

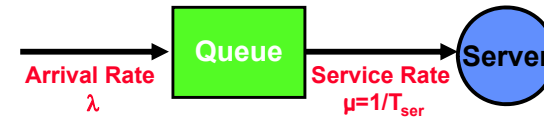
# CS162 Operating Systems and Systems Programming Lecture 19

## File Systems (Con't), MMAP, Buffer Cache

April 11<sup>th</sup>, 2019  
Prof. Ion Stoica  
<http://cs162.eecs.Berkeley.edu>

## Recall: A Little Queuing Theory: Some Results

- Assumptions:
  - System in equilibrium; No limit to the queue
  - Time between successive arrivals is random and memoryless



- Parameters that describe our system:
  - $\lambda$ : mean number of arriving customers/second
  - $T_{ser}$ : mean time to service a customer ("m1")
  - $C$ : squared coefficient of variance =  $\sigma^2/m^2$
  - $\mu$ : service rate =  $1/T_{ser}$
  - $u$ : server utilization ( $0 \leq u \leq 1$ ):  $u = \lambda/\mu = \lambda \times T_{ser}$
- Parameters we wish to compute:
  - $T_q$ : Time spent in queue
  - $L_q$ : Length of queue =  $\lambda \times T_q$  (by Little's law)
- Results:
  - Memoryless service distribution ( $C = 1$ ):
    - Called M/M/1 queue:  $T_q = T_{ser} \times u / (1 - u)$
  - General service distribution (no restrictions), 1 server:
    - Called M/G/1 queue:  $T_q = T_{ser} \times \frac{1}{2}(1+C) \times u / (1 - u)$

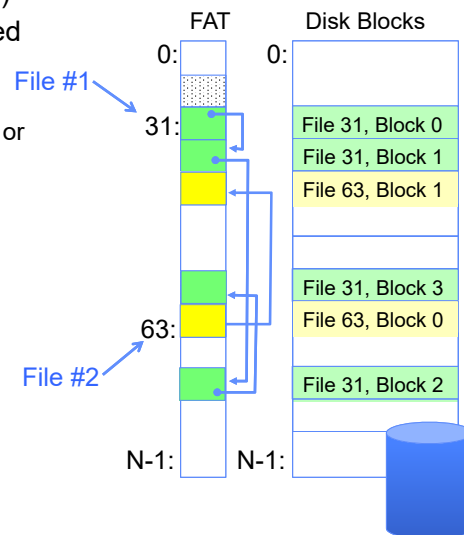
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.2

## Recall: FAT Properties

- File is collection of disk blocks (Microsoft calls them "clusters")
- FAT is array of integers mapped 1-1 with disk blocks
  - Each integer is either:
    - » Pointer to next block in file; or
    - » End of file flag; or
    - » Free block flag
- File Number is index of root of block list for the file
  - Follow list to get block #
  - Directory must map name to block number at start of file
- But: Where is FAT stored?
  - Beginning of disk, before the data blocks
  - Usually 2 copies (to handle errors)



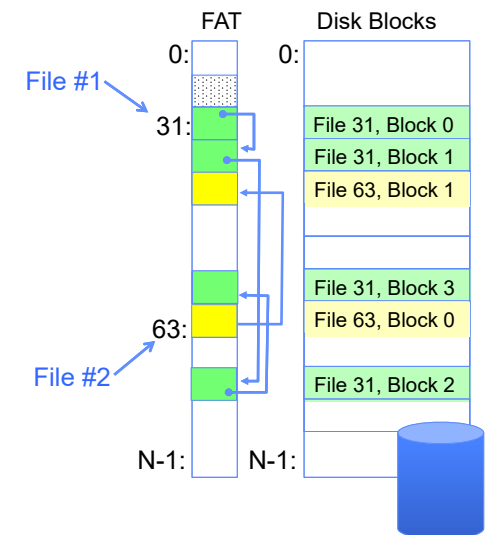
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.3

## FAT Assessment – Issues

- Time to find block (large files) ??
- Block layout for file ???
- Sequential Access ???
- Random Access ???
- Fragmentation ???
  - MSDOS defrag tool
- Small files ???
- Big files ???



4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.4

## What about the Directory?

- Essentially a file containing  
    <file\_name: file\_number> mappings
- Free space for new/deleted entries
- In FAT: file attributes are kept in directory (!!!)
- Each directory a linked list of entries
- Where do you find root directory ( "/" )?
  - At well-defined place on disk
  - For FAT, this is at block 2 (there are no blocks 0 or 1)
  - Remaining directories

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.5

## Directory Structure (cont'd)

- 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")

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.6

## Many Huge FAT Security Holes!

- FAT has no access rights
  - No way, even in principle, to track ownership of data
- FAT has no header in the file blocks
  - No way to enforce control over data, since all processes have access of FAT table
  - Just follow pointer to disk blocks
- Just gives an index into the FAT to read data
  - (file number = block number)
  - Could start in middle of file or access deleted data

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.7

## Characteristics of Files

### A Five-Year Study of File-System Metadata

NITIN AGRAWAL  
University of Wisconsin, Madison  
and

WILLIAM J. BOLOSKY, JOHN R. DOUCEUR, and JACOB R. LORCH  
Microsoft Research

9.9

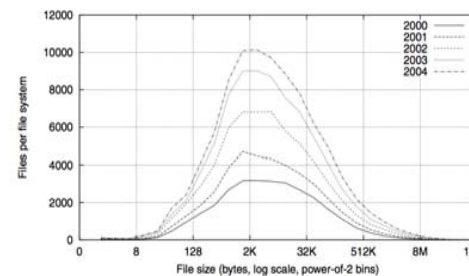


Fig. 2. Histograms of files by size.

