# CS162
# Operating Systems and Systems Programming
# Lecture 19

## File Systems (Con't), MMAP, Buffer Cache
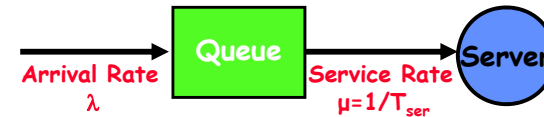
November 4th, 2015

Prof. John Kubiatowicz

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 = $\sigma^2/m1^2$
  - $\mu$:     service rate = $1/T_{ser}$
  - $u$:     server utilization ($0 \le u \le 1$): $u = \lambda/\mu = \lambda \times T_{ser}$
- **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:**
  - **M**emoryless service distribution ($C = 1$):
    - » Called M/M/1 queue: $T_q = T_{ser} \times u/(1 - u)$
  - **G**eneral service distribution (no restrictions), 1 server:
    - » Called M/G/1 queue: $T_q = T_{ser} \times \frac{1}{2}(1+C) \times u/(1 - u))$

---

## Recall: 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
  - Disk Management: collecting disk blocks into files
  - Naming: Interface to find files by name, not by blocks
  - Protection: Layers to keep data secure
  - Reliability/Durability: Keeping of files durable despite crashes, media failures, attacks, etc
- 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 ≥ sector size; in UNIX, block size is 4KB

---

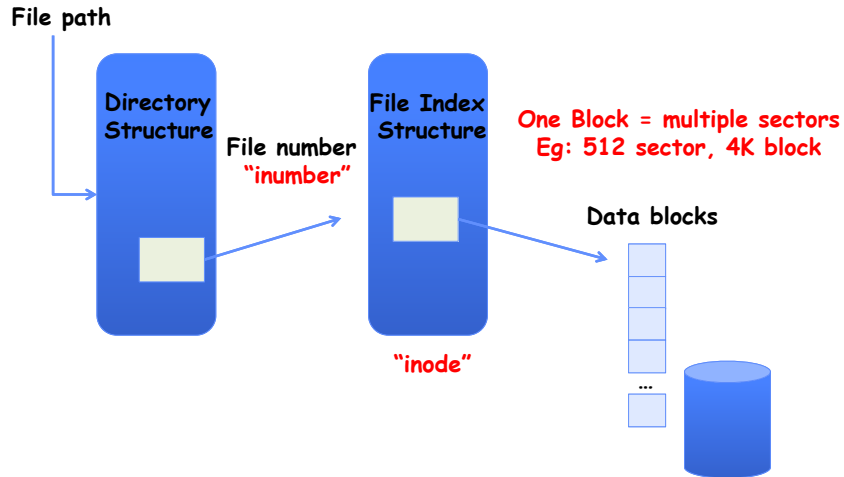## Disk Management Policies

- **Basic entities on a disk:**
  - **File:** user-visible group of blocks arranged sequentially in logical space
  - **Directory:** user-visible index mapping names to files (next lecture)
- **Access disk as linear array of sectors. Two Options:**
  - Identify sectors as vectors [cylinder, surface, sector]. Sort in cylinder-major order. Not used much anymore.
  - **Logical Block Addressing (LBA).** Every sector has integer address from zero up to max number of sectors.
  - Controller translates from address ⇒ physical position
    - » First case: OS/BIOS must deal with bad sectors
    - » Second case: hardware shields OS from structure of disk
- **Need way to track free disk blocks**
  - Link free blocks together ⇒ too slow today
  - Use bitmap to represent free space on disk
- **Need way to structure files: File Header ("Inode")**
  - Track which blocks belong at which offsets within the logical file structure
  - Optimize placement of files' disk blocks to match access and usage patterns

## Components of a File System

File path

Directory Structure

File number "inumber"

File Index Structure

**One Block = multiple sectors
Eg: 512 sector, 4K block**

Data blocks

"inode"

...

