

Operating Systems

Lecture 14

fs design

Prof. Mengwei Xu

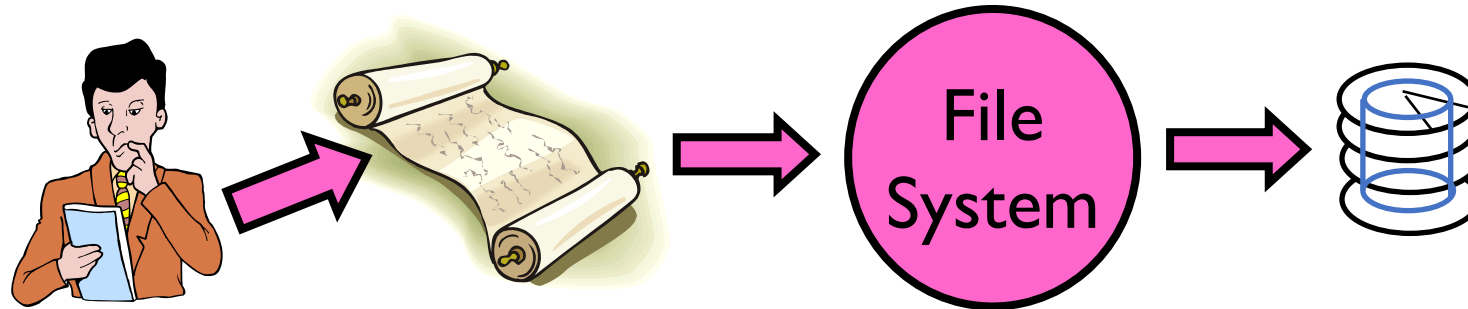
Recap: Building a File System

- **File System:** Layer of OS that transforms block interface of disks (or other block devices) into Files, Directories, etc.
- File System Components
 - **Naming:** Interface to find files by name, not by blocks
 - **Disk Management:** collecting disk blocks into files
 - **Protection:** Layers to keep data secure
 - **Reliability/Durability:** Keeping of files durable despite crashes, media failures, attacks, etc.

Recap: User vs. System View of a File

- User's view:
 - Durable Data Structures
- System's view (system call interface):
 - Collection of Bytes (UNIX)
 - Doesn't matter to system what kind of data structures you want to store on disk!
- System's view (inside OS):
 - Collection of blocks (a block is a logical transfer unit, while a sector is the physical transfer unit)
 - Block size \geq sector size; in UNIX, block size is 4KB

Translating from User to System View



- What happens if user says: give me bytes 2—12?
 - Fetch block corresponding to those bytes
 - Return just the correct portion of the block
- What about: write bytes 2—12?
 - Fetch block
 - Modify portion
 - Write out block
- Everything inside File System is in whole size blocks
 - For example, `getc()`, `putc()` \Rightarrow buffers something like 4096 bytes, even if interface is one byte at a time
- From now on, file is a collection of blocks

Disk Management Policies (1/2)

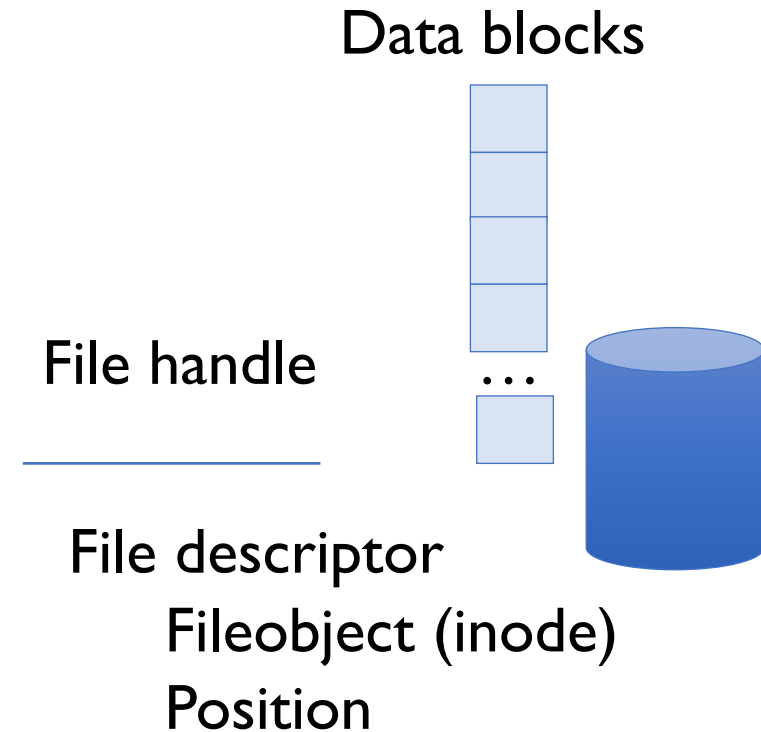
- Basic entities on a disk:
 - **File**: user-visible group of blocks arranged sequentially in logical space
 - **Directory**: user-visible index mapping names to files
- Access disk as linear array of sectors. Two Options:
 - Identify sectors as vectors [cylinder, surface, sector], sort in cylinder-major order
 - ❑ Used in BIOS, but not in OSes anymore
 - **Logical Block Addressing (LBA, 逻辑块寻址)**: Every sector has integer address from zero up to max number of sectors
 - Controller translates from address \Rightarrow physical position
 - ❑ First case: OS/BIOS must deal with bad sectors
 - ❑ Second case: hardware shields OS from structure of disk

Recap: Disk Management Policies (2/2)

- Need way to track free disk blocks
 - Link free blocks together \Rightarrow too slow today
 - Use bitmap to represent free space on disk
- Need way to structure files: [File Header](#)
 - Track which blocks belong at which offsets within the logical file structure
 - [Optimize placement of files' disk blocks to match access and usage patterns](#)

Recap: File

- Named permanent storage
- Contains
 - Data
 - ☐ Blocks on disk somewhere
 - Metadata (Attributes)
 - ☐ Owner, size, last opened, ...
 - ☐ Access rights
 - R, W, X
 - Owner, Group, Other (in Unix systems)
 - Access control list in Windows system



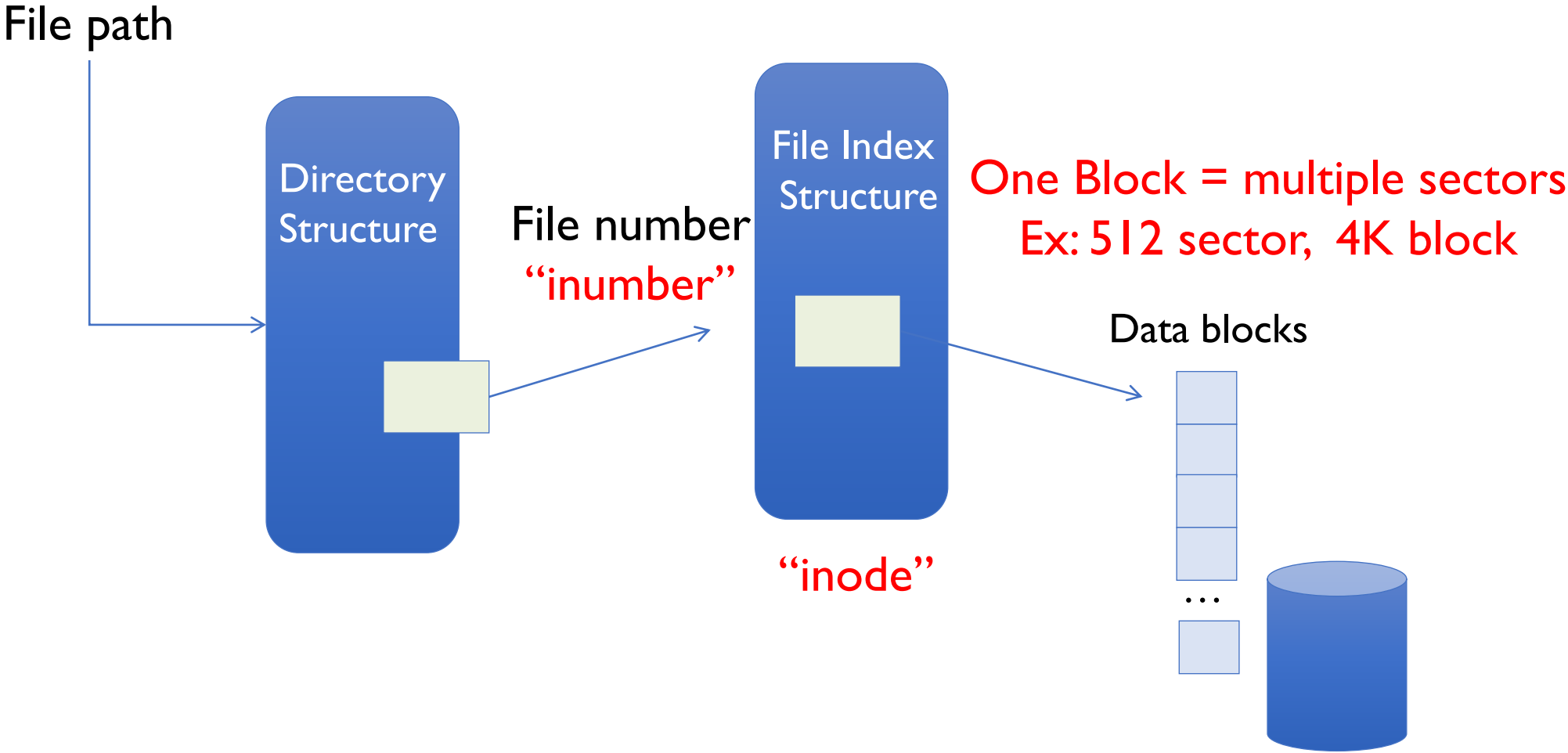
Recap: Directory

- Basically a hierarchical structure
- Each directory entry is a collection of
 - Files
 - Directories
 - A link to another entries
- Each has a name and attributes
 - Files have data
- Links (hard links) make it a DAG, not just a tree
 - Softlinks (aliases) are another name for an entry

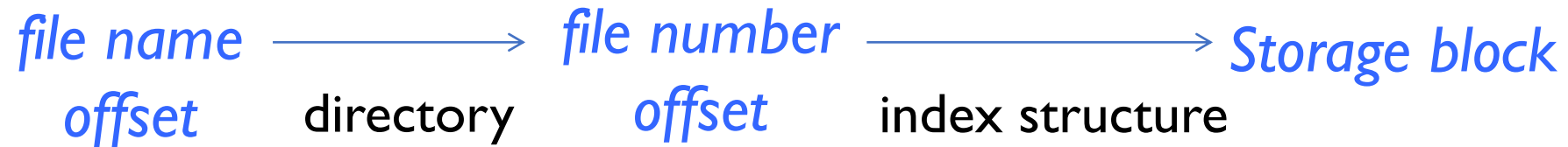
Recap: Designing a File System ...

- What factors are critical to the design choices?
- Durable data store => it's all on disk
- (Hard) Disks Performance !!!
 - Maximize sequential access, minimize seeks
- Open before Read/Write
 - Can perform protection checks and look up where the actual file resource are, in advance
- Size is determined as they are used !!!
 - Can write (or read zeros) to expand the file
 - Start small and grow, need to make room
- Organized into directories
 - What data structure (on disk) for that?
- Need to allocate / free blocks
 - Such that access remains efficient

Components of a File System

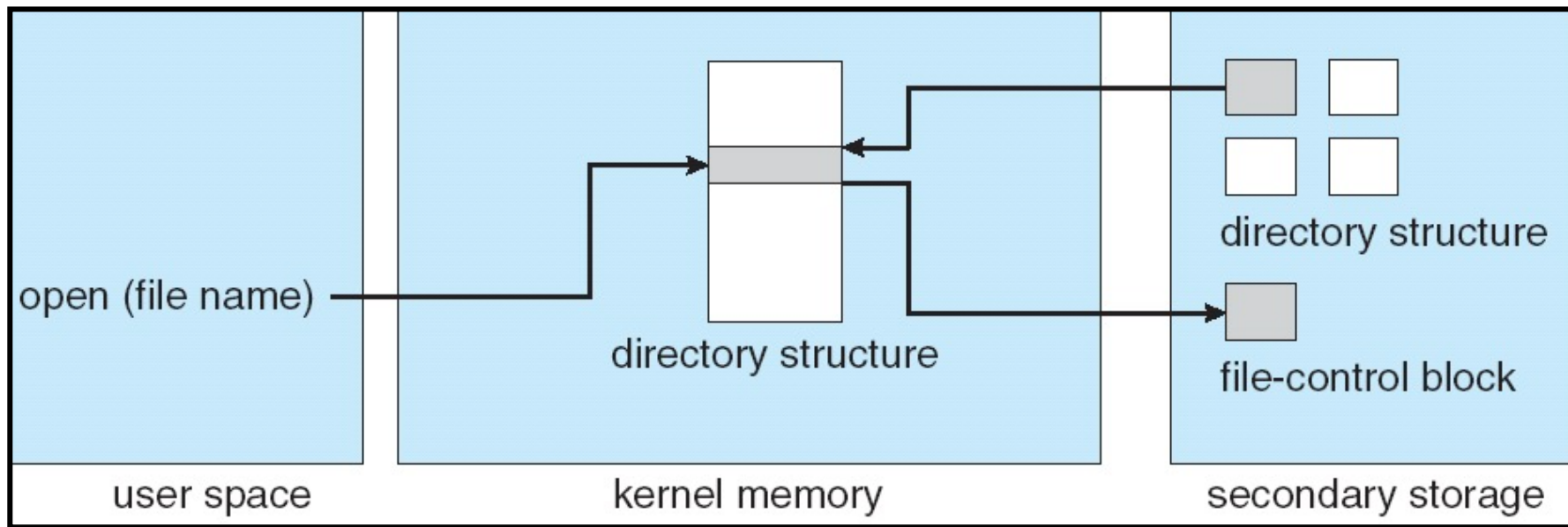


Components of a file system



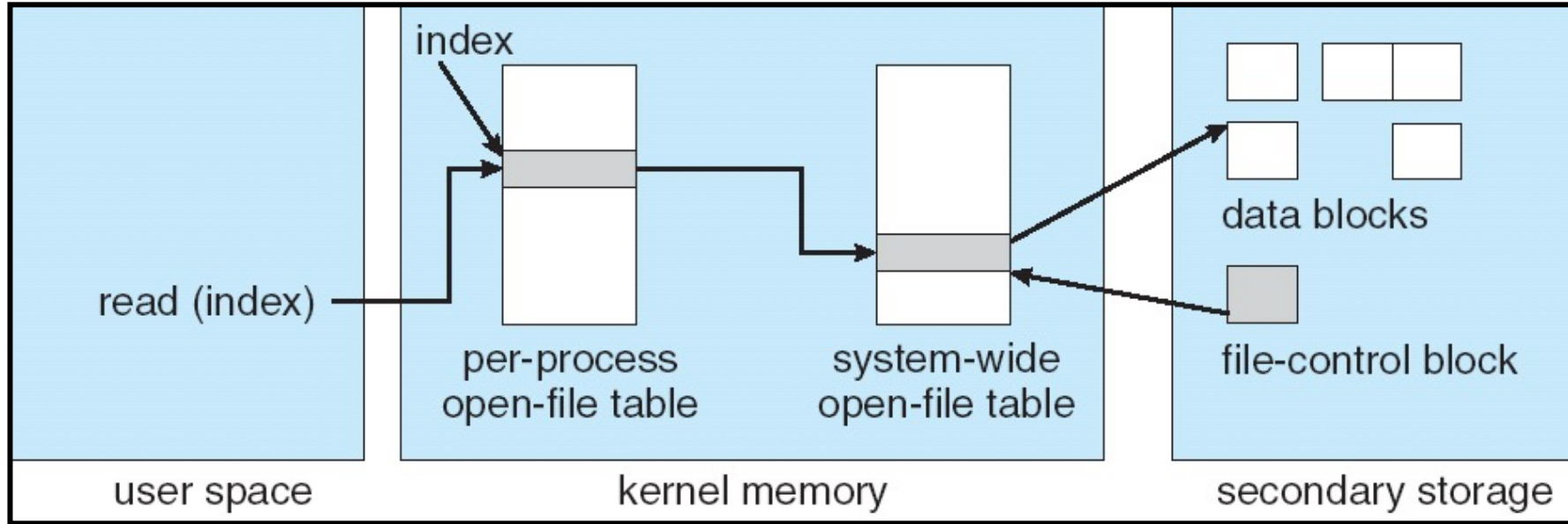
- Open performs *Name Resolution*
 - Translates pathname into a “file number”
 - Used as an “index” to locate the blocks
 - Creates a file descriptor in PCB within kernel
 - Returns a “handle” (another integer) to user process
- Read, Write, Seek, and Sync operate on handle
 - Mapped to file descriptor and to blocks

In-Memory File System Structures



- Open system call:
 - Resolves file name, finds file control block (inode)
 - Makes entries in per-process and system-wide tables
 - Returns index (called “file handle”) in open-file table

In-Memory File System Structures



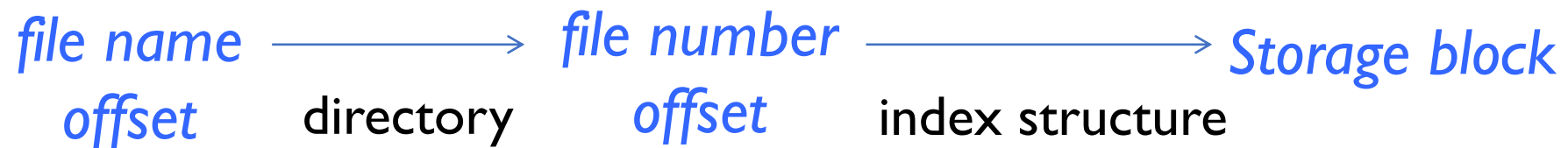
- Read/write system calls:
 - Use file handle to locate inode
 - Perform appropriate reads or writes

Typical File Systems

- FAT (Microsoft File Allocation Table), 1970s.
 - Extremely simple index structure: a linked list.
 - Still widely used in devices like flash memory sticks and digital cameras
- FFS (Unix Fast File System), 1980s.
 - Tree-based multilevel index to improve random access efficiency.
 - Uses a collection of locality heuristics to get good spatial locality.
 - EXT2 and EXT3 are based on FFS.
- NTFS (Microsoft New Technology File System): 1990s.
 - More flexible tree structure.
 - Mainstream file system on MS.
 - It's representative to EXT4, XFS, and Apple's Hierarchical File Systems (HFS and HFS+).