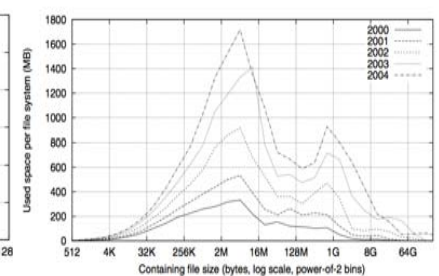


Fig. 4. Histograms of bytes by containing file size.

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.8

## Characteristics of Files

- Most files are small, growing numbers of files

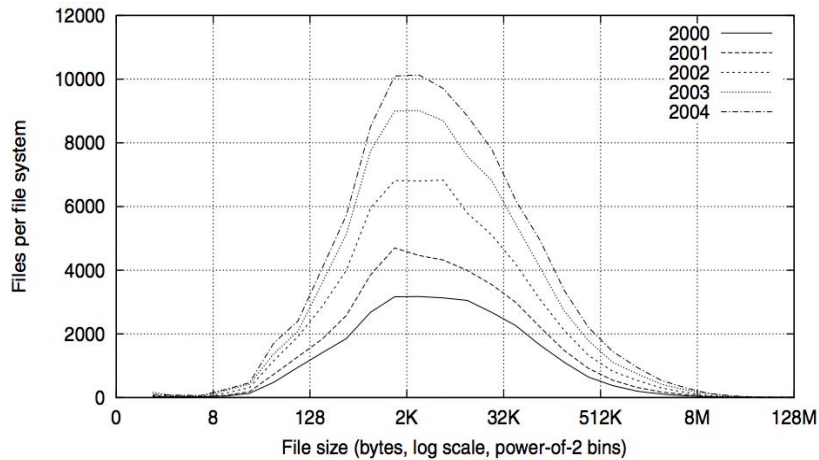


Fig. 2. Histograms of files by size.

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.9

## Characteristics of Files

- Most of the space is occupied by the rare big ones

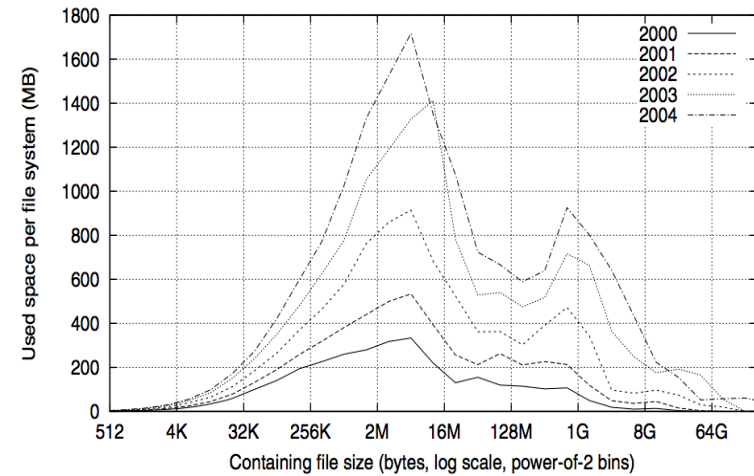


Fig. 4. Histograms of bytes by containing file size.

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.10

## Unix File System (1/2)

- Original inode format appeared in BSD 4.1
  - Berkeley Standard Distribution Unix
  - Part of your heritage!
  - 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

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.11

## 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

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.12

## Administrivia

- Watch for group evaluations
  - We are asking for additional information about both project 1 and project 2. Tell about your group dynamics
    - » Not about getting points for yourself, but about evaluating partners
  - Please collaborate with your group! (potentially a zero-sum game if you are not contributing to your group's success!)
- Final Midterm: 5/2
  - All topics up to and including that last lecture before then
  - You can bring 3 sheets of notes (both sides)
- Help to name new planet: <https://2007or10.name>
  - You have until May 10<sup>th</sup>, 2019
  - In the Kuiper Belt, beyond Neptune, smaller than Pluto and Eris
  - Currently named “2007 OR10”
    - » **Gonggong**: A Chinese water god with red hair and a serpent-like tail
    - » **Holle**: A European winter goddess of fertility, rebirth and women
    - » **Vili**: A Nordic deity. Together with his brothers Odin and Ve, he defeated frost giant Ymir and used Ymir's body to create universe