## Components of a file system

file name offset → directory → file number offset → index structure → Storage block
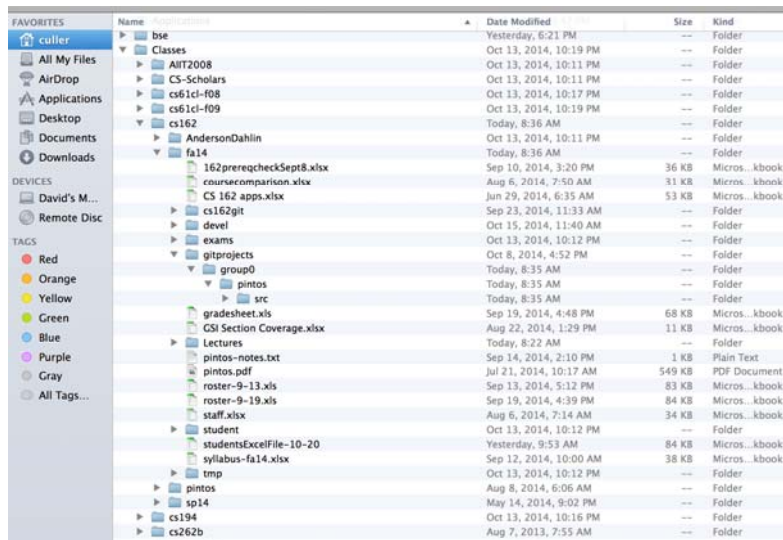
- **Open performs *name resolution***
  - **Translates pathname into a "file number"**
    - » Used as an "index" to locate the blocks
  - **Creates a file descriptor in PCB within kernel**
  - **Returns a "handle" (another int) to user process**
- **Read, Write, Seek, and Sync operate on handle**
  - **Mapped to descriptor and to blocks**

## Directories

| FAVORITES | Name | Date Modified | Size | Kind |
|---|---|---|---|---|
| culler | ▶ bse | Yesterday, 6:21 PM | -- | Folder |
| All My Files | ▼ Classes | Oct 13, 2014, 10:19 PM | -- | Folder |
| AirDrop | ▶ AIIT2008 | Oct 13, 2014, 10:11 PM | -- | Folder |
| Applications | ▶ CS-Scholars | Oct 13, 2014, 10:11 PM | -- | Folder |
| Desktop | ▶ cs61cl-f08 | Oct 13, 2014, 10:17 PM | -- | Folder |
| Documents | ▶ cs61cl-f09 | Oct 13, 2014, 10:19 PM | -- | Folder |
| Downloads | ▼ cs162 | Today, 8:36 AM | -- | Folder |
| | ▶ AndersonDahlin | Oct 13, 2014, 10:11 PM | -- | Folder |
| DEVICES | ▼ fa14 | Today, 8:36 AM | -- | Folder |
| David's M... | 162prereqcheckSept8.xlsx | Sep 10, 2014, 3:20 PM | 36 KB | Micros...kbook |
| Remote Disc | coursecomparison.xlsx | Aug 6, 2014, 7:50 AM | 31 KB | Micros...kbook |
| | CS 162 apps.xlsx | Jun 29, 2014, 6:35 AM | 53 KB | Micros...kbook |
| TAGS | ▶ cs162git | Sep 23, 2014, 11:33 AM | -- | Folder |
| Red | ▶ devel | Oct 15, 2014, 11:40 AM | -- | Folder |
| Orange | ▶ exams | Oct 13, 2014, 10:12 PM | -- | Folder |
| Yellow | ▼ gitprojects | Oct 8, 2014, 4:52 PM | -- | Folder |
| Green | ▼ group0 | Today, 8:35 AM | -- | Folder |
| Blue | ▼ pintos | Today, 8:35 AM | -- | Folder |
| Purple | ▶ src | Today, 8:35 AM | -- | Folder |
| Gray | gradesheet.xls | Sep 19, 2014, 4:48 PM | 68 KB | Micros...kbook |
| All Tags... | GSI Section Coverage.xlsx | Aug 22, 2014, 1:29 PM | 11 KB | Micros...kbook |
| | ▶ Lectures | Today, 8:22 AM | -- | Folder |
| | pintos-notes.txt | Sep 14, 2014, 2:10 PM | 1 KB | Plain Text |
| | pintos.pdf | Jul 21, 2014, 10:17 AM | 549 KB | PDF Document |
| | roster-9-13.xls | Sep 13, 2014, 5:12 PM | 83 KB | Micros...kbook |
| | roster-9-19.xls | Sep 19, 2014, 4:39 PM | 84 KB | Micros...kbook |
| | staff.xlsx | Aug 6, 2014, 7:14 AM | 34 KB | Micros...kbook |
| | ▶ student | Oct 13, 2014, 10:12 PM | -- | Folder |
| | studentsExcelFile-10-20 | Yesterday, 9:53 AM | 84 KB | Micros...kbook |
| | syllabus-fa14.xlsx | Sep 12, 2014, 10:00 AM | 38 KB | Micros...kbook |
| | ▶ tmp | Oct 13, 2014, 10:12 PM | -- | Folder |
| | ▶ pintos | Aug 8, 2014, 6:06 AM | -- | Folder |
| | ▶ sp14 | May 14, 2014, 9:02 PM | -- | Folder |
| | ▶ cs194 | Oct 13, 2014, 10:16 PM | -- | Folder |
| | ▶ cs262b | Aug 7, 2013, 7:55 AM | -- | Folder |