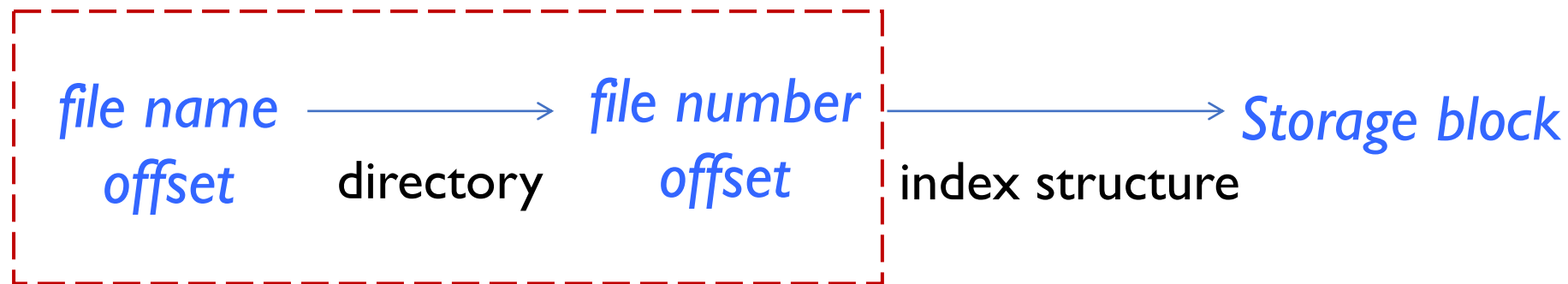
Goals for Today

- Directories: naming data
 - How do we convert a file name to the file number?
- Files: finding data
 - How do we locate storage block based on file number?
- Virtual file systems (VFS)
 - How do we make different FSs work together easily?



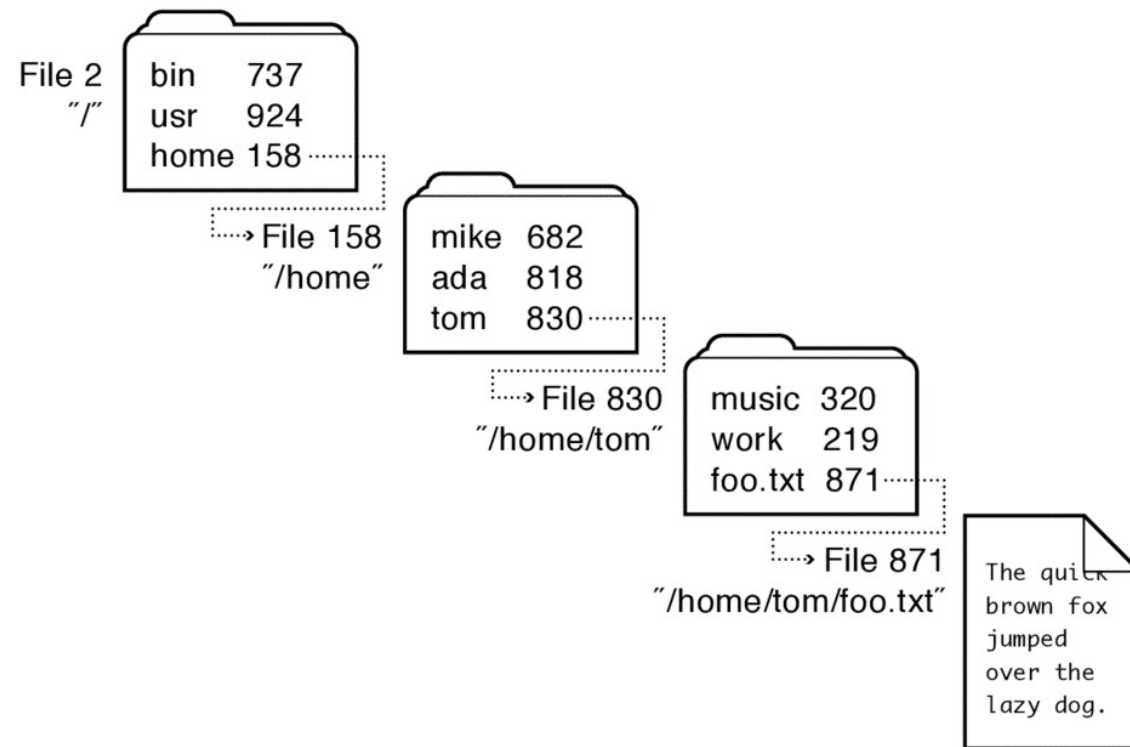
Goals for Today

- Directories: naming data
 - How do we convert a file name to the file number?
- Files: finding data
 - How do we locate storage block based on file number?
- Virtual file systems (VFS)
 - How do we make different FSs work together easily?



Directory Structure

- Directory is treated as a file with a list of <file name: file number> mappings
- The file number of the root directory is agreed ahead of time
 - In many Unix FSs, it's 2.



Directory Operations

- Stored in files, can be read, but typically don't

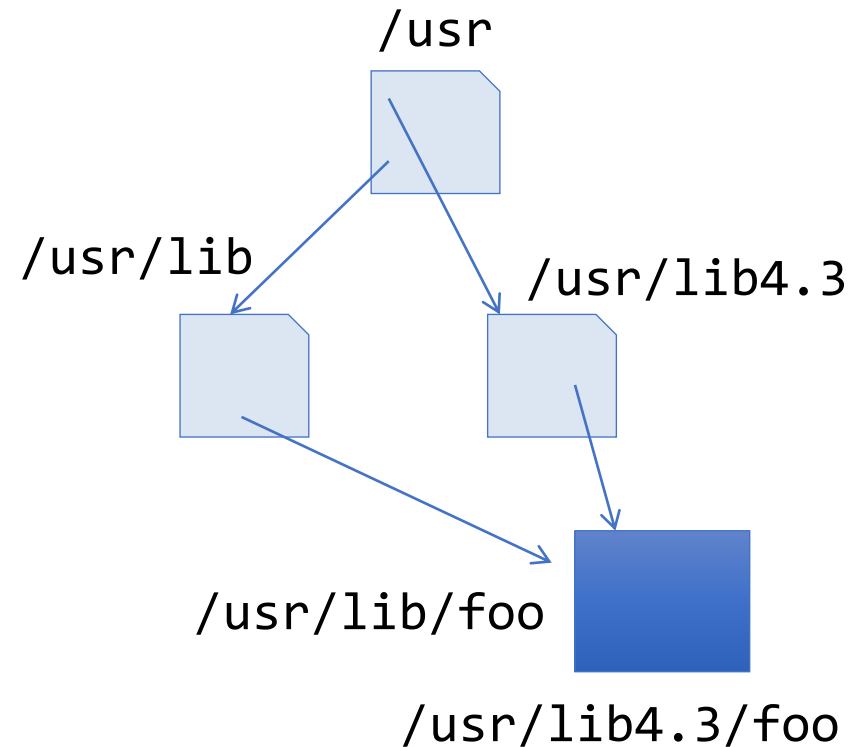
- System calls to access directories
- **open** / **creat** traverse the structure
- **mkdir** / **rmdir** add/remove entries
- **link** / **unlink (rm)**
 - ❑ Link existing file to a directory
 - Not in FAT !
 - ❑ Forms a DAG

- When can file be deleted?

- Maintain ref-count of links to the file
- Delete after the last reference is gone

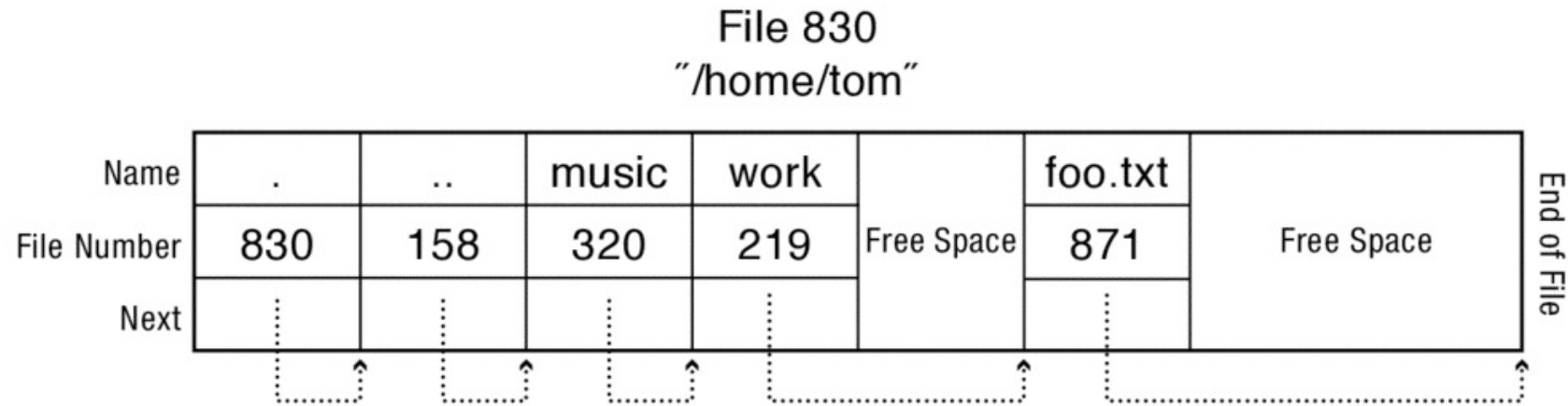
- libc support

- `DIR * opendir (const char *dirname)`
- `struct dirent * readdir (DIR *dirstream)`
- `int readdir_r (DIR *dirstream, struct dirent *entry, struct dirent **result)`



Directory Internals

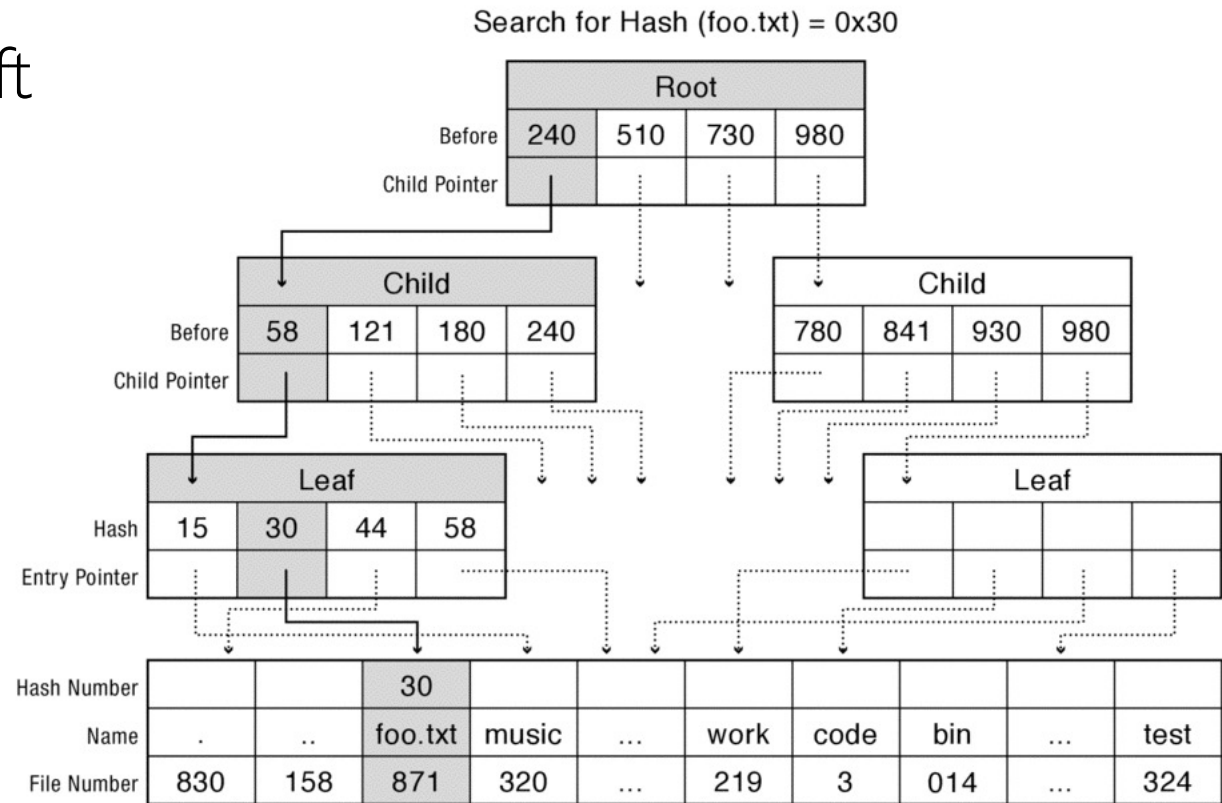
- Early implementations simply stored linear lists of <file name, file number> in directory files.
 - Free spaces are for new entries. Note: files can be added/deleted.



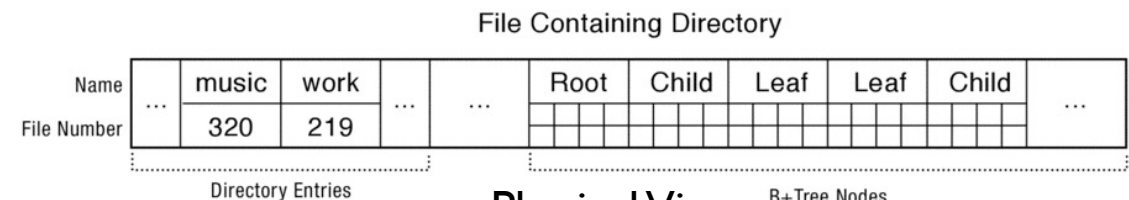
- Works fine in most cases. But when there are thousands of files in a directory..The access could be slow!

