# Recap

File System Interface

- File & Directory

- POSIX API: open/close, read/write, random access, hard/soft link

# File System Implementation

Based on Ch. 40

# File System Implementation

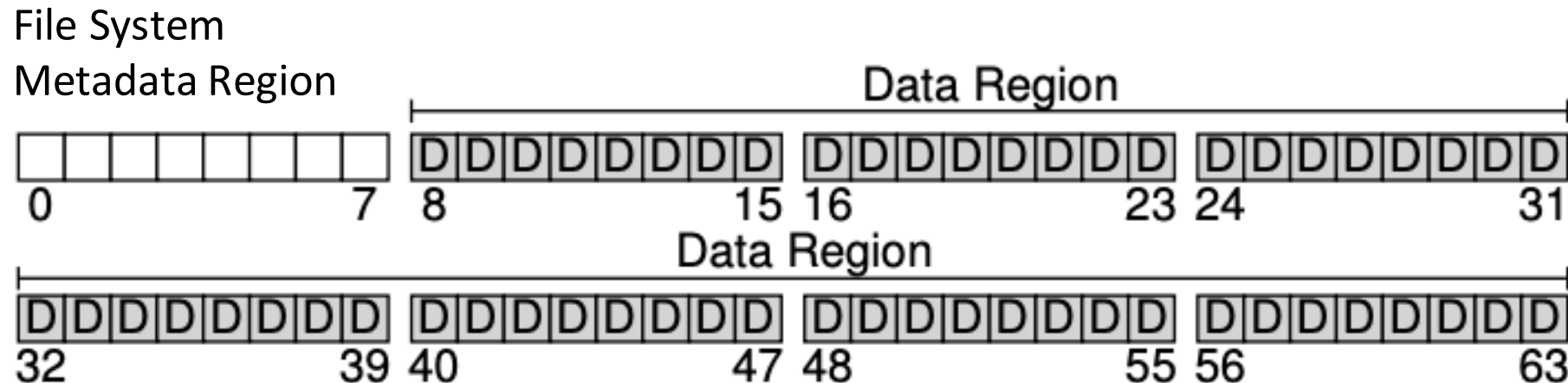Many possible approaches to making a file system

Need to understand the data structures for files, directories and free space and how access (e.g., read and write) is implemented

How to implement a simple file system?

# A Very Simple File System (VSFS)

Disk is divided into fixed sized blocks (e.g., block size = 4KB)

A few blocks are reserved for file system metadata (information about the files), the rest stores the data of the files

# Contiguous vs Block Allocation

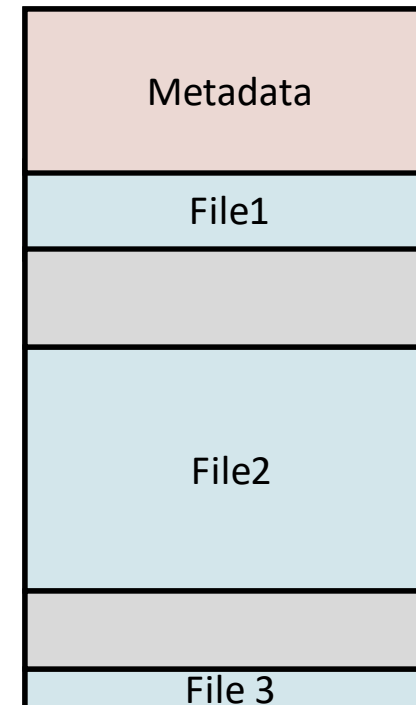With memory we explored base and bounds (contiguous allocation) or paging

File system have a similar choice (blocks are like pages)

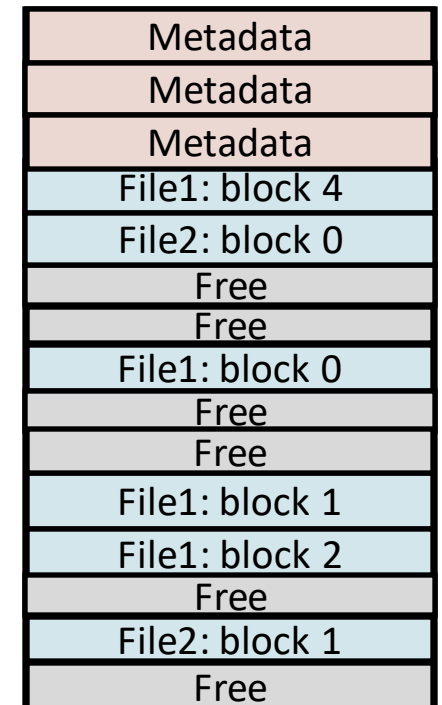**Contiguous allocation** has problem of **external fragmentation**

**Block allocation** has problem of **internal fragmentation**

Block allocation usually wins because internal fragmentation is less of a problem and block allocation is much more efficient in allocating and resizing files

