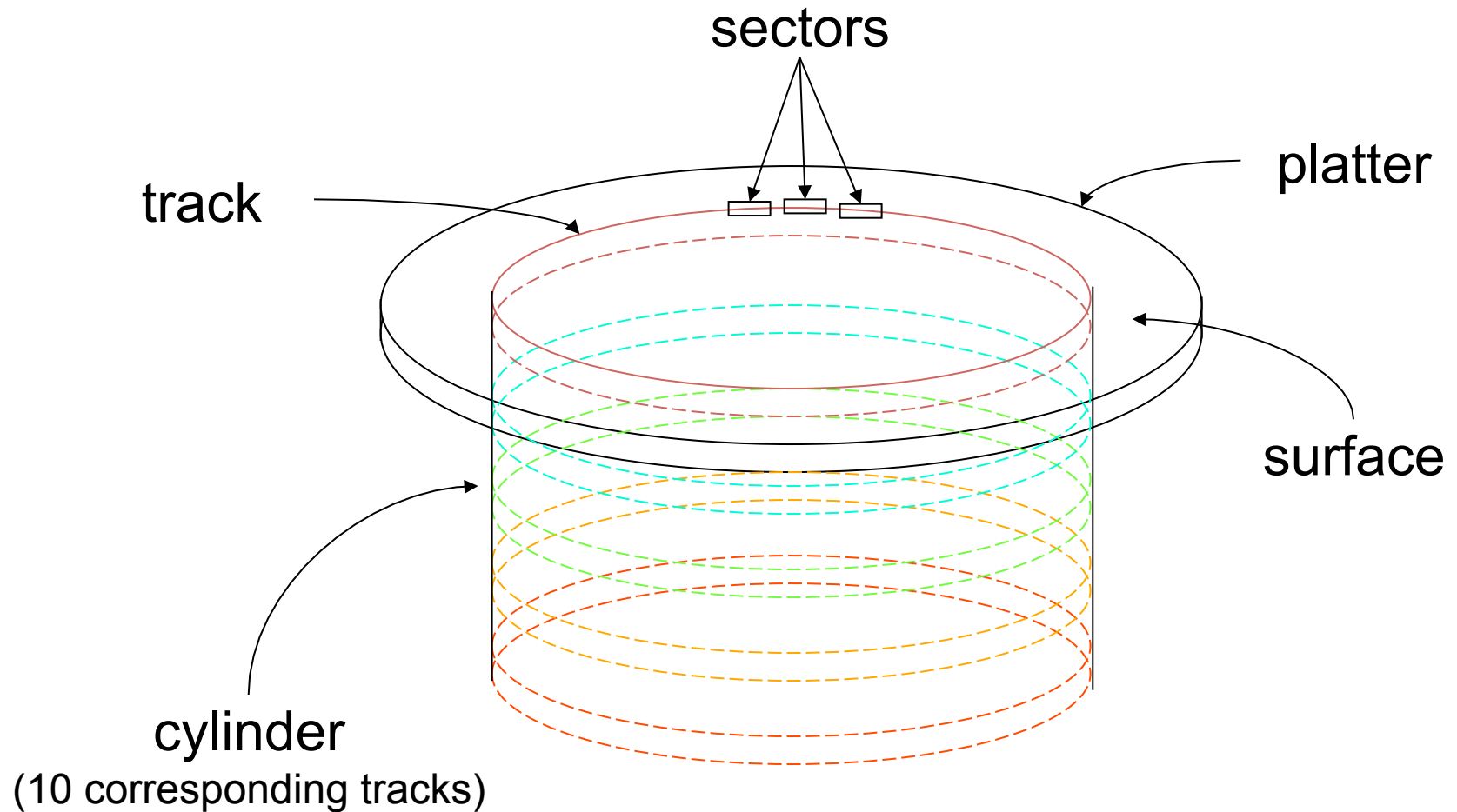
  - Somewhat uncomfortably

# Some Important Disk Characteristics

- Disks are random access devices (mostly . . .)
  - With complex usage, performance, and scheduling