Directory Internals

- Modern FSs (Linux XFS, Microsoft NTFS, and Oracle ZFS) organize directory's contents as a tree.
 - B/B+ tree: fast lookup, insert, and removal
 - Names are first hashed into a key, which is used to find the file number in the tree



Logical View



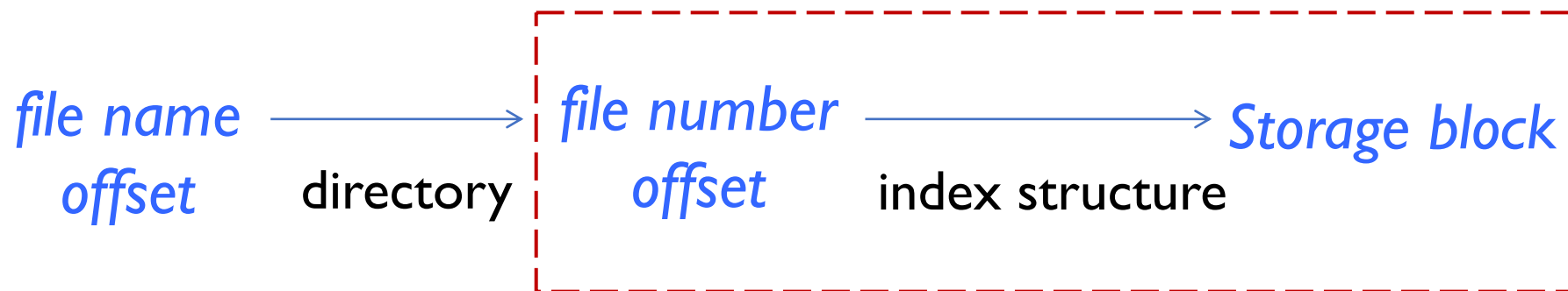
Physical View

Directory Structure Access Cost

- How many disk accesses to resolve “/my/book/count”?
 - Read in file header for root (fixed spot on disk)
 - Read in first data block for root
 - ❑ Table of file name/index pairs. Search linearly – ok since directories typically very small
 - Read in file header for “my”
 - Read in first data block for “my”; search for “book”
 - Read in file header for “book”
 - Read in first data block for “book”; search for “count”
 - Read in file header for “count”
- **Current working directory:** Per-address-space pointer to a directory (inode) used for resolving file names
 - Allows user to specify relative filename instead of absolute path (say CWD=“/my/book” can resolve “count”)

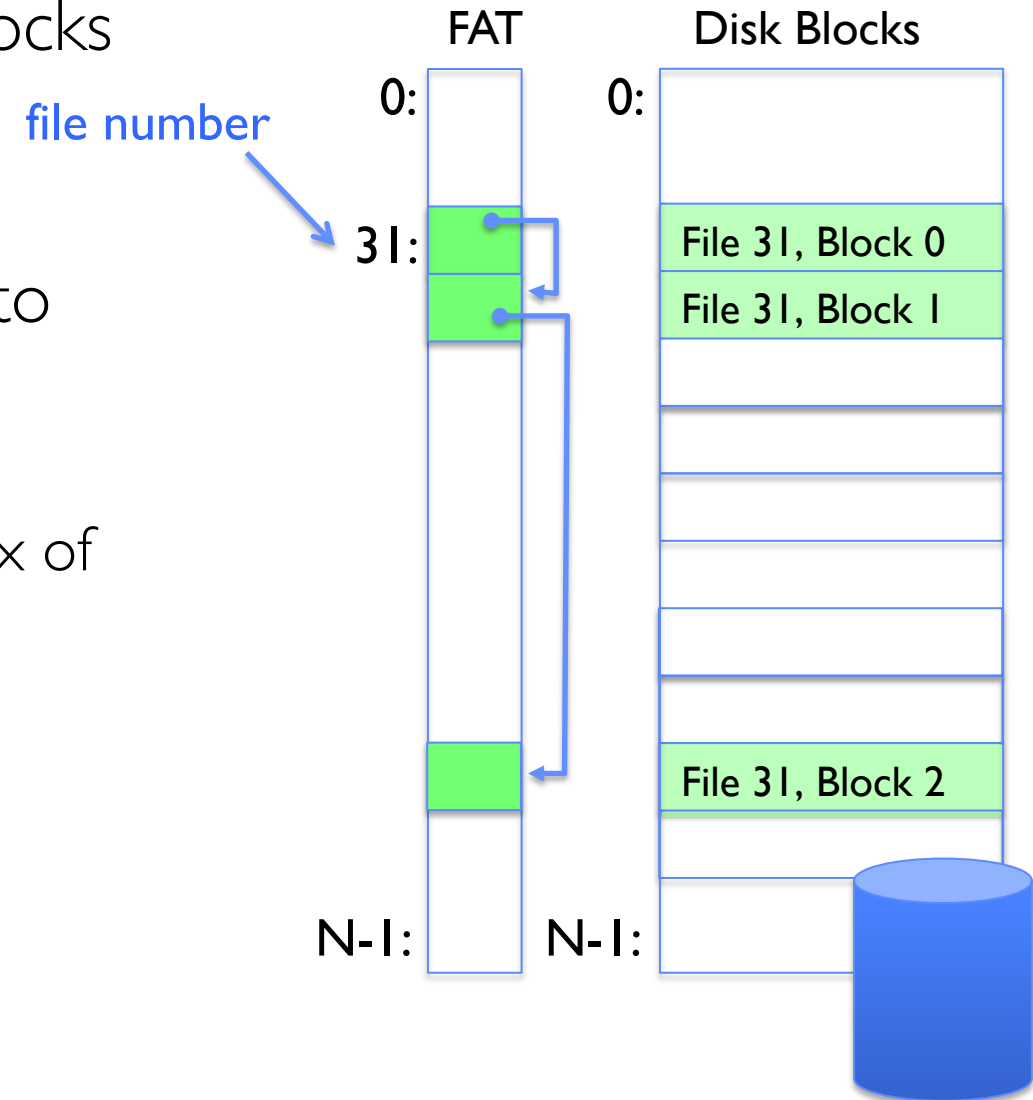
Goals for Today

- Directories: naming data
 - How do we convert a file name to the file number?
- Files: finding data
 - How do we locate storage block based on file number?
- Virtual file systems (VFS)
 - How do we make different FSs work together easily?



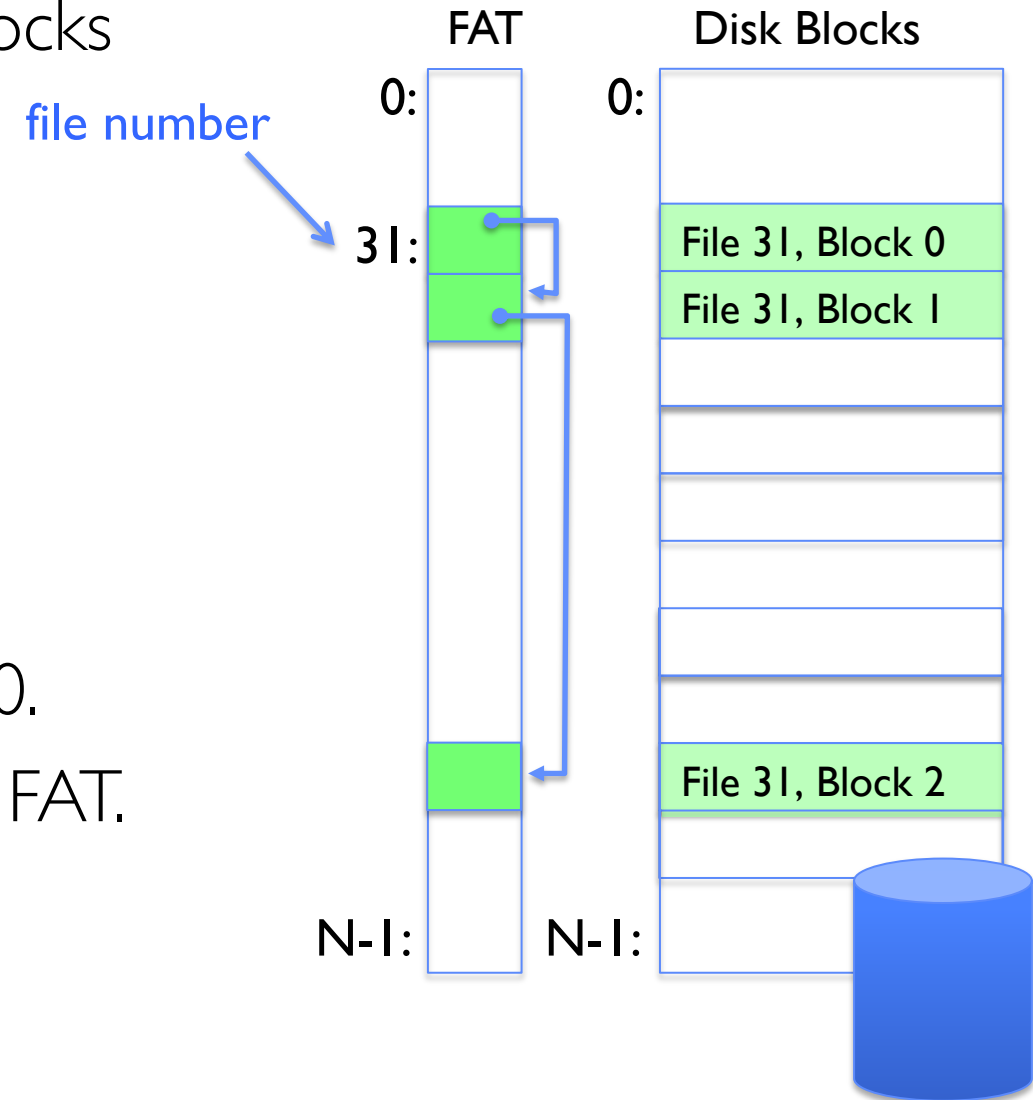
FAT (File Allocation Table)

- FAT is a linked list as I-I map with blocks
 - Represented as a list of 32-bit entries
 - Older versions use fewer bits
- Each entry in FAT contains a pointer to the next FAT entry of the same file
 - Or a special END_OF_FILE value.
 - The file number is the 1st (or root) index of the block list for the file
- For File No. #i, its
 - 1st data block index: i
 - 2nd data block index: $*(FAT[i])$
 - 3rd data block index: $*(*(FAT[i]))$
 - ..



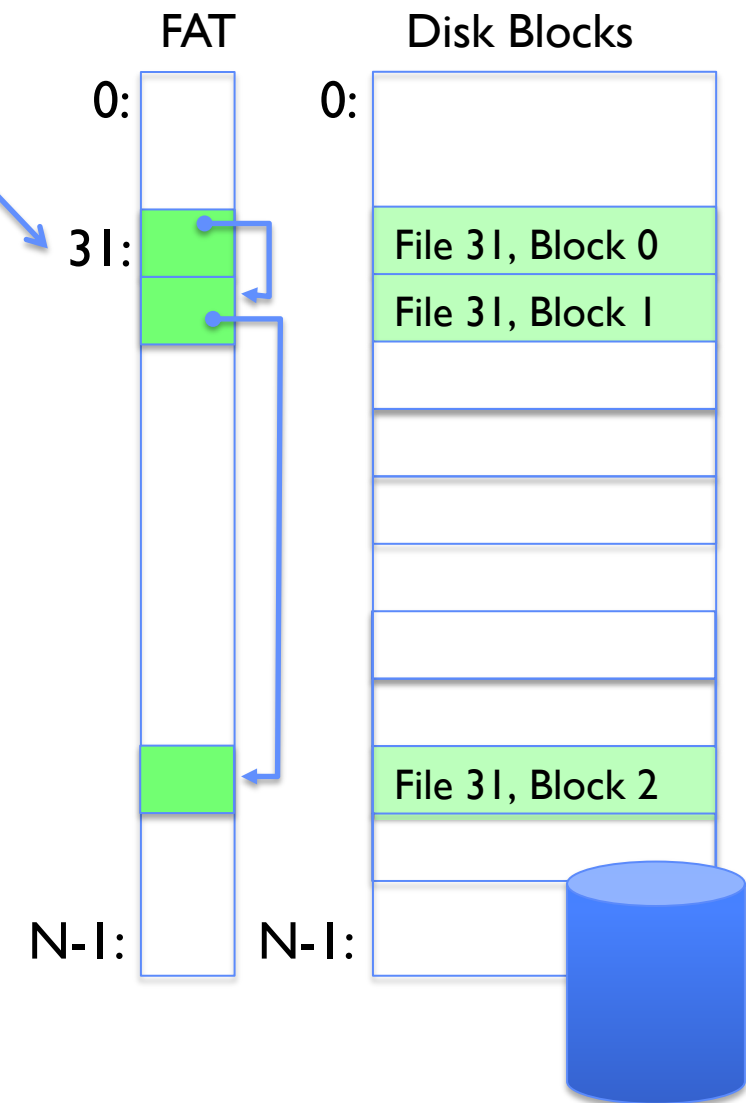
FAT (File Allocation Table)

- FAT is a linked list as I-I map with blocks
 - Represented as a list of 32-bit entries
- Where is FAT stored?
 - On Disk, on boot cache in memory, second (backup) copy on disk
- Free space: FAT free list, i.e., $FAT[i] = 0$.
- To find a free block: scanning through FAT.



FAT (File Allocation Table)

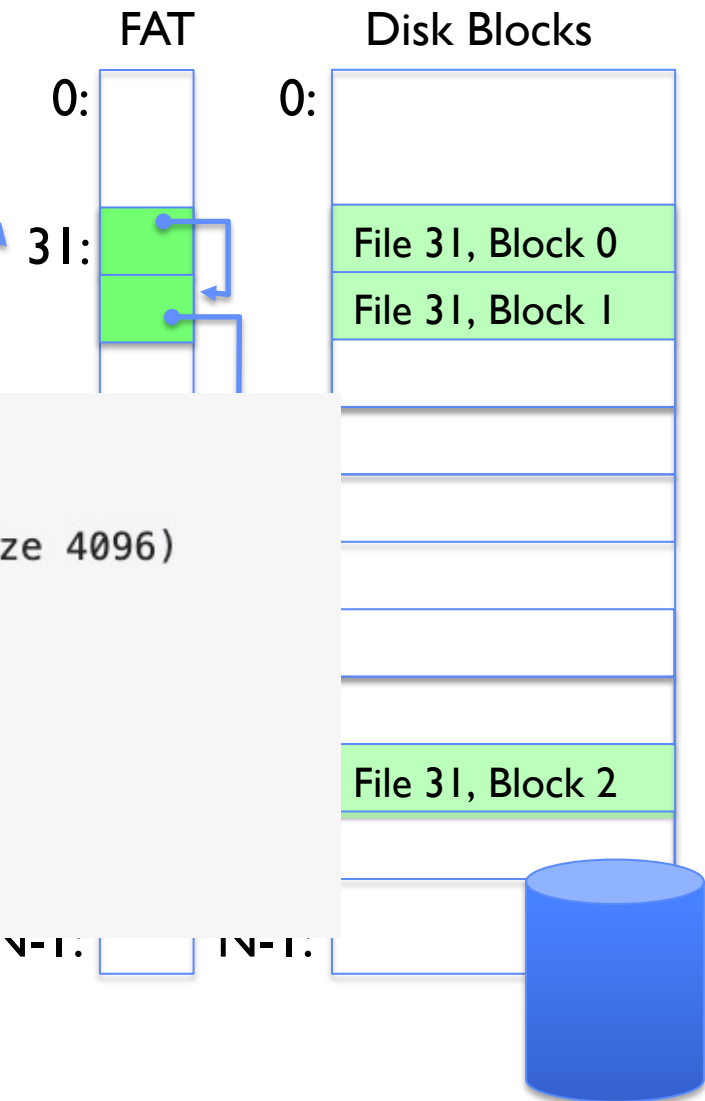
- FAT is a linked list as I-I map with blocks
 - Represented as a list of 32-bit entries
- **Locality** (storing a file in sequential blocks) is important for fast I/O
 - Sequential I/O is much faster than random I/O
 - Imagine you want to write 100MB to a 200MB file.. FS cannot guarantee they are stored sequential
- How to ensure good locality heuristics in FAT?
 - Simple strategy: *next fit*, i.e., scans sequentially through the FAT starting from the last entry that was allocated and return the next free entry
 - Still, there will be increasing fragment
 - Defragmentation tool: read files and rewrite them to new locations with better locality



FAT (File Allocation Table)

- FAT is a linked list as I-I map with blocks
 - Represented as a list of 32-bit entries
- Locality** (storing a file in sequential blocks) is important for fast I/O
 - Sequential I/O is much faster than random I/O

file number

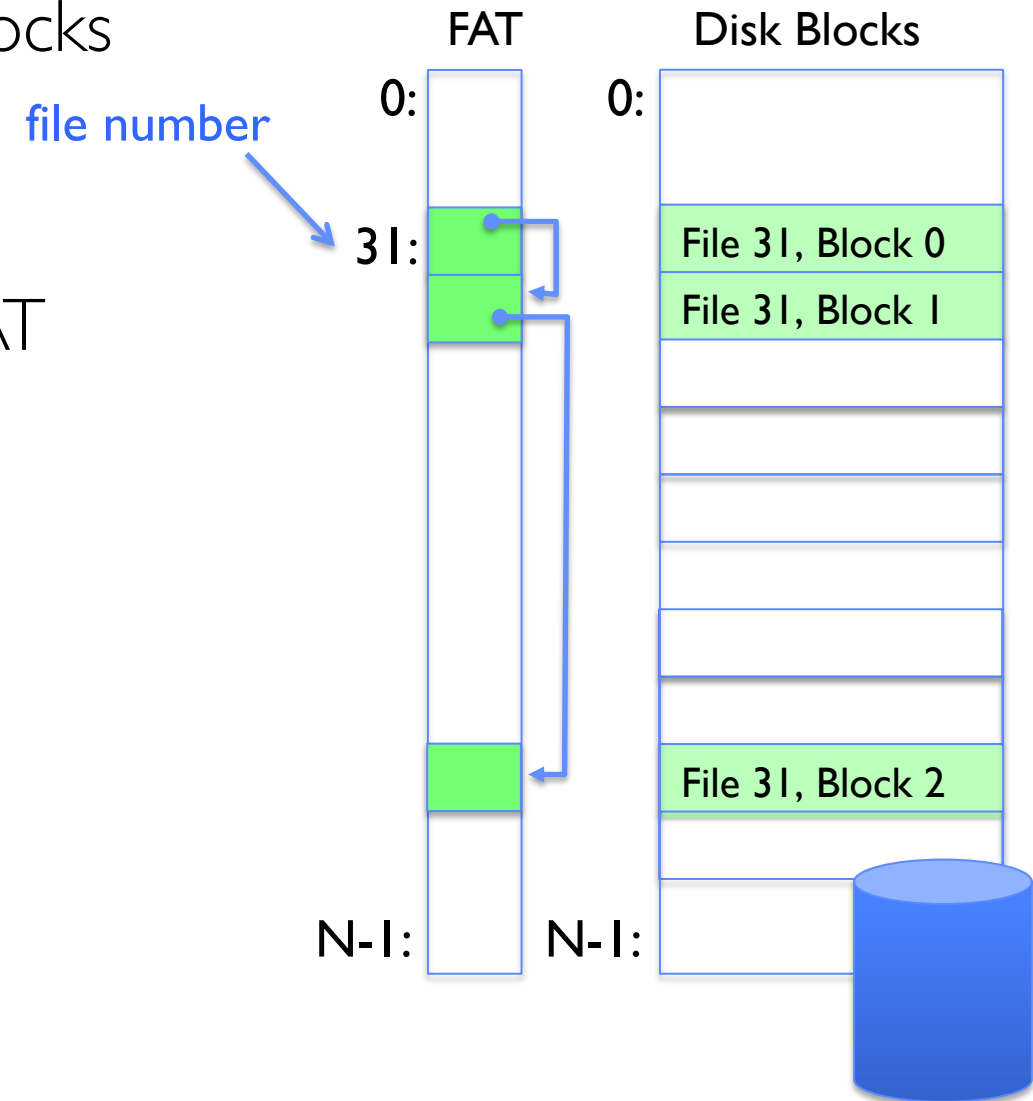


```
# filefrag -v /var/log/messages.1
Filesystem type is: ef53
File size of /var/log/messages.1 is 41733 (11 blocks, blocksize 4096)
ext logical physical expected length flags
0      0  2130567          1
1      1 15907576  2130568          1
2      2 15910400 15907577          1
3      3 15902720 15910401          7
4     10 2838546 15902727          1 eof
/var/log/messages.1: 5 extents found
```

- Still, there will be increasing fragment
- Defragmentation tool: read files and rewrite them to new locations with better locality

FAT (File Allocation Table)

- FAT is a linked list as I-I map with blocks
 - Represented as a list of 32-bit entries
- READ: just get block by block with FAT
- WRITE
 - Get blocks from free list
 - Linking them into a file
- Format a disk
 - Zero the blocks, link up the FAT free list
- Quick format
 - Link up the FAT free-list



FAT Issues

- Poor locality: there will be fragmentations
- Poor random access: needs to traverse the file's FAT entries till the block is reached
- Limited file metadata and access control: only has file's name, size, and creation time, but cannot specify the file's owner or group.
- No support for hard links: no room for any other file metadata.
- Limitations on volume and file size
 - With top 4 bits reserved.
 - 2^{28} blocks * 4KB block size = 1TB.
 - Larger block size (up to 256KB)?
 - File size is encoded in 32 bits, so less than 4GB.

