

CS111, Lecture 2

Introduction to Filesystems

Optional reading:

Operating Systems: Principles and Practice (2nd Edition): Chapter 11,
Section 12.1, 12.2 and Section 13.3 (up through page 567)

***While you're waiting – get set up with PollEverywhere! Visit pollev.stanford.edu to set up your account.
When responding, make sure to allow location permissions, as location is required to confirm credit.***

PollEverywhere

- Today we're doing a "trial run" of using PollEverywhere for poll questions
 - Not counted for attendance (that starts next lecture), just a chance to try it out
 - Canvas Gradebook records lecture scores
 - Responses not anonymized, but we look only at aggregated results and totals
 - Poll responses require being present in the lecture room to receive credit
 - Alternative option instead of lecture credit is to shift the full 5% weight to the final exam
- Visit pollev.stanford.edu to log in (or use the PollEverywhere app) and sign in with your **@stanford.edu email – NOT your personal email!**
- Compatible with any device with a web browser, mobile app also available.
- Poll questions in slides will automatically activate the poll and respond at pollev.com/cs111.
- Polls will prompt for location check-in; check in **just prior** to responding.

Announcements

- Assign0 released - see course website for more information
 - No late submissions accepted (except for OAE/Head TA accommodations – during the quarter, extension requests must be received in advance of the assignment on-time deadline, or as soon as possible if extenuating circumstances occur later, or extenuating circumstances prevent reaching out prior to the deadline.)
- Remember to input your section preferences by 11:59PM Thurs! Link is on the course website (under “Sections”).
- Helper Hours now scheduled, starting this week! Tea hours also starting!
- Please let us know about OAE accommodations or midterm conflicts as soon as possible – see our midterm conflict form on the Midterm page of the course website, and our OAE letter submission form at the bottom of the course homepage.

Topic 1: Filesystems - How can we design filesystems to manage files on disk, and what are the tradeoffs inherent in designing them? How can we interact with the filesystem in our programs?

CS111 Topic 1: Filesystems

Filesystems - *How can we design filesystems to manage files on disk, and what are the tradeoffs inherent in designing them? How can we interact with the filesystem in our programs?*

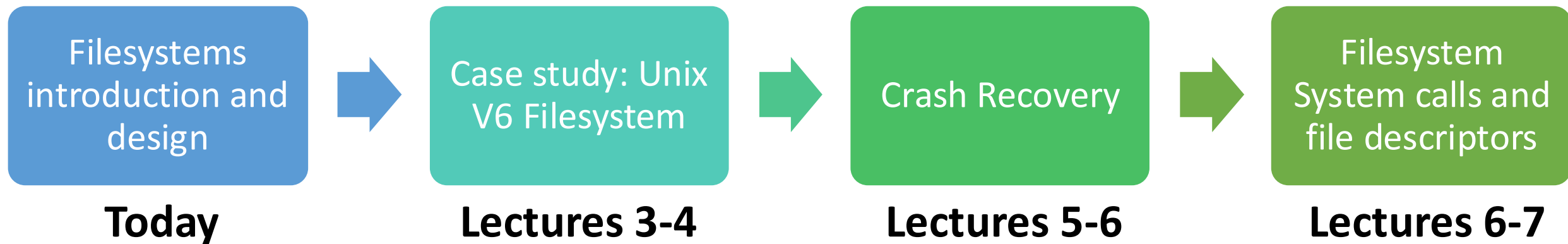
Why is answering this question important?

- Helps us understand what filesystems do (today and next time)
- Provides insight into the challenges and tradeoffs in designing large systems (next few lectures)
- Shows us how we can directly manipulate files in our programs (next week)

assign1: implement layers of the Unix v6 filesystem to read a file from disk given its path.

CS111 Topic 1: Filesystems

Filesystems - *How can we design filesystems to manage files on disk, and what are the tradeoffs inherent in designing them? How can we interact with the filesystem in our programs?*



assign1: implement portions of the Unix v6 filesystem!

Learning Goals

- Understand the key responsibilities and requirements of a filesystem
- Get practice identifying tradeoffs in different filesystem designs
- Explore the design of the Unix V6 filesystem

Plan For Today

- Filesystems Introduction
- Methods for Storing Files
 - Contiguous Allocation
 - Linked Files
 - Windows FAT
 - Multi-level indexes
- The Unix V6 Filesystem
 - Inodes

Plan For Today

- **Filesystems Introduction**
- Methods for Storing Files
 - Contiguous Allocation
 - Linked Files
 - Windows FAT
 - Multi-level indexes
- The Unix V6 Filesystem
 - Inodes

Filesystems

A **filesystem** is the portion of the OS that manages the disk.

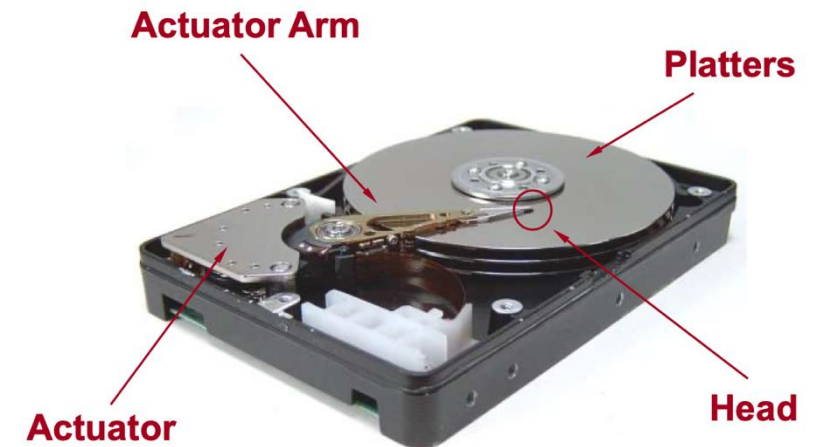
- A hard drive (or, more commonly these days, flash storage) is persistent storage – it can store data between power-offs.