- Key OS services depend on disk I/O
  - Program loading, file I/O, paging
  - Disk performance drives overall performance

- Disk I/O operations are subject to overhead
  - Higher overhead means fewer operations/second
  - Careful scheduling can reduce overhead
  - Clever scheduling can improve throughput, delay

# Disk Drives – A Physical View

Spindle          10 heads

5 platters
10 surfaces

0
1

8
9

head
positioning
assembly

Motor

# Disk Drives – A Logical View

sectors

track

platter

surface

cylinder
(10 corresponding tracks)

# Disk Drive Terms

- *Spindle*
  - A mounted assembly of circular platters
- *Head assembly*
  - Read/write head per surface, all moving in unison
- *Track*
  - Ring of data readable by one head in one position
- *Cylinder*
  - Corresponding tracks on all platters
- *Sector*
  - Logical records written within tracks
- *Disk address* = <cylinder / head / sector >

# Disk Overheads

- ## Seek time
  - Time to move heads from current track to the right track
  - Not constant

- ## Rotational delay
  - Time for the right sector to rotate under the head
  - Not constant

- ## Transfer time
  - Time to read all the bytes of a sector
  - Constant

# Typical Disk Drive Performance

| heads | 10 | platters | 5 |
|---|---|---|---|
| cylinders | 17,000 | tracks/inch | 18,000 |
| sectors/track | 400 | bytes/sector | 512 |
| RPM | 7200 | speed | 196Mb/sec |
| seek time | 0-15 ms | latency | 0-8ms |

## Time to read one 8192 byte block

| | seek | rotate | transfer | total |
|---|---|---|---|---|
| best case | 0ms | 0ms | 333us | 333us |
| worst case | 15ms | 8ms | 333us | 23.3ms (70X) |
| average | 9ms | 4ms | 333us | 13.3ms (40X) |

# Why Is This Problematic For the OS?

- When you go to disk, it could be fast or slow
  - If you go to disk a lot, that matters
- The OS can make choices that make it faster or slower
  - Deciding where to put a piece of data on disk
  - Deciding when to perform an I/O
  - Reordering multiple I/Os to minimize seek time and latency
  - Perhaps optimistically performing I/Os that haven't been requested

# File Systems Control Structures

- A file is a named collection of information
- Primary roles of file system:
  - To store and retrieve data
  - To manage the media/space where data is stored
- Typical operations:
  - Where is the first block of this file?
  - Where is the next block of this file?
  - Where is block 35 of this file?
  - Allocate a new block to the end of this file
  - Free all blocks associated with this file

# Finding Data On Disks

- Essentially a question of how you managed the space on your disk

- Space management on disk is complex
  - There are millions of blocks and thousands of files
  - Files are continuously created and destroyed
  - Files can be extended after they have been written
  - Data placement on disk has performance effects
  - Poor management leads to poor performance

- Must track the space assigned to each file
  - On-disk, master data structure for each file

# On-Disk File Control Structures

- On-disk description of important attributes of a file
  - Particularly where its data is located

- Virtually all file systems have such data structures
  - Different implementations, performance & abilities
  - Implementation can have profound effects on what the file system can do (well or at all)

- A core design element of a file system

- Paired with some kind of in-memory representation of the same information

# The Basic File Control Structure Problem

- A file typically consists of multiple data blocks

- The control structure must be able to find them

- Preferably able to find any of them quickly

  - I.e., shouldn't need to read the entire file to find a block near the end

- Blocks can be changed

- New data can be added to the file

  - Or old data deleted

- Files can be sparsely populated

# The In-Memory Representation

- On file open, create an in-memory structure
- Not an exact copy of the disk version
    - The disk version points to disk sectors
    - The in-memory version points to RAM pages
        - Or indicates that the block isn't in memory
    - Also keeps track of which blocks are dirty and which aren't
- Handles issues of multiple processes sharing an open file simultaneously

# File System Structure

- How do I organize a disk into a file system?
  - Linked extents
    - The DOS FAT file system
  - File index blocks
    - Unix System V file system