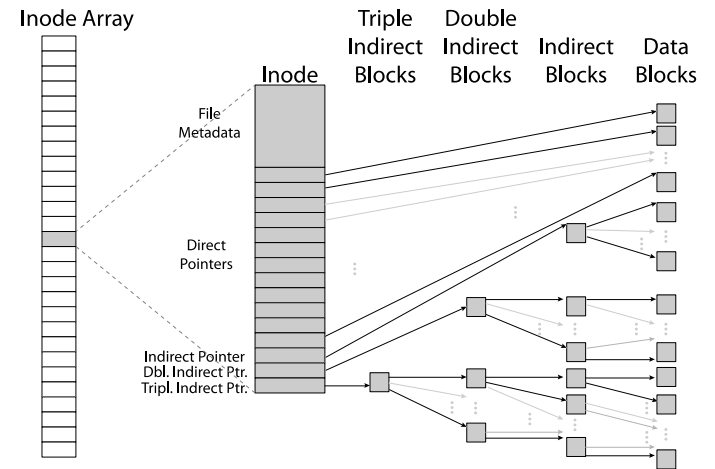
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.13

## Inode Structure

- inode metadata



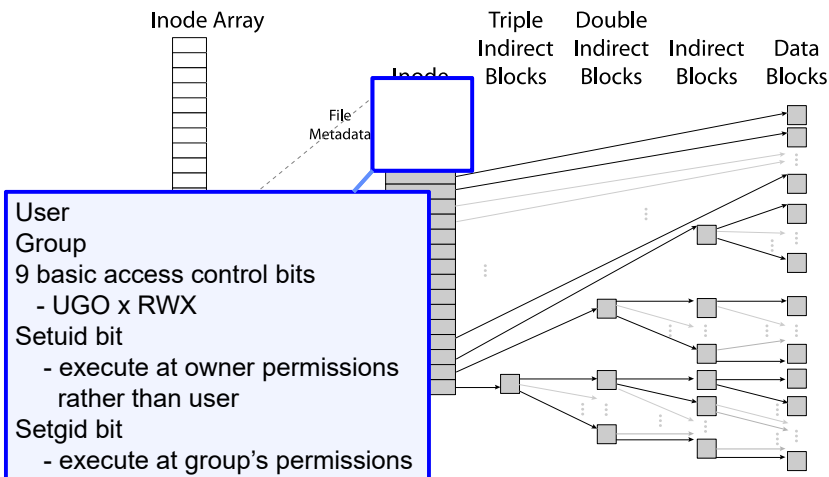
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.14

## File Attributes

- inode metadata



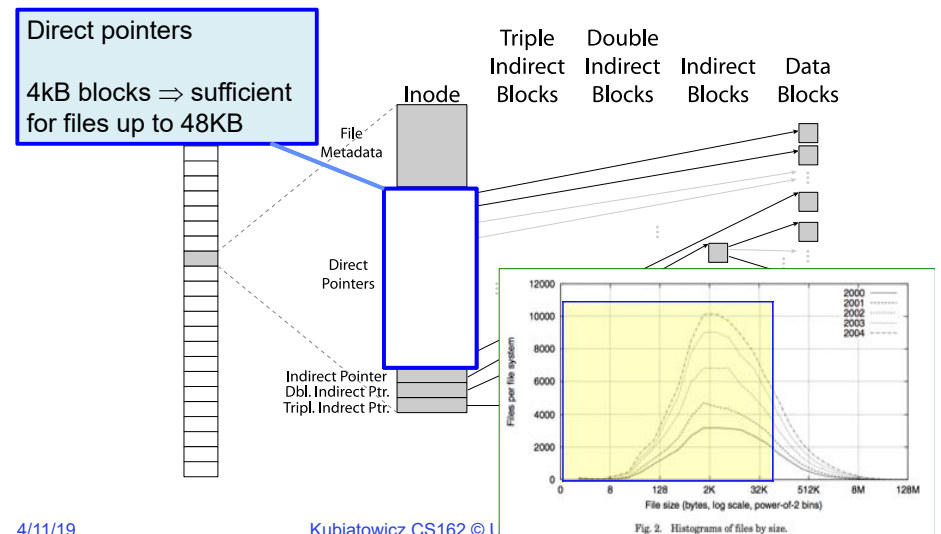
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.15

## Data Storage

- Small files: 12 pointers direct to data blocks



4/11/19

Kubiatowicz CS162 © UCB Spring 2019

## Data Storage

- Large files: 1,2,3 level indirect pointers

Indirect pointers  
 - point to a disk block containing only pointers  
 - 4 kB blocks => 1024 ptrs  
 => 4 MB @ level 2  
 => 4 GB @ level 3  
 => 4 TB @ level 4