Memory (RAM)	Disk
<ul style="list-style-type: none">• Fast, less space (e.g. 8-32GB), more expensive to buy per GB• <i>Byte-addressable</i>: can quickly access any byte of data by address, but not individual bits by address• Not persistent: cannot store data between power-offs	<ul style="list-style-type: none">• Slower, more space (e.g. 256GB – 1024GB), cheaper to buy per GB• <i>Sector-addressable</i>: cannot read/write just one byte of data – can only read/write “sectors” of data at a time• Persistent: stores data between power-offs

Hard Drives

Magnetic disks (hard drives) have been the standard storage mechanism for files.

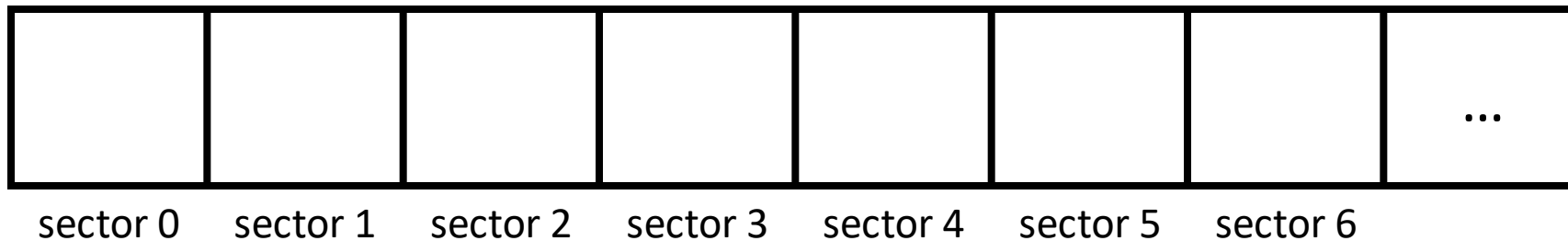
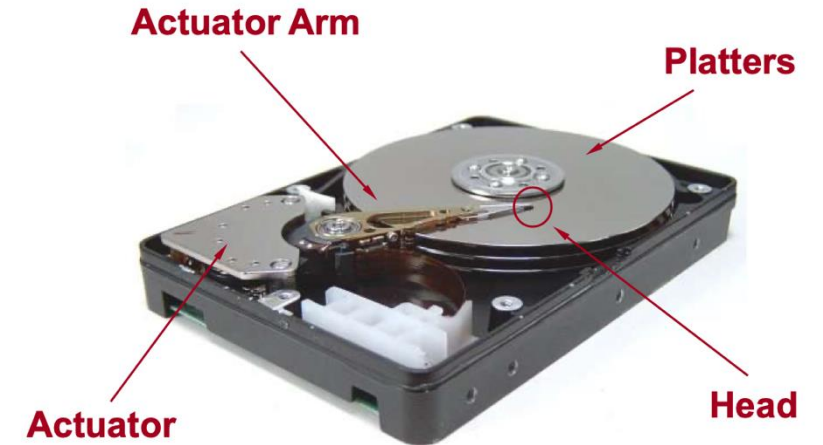
- Spinning, magnetically-coated platters
- Actuator arm positions *heads*, which can read and write data on the magnetic surfaces
- Moving parts means risk of damage from sudden movement, dust, etc.



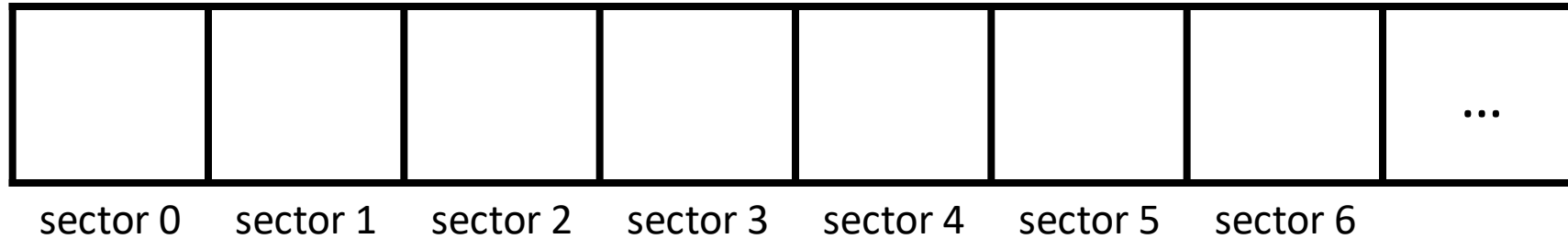
Hard Drives

Hard drives have peculiar performance characteristics that have a big impact on how we build filesystems.

- Reading and writing requires *seeking* (moving arm to position heads over desired track) and waiting for desired location to pass underneath. Want to minimize this time.
- We can only read data in chunks of *sectors*. Example of *virtualization*; making one thing look like another.