# Basics of File System Structure

- Most file systems live on disks
- Disk volumes are divided into fixed-sized blocks
  - Many sizes are used: 512, 1024, 2048, 4096, 8192 ...
- Most blocks will be used to store user data
- Some will be used to store organizing "meta-data"
  - Description of the file system (e.g., layout and state)
  - File control blocks to describe individual files
  - Lists of free blocks (not yet allocated to any file)
- All operating systems have such data structures
  - Different OSes and file systems have very different goals
  - These result in very different implementations

# The Boot Block

- The $0^{th}$ block of a disk is usually reserved for the boot block

  – Code allowing the machine to boot an OS

- Not usually under the control of a file system

  – It typically ignores the boot block entirely

- Not all disks are bootable

  – But the $0^{th}$ block is usually reserved, "just in case"
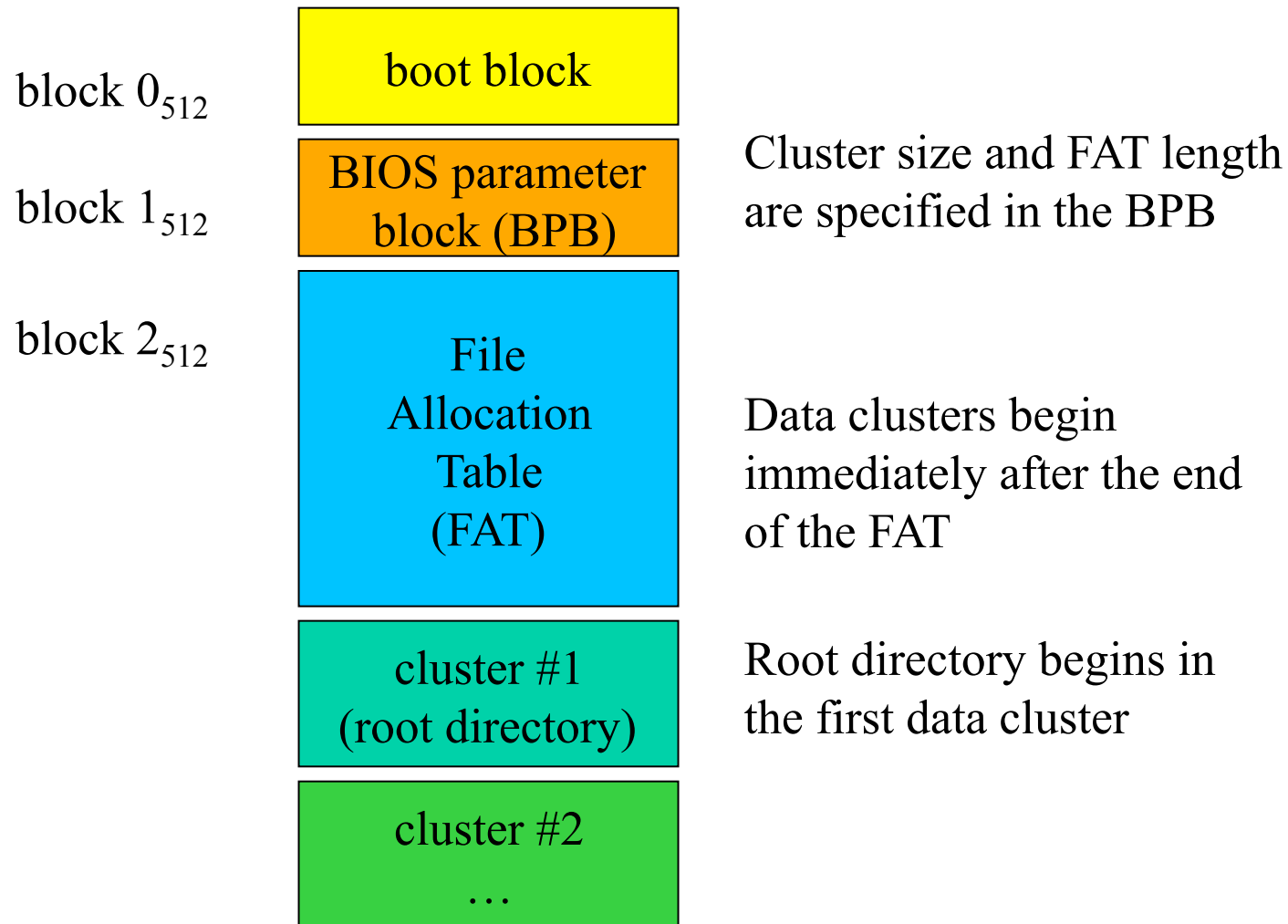
- So file systems start work at block 1