Unix File System (1/2)

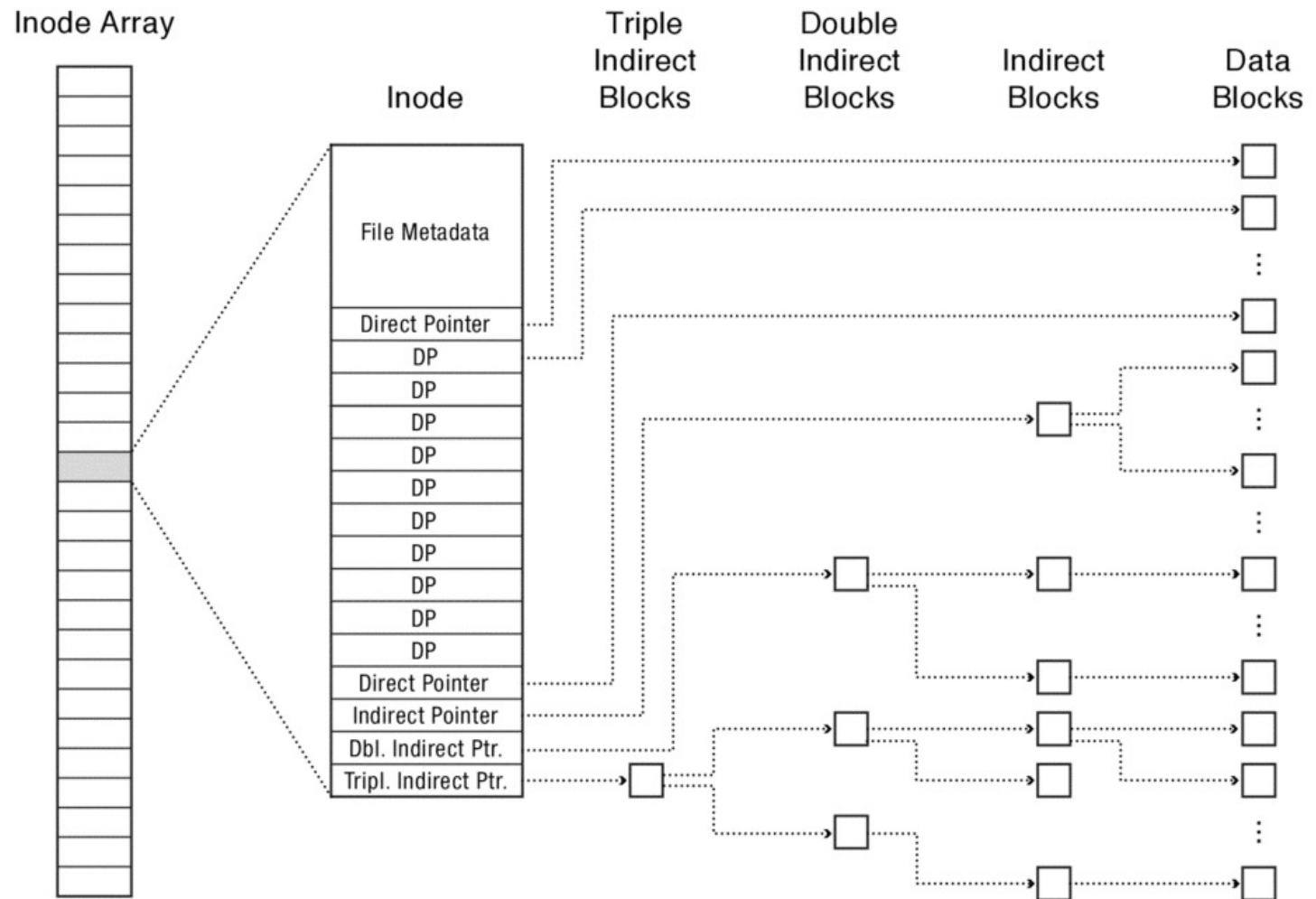
- Original inode format appeared in BSD 4.1
 - Berkeley Standard Distribution Unix
 - Similar structure for Linux Ext2/3
- File Number is index into inode arrays
- Multi-level index structure
 - Great for little and large files
 - Asymmetric tree with fixed sized blocks

Unix File System (2/2)

- Metadata associated with the file
 - Rather than in the directory that points to it
- UNIX Fast File System (FFS) BSD 4.2 Locality Heuristics:
 - Block group placement
 - Reserve space
- Scalable directory structure

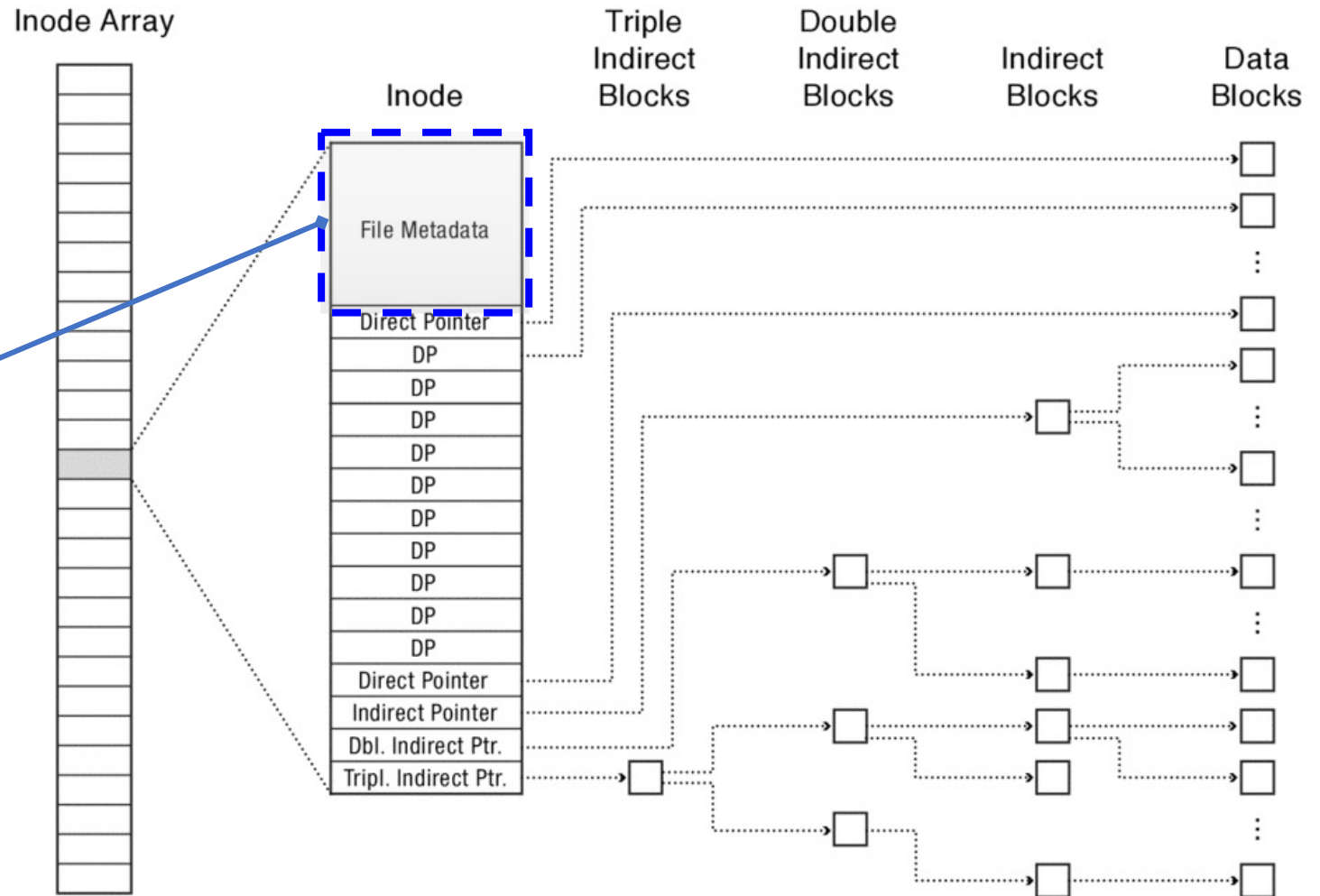
File Attributes

- Multi-level index
 - Fixed, asymmetric tree



File Attributes

- User
- Group
- 9 basic access control bits
 - UGO x RWX
- Setuid bit
 - execute at owner permissions rather than user
- Setgid bit
 - execute at group's permissions



File Attributes

```

echo:homepage echo$ ls -la
total 176
drwxr-xr-x@ 11 echo staff 352 Nov 2 13:55 .
drwxr-xr-x@ 10 echo staff 320 Nov 2 13:55 ..
-rw-r--r--@ 1 echo staff 6148 Nov 2 13:45 .DS_Store
drwxr-xr-x@ 12 echo staff 384 Nov 2 13:55 .git
-rw-r--r--@ 1 echo staff 1374 Jul 5 09:55 awards.html
drwxr-xr-x@ 48 echo staff 1536 Nov 2 13:47 files
drwxr-xr-x@ 8 echo staff 256 Dec 9 2021 image
-rwxr-xr-x@ 1 echo staff 55677 Nov 2 13:48 index.html
-rw-r--r--@ 1 echo staff 18233 Jun 22 2021 index.old.html
drwxr-xr-x@ 4 echo staff 128 Jan 8 2022 materials
drwxr-xr-x@ 6 echo staff 192 Dec 9 2021 projects

```

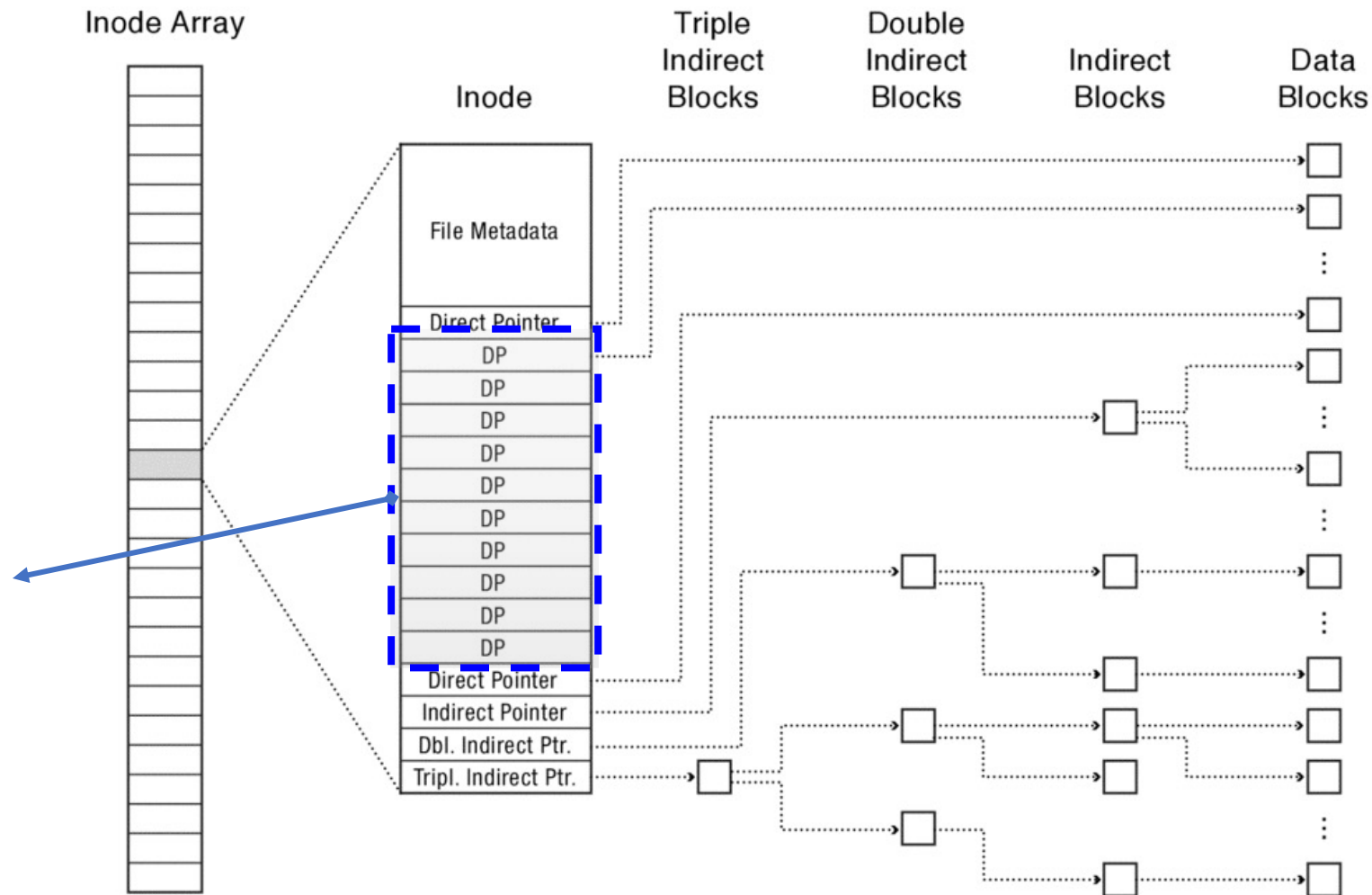
User
 Group
 9 basic
 - UC
 Setuid
 - ex
 ra
 Setgid
 - execute at group's permissions



File Attributes

12x Direct Pointers (直接索引)

4kB blocks \Rightarrow
sufficient for files up to 48KB



File Attributes

Indirect pointers

- point to a disk block containing only pointers

4 kB blocks => 1024 ptrs

Indirect Pointer (一级间接索引)

=> 4 MB

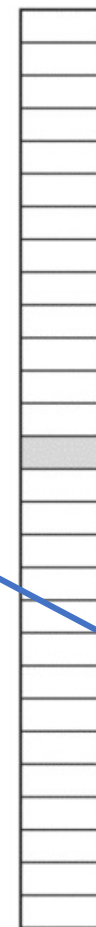
Double Indirect Pointer (二级..)

=> 4 GB

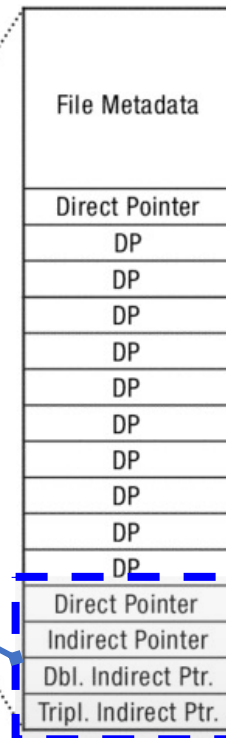
Triple Indirect Pointer (三级..)