Hard Disks are Sector-Addressable



If we are the OS, the hard disk creators might provide this API (“application programming interface”) – a set of public functions - to interface with the disk:

```
void readSector(size_t sectorNumber, void *data);  
void writeSector(size_t sectorNumber, const void *data);
```

This is all we get! We (the OS) must build a filesystem by layering functions on top of these to ultimately allow us to read, write, lookup, and modify entire files.₁₄

Filesystem Functionality

We want to read/write file on disk and have them persist even when the device is off. This may include operations like:

- creating a new file on disk
- looking up the location of a file on disk
- Reading/editing all or part of an existing file from disk – e.g., sequential/random access
- creating folders on disk
- getting the contents of folders on disk
- ...

Filesystems

Functions for user programs to read/write files

Filesystem

readSector and writeSector

Filesystem Challenges

Problems addressed by modern file systems:

- **Disk space management:**
 - Fast access to files (minimize seeks)
 - Sharing space between users
 - Efficient use of disk space
- **Naming:** how do users select files?
- **Reliability:** information must survive OS crashes and hardware failures.
- **Protection:** isolation between users, controlled sharing.

Flash Storage

Recently, flash storage (“SSD”) has become more popular and commonplace, especially with the growth in mobile devices.

- Much faster (100x faster access), but more expensive
- No moving parts, so more reliable
- Issues with *wear-out*; once a chunk of the drive has been erased many times (~100k), it no longer stores info reliably.
- Typically, still use an interface of reading/writing in units of sectors.



<https://www.samsung.com/us/computing/memory-storage/solid-state-drives/980-pro-pcie-4-0-nvme-ssd-1tb-mz-v8p1t0b-am/>

Plan For Today

- Filesystems Introduction
- **Methods for Storing Files**
 - Contiguous Allocation
 - Linked Files
 - Windows FAT
 - Multi-level indexes
- The Unix V6 Filesystem
 - Inodes

Sectors and Blocks

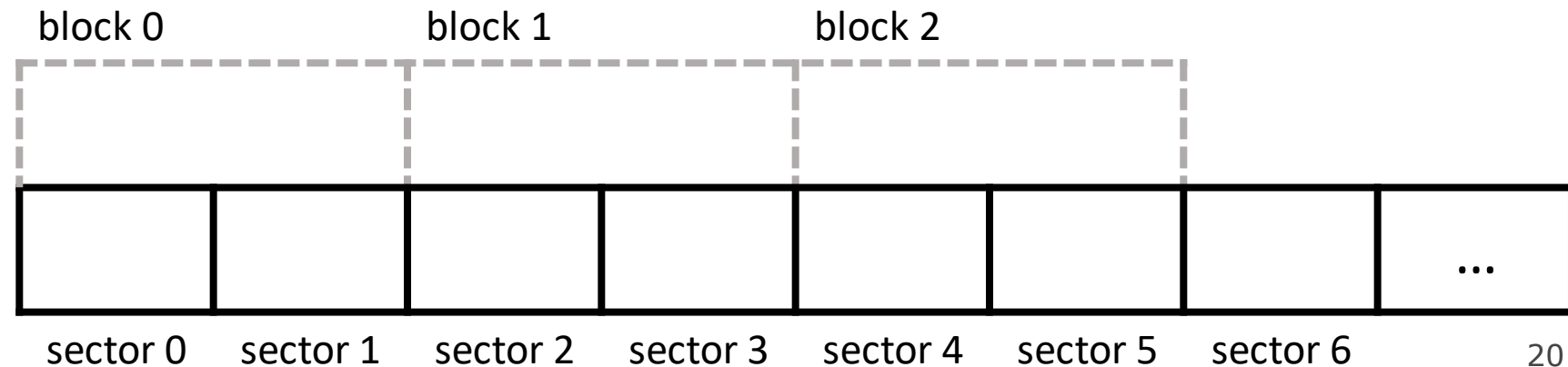
A filesystem generally defines its own unit of data, a "block," that it reads/writes at a time.

- "Sector" = hard disk storage unit
- "Block" = filesystem storage unit (1 or more sectors) - software abstraction

Pros of larger block size? Smaller block size?

- more efficient I/O if larger, but less internal fragmentation if smaller

Example: the block size could be defined as two sectors



Storing Files on Disk

Two types of data we will be working with:

1. file payload data - contents of files (e.g. text in documents, pixels in images)
2. file metadata - information about files (e.g. name, size)

Key insight: *both* must be stored on the hard disk. Otherwise, we will not have it across power-offs! (E.g. without storing metadata we would lose all filenames after shutdown). *This means some blocks must store data other than payload data.*

Storing Files on Disk

Two types of data we will be working with:

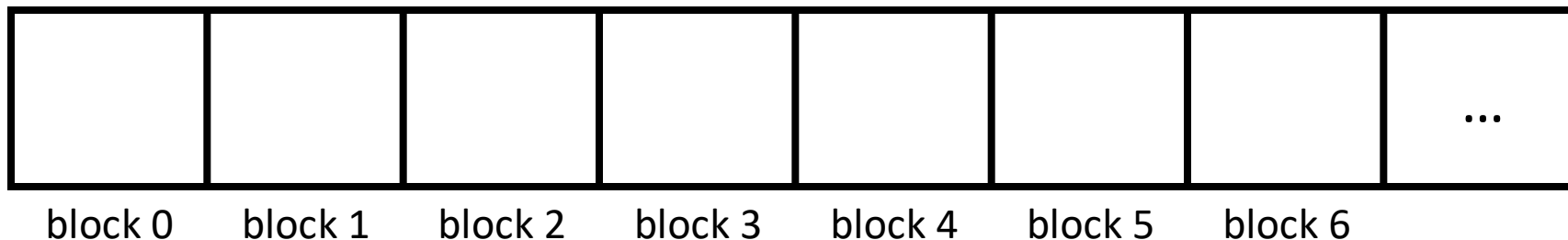
1. file payload data - contents of files (e.g. text in documents, pixels in images)
2. file metadata - information about files (e.g. name, size)

Key insight: *both* must be stored on the hard disk. Otherwise, we will not have it across power-offs! (E.g. without storing metadata we would lose all filenames after shutdown). *This means some blocks must store data other than payload data.*

Contiguous Allocation

First key question: should we store files contiguously on disk? What would it look like if we did?