## 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**
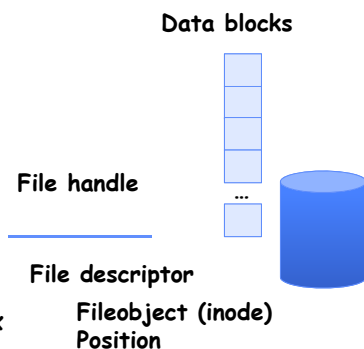  - **Softlinks (aliases) are another name for an entry**

## 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)**
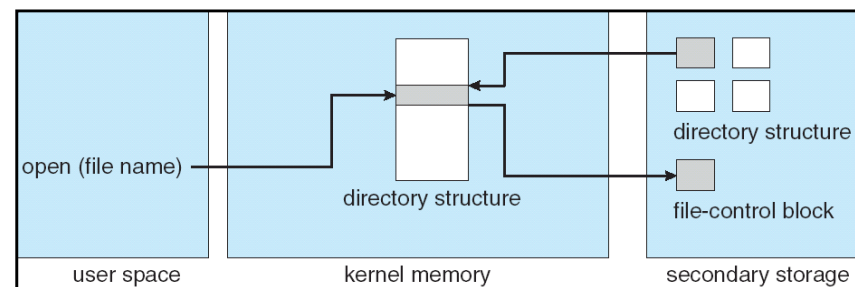      - · **Access control list in Windows system**

**Data blocks**

**File handle**

...

**File descriptor**

**Fileobject (inode)**
**Position**

## In-Memory File System Structures

open (file name)

directory structure

directory structure

file-control block

user space          kernel memory          secondary storage

- **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

index

read (index)

per-process
open-file table

system-wide
open-file table

data blocks

file-control block

user space          kernel memory          secondary storage

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

## Administrivia

- **HW3 – Moved deadline to Wednesday (11/04)**
  - **Sorry about fact that server was down!**
- **Project 2 code due this Friday!**
- **Don't forget Peer-review for project 1**
  - **Everyone must submit one!**
- **Midterm I Regrade requests: Due Today!**
- **Midterm II: Coming up in 3 weeks! (11/23)**
  - **7-10PM, "here" (2040, 2050, 2060 VLSB)**
  - **Topics up to and including previous Wednesday**
  - **1 page of hand-written notes, both sides**
- **Only 30 people filled out course survey**
  - **Please give us more feedback!**

## Our first filesystem: FAT (File Allocation Table)

- **Assume (for now) we have a way to translate a path to a "file number"**
  - i.e., a directory structure
- **Disk Storage is a collection of Blocks**
  - Just hold file data
- **Example: file_read 31, < 2, x >**
  - Index into FAT with file number
  - Follow linked list to block
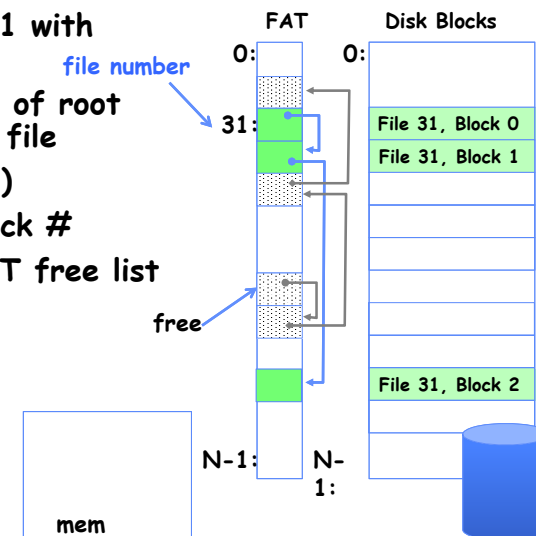  - Read the block from disk into mem

*FAT*   *Disk Blocks*

*file number*

0:

31:  File 31, Block 0
     File 31, Block 1

     File 31, Block 2

N-1:  N-1:

**mem**

## FAT Properties

- **File is collection of disk blocks**
- **FAT is linked list 1-1 with blocks**
- **File Number is index of root of block list for the file**
- **File offset (o = B:x )**
- **Follow list to get block #**
- **Unused blocks ⇔ FAT free list**

*FAT*   *Disk Blocks*

*file number*

0:   0:

31:  File 31, Block 0
     File 31, Block 1

free

     File 31, Block 2

N-1:  N-1:

**mem**

## FAT Properties

- **File is collection of disk blocks**
- **FAT is linked list 1-1 with blocks**
- **File Number is index of root of block list for the file**
- **File offset (o = B:x )**
- **Follow list to get block #**
- **Unused blocks ⇔ FAT free list**
- **Ex: file_write(51, <3, y> )**
  - Grab blocks from free list
  - Linking them into file

*FAT*   *Disk Blocks*

*file number*

0:   0:

31:  File 31, Block 0
     File 31, Block 1

free

     File 31, Block 3

     File 31, Block 2

N-1:  N-1:

**mem**

## FAT Properties

- **File is collection of disk blocks**
- **FAT is linked list 1-1 with blocks**
- **File Number is index of root of block list for the file**
- **Grow file by allocating free blocks and linking them in**
- **Ex: Create file, write, write**

*FAT*   *Disk Blocks*

*file number*

0:   0:

31:  File 31, Block 0
     File 31, Block 1
     File 63, Block 1

     File 31, Block 3
free  63:  File 63, Block 0

*File 2 number*

     File 31, Block 2

N-1:  N-1:

**mem**

## FAT Assessment

- **Used in DOS, Windows, thumb drives, …**
- **Where is FAT stored?**
  - On Disk, restore on boot, copy in memory
- **What happens when you format a disk?**
  - Zero the blocks, link up the FAT free-list
- **Simple**

file number

FAT     Disk Blocks

0:     0:

31:   File 31, Block 0
File 31, Block 1
File 63, Block 1

File 31, Block 3
free 63:   File 63, Block 0

File 2 number

File 31, Block 2

N-1:   N-1:

mem

---

## FAT Assessment

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

file number

FAT     Disk Blocks

0:     0:

31:   File 31, Block 0
File 31, Block 1
File 63, Block 1

File 31, Block 3
free 63:   File 63, Block 0

File 2 number

File 31, Block 2

N-1:   N-1:

mem

---

## What about the Directory?

file 5268830
"/home/tom"

end of file

| Name | . | .. | Music | Work | Free Space | foo.txt | Free Space |
| File Number | 5268830 | 88026158 | 35002320 | 85200219 | | 66212871 | |
| Next | | | | | | | |

- **Essentially a file containing <file_name: file_number> mappings**
- **Free space for new entries**
- **In FAT: attributes kept in directory (!!!)**
- **Each directory a linked list of entries**
- **Where do you find root directory ( "/" )?**

---

## Directory Structure (Con't)

- **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")

## Big FAT security holes

- **FAT has no access rights**
- **FAT has no header in the file blocks**
- **Just gives and index into the FAT**
  - **– (file number = block number)**

## Characteristics of Files

- **Most files are small**
- **Most of the space is occupied by the rare big ones**

A Five-Year Study of File-System Metadata

NITIN AGRAWAL
University of Wisconsin, Madison
and
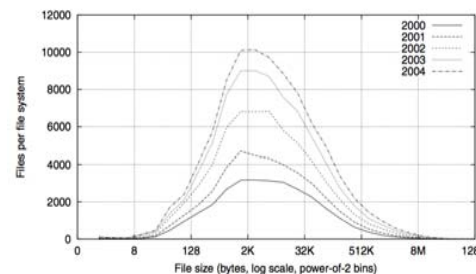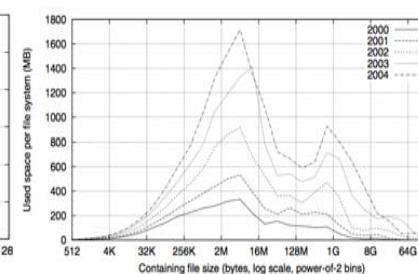WILLIAM J. BOLOSKY, JOHN R. DOUCEUR, and JACOB R. LORCH
Microsoft Research



A Five-Year Study of File-System Metadata    •    9:9

Fig. 2. Histograms of files by size.

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

## So what about a "real" file system

- **Meet the inode:**

## Unix File System

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

## An "almost real" file system

• **Pintos: src/filesys/file.c, inode.c**

```
/* An open file. */
struct file
  {
    struct inode *inode;      /* File's inode. */
    off_t pos;                /* Current position. */
    bool deny_write;          /* Has file_deny_write() been called? */
  };
```

```
/* In-memory inode. */
struct inode
  {
    struct list_elem elem;           /* Element in inode list. */
    block_sector_t sector;           /* Sector number of disk location. */
    int open_cnt;                    /* Number of openers. */
    bool removed;                    /* True if deleted, false otherwise. */
    int deny_write_cnt;              /* 0: writes ok, >0: deny writes. */
    struct inode_disk data;          /* Inode content. */
  };
```

```
/* On-disk inode.
   Must be exactly BLOCK_SECTOR_SIZE bytes long. */
struct inode_disk
  {
    block_sector_t start;            /* First data sector. */
    off_t length;                    /* File size in bytes. */
    unsigned magic;                  /* Magic number. */
    uint32_t unused[125];            /* Not used. */
  };
```

---

## File Attributes

• **Inode metadata**



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

---

## Data Storage

• **Small files: 12 pointers direct to data blocks**

```
Direct pointers

4kB blocks ⇒ sufficient
For files up to 48KB
```



Fig. 2.  Histograms of files by size.