# Managing Allocated Space

- A core activity for a file system, with various choices
- What if we give each file same amount of space?
  - Internal fragmentation ... just like memory
- What if we allocate just as much as file needs?
  - External fragmentation, compaction ... just like memory
- Perhaps we should allocate space in "pages"
  - How many chunks can a file contain?
- The file control data structure determines this
  - It only has room for so many pointers, then file is "full"
- So how do we want to organize the space in a file?

# Linked Extents

- A simple answer
- File control block contains exactly one pointer
  - To the first chunk of the file
  - Each chunk contains a pointer to the next chunk
  - Allows us to add arbitrarily many chunks to each file
- Pointers can be in the chunks themselves
  - This takes away a little of every chunk
  - To find chunk N, you have to read the first N-1 chunks
- Pointers can be in auxiliary "chunk linkage" table
  - Faster searches, especially if table kept in memory

# The DOS File System

block $0_{512}$

block $1_{512}$

block $2_{512}$

boot block

BIOS parameter
block (BPB)

File
Allocation
Table
(FAT)

cluster #1
(root directory)

cluster #2
...

Cluster size and FAT length
are specified in the BPB

Data clusters begin
immediately after the end
of the FAT

Root directory begins in
the first data cluster

# DOS File System Overview

- DOS file systems divide space into "clusters"
  - Cluster size (multiple of 512) fixed for each file system
  - Clusters are numbered 1 though N

- File control structure points to first cluster of a file

- File Allocation Table (FAT), one entry per cluster
  - Contains the number of the next cluster in file
  - A 0 entry means that the cluster is not allocated
  - A -1 entry means "end of file"

- File system is sometimes called "FAT," after the name of this key data structure

# DOS FAT Clusters

**directory entry**

| name: | myfile.txt |
|---|---|
| length: | 1500 bytes |
| 1st cluster: | 3 |

**File Allocation Table**

| | |
|---|---|
| 1 | x |
| 2 | x |
| 3 | 4 |
| 4 | 5 |
| 5 | -1 |
| 6 | 0 |

Each FAT entry corresponds to a cluster, and contains the number of the next cluster.

-1 = End of File
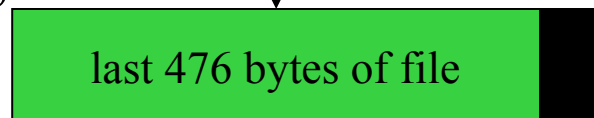
0 = free cluster

cluster #3

first 512 bytes of file

cluster #4

second 512 bytes of file

cluster #5

last 476 bytes of file

# DOS File System Characteristics

- To find a particular block of a file
  - Get number of first cluster from directory entry
  - Follow chain of pointers through File Allocation Table

- Entire File Allocation Table is kept in memory
  - No disk I/O is required to find a cluster
  - For very large files the search can still be long

- No support for "sparse" files
  - Of a file has a block $n$, it must have all blocks $< n$