- Called *contiguous allocation* – allocate a file in one contiguous group of blocks
- For each file, keep track of the number of its first block and its length
- Keep a free list of unused areas of the disk
- *For our purposes, we won't allocate partial blocks (only give out full blocks)*
- Example: IBM OS/360
- Advantages?



Contiguous Allocation

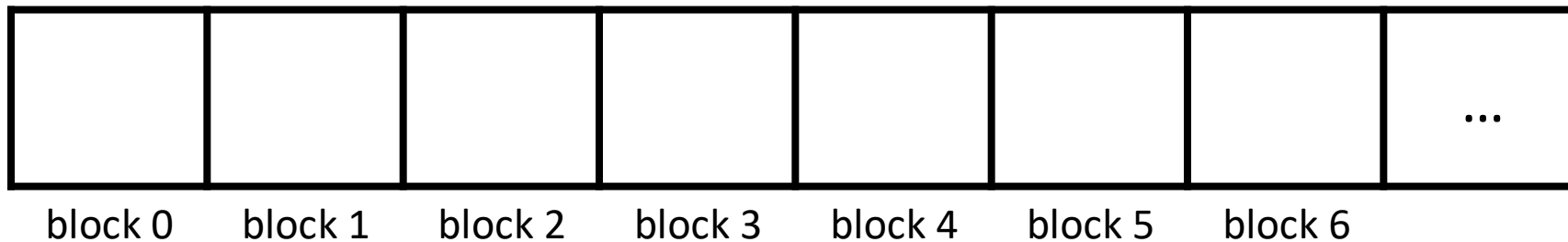
First key question: should we store files contiguously on disk? What would it look like if we did?

- Called *contiguous allocation* – allocate a file in one contiguous group of blocks

Advantages:

- simple
- can read sequentially or easily jump to any location in file (“random access”)
- all data in one place (few seeks)

What about disadvantages?



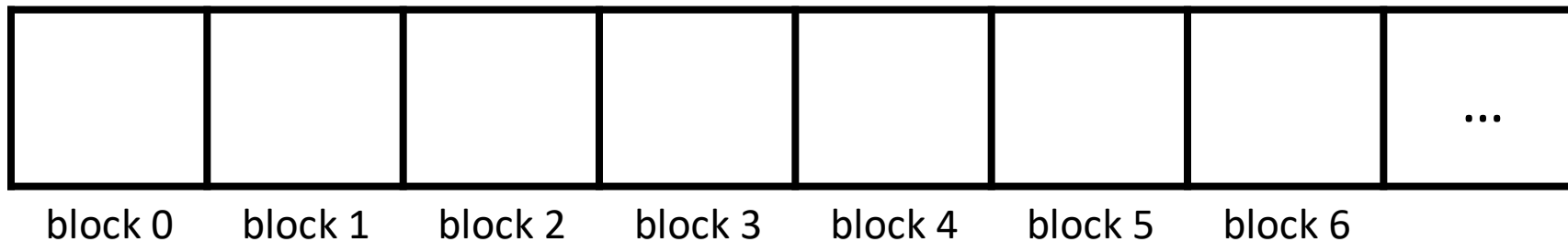
Contiguous Allocation

First key question: should we store files contiguously on disk? What would it look like if we did?

- Called *contiguous allocation* – allocate a file in one contiguous group of blocks

Disadvantages:

- hard to grow files
- hard to lay out files on disk – we may not be able to squeeze a new file in a block of free space (*external fragmentation* – no single space is large enough, but enough aggregate space is available)
- Also *internal fragmentation* (space allocated is larger than needed) if we allocate in units of blocks



Fragmentation

- **Internal Fragmentation:** space allocated for a file is larger than what is needed. A file may not take up all the space in the blocks it's using. E.g. block = 512 bytes, but file is only 300 bytes. (you could share blocks between multiple files, but this gets complex)
- **External Fragmentation (issue with contiguous allocation):** no single space is large enough to satisfy an allocation request, even though enough aggregate free disk space is available

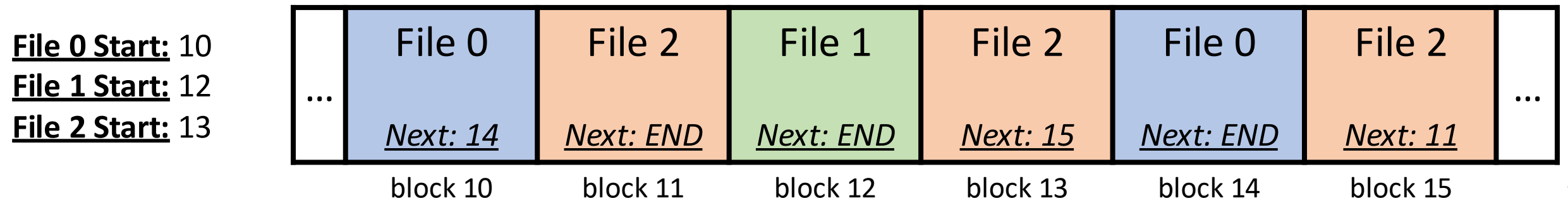
Linked Files

First key question: should we store files contiguously on disk? What would it look like if we *didn't*?