---

## 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
```



**48 KB**

**+4 MB**

**+4 GB**

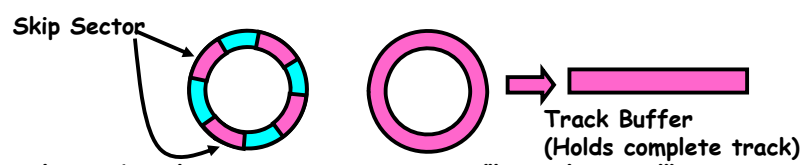**+4 TB**

## UNIX BSD 4.2

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



Skip Sector

Track Buffer
(Holds complete track)

  - **Solution1: Skip sector positioning ("interleaving")**
    - » **Place the blocks from one file on every other block of a track: give time for processing to overlap rotation**
  - **Solution2: 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**
- **Important Aside: Modern disks+controllers do many complex 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 were created (They were each 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 a portion of the file header array on each of 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)**
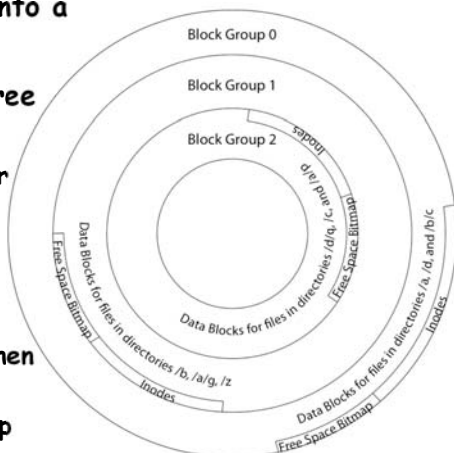    - » **General optimization to avoid seeks**

## 4.2 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
- 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!**
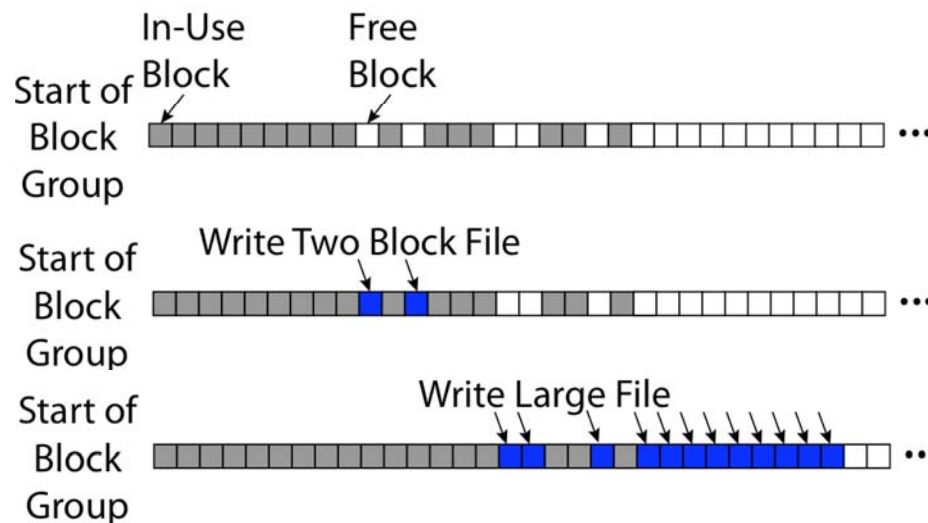  - **Reserve space in the BG**

---

## FFS First Fit Block Allocation

---

## FFS

- Pros
  - Efficient storage for both small and large files
  - Locality for both small and large files
  - Locality for metadata and data
- 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 (no equivalent to superpages)
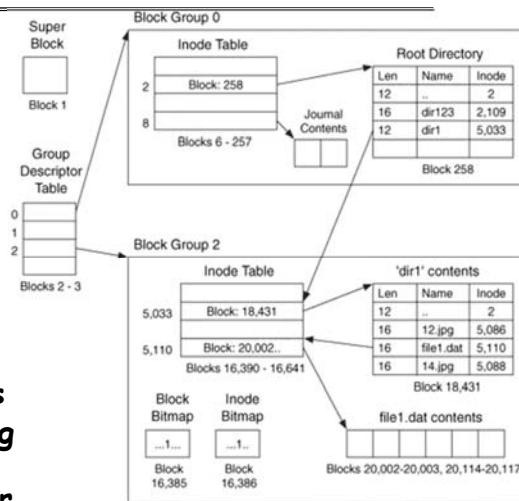  - Need to reserve 10-20% of free space to prevent fragmentation

---

## 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 4.2BSD
  - with 12 direct pointers
- Ext3: Ext2 w/Journaling
  - Several degrees of protection with more or less cost