Contiguous

| |
|---|
| Metadata |
| File1 |
| |
| File2 |
| |
| File 3 |

Blocks

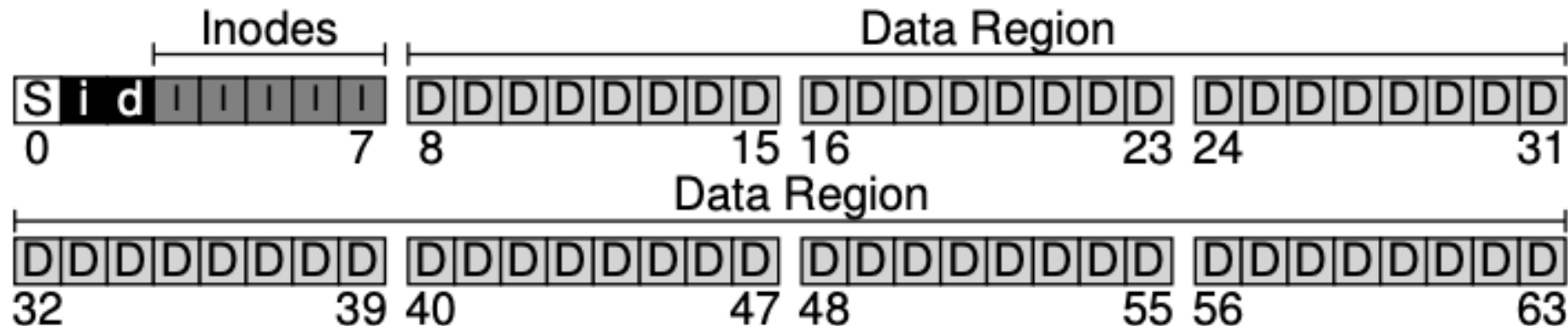| |
|---|
| Metadata |
| Metadata |
| Metadata |
| File1: block 4 |
| File2: block 0 |
| Free |
| Free |
| File1: block 0 |
| Free |
| Free |
| File1: block 1 |
| File1: block 2 |
| Free |
| File2: block 1 |
| Free |

# File System Metadata

**inodes** (I) contains information about a particular file

**data bitmap** (d) stores which blocks are free/used in the data region

**inode bitmap** (i) stores which blocks are free/used in the metadata region

**Super block** (S) contains information about the file system

# inode

inode (index node) contains metadata (information) about a file

Example: Ext2 file system  inode

| 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 |

← points to other inode or data blocks

# Direct Indexing

In **direct indexing** the inode for a file has pointers to the data blocks of the file

Suppose an inode has 15 block pointers and blocks are 4KB, what is the largest possible file size?

15 * 4KB = 60KB maximum file size

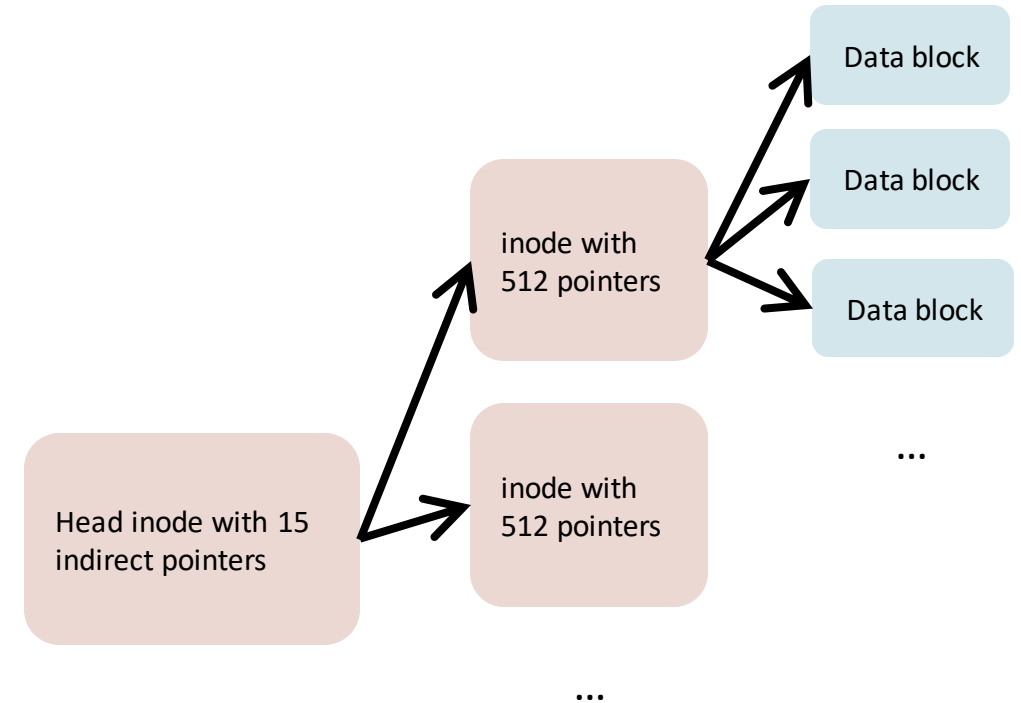# Indirect indexing

An **indirect pointer** points to an inode with more pointers

Suppose pointers are 8 bytes and block size is 4KB, then a block can hold 512 pointers