=> 4 TB

Inode Array



Inode

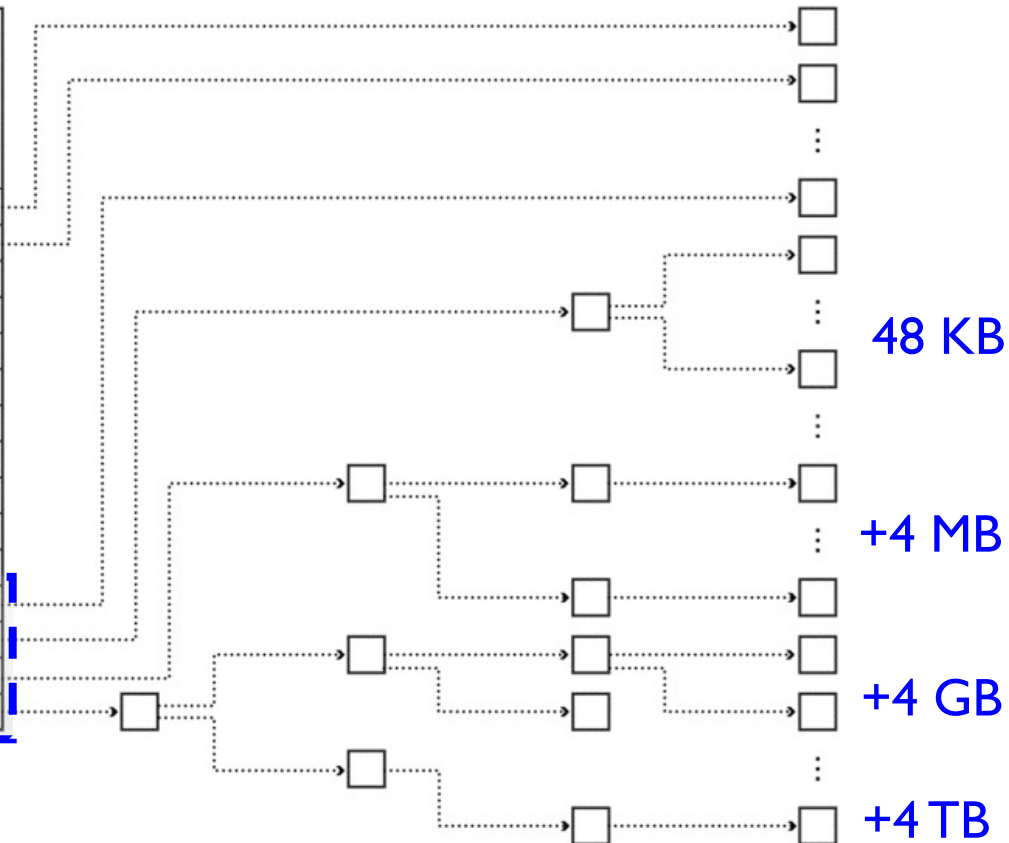


Triple Indirect Blocks

Double Indirect Blocks

Indirect Blocks

Data Blocks



FFS Characteristics

- **Tree structure.** Each file is represented as a tree, which allows the file system to efficiently find any block of a file.
- **High degree.** The FFS tree uses internal nodes with many children.
 - A 4KB file block contains 1024x blockID in 4 bytes.
 - Improves sequential reads and writes. Why?
- **Fixed structure.** The FFS tree has a fixed structure.
 - For a given configuration of FFS, the first set of d pointers always point to the first d blocks of a file; etc.
 - Make implementation easier.
- **Asymmetric.** FFS's tree structure is asymmetric, i.e., different depths.
 - Small files can be stored with low cost (size and access speed).
 - While we still support very large files.

Asymmetric vs. Symmetric

- In a symmetric tree with each entry to be triple indirect pointers

⇒ To store a 4B file, how much space we need?

Asymmetric vs. Symmetric

- In a symmetric tree with each entry to be triple indirect pointers

⇒ To store a 4B file, how much space we need?

- 4B data + small inode + 3x 4KB indirect blocks
- How about our asymmetric tree?

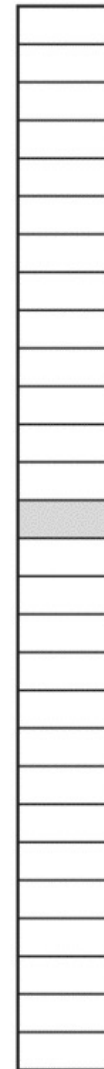
Asymmetric vs. Symmetric

- In a symmetric tree with each

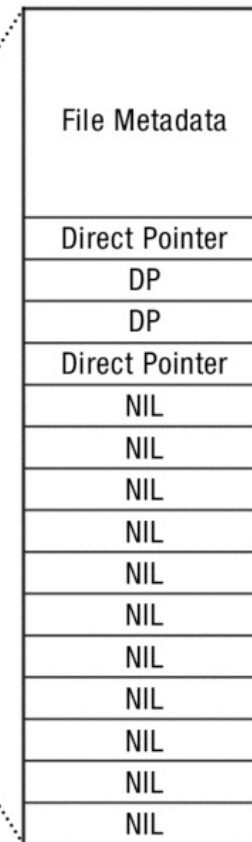
⇒ To store a 4B file, how much

- 4B data + small inode + 3x 4B
- How about our asymmetric tree

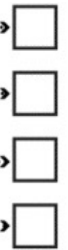
Inode Array



Inode



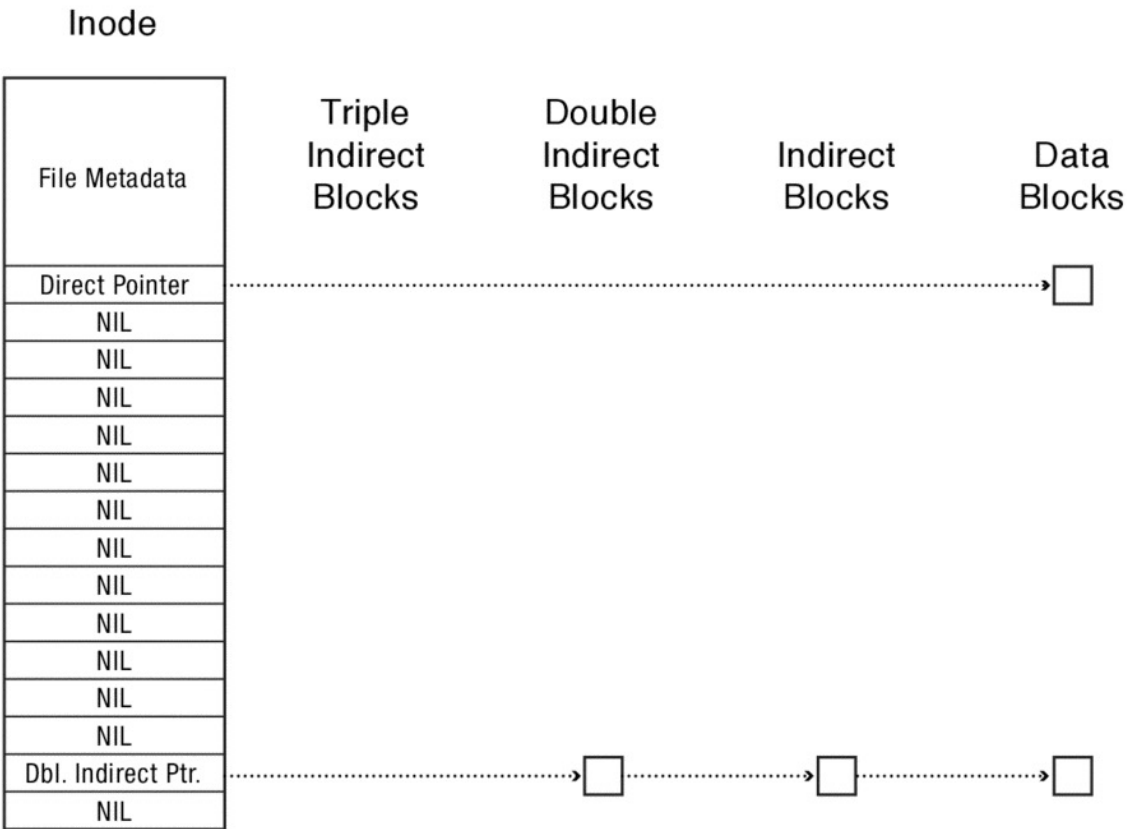
Data Blocks



A file with 4 blocks, each block accessed directly

Sparse Files

- FFS can support sparse files, in which one or more ranges of empty space are surrounded by file data.
 - Those empty space shall not consume disk space.



<= A sparse file with 4KB data at offset 0, and 4KB data at offset 2^{30} .

Command `ls` shows it takes 1.1GB.
 Command `du` shows it takes 16KB.

Free Space Management

- FFS allocates a bitmap with one bit per storage block. The i -th bit in the bitmap indicates whether the i th block is free or in use.
- The position of FFS's bitmap is fixed when the file system is formatted.
 - So it is easy to find the part of the bitmap that identifies free blocks near any location of interest.

Where are inodes Stored?

- In early UNIX and DOS/Windows' FAT file system, headers stored in special array in outermost cylinders
- Header not stored anywhere near the data blocks
 - To read a small file, seek to get header, seek back to data
- Fixed size, set when disk is formatted
 - At formatting time, a fixed number of inodes are created
 - Each is given a unique number, called an "inumber"

Where are inodes Stored?

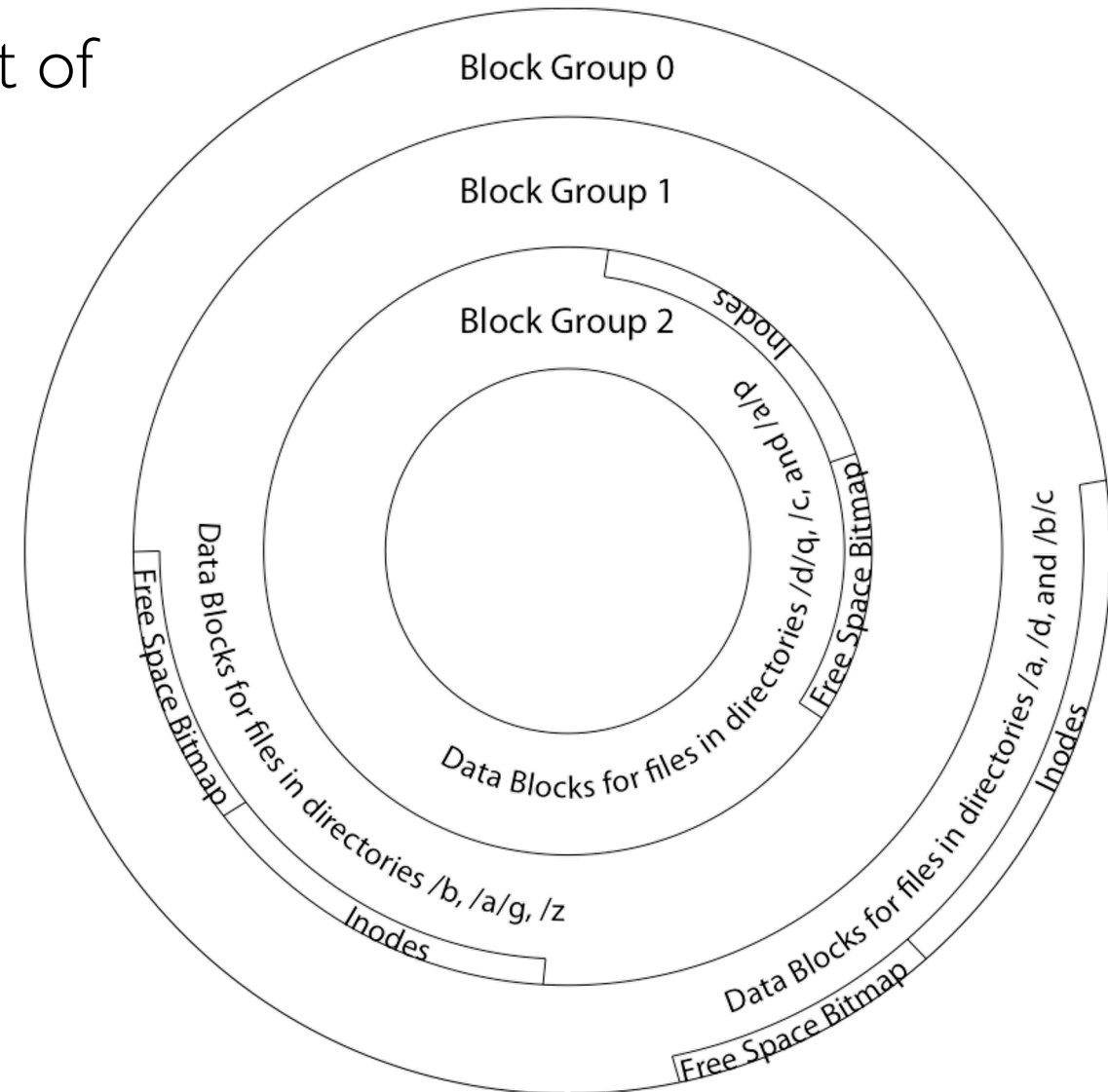
- Later versions of UNIX moved the header information to be closer to the data blocks
 - Often, inode for file stored in same “cylinder group” as parent directory of the file (makes an **ls** of that directory run fast)
- Pros:
 - UNIX BSD 4.2 puts bits of file header array on many cylinders
 - For small directories, can fit all data, file headers, etc. in same cylinder \Rightarrow no seeks!
 - File headers much smaller than whole block (a few hundred bytes), so multiple headers fetched from disk at same time
 - Reliability: whatever happens to the disk, you can find many of the files (even if directories disconnected)
- Part of the Fast File System (FFS)
 - General optimization to avoid seeks

Locality Heuristics

- **Block group placement:** FFS places data to optimize for the common case where a file's data blocks, a file's data and metadata, and different files from the same directory are accessed together.
- **Reserved space:** FFS reserves some fraction of the disk's space (e.g., 10%) and presents a slightly reduced disk size to applications.
 - When disk is full, there's little opportunity for file system to optimize locality.
 - Sacrifices a little disk capacity for better locality thus reduced seek times.

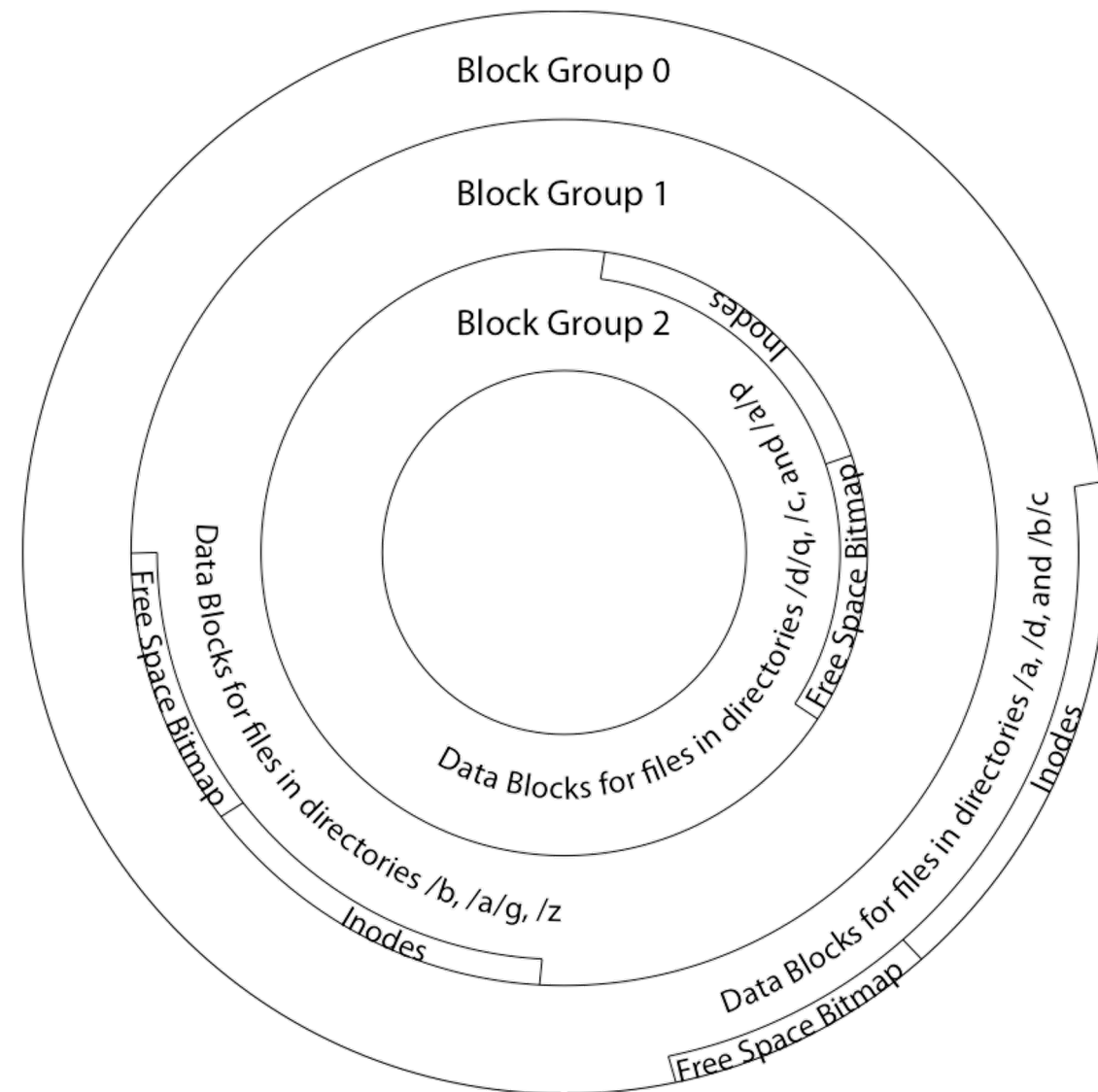
Block Group Placement

- File system volume is divided into a set of block groups
 - Small seek time
- Data blocks, metadata, and free space are distributed to different block
 - Avoid huge seeks between user data and system structure