- Problem: we need to know what blocks are associated with what files

One idea: *linked files* – like a linked list

- Each block contains file data *as well as* the location of the next block
- For each file, keep track of the number of its first block in separate location
- Approximate examples: TOPS-10, Xerox Alto
- Advantages?



Linked Files

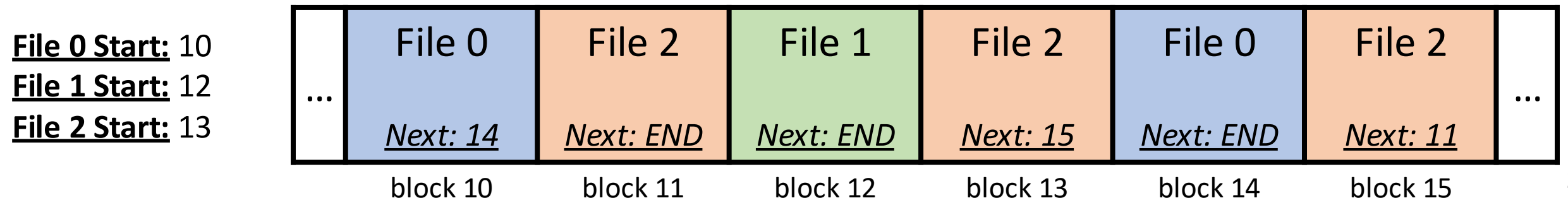
First key question: should we store files contiguously on disk? What would it look like if we *didn't*? One idea: *linked files* – like a linked list

- Each block contains file data *as well as* the location of the next block

Advantages:

- Easy to grow files
- Easier to fit files in available space – less *fragmentation*
- Still supports simple sequential access

What about disadvantages?



Linked Files

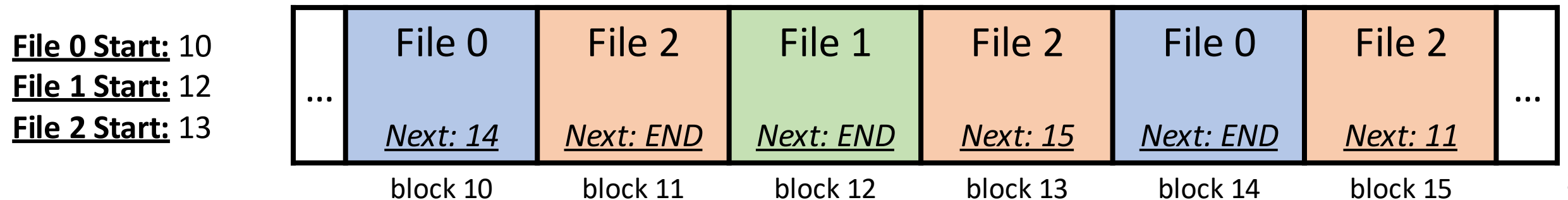
First key question: should we store files contiguously on disk? What would it look like if we *didn't*?

One idea: *linked files* – like a linked list

- Each block contains file data *as well as* the location of the next block

Disadvantages:

- Can't easily jump to any arbitrary location in the file
- Data scattered throughout disk (more seeks vs. contiguous allocation)



Linked Files

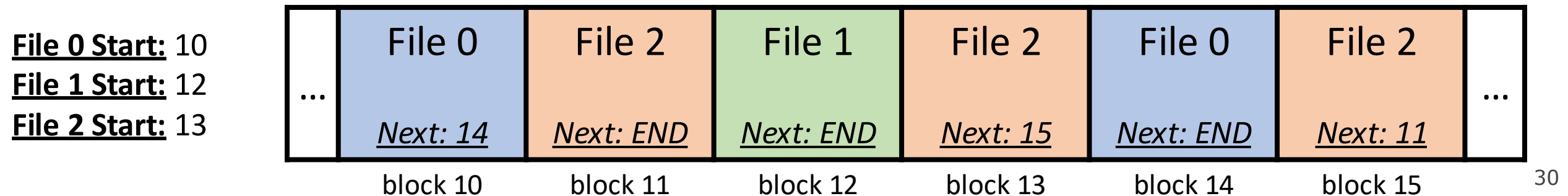
First key question: should we store files contiguously on disk? What would it look like if we *didn't*?

One idea: *linked files* – like a linked list

- Each block contains file data *as well as* the location of the next block

Disadvantages:

- **Can't easily jump to any arbitrary location in the file**
- Data scattered throughout disk (more seeks)



Windows FAT

First key question: should we store files contiguously on disk? What would it look like if we *didn't*?

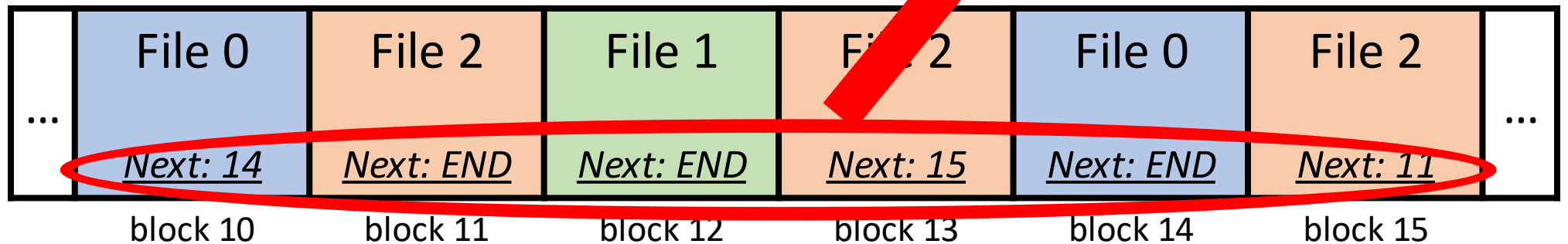
Interesting idea: what if we instead stored the links in one big table in memory?