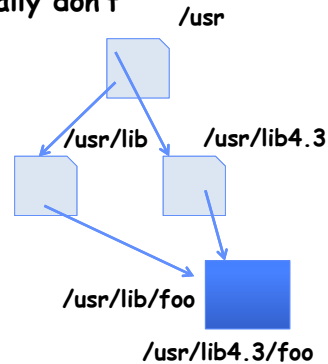- **Example: create a *file1.dat* under */dir1/* in Ext3**

## 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**
    - » **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)`**

/usr

/usr/lib    /usr/lib4.3

/usr/lib/foo

/usr/lib4.3/foo

## Links

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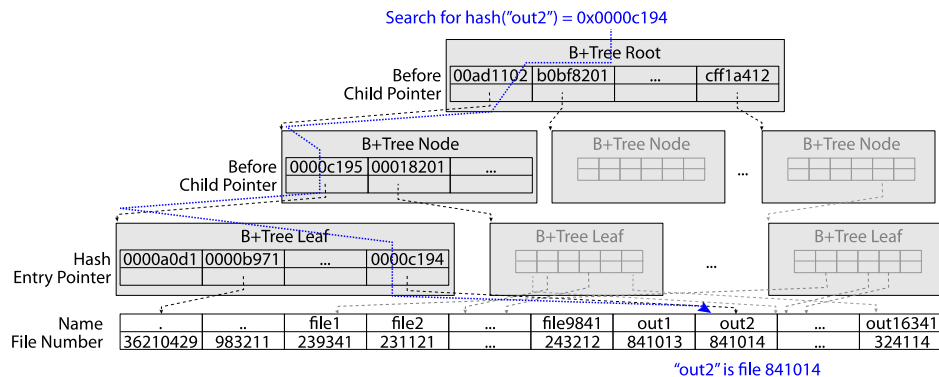  - **Each is "first class"**
- **Soft link or Symbolic Link**
  - **Directory entry contains the name of the file**
  - **Map one name to another name**

## Large Directories: B-Trees (dirhash)

Search for hash("out2") = 0x0000c194

B+Tree Root
| 00ad1102 | b0bf8201 | … | cff1a412 |

Before Child Pointer

B+Tree Node
| 0000c195 | 00018201 | … |

Before Child Pointer

B+Tree Node   …   B+Tree Node

B+Tree Leaf
| 0000a0d1 | 0000b971 | … | 0000c194 |

Hash Entry Pointer

B+Tree Leaf   …   B+Tree Leaf

| Name | . | .. | file1 | file2 | … | file9841 | out1 | out2 | … | out16341 |
|---|---|---|---|---|---|---|---|---|---|---|
| File Number | 36210429 | 983211 | 239341 | 231121 | … | 243212 | 841013 | 841014 | … | 324114 |

"out2" is file 841014

## NTFS

- **New Technology File System (NTFS)**
  - **Common 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 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
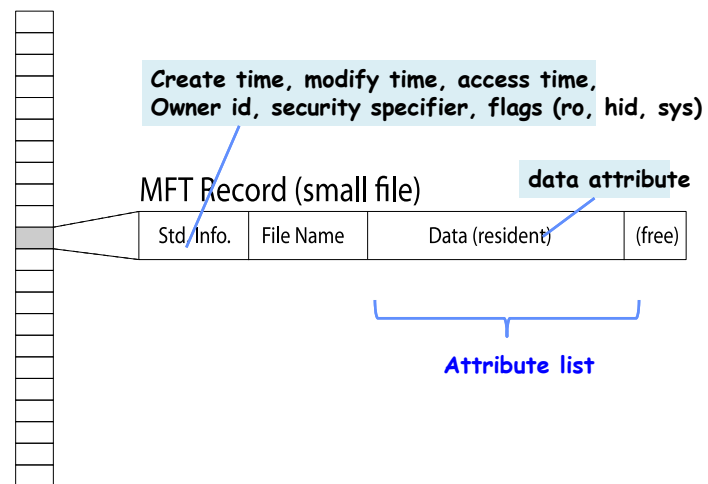- **Journalling for reliability**
  - Discussed later

## NTFS Small File

Master File Table

Create time, modify time, access time, Owner id, security specifier, flags (ro, hid, sys)

MFT Record (small file)          data attribute

| Std. Info. | File Name | Data (resident) | (free) |

**Attribute list**

## NTFS Medium File



Master File Table