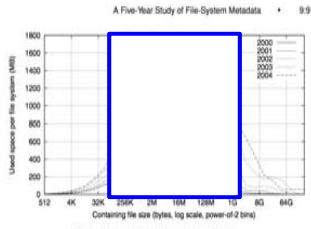
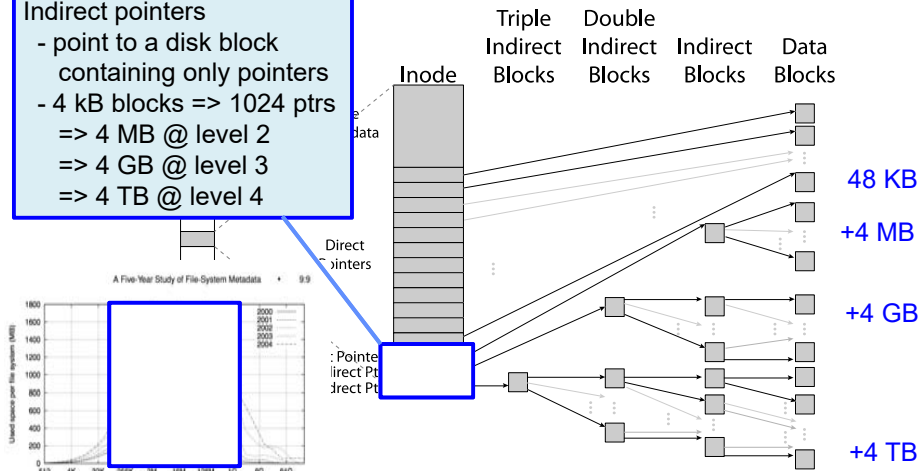


Fig. 4. Histogram of bytes by containing file size.

## UNIX BSD 4.2 (1984) (1/2)

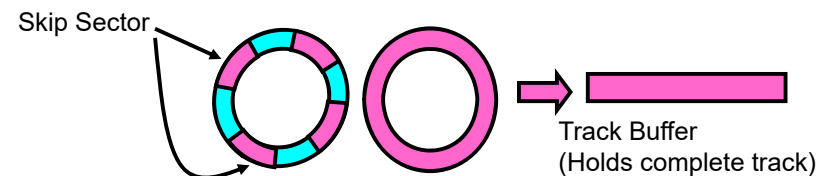
- Same as BSD 4.1 (same file header and triply indirect blocks), except incorporated ideas from Cray Operating System:
  - Uses bitmap allocation in place of freelist
  - Attempt to allocate files contiguously
  - 10% reserved disk space
  - Skip-sector positioning (mentioned later)

## UNIX BSD 4.2 (1984) (2/2)

- Problem: When create a file, don't know how big it will become (in UNIX, most writes are by appending)
  - How much contiguous space do you allocate for a file?
  - In BSD 4.2, just find some range of free blocks
    - » Put each new file at the front of different range
    - » To expand a file, you first try successive blocks in bitmap, then choose new range of blocks
  - Also in BSD 4.2: store files from same directory near each other
- Fast File System (FFS)
  - Allocation and placement policies for BSD 4.2

## Attack of the Rotational Delay

- Problem 2: Missing blocks due to rotational delay
  - Issue: Read one block, do processing, and read next block. In meantime, disk has continued turning: missed next block! Need 1 revolution/block!



- Solution 1: Skip sector positioning (“interleaving”)
  - » Place the blocks from one file on every other block of a track: give time for processing to overlap rotation
  - » Can be done by OS or in modern drives by the disk controller
- Solution 2: Read ahead: read next block right after first, even if application hasn't asked for it yet
  - » This can be done either by OS (read ahead)
  - » By disk itself (track buffers) - many disk controllers have internal RAM that allows them to read a complete track
- Modern disks + controllers do many things “under the covers”
  - Track buffers, elevator algorithms, bad block filtering

## 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"

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.21

## 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 ⇒ 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

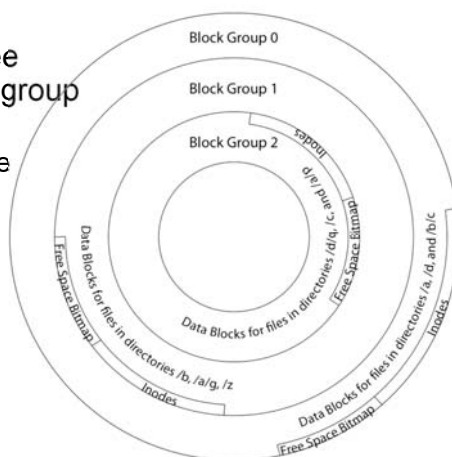
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.22

## 4.2 BSD Locality: Block Groups

- File system volume is divided into a set of block groups
  - Close set of tracks
- Data blocks, metadata, and free space interleaved within block group
  - Avoid huge seeks between user data and system structure
- Put directory and its files in common block group



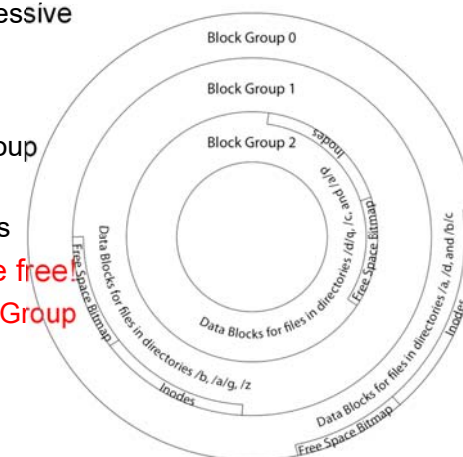
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.23

## 4.2 BSD Locality: Block Groups

- 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
- Important: keep 10% or more free
  - Reserve space in the Block Group