- *Windows (DOS) FAT*: like linked allocation, except links aren't in blocks, they are in a "file allocation table" in memory and disk (originally 16 bits per entry)

In-Memory File Allocation Table

	...
10	14
11	END
12	END
13	15
14	END
15	11
	...

Disk



Windows FAT

First key question: should we store files contiguously on disk? What would it look like if we *didn't*?

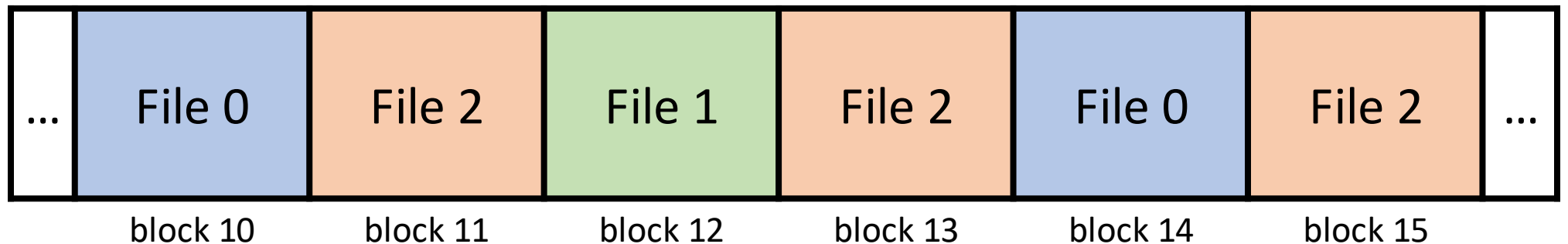
Interesting idea: what if we instead stored the links in one big table in memory?

- *Windows (DOS) FAT*: like linked allocation, except links aren't in blocks, they are in a "file allocation table" in memory and disk (originally 16 bits per entry)

In-Memory File Allocation Table

	...
10	14
11	END
12	END
13	15
14	END
15	11
	...

Disk



File 0 Start: 10

File 1 Start: 12

File 2 Start: 13

Windows FAT

First key question: should we store files contiguously on disk? What would it look like if we *didn't*?

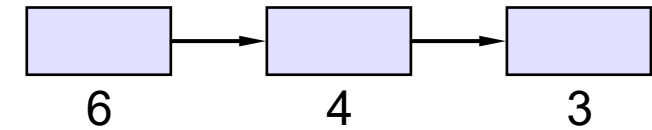
Interesting idea: what if we instead stored the links in one big table in memory?

- *Windows (DOS) FAT*: like linked allocation, except links aren't in blocks, they are in a "file allocation table" in memory and disk (originally 16 bits per entry)
- Still keep track of each file's first block
- (Still used today for flash sticks, digital cameras, many embedded devices)
- Advantages?

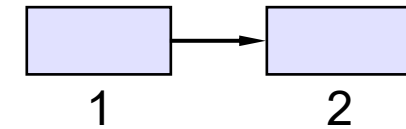
File Allocation Table

0	free
1	2
2	end
3	end
4	3
5	end
6	4
7	free
	...

File A:



File B:



Windows FAT

First key question: should we store files contiguously on disk? What would it look like if we *didn't*?

- *Windows (DOS) FAT*: like linked allocation, except links aren't in blocks, they are in a "file allocation table" in memory

Advantages:

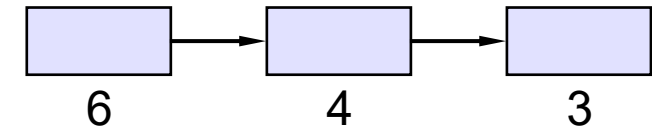
- Can more quickly jump to various locations in a file
- Still supports easy sequential access

What about disadvantages?

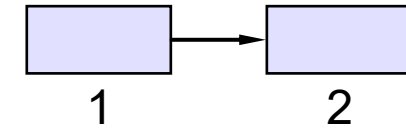
File Allocation Table

0	free
1	2
2	end
3	end
4	3
5	end
6	4
7	free
	...

File A:



File B:



Windows FAT

First key question: should we store files contiguously on disk? What would it look like if we *didn't*?

- *Windows (DOS) FAT*: like linked allocation, except links aren't in blocks, they are in a “file allocation table” in memory (and also stored on disk)

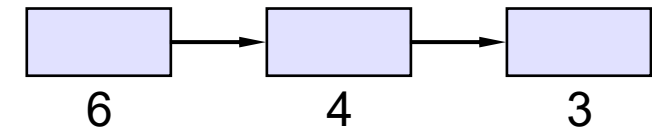
Disadvantages:

- Data scattered throughout disk (more seeks vs. contiguous allocation)
- Still need to jump through table to get to an arbitrary location in the file
- Must store table in memory

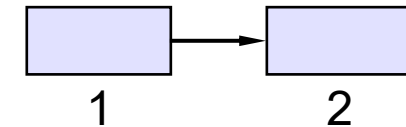
File Allocation Table

0	free
1	2
2	end
3	end
4	3
5	end
6	4
7	free
	...