MFT Record

| Std. Info. | File Name | Data (nonresident) | (free) |

Start
Length
Data Extent
Start + Length

Start
Length
Data Extent
Start + Length

## NTFS Multiple Indirect Blocks



MFT Record (big/fragmented file)

| Std. Info. | Attr. List | ... | Data (nonresident) |

Data (nonresident)

Data (nonresident)

Data (nonresident)

## Slide 19.45

Master File Table

MFT Record
(huge/badly-fragmented file)

| Std. Info. | Attr. List (nonresident) | ... |

... Extent with part of attribute list

Data (nonresident)
...

Data (nonresident)
...

Data (nonresident)
...

... Extent with part of attribute list

Data (nonresident)
...

Data (nonresident)
...

... Extent with part of attribute list

Data (nonresident)
...

Data (nonresident)
...
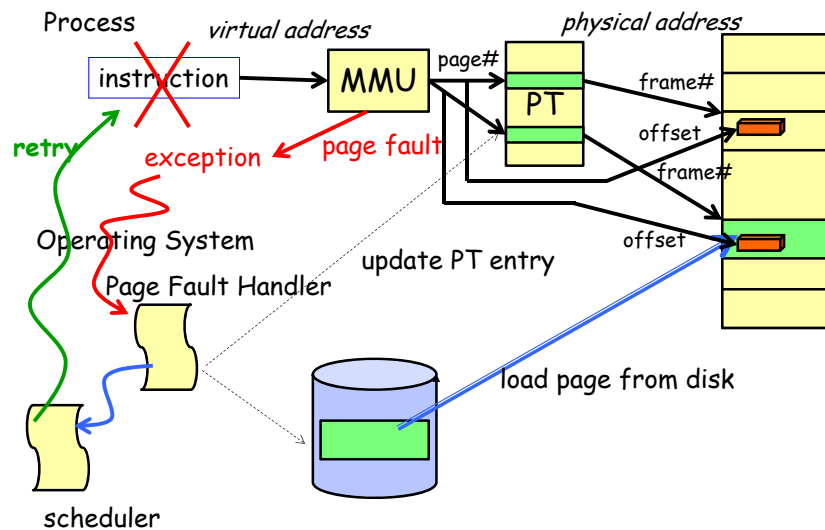
---

## Memory Mapped Files

- **Traditional I/O involves explicit transfers between buffers in process address space to 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 file is treated this way when we exec the process !!**

---

## Recall: Who does what, when?

Process

virtual address

instruction

retry

exception

MMU

page#

page fault

PT

physical address

frame#

offset

frame#

offset

Operating System

update PT entry

Page Fault Handler

load page from disk

scheduler

---

## Using Paging to mmap files

Process

virtual address

instruction

retry

excep

MMU

page#

page fault

PT

physical address

frame#

offset

**Read File contents from memory!**

tries
region
as "backed" by file

Operating System

Page Fault Handler

File

**mmap file to region of VAS**

scheduler

## 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**

## An example

