

Chapter 9 Main Memory

Cache sits between main memory and CPU registers

Base and limit registers define logical address space usable by a process

Compiled code addresses bind to relocatable addresses

Can happen at three different stages

- Compile time: If memory location known a priori, absolute code can be generated
- Load time: Must generate relocatable code if memory location not known at compile time
- Execution time: Binding delayed until run time if the process can be moved during its execution

Memory Management Unit (MMU) device that maps virtual to physical address

Simple scheme uses a relocation register which adds a base value to address

Swapping allows total physical memory space of processes to exceed physical memory

- Def: process swapped out temporarily to backing store then brought back in for continued execution
- Backing store: fast disk large enough to accommodate copies of all memory images

Roll out, roll in: swapping variant for priority-based scheduling

- Lower priority process swapped out so that higher priority process can be loaded

Solutions to Dynamic Storage-Allocation Problem:

- First-fit: allocate the first hole that is big enough
- Best fit: allocate the smallest hole that is big enough (must search entire list) → smallest leftover hole
- Worst fit: allocate the largest hole (search entire list) → largest leftover hole

External Fragmentation: total memory space exists to satisfy request, but is not contiguous

- Reduced by compaction: relocate free memory to be together in one block
- Only possible if relocation is dynamic

Internal Fragmentation: allocated memory may be slightly larger than requested memory

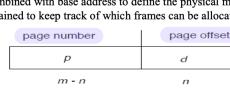
Physical memory divided into fixed-sized frames: size is power of 2, between 512 bytes and 16 MB

Logical memory divided into same sized blocks: pages

Page table used to translate logical to physical addresses

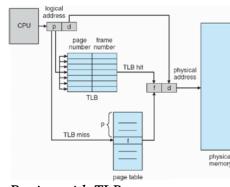
- Page number (p), used as an index into a page table
- Page offset (d): combined with base address to define the physical memory address

Free-frame list is maintained to keep track of which frames can be allocated



Transition Look-aside Buffer (TLB) is a CPU cache that memory management hardware uses to improve virtual address translation speed

- Typically small – 64 to 1024 entries
- On TLB miss, value loaded to TLB for faster access next time
- TLB is associative – searched in parallel



Paging with TLB



Paging without TLB

Effective Access Time: $EAT = (1 + \epsilon) \alpha + (2 + \epsilon)(1 - \alpha)$

ϵ = time unit, α = hit ratio

Valid and invalid bits can be used to protect memory

- "Valid" if the associated page is in the process' logical address space, so it is a legal page

Can have multilevel page tables (paged page tables)

Hashed Page Tables: virtual page number hashed into page table

- Page table has chain of elements hashing to the same location

Each element has (1) virtual page number, (2) value of mapped page frame, (3) a pointer to the next element

Search through the chain for virtual page number

Segment table – maps two-dimensional physical addresses

- Entries protected with valid bits and r/w/x privileges

OSC 9.15: Compare the memory organization schemes or contiguous memory allocation and paging with respect to the following issues:

a. External fragmentation

Answer: Contiguous allocation can result in external fragmentation, while paging essentially eliminates it.

b. Internal fragmentation

Answer: In contrast, contiguous allocation has no internal fragmentation, while it is possible for paging to have unused space in a page up to the page size - 1.

c. Ability to share code across processes

Answer: It is possible to map code to the same logical addresses across processes such that the pages references are the same. Then there needs only be 1 set of frames used to hold the code in memory. However, contiguous allocation requires all addresses reference by a process to be in its own process-specific memory region. Thus, it is not possible to share code in the same way. A copy of the code would be in the memory allocated to each process.

Swapping

Physical memory space is finite

Suppose we allocate physical memory to meet the needs of a process

- What happens when the total physical memory space requested by processes exceeds physical memory?

A process can be swapped temporarily out of memory

- Moved to a backing store– fast disk large enough to accommodate copies of all memory images for all users
- Later brought back into memory for continued execution

Swapping is used along with process scheduling

- Roll out, roll in – swapping variant used for priority-based scheduling algorithms

Swapping Methods

Does the swapped out process need to swap back in to same physical addresses?