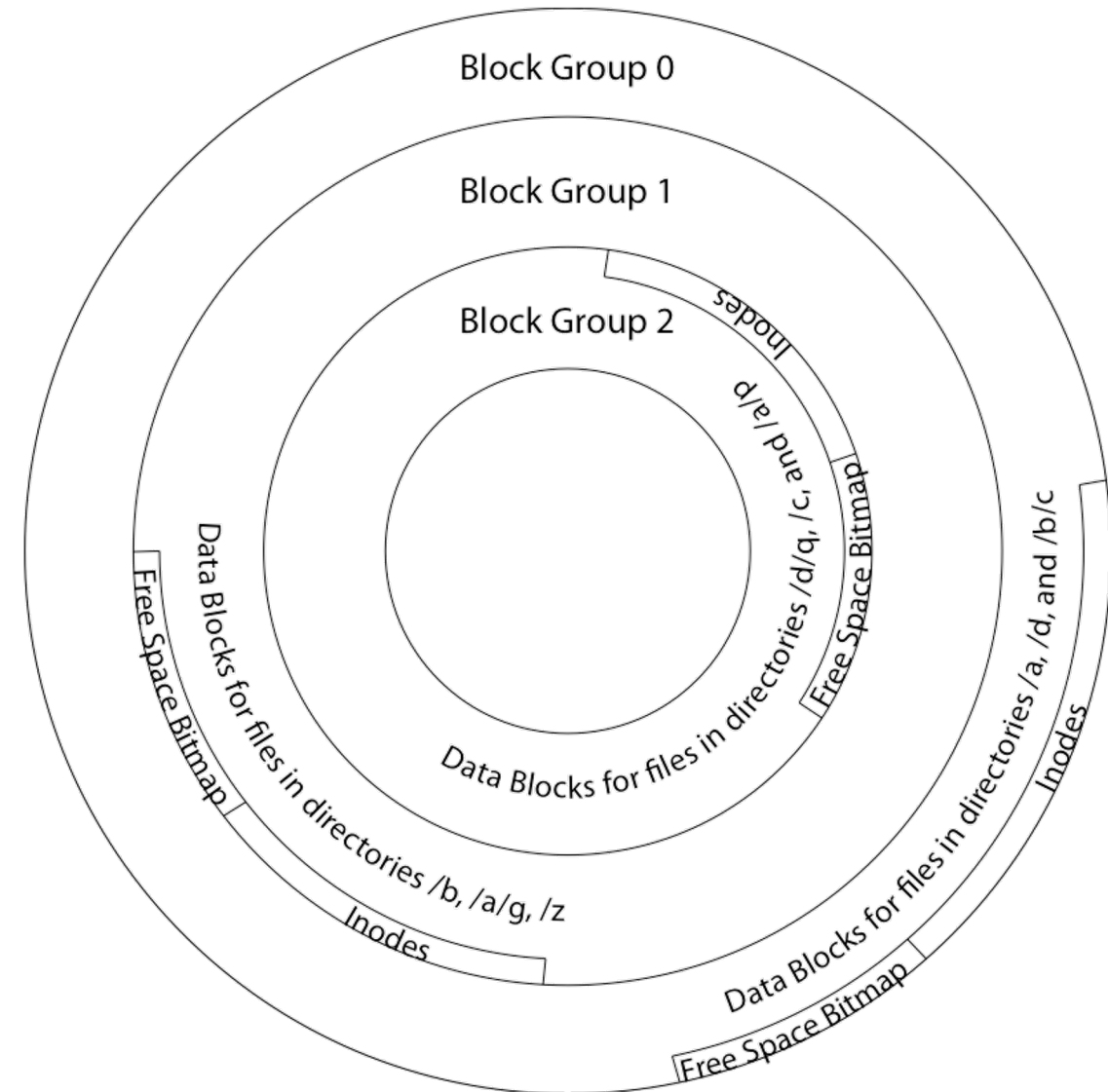
Block Group Placement

- Files in the same directory are placed in the same block group
 - The same for the “directory file” as well
 - i.e., when a new file is created, find an inode number within the block where its directory resides and give it to the file.
 - ❑ Unless there's no free inode number in that block
- But don't put the directory and its sub-directory together
 - Though they might have locality, it will easily fill the block.

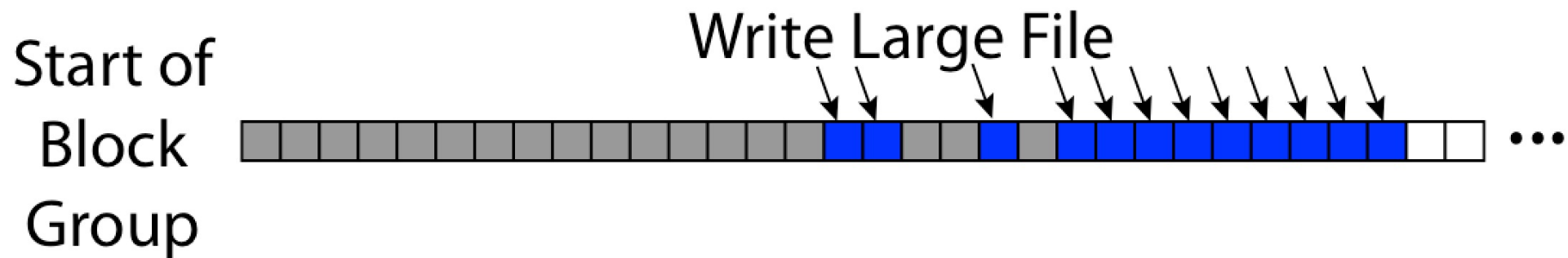
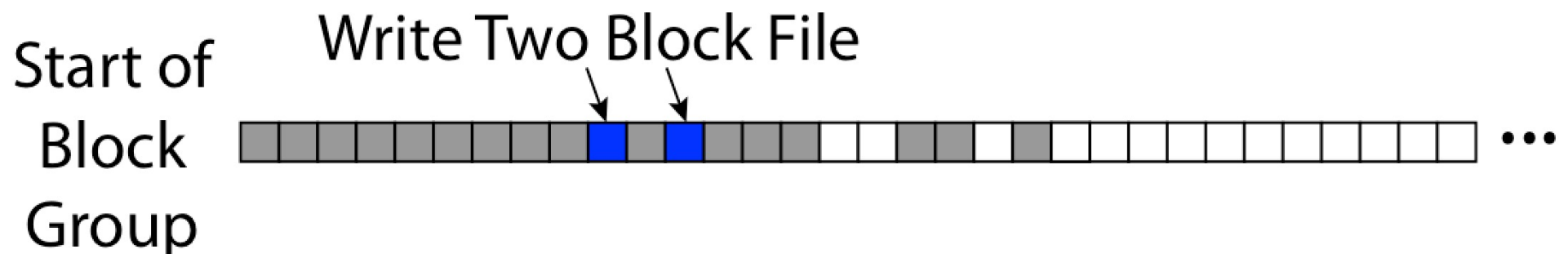
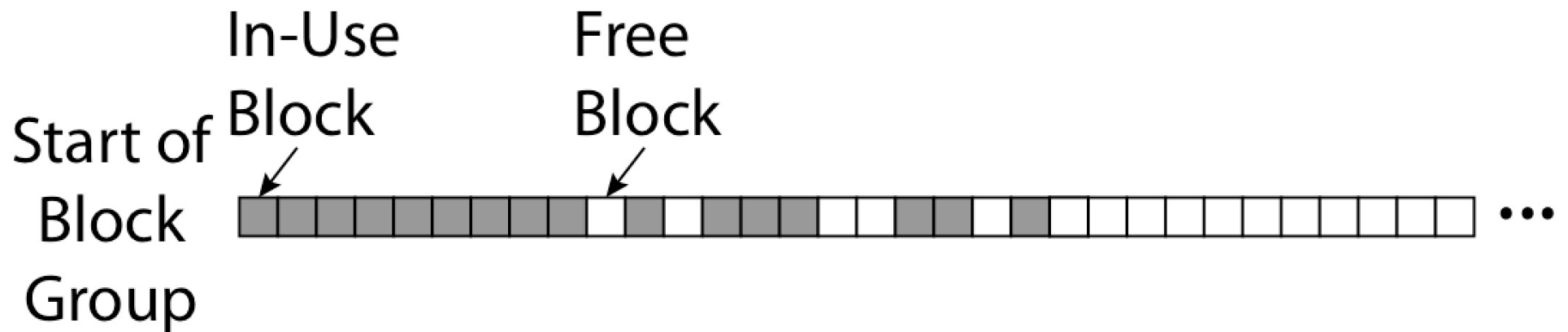


Block Group Placement

- First-Free allocation of new file blocks
 - To expand file, first try successive blocks in bitmap, then choose new range of blocks
 - Few little holes at start, big sequential runs at end of group
 - Avoids fragmentation
 - Sequential layout for big files



UNIX 4.2 BSD FFS First Fit Block Allocation



FFS Summary

- Pros
 - Efficient storage for both small and large files
 - Locality for both small and large files
 - Locality for metadata and data
 - No defragmentation necessary!

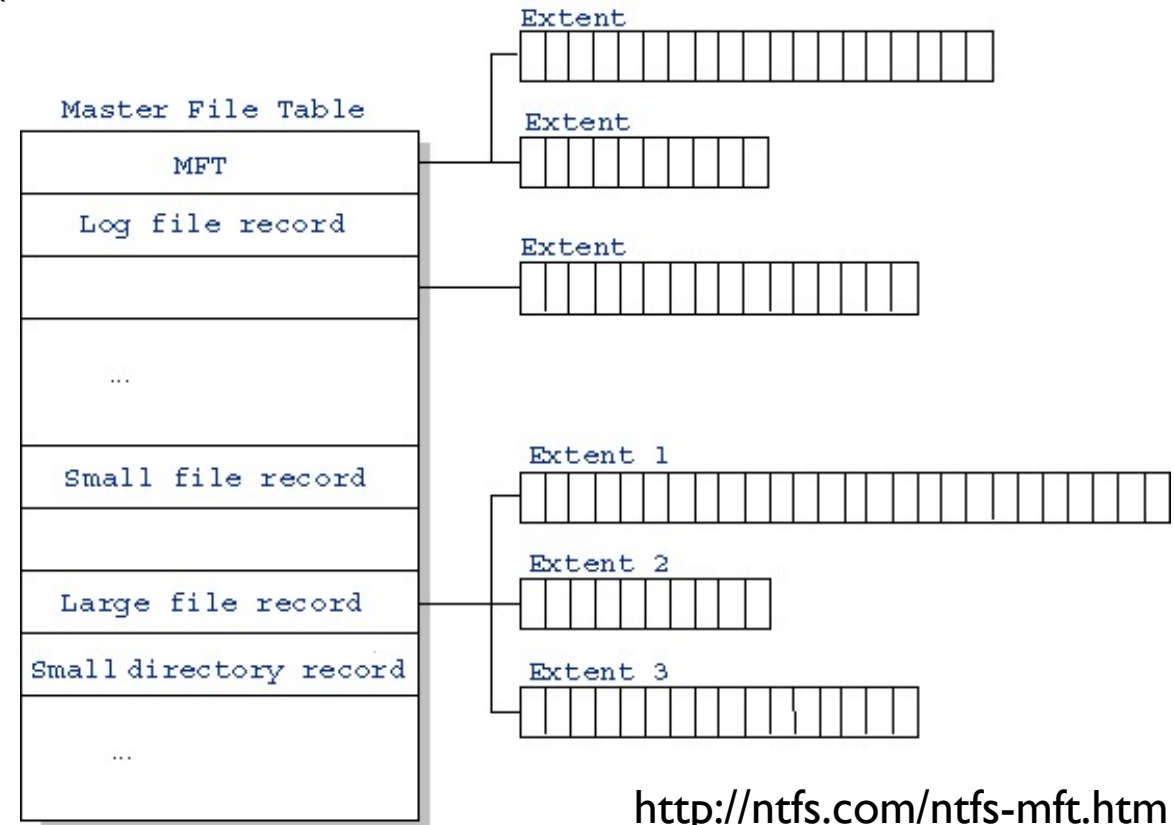
- Cons
 - Inefficient for tiny files (a 1 byte file requires both an inode and a data block)
 - Inefficient encoding when file is mostly contiguous on disk
 - Need to reserve 10-20% of free space to prevent fragmentation

NTFS

- New Technology File System (NTFS)
 - Default on Microsoft Windows systems
- Variable length extents
 - Rather than fixed blocks
- Everything (almost) is a sequence of <attribute:value> pairs
 - Meta-data and data
- Mix direct and indirect freely
- Directories organized in B-tree structure by default

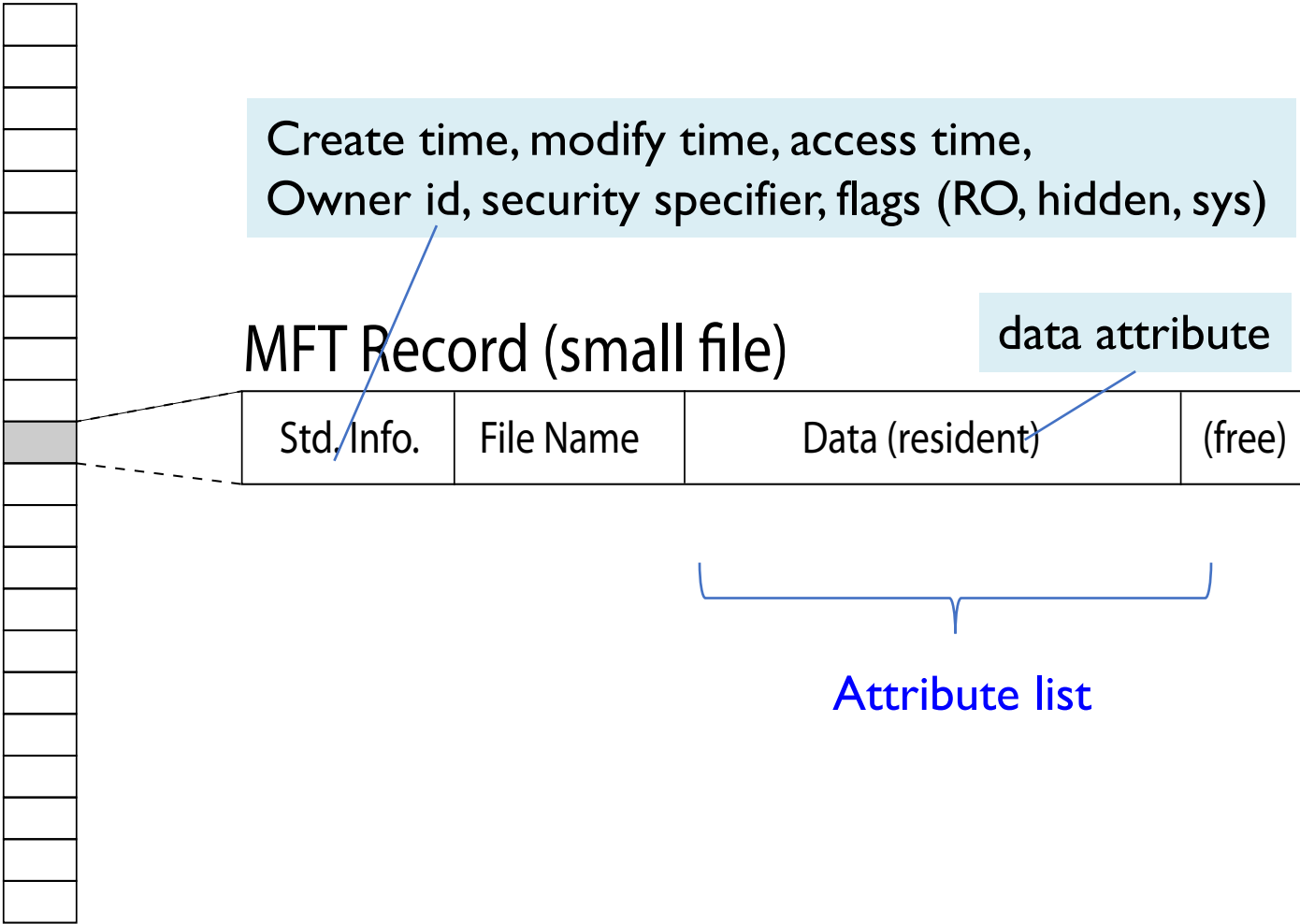
NTFS

- Master File Table
 - Database with Flexible 1 KB entries for metadata/data
 - Variable-sized attribute records (data or metadata)
 - Extend with variable depth tree (non-resident)
- Extents – variable length contiguous regions
 - Block pointers cover runs of blocks
 - Similar approach in Linux (ext4)
 - File create can provide hint as to size of file
- Journaling for reliability
 - Discussed later