With one level of indirection the previous example has:

15 * 512 * 4KB = about 30MB maximum file size

Data block

Data block

Data block

inode with 512 pointers

…

Head inode with 15 indirect pointers

inode with 512 pointers

…

# Multiple Levels of Indirection

**Double indirect** means two levels of indirection

# Multiple Levels of Indirection

**Triple indirect** means three levels of indirection

triple indirect

Data block

Data block

Data block

double indirect

inode with
512 pointers

inode with
512 pointers

...

single indirect

inode with
512 pointers

inode with
512 pointers

...

inode with
512 pointers

Head inode with 15
indirect pointers

inode with
512 pointers

...

...

# Multi-Level Indexing

Head inode may combine multiple levels of indirection



12 * 4KB = 48KB

1 * 512 * 4KB = 24MB

1 * 512 * 512 * 4KB = 12GB

Head inode
12 directory pointers
1 single indirect pointer
1 double indirect pointer

Data block
Data block
Data block
...

512 pointers

Data block
Data block
Data block
...

512 pointers

512 pointers
512 pointers
512 pointers
...

Data block
Data block
Data block
...

# Common Observations of File Systems

| | |
|---|---|
| **Most files are small** | ˜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 |

# Example Reading a File

How to read from the file /foo/bar? First, fd=open("/foo/bar").

1. inode number for root directory is well known; read "root inode" and get the location for "root data".

2. Read "root data" to get inode number for "/foo".

3. Read "foo inode" to get location of "foo data".

4. Read "foo data" to get inode number for "/foo/bar".

5. Read "bar inode" and bring the metadata of /foo/bar to main memory.

| | data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data | bar data [0] | bar data [1] | bar data [2] |
|---|---|---|---|---|---|---|---|---|---|---|
| open(bar) | | | read | read | read | read | read | | | |
| read() | | | | | read | | | read | | |
| | | | | | write | | | | | |
| read() | | | | | read | | | | read | |
| | | | | | write | | | | | |
| read() | | | | | read | | | | | read |
| | | | | | write | | | | | |

# Example Reading a File

How to read from the file /foo/bar?
Second, read(fd,...).

1. Read "bar inode" to get the location of the first data block of /foo/bar, i.e., "bar data [0]".

2. Read "bar data [0]".

3. Update last access time at "bar inode".

4. Read "bar inode" to get the location of the second block of /foo/bar, i.e., "bar data[1]".

5. Read "bar data [1]".

6. Update last access time at "bar inode".

7. ... ...

| | data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data | bar data [0] | bar data [1] | bar data [2] |
|---|---|---|---|---|---|---|---|---|---|---|
| open(bar) | | | read | read | read | read | read | read | | |
| read() | | | | | read<br>write | | | read | | |
| read() | | | | | read<br>write | | | | read | |
| read() | | | | | read<br>write | | | | | read |

# Example: Creating and Writing File

How to create file /foo/bar?

1. Read "root inode" (i.e., inode for "/") to get location of "/".

2. Read data block of "/" (i.e., "root data") to get inode number of "/foo".

3. Read "foo inode" to get location of "/foo".

4. Read "/foo" (i.e., "foo data"); read inode bitmap to find an empty inode, denoted as x, and update the bitmap; add a pair ("bar", x) to "foo data".

5. Read "bar inode" (i.e., inode with number x) and initialize it.

6. Update last access time at "foo inode".

| | data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data | bar data [0] | bar data [1] | bar data [2] |
|---|---|---|---|---|---|---|---|---|---|---|
| create (/foo/bar) | | read write | read | read | read write write | read | read write | read | | |
| write() | | read write | | | read write | | | write | | |
| write() | | read write | | | read write | | | | write | |
| write() | | read write | | | read write | | | | | write |

# Example: Creating and Writing File

How to <span style="color:red">write</span> to new file /foo/bar?

1. Read "bar inode".
2. Read data bitmap to find an empty data block x and update the bitmap.
3. Write to data block x (i.e., "bar data [0]").
4. Add block x to "bar inode".
5. Update the last access time at "bar inode"
6. Read "bar inode".
7. Read data bitmap to find an empty data block y and update the bitmap.
8. Write to data block y (i.e., "bar data [1]").
9. Add block y to "bar inode".
10. … ...

| | data bitmap | inode bitmap | root inode | foo inode | bar inode | root data | foo data | bar data [0] | bar data [1] | bar data [2] |
|---|---|---|---|---|---|---|---|---|---|---|
| create (/foo/bar) | | read write | read, write | read | read write, write | read | read, write | | | |
| write() | read write | | | | read, write | | | write | | |
| write() | read write | | | | read, write | | | | write | |
| write() | read write | | | | read, write | | | | | write |

# Basic Performance Improvements

**Caching** – holds popular blocks to decrease number of times blocks are read from disk

**Write buffering** - batch multiple updates into a smaller set of I/O operations