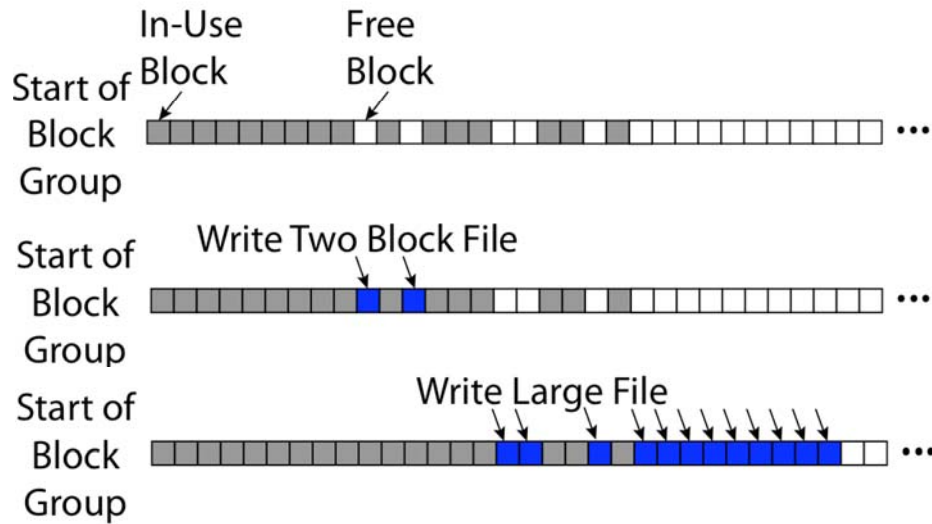
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.24



## UNIX 4.2 BSD FFS First Fit Block Allocation



4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.25

## UNIX 4.2 BSD FFS

- 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

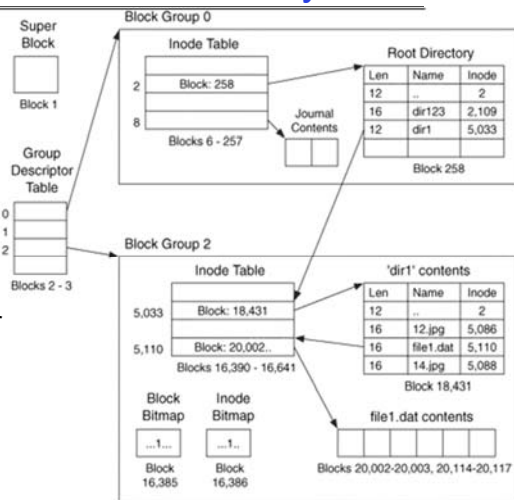
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.26

## Linux Example: Ext2/3 Disk Layout

- Disk divided into block groups
  - Provides locality
  - Each group has two block-sized bitmaps (free blocks/inodes)
  - Block sizes settable at format time: 1K, 2K, 4K, 8K...
- Actual inode structure similar to 4.2 BSD
  - with 12 direct pointers
- Ext3: Ext2 with Journaling
  - Several degrees of protection with comparable overhead



- Example: create a file1.dat under /dir1/ in Ext3

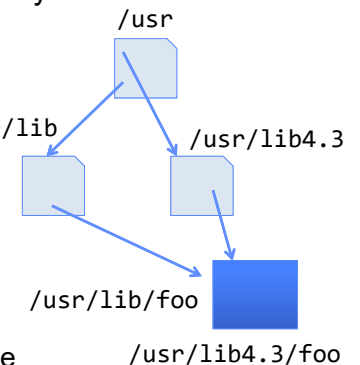
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.27

## A bit more on directories

- 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)



4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.28

## Links

- Hard link
  - Sets another directory entry to contain the file number for the file
  - Creates another name (path) for the file
  - Each is “first class”
- Soft link or Symbolic Link or Shortcut
  - Directory entry contains the path and name of the file
  - Map one name to another name

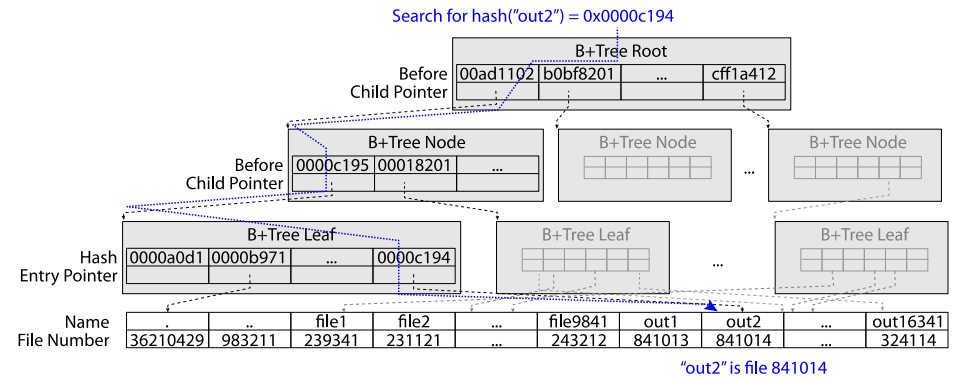
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.29