File A:



File B:



Windows FAT

First key question: should we store files contiguously on disk? What would it look like if we *didn't*?

- *Windows (DOS) FAT*: like linked allocation, except links aren't in blocks, they are in a “file allocation table” in memory (and also stored on disk)

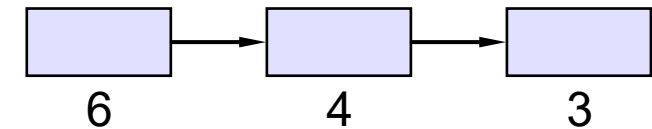
Disadvantages:

- Data scattered throughout disk (more seeks vs. contiguous allocation)
- **Still need to jump through table to get to an arbitrary location in the file**
- Must store table in memory

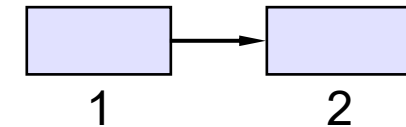
File Allocation Table

0	free
1	2
2	end
3	end
4	3
5	end
6	4
7	free
	...

File A:



File B:



File Payload Data

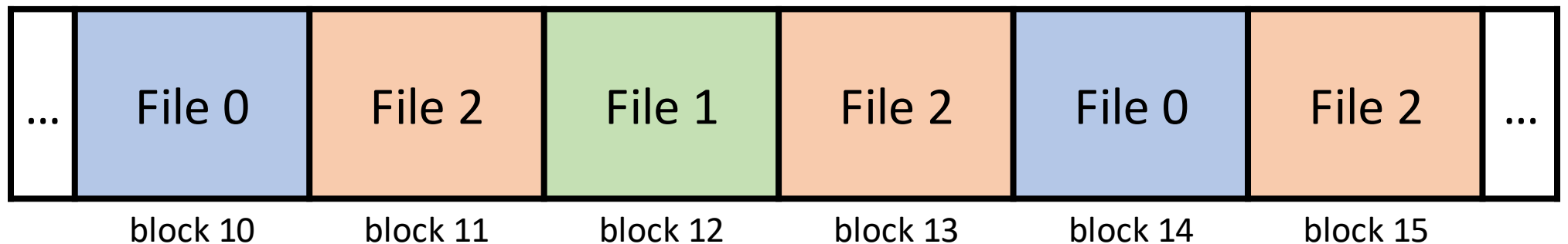
First key question: should we store files contiguously on disk? What would it look like if we *didn't*?

Interesting idea: what if did not have a file allocation table or links, and instead we stored all the block numbers for a file in order? That way we could quickly jump to any point in the file.

File 0: 10, 14

File 1: 12

File 2: 13, 15, 11



File Payload Data

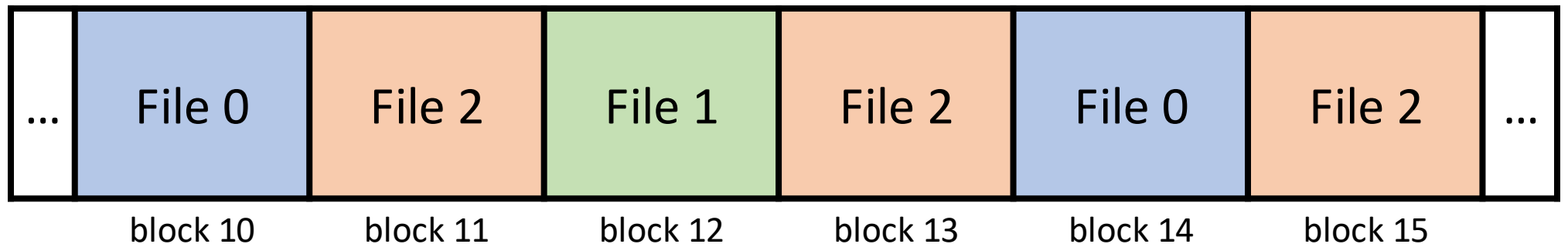
First key question: should we store files contiguously on disk? What would it look like if we *didn't*?

- *Multi-level indexes*: store all block numbers for a given file (but how?)
 - Example: 4.3 BSD Unix, Unix V6 Filesystem (~1975)
 - More modern ext2 and ext3 Linux file systems based on this idea; Windows NTFS also uses a tree-based structure, though slightly different