Depends on address binding method

- Also consider pending I/O since it might be associated with particular addresses

Modified versions of swapping are found on many systems (i.e., UNIX, Linux, and Windows)

Swapping normally disabled

Started if greater then memory allocation threshold

Disabled again once memory demand reduced

8.3 Why are page sizes always powers of 2?

Answer: Recall that pages are implemented by breaking up an address into a page offset and page number. It is most efficient to break the address into X page bits and Y offset bits, rather than perform arithmetic on the address to calculate the page number and offset. Because each bit position represents a power of 2, splitting an address between bits results in a page size that is a power of 2.

Simulations have shown that both first fit and best fit are better than worst fit in terms of decreasing time and storage utilization. Neither first fit nor best fit is clearly better than the other in terms of storage utilization, but first fit is generally faster.

Segmentation

- Memory-management scheme that supports a user view of memory
- Allows different parts to be allocated separately
- A program is a collection of segments
- A segment is a logical unit such as:
 - main program
 - procedure
 - function
 - method
 - object
 - local variables, global variables
 - common block
 - stack
 - symbol table
 - arrays



memory

MMU

physical address

logical address

segment table

base address

limit

offset

segment offset

segment selector

segment limit

segment base

segment

program

function

method

object

local

global

common

stack

symbol

table

arrays

logical address

segment

offset

segment

base

segment

limit

segment

selector

segment

offset

segment

base

segment

limit

segment

offset

segment

base

segment</

Chapter 10 - Virtual Memory

| Rank | Algorithm | Suffer from Belady's anomaly |
|------|---------------|------------------------------|
| 1 | Optimal | no |
| 2 | LRU | no |
| 3 | Second-chance | yes |
| 4 | FIFO | yes |

9.1 Under what circumstances do page faults occur? Describe the actions taken by the operating system when a page fault occurs.

Answer:

A page fault occurs when an access to a page that has not been brought into main memory takes place. The operating system verifies the memory access, aborting the program if it is invalid. If it is valid, a free frame is located and I/O is requested to read the needed page into the free frame. Upon completion of I/O, the process table and page table are updated and the instruction is restarted.

9.5 Discuss the hardware support required to support demand paging.

Answer:

For every memory-access operation, the page table needs to be consulted to check whether the corresponding page is resident or not and whether the program has read/write privileges for accessing the page. These checks have to be performed in hardware. A TLB could serve as a cache and improve the performance of the lookup operation.

OSC 10.21: Assume that we have a demand-paged memory. The page table is held in registers. It takes 8 milliseconds to service a page fault if an empty frame is available or if the replaced page is not modified. 20 milliseconds if the replaced page is modified. Memory-access time is 100 nanoseconds. Assume that the page to be replaced is modified 70 percent of the time. What is the maximum acceptable page-fault rate for an effective access time of no more than 200 nanoseconds?

$$\text{effective access time} = (1 - p) \times \text{ma} + p \times \text{page fault time.}$$

OSC 10.37: What is the cause of thrashing? How does the system detect thrashing? Once it detects thrashing, what can the system do to eliminate this problem?

Answer: Thrashing is occurs in a virtual memory system when there is an increase in the number of page faults in the system caused by the inability to keep pages that processes need in memory. Thrashing can be attributed to memory allocation difficulties brought on by increasing the degree of multiprogramming. Essentially, the higher the degree of multiprogramming, the smaller the amount of memory that can be allocated to any process. Consequently, the OS will spend a relatively larger amount of time servicing page faults versus performing useful work. Thrashing can be detected when the page fault rates increases and the CPU utilization falls significantly. The only real way to eliminate the problem is to reduce the degree of multiprogramming to an acceptable level.