## Large Directories: B-Trees (dirhash)

in FreeBSD, NetBSD, OpenBSD



4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.30

## 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

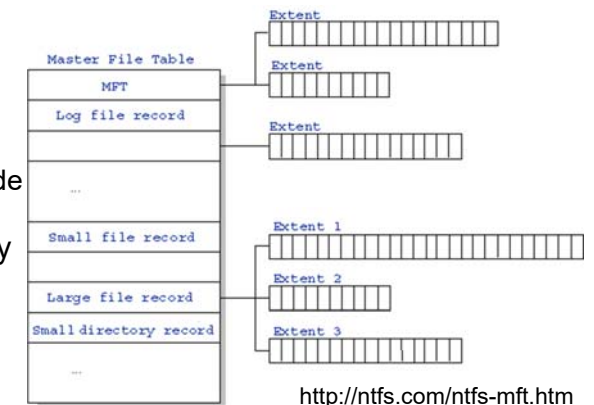
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.31

## NTFS

- Master File Table
  - Database with Flexible 1KB 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



4/11/19

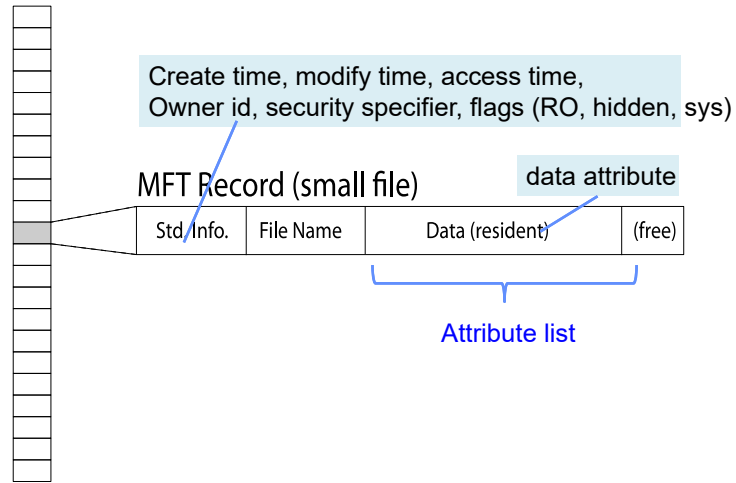
Kubiatowicz CS162 © UCB Spring 2019

Lec 19.32



## NTFS Small File

Master File Table



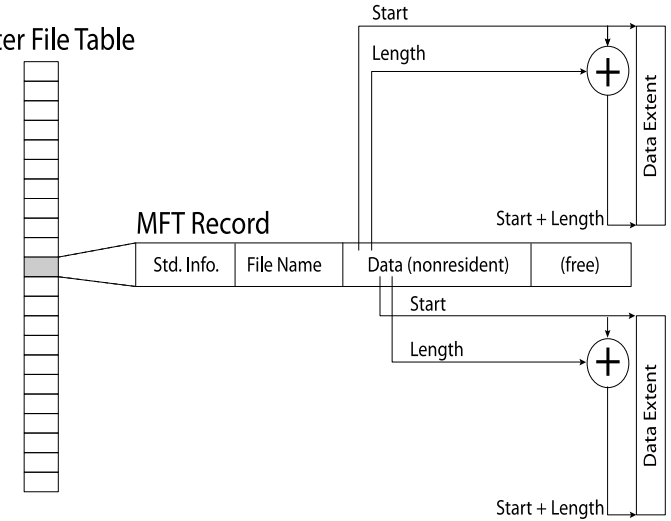
4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.33