File 0: 10, 14

File 1: 12

File 2: 13, 15, 11



Plan For Today

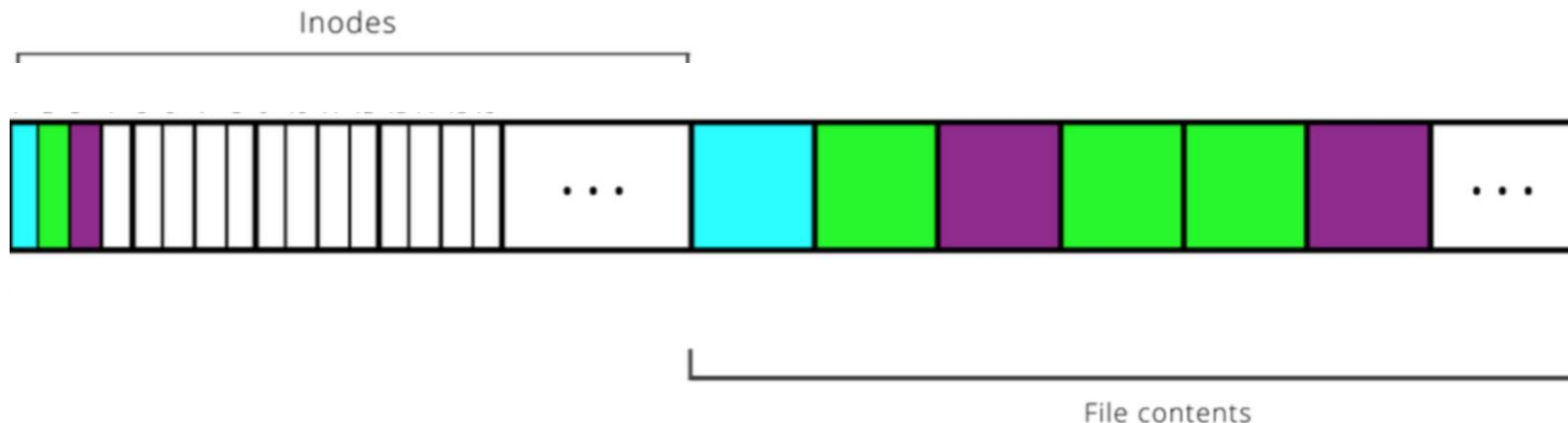
- Filesystems Introduction
- Methods for Storing Files
 - Contiguous Allocation
 - Linked Files
 - Windows FAT
 - Multi-level indexes
- **The Unix V6 Filesystem**
 - Inodes

Unix V6 Filesystem

Key Idea: files don't need to be stored contiguously on disk, but we want to store all the block numbers in order that make up the data for a file.

Where could we store this information for each file for easy lookup?

Let's reserve some space on disk to store this information for each file, separately from its payload data. This per-file space is called an **inode**.



Inodes

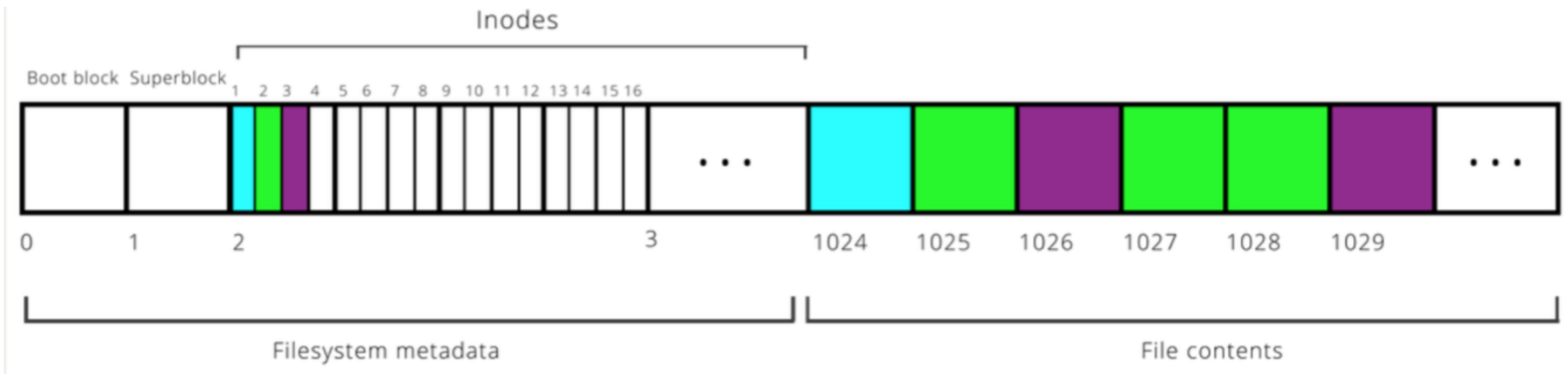
An **inode** ("index node") is a grouping of data about a single file, stored on disk.

- Various filesystems (e.g. Unix v6, contiguous allocation/linked files, but not FAT) store file metadata in inodes
- In Unix v6, an inode contains an ordered list of block numbers that store the file's payload data, and also stores other metadata like file size.
- Unix v6 stores inodes contiguously in a reserved space called the **inode table**, which starts at block **2**. (block 0 is "boot block" containing hard drive info, block 1 is "superblock" containing filesystem info).
- Inodes can be read into memory when used for quicker access

Unix V6 Inodes

The Unix v6 filesystem stores inodes on disk together in the **inode table** for quick access.

- Unix v6 Inodes are 32 bytes big, and 1 block = 1 sector = 512 bytes, so 16 inodes/block.
- Typically, at most 10% of the drive stores metadata.
- Filesystem goes from **filename** to **inode number** ("inumber") to **file data**.



Unix V6 Inodes

We need inodes to be a fixed size, and not too large. So how should we store the block numbers? How many should there be?

1. if variable number, there's no fixed inode size
2. if fixed number, this limits maximum file size

The inode design here has space for 8 block numbers, which are stored in order. (i.e. first block number stores first chunk of file, etc.). But we will see later how we can build on this to support very large files.

Recap

- Filesystems Introduction
- Methods for Storing Files
 - Contiguous Allocation
 - Linked Files
 - Windows FAT
 - Multi-level indexes
- The Unix V6 Filesystem
 - Inodes

Next time: more about the Unix v6 Filesystem

Lecture 2 takeaway:

Filesystems need to store both file metadata and payload data. There are various ways to store payload data, each with different pros/cons. The Unix V6 filesystem uses inodes to store file data, including block numbers.