```c
#include <sys/mman.h>

int something = 162;

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

  printf("Data  at: %16lx\n", (long unsigned int) &something);
  printf("Heap at : %16lx\n", (long unsigned int) malloc(1));
  printf("Stack at: %16lx\n", (long unsigned int) &mfile);

  /* Open the file */
  myfd = open(argv[1], O_RDWR | O_CREATE);
  if (myfd < 0) { perror(("open failed!");exit(1); }

  /* map the file */
  mfile = mmap(0, 10000, PROT_READ|PROT_WRITE, MAP_FILE|MAP_SHARED, myfd, 0);
  if (mfile == MAP_FAILED) {perror("mmap failed"); exit(1);}

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

  puts(mfile);
  strcpy(mfile+20,"Let's write over it");
  close(myfd);
  return 0;
}
```
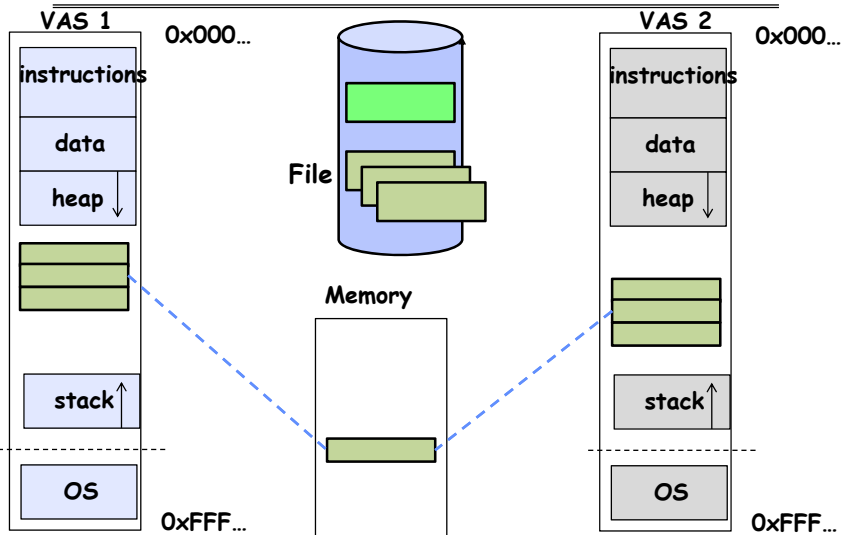
## Sharing through Mapped Files



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

## System-V-style Shared Memory
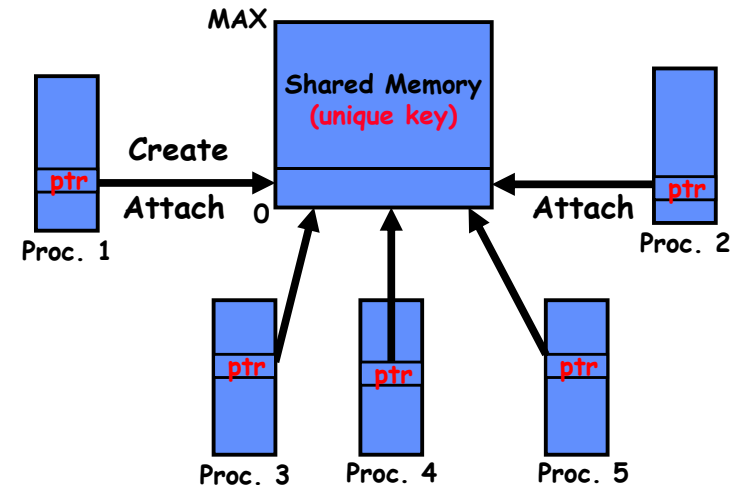
**Common chunk of read/write memory among processes**

## Creating Shared Memory

```
// Create new segment
int shmget(key_t key, size_t size, int shmflg);
```

Example:

```
key_t key;
int shmid;
```

> Filename only used to generate key – not for storage

```
key = ftok("<somefile>", 'A');

shmid = shmget(key, 1024, 0644 | IPC_CREAT);

Special key: IPC_PRIVATE (create new segment)
Flags: IPC_CREAT (Create new segment)
       IPC_EXCL (Fail if segment with key already exists)
       lower 9 bits – permissions use on new segment
```

## Attach and Detach Shared Memory

```
// Attach
void *shmat(int shmid, void *shmaddr, int shmflg);
       Flags: SHM_RDONLY, SHM_REMAP
// Detach
int shmdt(void *shmaddr);
```

Example:

```
key_t key;
int shmid;
char *sharedmem;

key = ftok("<somefile>", 'A');
shmid = shmget(key, 1024, 0644);
sharedmem = shmat(shmid, (void *)0, 0);  // Attach smem
// Use shared memory segment (address is in sharedmem)
…
shmdt(sharedmem); // Detach smem (all finished)
```

## 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!)

## 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

## 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: act of translating from user-visible names to actual system 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 uses variable extents, rather than fixed blocks, and tiny files data is in the header
- 4.2 BSD Multilevel index files
  - Inode contains pointers 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 System Summary (2/2)

- **File layout driven by freespace management**
  - Integrate freespace, inode table, file blocks and directories into block group
- **Deep interactions between memory management, file system, and sharing**
  - mmap(): map file or anonymous segment to memory
  - ftok/shmget/shmat: Map (anon) shared-memory segments
- **Buffer Cache: Memory used to cache kernel resources, including disk blocks and name translations**
  - Can contain "dirty" blocks (blocks yet on disk)
- **Important system properties**
  - Availability: how often is the resource available?
  - Durability: how well is data preserved against faults?
  - Reliability: how often is resource performing correctly?