- Width of FAT determines max file system size
  - How many bits describe a cluster address
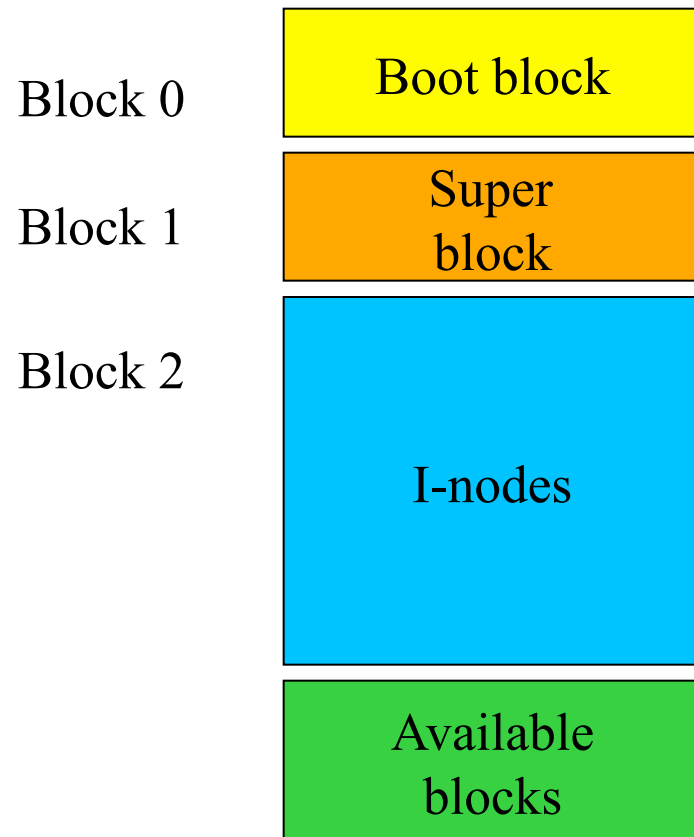  - Originally 8 bits, eventually expanded to 32

# File Index Blocks

- A different way to keep track of where a file's data blocks are on the disk

- A file control block points to all blocks in file
  - Very fast access to any desired block
  - But how many pointers can the file control block hold?

- File control block could point at extent descriptors
  - But this still gives us a fixed number of extents

# Hierarchically Structured File Index Blocks

- To solve the problem of file size being limited by entries in file index block

- The basic file index block points to blocks

- Some of those contain pointers which in turn point to blocks

- Can point to many extents, but still a limit to how many

  - But that limit might be a very large number

  - Has potential to adapt to wide range of file sizes
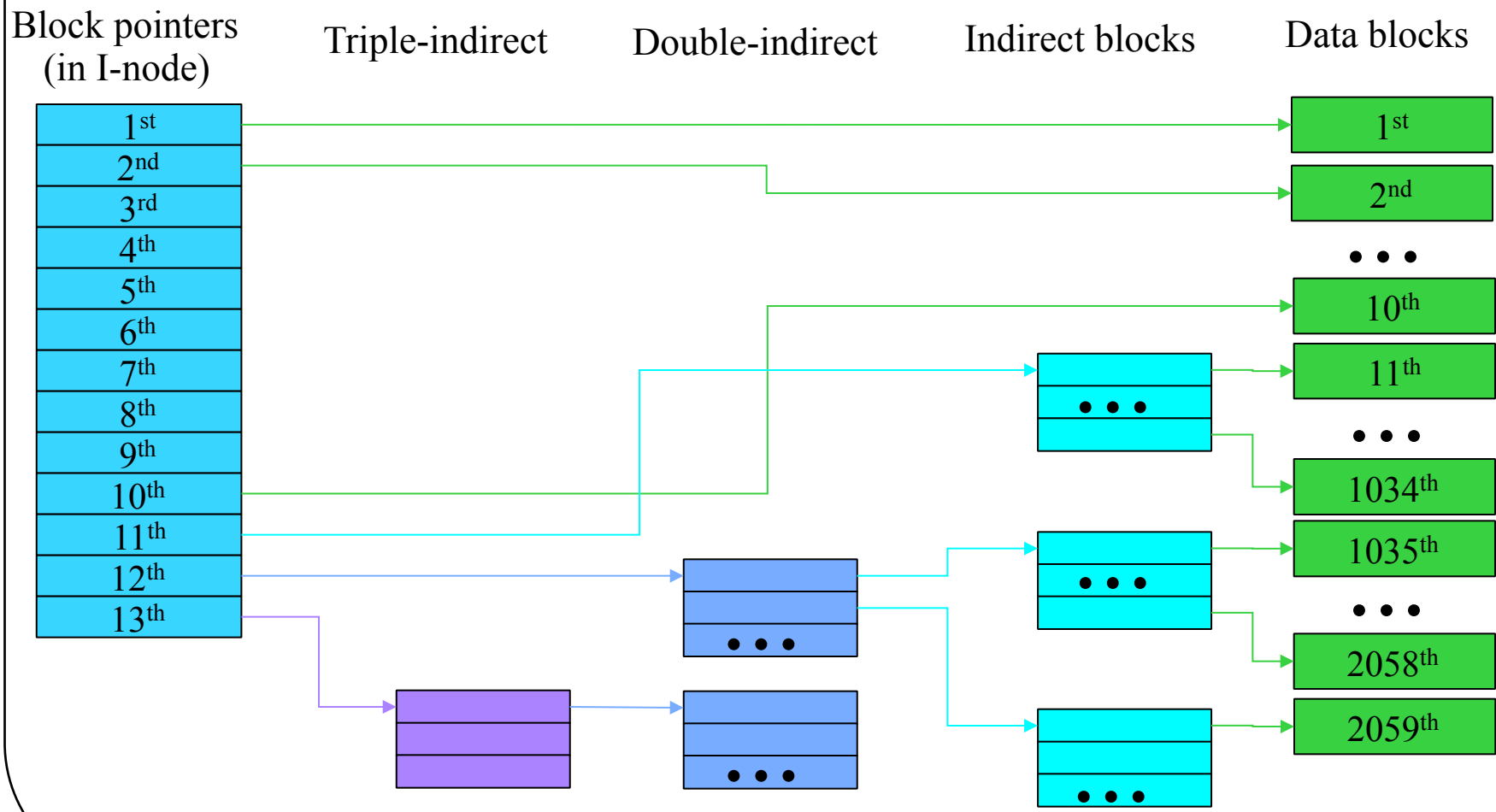
# Unix System V File System