## NTFS Medium File

Master File Table

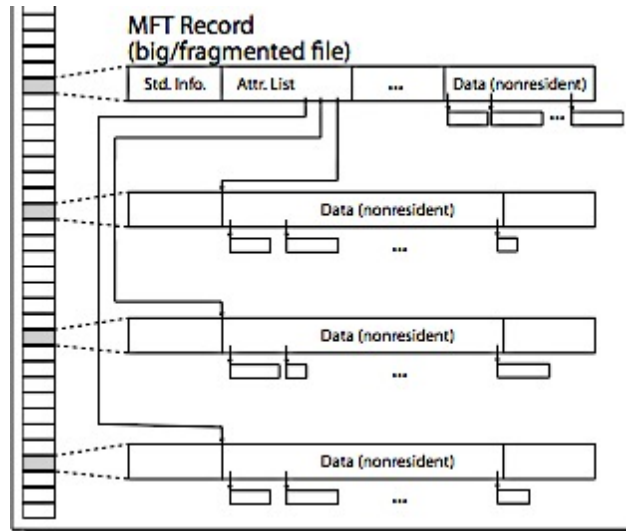


4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.34

## NTFS Multiple Indirect Blocks

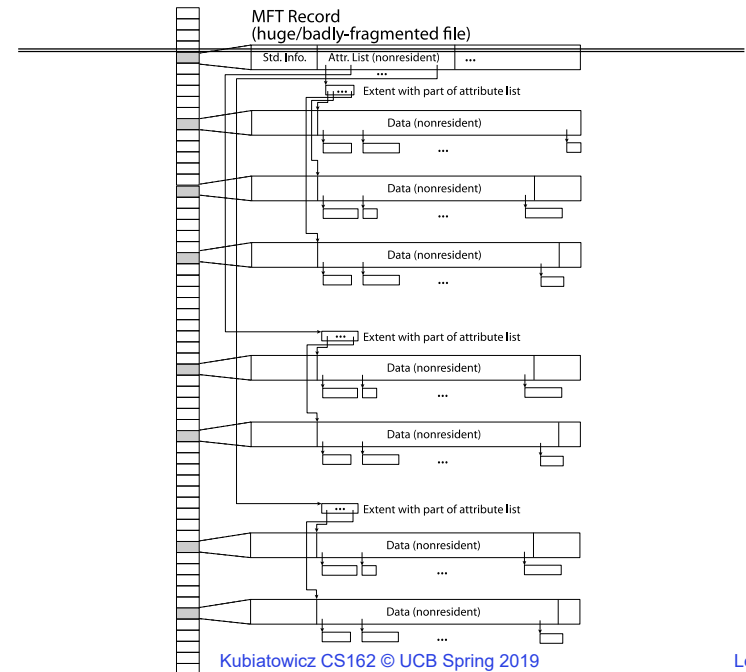


4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.35

Master File Table



4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.36

## Memory Mapped Files

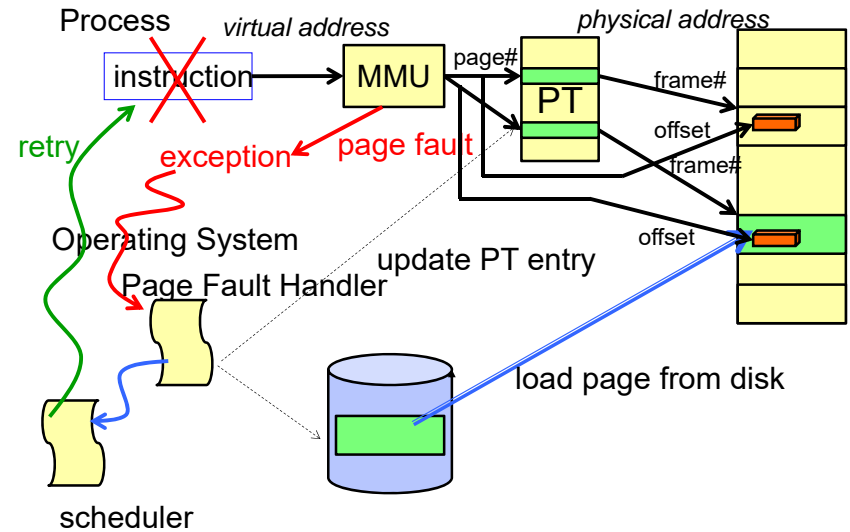
- Traditional I/O involves explicit transfers between buffers in process address space to/from regions of a file
  - This involves multiple copies into caches in memory, plus system calls
- What if we could “map” the file directly into an empty region of our address space
  - Implicitly “page it in”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