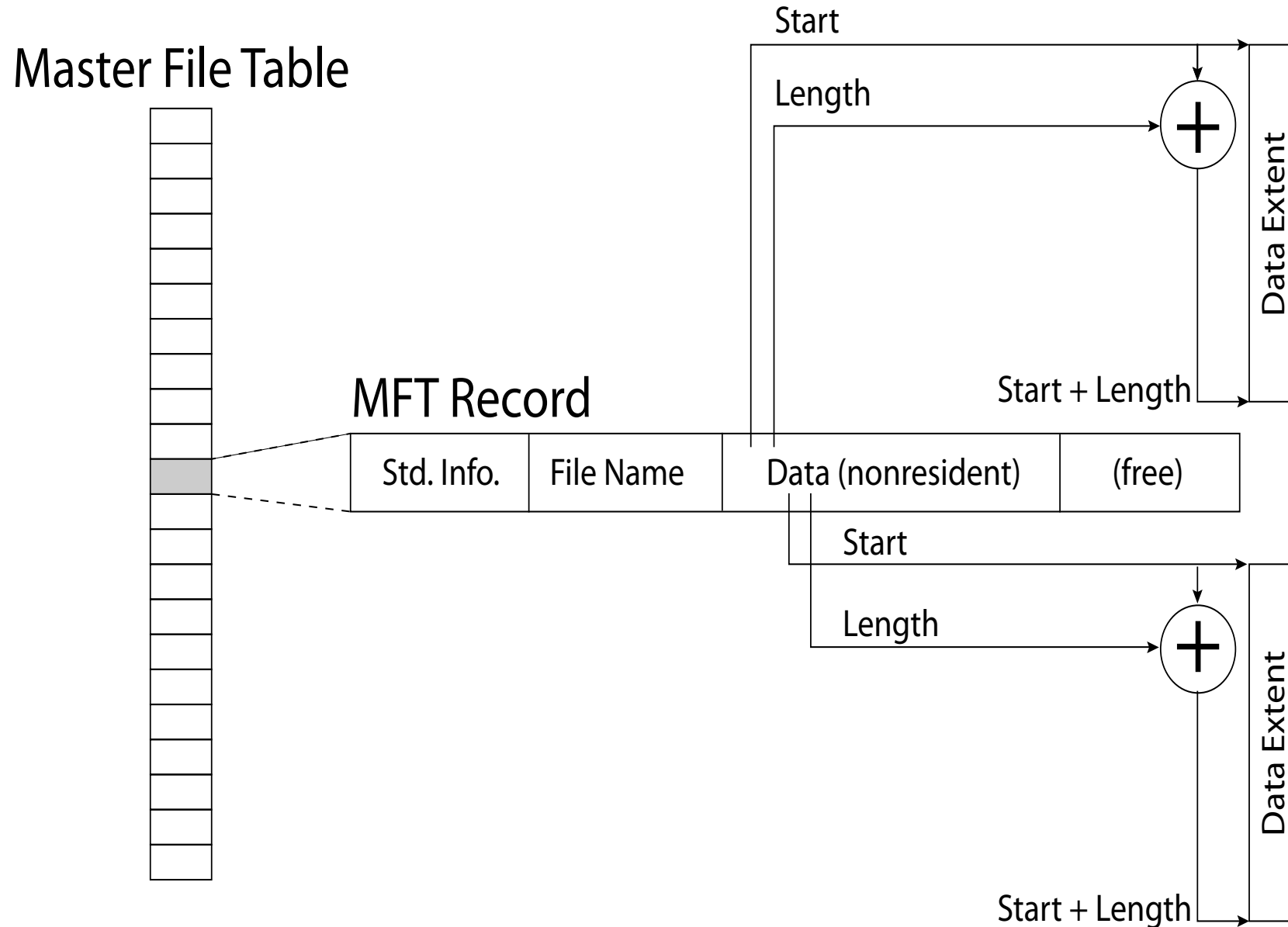


NTFS Small File

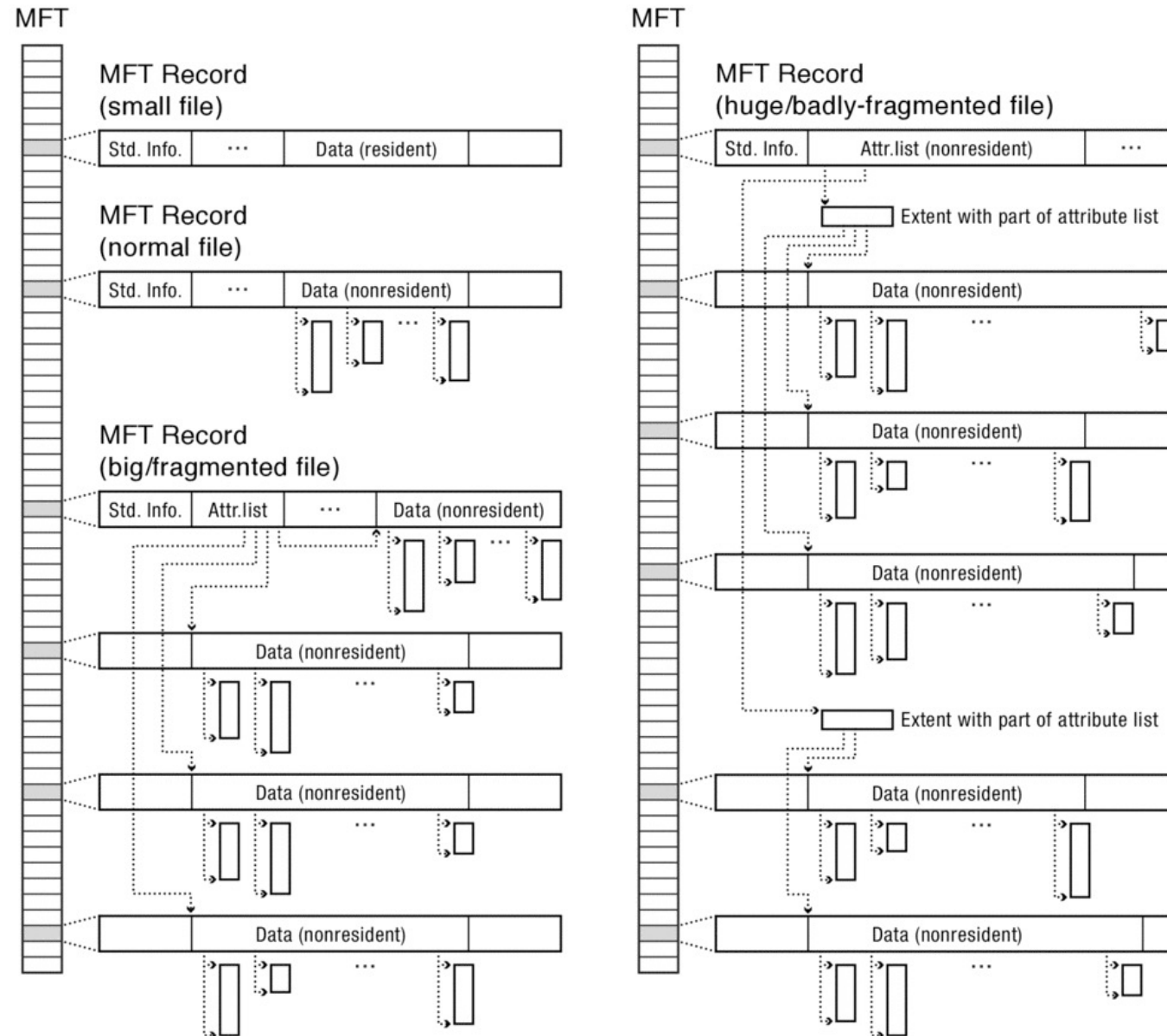
Master File Table



NTFS Medium File



NTFS Multiple Indirect Blocks



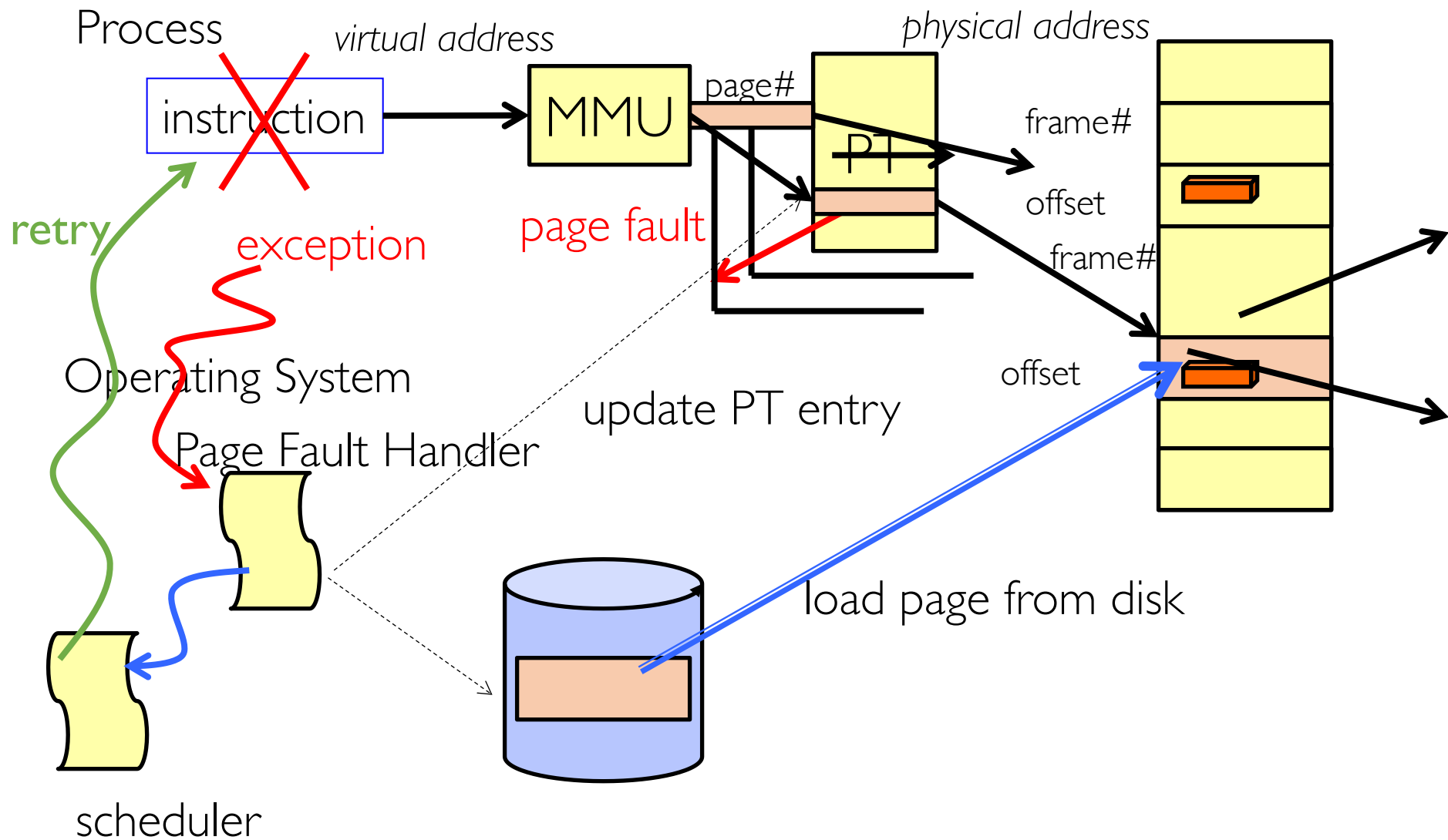
NTFS Locality Heuristics

- Best fit: where the system tries to place a newly allocated file in the smallest free region that is large enough to hold it.
 - In most implementations
- A variation of NTFS: rather than trying to keep the allocation bitmap for the entire disk in memory, the system caches the allocation status for a smaller region of the disk and searches that region first.
 - If the bitmap cache holds information for areas where writes recently occurred, then writes that occur together in time will tend to be clustered together.
- An important NTFS feature: `SetEndOfFile()` to specify the expected size of a file at creation time.
 - Why it is useful?

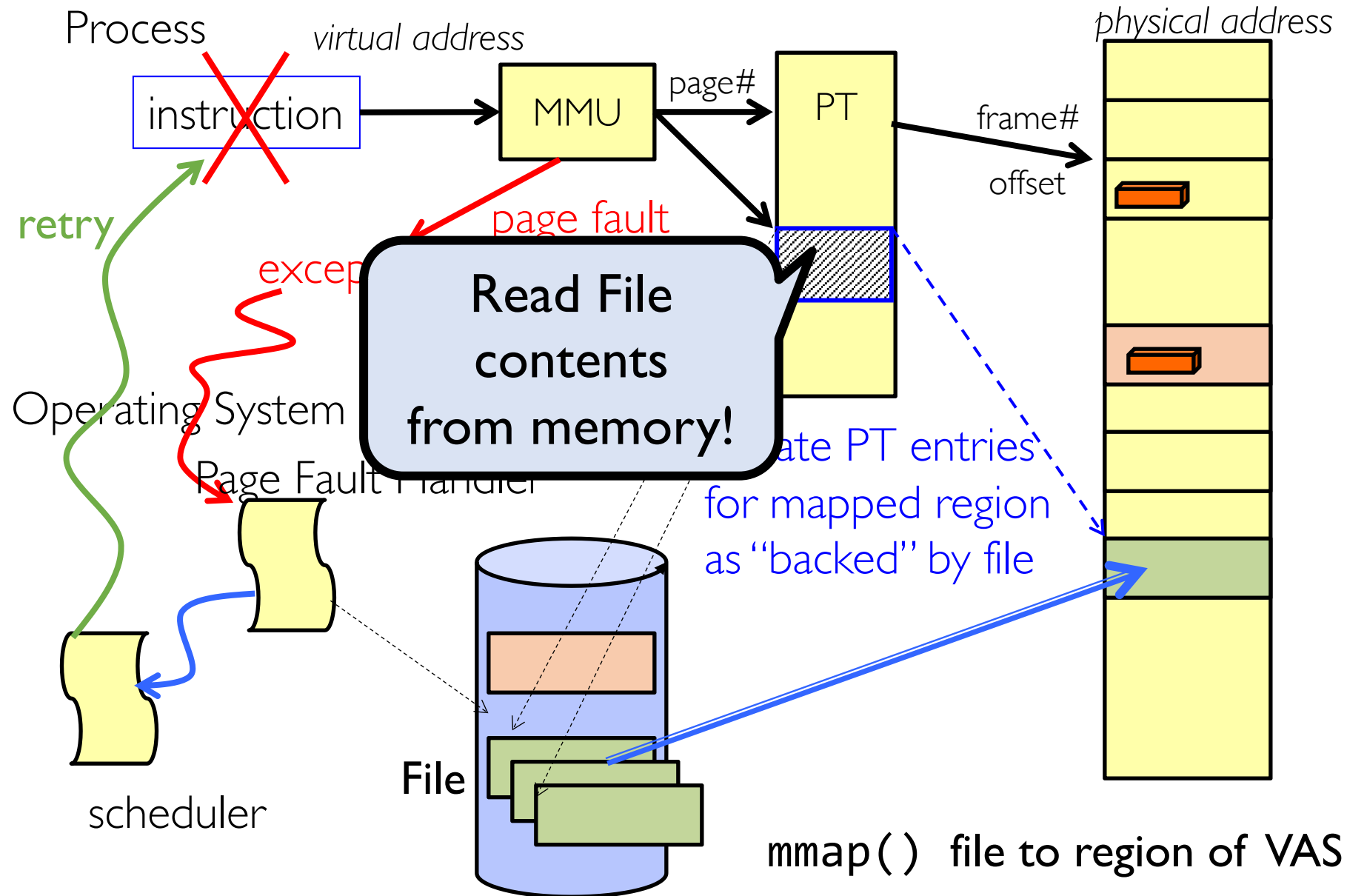
Memory Mapped Files

- Traditional I/O involves explicit transfers between buffers in process address space to/from regions of a file
 - This involves multiple copies into caches in memory, plus system calls
- What if we could “map” the file directly into an empty region of our address space
 - Implicitly “page it in” when we read it
 - Write it and “eventually” page it out
- Executable files are treated this way when we exec the process!!

Recall: Who Does What, When?



Using Paging to mmap() Files



mmap () system call

MMAP(2)	BSD System Calls Manual	MMAP(2)
NAME		
mmap -- allocate memory, or map files or devices into memory		
LIBRARY		
Standard C Library (libc, -lc)		
SYNOPSIS		
#include <sys/mman.h>		
 void * mmap(void *addr, size_t len, int prot, int flags, int fd, off_t offset);		
DESCRIPTION		
The mmap() system call causes the pages starting at <u>addr</u> and continuing for at most <u>len</u> bytes to be mapped from the object described by <u>fd</u> , starting at byte offset <u>offset</u> . If <u>offset</u> or <u>len</u> is not a multiple of the page size, the mapped region may extend past the specified range.		