Block 0

**Boot block**

Block 1

**Super block**

Block size and number of I-nodes are specified in super block

Block 2

**I-nodes**

I-node #1 (traditionally) describes the root directory

**Available blocks**

Data blocks begin immediately after the end of the I-nodes.

# Unix Inodes and Block Pointers

Block pointers
(in I-node)

Triple-indirect

Double-indirect

Indirect blocks

Data blocks

# Why Is This a Good Idea?

- The UNIX pointer structure seems ad hoc and complicated

- Why not something simpler?
  - E.g., all block pointers are triple indirect

- File sizes are not random
  - The majority of files are only a few thousand bytes long

- Unix approach allows us to access up to 40Kbytes (assuming 4K blocks) without extra I/Os
  - Remember, the double and triple indirect blocks must themselves be fetched off disk
  - Also remember, it's invisible to users

# How Big a File Can Unix Handle?

- The on-disk inode contains 13 block pointers
  - First 10 point to first 10 blocks of file
  - 11th points to an indirect block (which contains pointers to 1024 blocks)
  - 12th points to a double indirect block (pointing to 1024 indirect blocks)
  - 13th points to a triple indirect block (pointing to 1024 double indirect blocks)

- Assuming 4k bytes per block and 4-bytes per pointer
  - 10 direct blocks = 10 * 4K bytes = 40K bytes
  - Indirect block = 1K * 4K = 4M bytes
  - Double indirect = 1K * 4M = 4G bytes
  - Triple indirect = 1K * 4G = 4T bytes
  - At the time system was designed, that seemed impossibly large
  - But . . .

# Unix Inode Performance Issues

- The inode is in memory whenever file is open
- So the first ten blocks can be found with no extra I/O
- After that, we must read indirect blocks
  - The real pointers are in the indirect blocks
  - Sequential file processing will keep referencing it
  - Block I/O will keep it in the buffer cache
- 1-3 extra I/O operations per thousand pages
  - Any block can be found with 3 or fewer reads
- Index blocks can support "sparse" files
  - Not unlike page tables for sparse address spaces