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.37

## Recall: Who Does What, When?

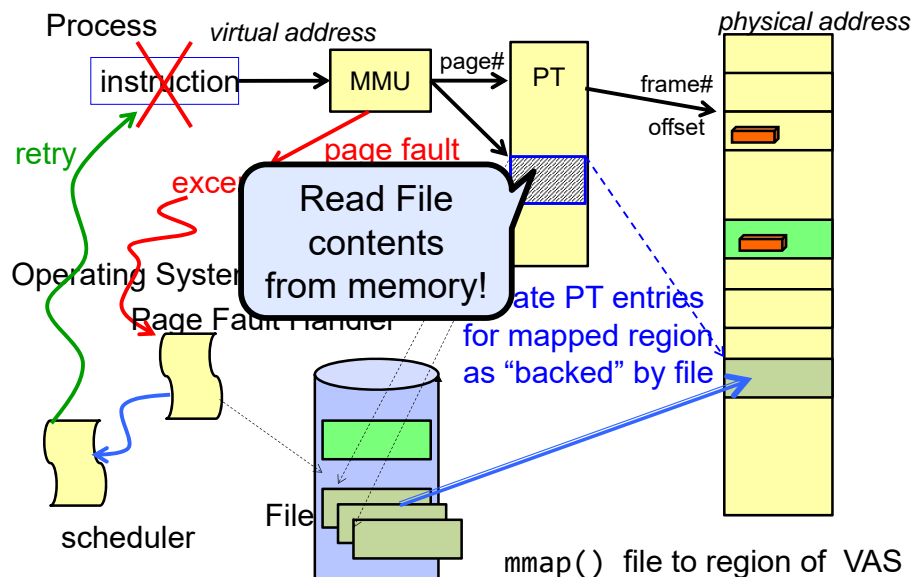


4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.38

## Using Paging to `mmap()` Files



4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.39

## `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 addr and continuing
  for at most len bytes to be mapped from the object described by fd,
  starting at byte offset offset. If offset or len is not a multiple of
  
```

- 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

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.40

## 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; }

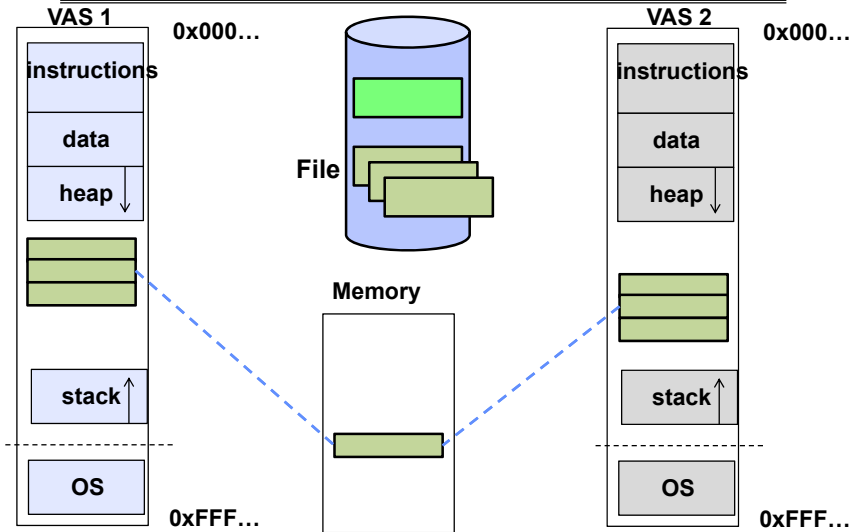
    printf("mmap at : %16lx\n", (long) mfile);

    puts(mfile);
    strcpy(mfile+20,"Let's write over its line three");
    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 its line three
This is line four
```

## Sharing through Mapped Files



- Also: anonymous memory between parents and children
  - no file backing – just swap space

## File System Caching

- Key Idea: Exploit locality by caching data in memory
  - Name translations: Mapping from paths→inodes
  - Disk blocks: Mapping from block address→disk content
- **Buffer Cache:** Memory used to cache kernel resources, including disk blocks and name translations
  - Can contain “dirty” blocks (blocks yet on disk)
- Replacement policy? LRU
  - Can afford overhead of timestamps for each disk block
  - Advantages:
    - » Works very well for name translation
    - » Works well in general as long as memory is big enough to accommodate a host’s working set of files.
  - Disadvantages:
    - » Fails when some application scans through file system, thereby flushing the cache with data used only once
    - » Example: `find . -exec grep foo {} \;`
- Other Replacement Policies?
  - Some systems allow applications to request other policies
  - Example, ‘Use Once’:
    - » File system can discard blocks as soon as they are used

## File System Caching (con’t)

- Cache Size: How much memory should the OS allocate to the buffer cache vs virtual memory?
  - Too much memory to the file system cache ⇒ won’t be able to run many applications at once
  - Too little memory to file system cache ⇒ many applications may run slowly (disk caching not effective)
  - Solution: adjust boundary dynamically so that the disk access rates for paging and file access are balanced
- **Read Ahead Prefetching:** fetch sequential blocks early
  - Key Idea: exploit fact that most common file access is sequential by prefetching subsequent disk blocks ahead of current read request (if they are not already in memory)
  - Elevator algorithm can efficiently interleave groups of prefetches from concurrent applications
  - How much to prefetch?
    - » Too many imposes delays on requests by other applications
    - » Too few causes many seeks (and rotational delays) among concurrent file requests

## File System Caching (con't)

- **Delayed Writes:** Writes to files not immediately sent out to disk
  - Instead, `write()` copies data from user space buffer to kernel buffer (in cache)
    - » Enabled by presence of buffer cache: can leave written file blocks in cache for a while
    - » If some other application tries to read data before written to disk, file system will read from cache
  - Flushed to disk periodically (e.g. in UNIX, every 30 sec)
  - Advantages:
    - » Disk scheduler can efficiently order lots of requests
    - » Disk allocation algorithm can be run with correct size value for a file
    - » Some files need never get written to disk! (e..g temporary scratch files written /tmp often don't exist for 30 sec)
  - Disadvantages
    - » What if system crashes before file has been written out?
    - » Worse yet, what if system crashes before a directory file has been written out? (lose pointer to inode!)

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.45

## Important “ilities”

- **Availability:** the probability that the system can accept and process requests
  - Often measured in “nines” of probability. So, a 99.9% probability is considered “3-nines of availability”
  - Key idea here is independence of failures
- **Durability:** the ability of a system to recover data despite faults
  - This idea is fault tolerance applied to data
  - Doesn't necessarily imply availability: information on pyramids was very durable, but could not be accessed until discovery of Rosetta Stone
- **Reliability:** the ability of a system or component to perform its required functions under stated conditions for a specified period of time (IEEE definition)
  - Usually stronger than simply availability: means that the system is not only “up”, but also working correctly
  - Includes availability, security, fault tolerance/durability
  - Must make sure data survives system crashes, disk crashes, other problems

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.46

## 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

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.47

## 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
- Buffer Cache: Memory used to cache kernel resources, including disk blocks and name translations
  - Can contain “dirty” blocks (blocks yet on disk)

4/11/19

Kubiatowicz CS162 © UCB Spring 2019

Lec 19.48