- May map a specific region or let the system find one for you
 - Tricky to know where the holes are
- Used both for manipulating files and for sharing between processes

An mmap() Example

```
#include <sys/mman.h> /* also stdio.h, stdlib.h, string.h,fcntl.h,unistd.h */
```

```
int something = 162;
```

```
int main (int argc, char *argv[]) {  
    int myfd;  
    char *mfile;
```

```
    printf("Data at: %16lx\n", (long) something);  
    printf("Heap at : %16lx\n", (long) something);  
    printf("Stack at: %16lx\n", (long) something);
```

```
    /* Open the file */
```

```
    myfd = open(argv[1], O_RDWR | O_CREAT, 0666);  
    if (myfd < 0) { perror("open failed!"); return 1; }
```

```
    /* map the file */
```

```
    mfile = mmap(0, 10000, PROT_READ|PROT_WRITE, MAP_SHARED, myfd, 0);  
    if (mfile == MAP_FAILED) { perror("mmap failed!"); return 1; }
```

```
    printf("mmap at : %16lx\n", (long) something);
```

```
    puts(mfile);
```

```
    strcpy(mfile+20, "Let's write over it");
```

```
    close(myfd);
```

```
    return 0;
```

```
}
```

```
$ ./mmap test
```

```
Data at:          105d63058
```

```
Heap at :          7f8a33c04b70
```

```
Stack at:          7fff59e9db10
```

```
mmap at :          105d97000
```

```
This is line one
```

```
This is line two
```

```
This is line three
```

```
This is line four
```

```
$ cat test
```

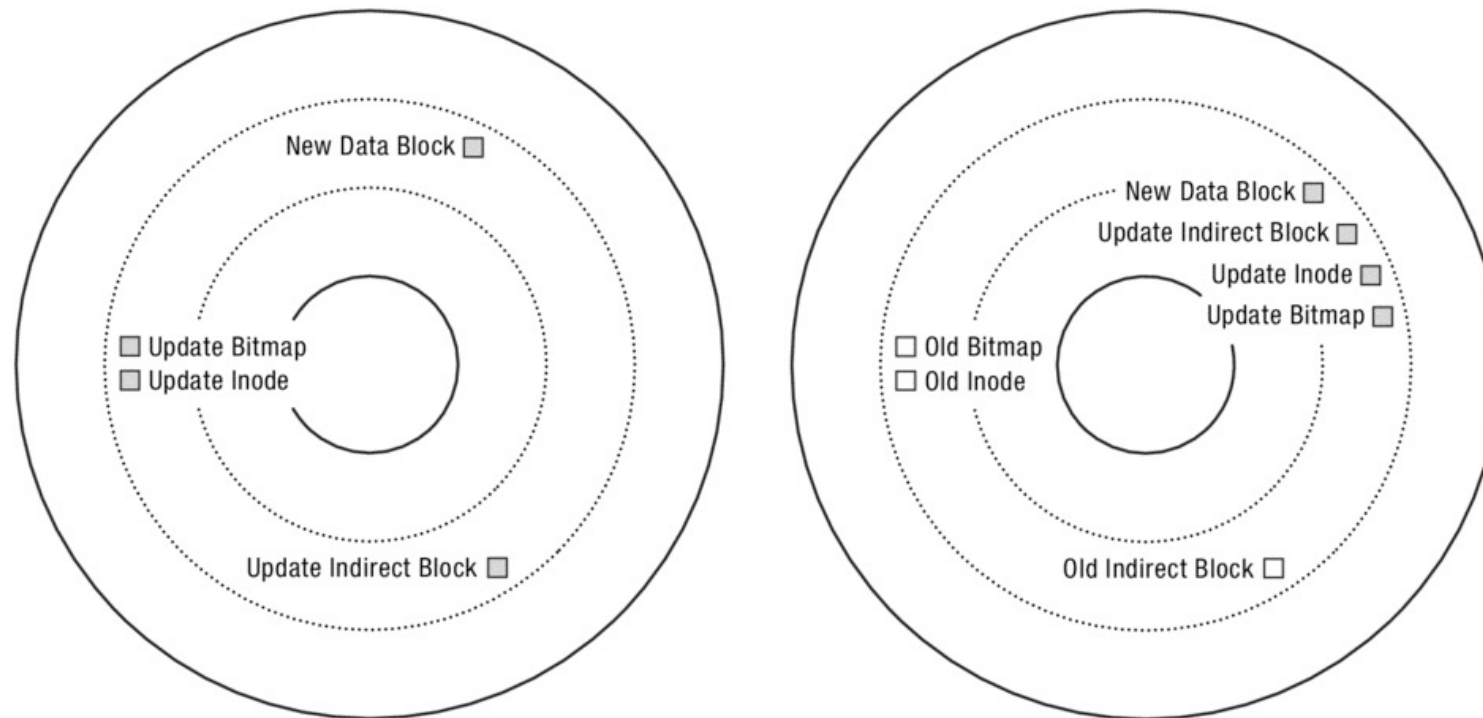
```
This is line one
```

```
This is line two  
Let's write over it  
This is line three
```

```
This is line four
```

Other File Systems..

- Copy-on-write (COW) file system: when updating an existing file, it does not overwrite the existing data or metadata; instead, it writes new versions to new locations
 - Turning random I/O updates to sequential ones.

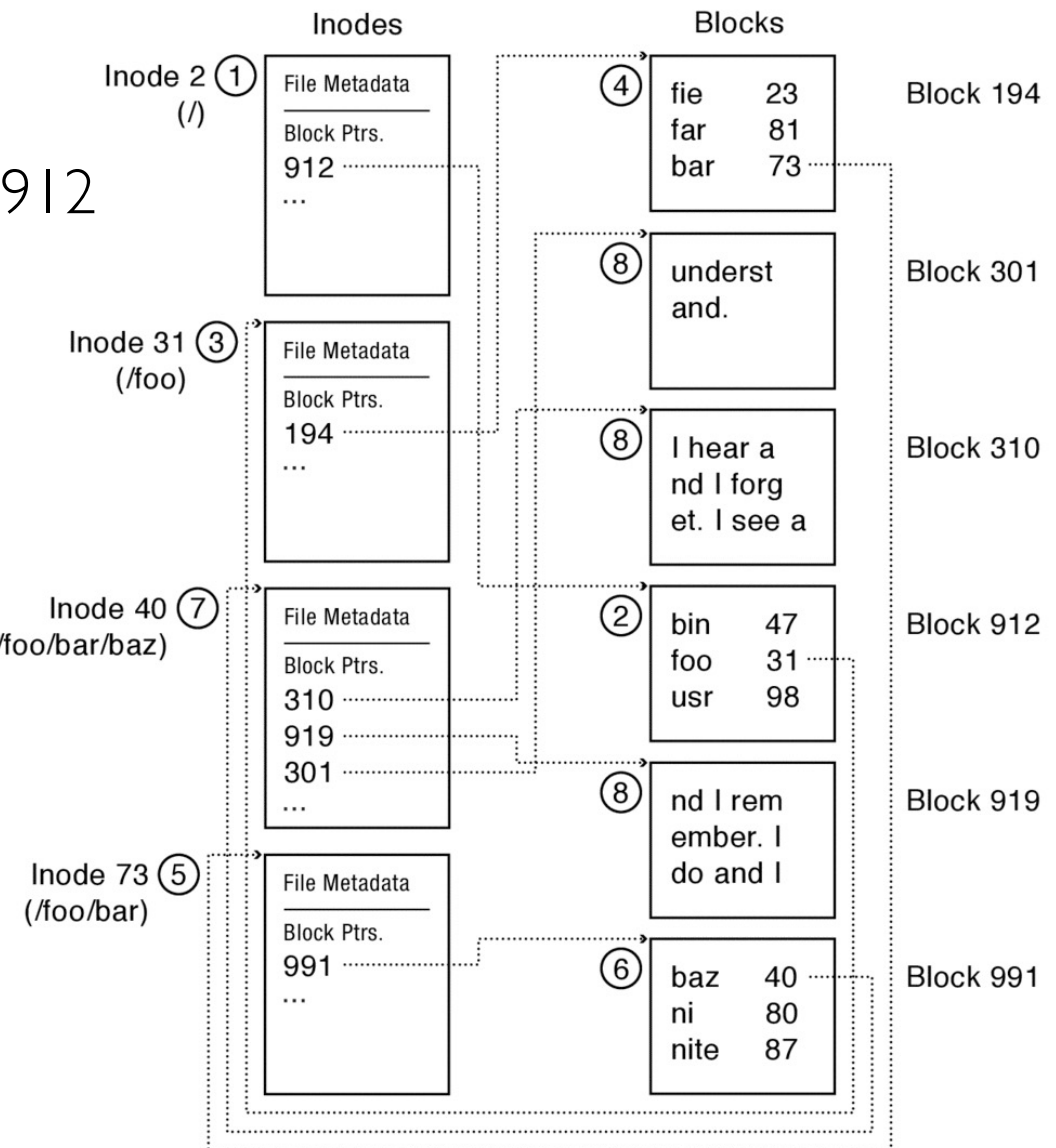


Read textbook for more information!

Put All Things Together (FFS)

Example: read the file /foo/bar/baz

1. Read "/" root inode #2's inode, get block #912
2. From block #912, get inode #31 for "foo"
3. From inode #31, get block #194
4. From block #194, get inode #73 for "bar"
5. From inode #73, get block #991
6. From block #991, get inode #40 for "baz"
7. From inode #40, get 3 data blocks
 - Block 310
 - Block 919
 - Block 301



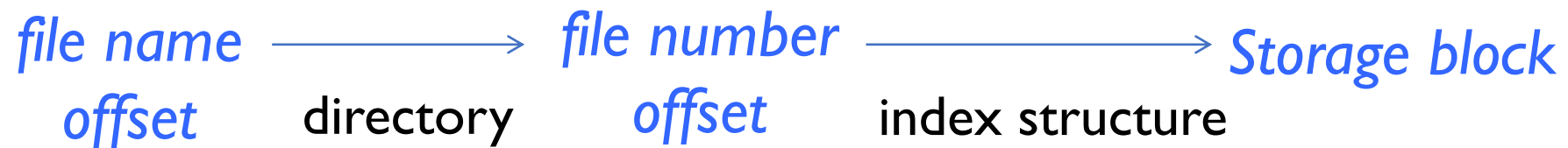
How Many Inodes in Linux

- For 32-bit inode number, it's 2^{32} (about 4 billions)
 - Max
- It's also configurable in many file systems
- *Out of inode error..*

```
echo:homepage echo$ df -i
Filesystem      512-blocks      Used Available Capacity  iused      ifree %iused
/dev/disk1s1    1953595632    21968928 991671656      3%   488378 9767489782      0%
devfs           387           387           0    100%      678           0    100%
/dev/disk1s2    1953595632    934163472 991671656     49%  4233888 9763744272      0%
/dev/disk1s5    1953595632     4194344 991671656      1%           2 9767978158      0%
map auto_home      0             0           0    100%           0           0    100%
```

Goals for Today

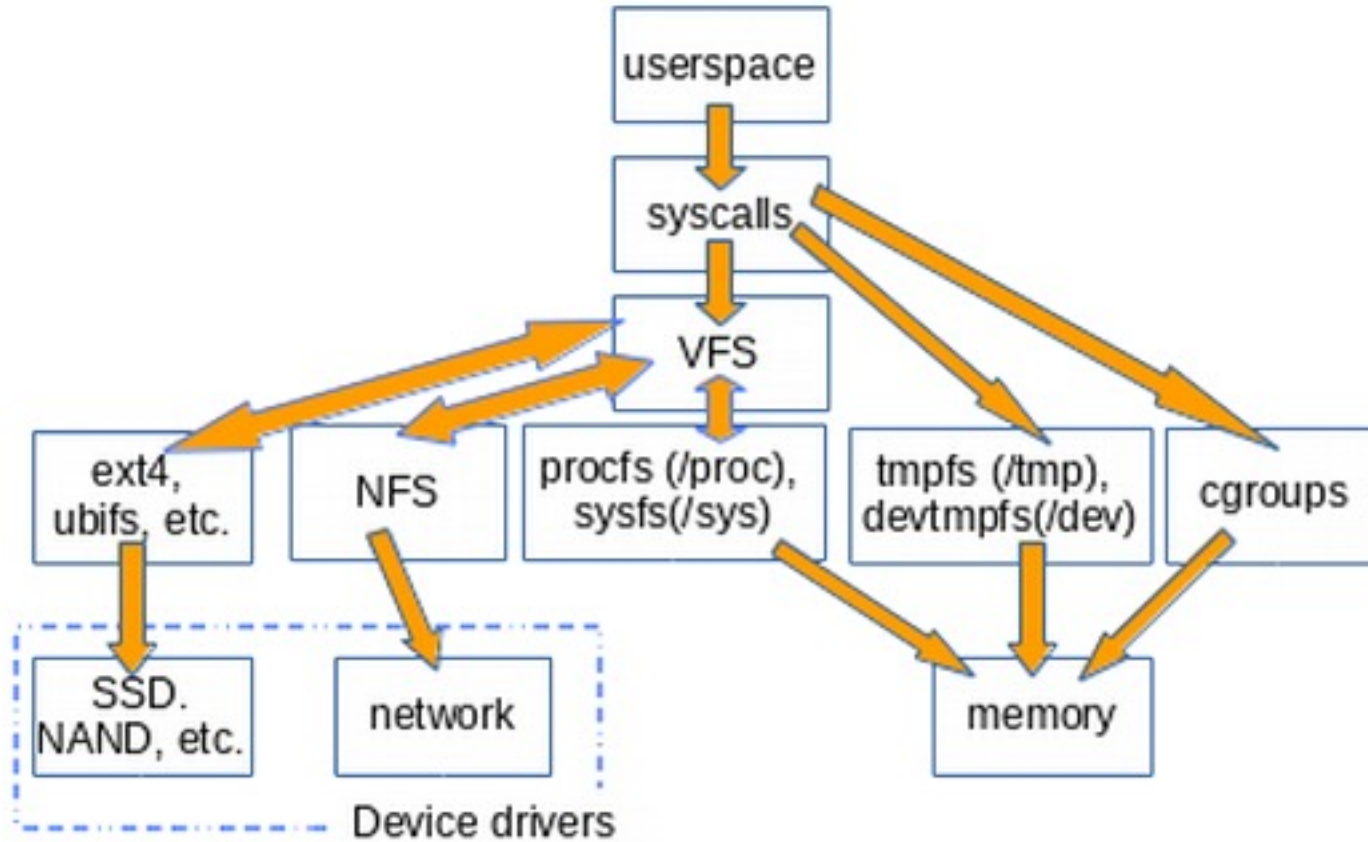
- Directories: naming data
 - How do we convert a file name to the file number?
- Files: finding data
 - How do we locate storage block based on file number?
- Virtual file systems (VFS)
 - How do we make different FSs work together easily?



History

- Early OSes provided a single file system
 - In general, system was tailored to target hardware
- people became interested in supporting more than one file system type on a single system
 - Any guesses why?
 - Networked file systems: sharing parts of a file system across a network of workstations

Virtual File System (VFS)



Modern VFS

- Dozens of supported file systems
 - Allows new features and designs transparent to apps
 - Interoperability with removable media and other Oses
- Independent layer from backing storage
 - In-memory file systems (ramdisks)
 - Pseudo file systems used for configuration
 - ❑ (/proc, /devtmps...) only backed by kernel data structures
- And, of course, networked file system support

What the VFS Does

- The VFS is a substantial piece of code
 - not just an API wrapper
- Caches file system metadata (e.g., names, attributes)
 - Coordinates data caching with the page cache
- Enforces a common access control model
- Implements complex, common routines
 - path lookup
 - opening files
 - file handle management

User's Perspective

- Single programming interface
 - (POSIX file system calls – open, read, write, etc.)
- Single file system tree
 - Remote FS can be transparently mounted (e.g., at /home)
- Alternative: Custom library for each file system
 - Much more trouble for the programmer

FS Developer's Perspective

- FS developer responsible for implementing standard objects/functions called by the VFS
 - Primarily populating in-memory objects
 - Typically from stable storage
 - Sometimes writing them back
- Can use block device interfaces to schedule disk I/O
 - And page cache functions
 - And some VFS helpers
- Analogous to implementing Java abstract classes

High-level FS dev. tasks

- Translate between VFS objects and backing storage (whether device, remote system, or other/none)
 - Potentially includes requesting I/O
- Read and write file pages
- VFS doesn't prescribe all aspects of FS design
 - More of a lowest common denominator
- Opportunities: (to name a few)
 - More optimal media usage/scheduling
 - Varying on-disk consistency guarantees
 - Features (e.g., encryption, virus scanning, snapshotting)

Core VFS Abstractions

- **super block:** FS-global data
 - Early/many file systems put this as first block of partition
- **inode:** (index node): metadata for one file
- **dentry:** (directory entry): name to inode mapping
- **file object:** pointer to dentry and cursor (file offset)
- SB and inodes are extended by file system developer

Embedded Inodes

- Many FSes embed VFS inode in FS-specific inode

```
struct myfs_inode {  
    int ondisk_blocks[];  
    /* other stuff */  
    struct inode vfs_inode;  
}
```

- Why? Finding the low-level from inode is simple
 - Compiler translates references to simple math

<https://compas.cs.stonybrook.edu/~nhonarmand/courses/fall4/cse506.2/slides/vfs.pdf>

File System Summary (1/2)

- File System:
 - Transforms blocks into Files and Directories
 - Optimize for size, access and usage patterns
 - Maximize sequential access, allow efficient random access
 - Projects the OS protection and security regime (UGO vs ACL)
- File defined by header, called “inode”
- Naming: translating from user-visible names to actual sys resources
 - Directories used for naming for local file systems
 - Linked or tree structure stored in files
- Multilevel Indexed Scheme
 - inode contains file info, direct pointers to blocks, indirect blocks, doubly indirect, etc..
 - NTFS: variable extents not fixed blocks, tiny files data is in header

File System Summary (2/2)

- 4.2 BSD Multilevel index files
 - Inode contains ptrs to actual blocks, indirect blocks, double indirect blocks, etc.
 - Optimizations for sequential access: start new files in open ranges of free blocks, rotational optimization
- File layout driven by freespace management
 - Integrate freespace, inode table, file blocks and dirs into block group
- Deep interactions between mem management, file system, sharing
 - `mmap()`: map file or anonymous segment to memory