Summary

- Virtual memory abstracts physical memory into an extremely large uniform array of storage.
- The benefits of virtual memory include the following: (1) a program can be larger than physical memory, (2) a program does not need to be entirely in memory, (3) processes can share memory, and (4) processes can be created more efficiently.
- Demand paging is a technique whereby pages are loaded only when they are demanded during program execution. Pages that are never demanded are thus never loaded into memory.
- A page fault occurs when a page that is currently not in memory is accessed. The page must be brought from the backing store into an available page frame in memory.
- Copy-on-write allows a child process to share the same address space as its parent. If either the child or the parent process writes (modifies) a page, a copy of the page is made.
- When available memory runs low, a page-replacement algorithm selects an existing page in memory to replace with a new page. Page-replacement algorithms include FIFO, optimal, and LRU. Pure LRU algorithms are impractical to implement, and most systems instead use LRU-approximation algorithms.
- Global page-replacement algorithms select a page from any process in the system for replacement, while local page-replacement algorithms select a page from the faulting process.
- Thrashing occurs when a system spends more time paging than executing.
- A locality represents a set of pages that are actively used together. As a process executes, it moves from locality to locality. A working set is based on locality and is defined as the set of pages currently in use by a process.
- Memory compression is a memory-management technique that compresses a number of pages into a single page. Compressed memory is an alternative to paging and is used on mobile systems that do not support paging.
- Kernel memory is allocated differently than user-mode processes; it is allocated in contiguous chunks of varying sizes. Two common techniques for allocating kernel memory are (1) the buddy system and (2) slab allocation.
- TLB reach refers to the amount of memory accessible from the TLB and is equal to the number of entries in the TLB multiplied by the page size. One technique for increasing TLB reach is to increase the size of pages.
- Linux, Windows, and Solaris manage virtual memory similarly, using demand paging and copy-on-write, among other features. Each system also uses a variation of LRU approximation known as the clock algorithm.

EXAMPLE:

$$\text{L.A.S.} = 4 \cdot 6 \rightarrow 2^2 \cdot 2^{30} = 2^{32} \rightarrow \text{LA: } \boxed{10 \ 12} = 32$$

$$\text{P.A.S.} = 64 \cdot 6 \rightarrow 2^6 \cdot 2^{20} = 2^{26} \rightarrow \text{PA: } \boxed{14 \ 12} = 26$$

$$\# \text{pages} = 2^{20}$$

$$\# \text{frames} = 2^{14}$$

$$\# \text{entries into page table} = 2^{20}$$

$$\text{Size of page table} = 2^{20} \times 14 \text{ bits}$$

PHENOMENON IN WHICH INCREASING THE NUMBER OF PAGE FRAMES RESULTS IN INCREASED NUMBER OF PAGE FAULTS FOR CERTAIN PAGES (COMMON W/ FIFO)

9.10 You have devised a new page-replacement algorithm that you think may be optimal. In some contested test cases, Belady's anomaly occurs. Is your new algorithm optimal? Explain your answer.

Answer: No. An optimal algorithm will not suffer from Belady's anomaly because by definition—an optimal algorithm replaces the page that will not be used for the longest time. Belady's anomaly occurs when a page-replacement algorithm evicts a page that will be needed in the immediate future. An optimal algorithm would not have selected such a page.

9.9 Suppose that you want to use a paging algorithm that requires a reference bit (such as second-chance replacement or working-set model), but the hardware does not support it. Sketch how you would simulate a reference bit even if one were not provided by the hardware, or explain why it is not possible to do so. If it is possible, calculate what the cost would be.

Answer: You can use the valid/invalid bit supported in hardware to simulate the reference bit. Initially set the bit to invalid. On first reference to the operating system is generated. The operating system will set a software bit to 1 and reuse the valid/invalid bit to valid.

9.11 Segmentation is similar to paging but uses variable-sized "pages." Define two segment-replacement algorithms based on FIFO and LRU page-replacement schemes. Recall that segments are contiguous in memory, so the segment that is chosen to be replaced may not be big enough to leave enough consecutive locations for the needed segment. Consider strategies for systems where segments cannot be relocated, and those for systems where they can.

a. **FIFO:** Find the first segment large enough to accommodate the incoming segment. If relocation is not possible and no one segment is large enough, select a combination of segments whose memories are contiguous, which are "closest" to the first of the list and which can accommodate the new segment. If relocation is possible, rearrange the memory so that the first N segments large enough for the incoming segment are contiguous in memory. Add any leftover space to the free-space list in both cases.

b. **LRU:** Select the segment that has not been used for the longest period of time and that is large enough, adding any leftover space to the free-space list. If no one segment is large enough, select a combination of segments that are contiguous and are the oldest N segments. If relocation is available, rearrange the oldest N segments to be contiguous in memory and replace those with the new segment.

Answer: The effective access time formula we are interested in is:

$$r * 3 * 8 + r * 7 * 20 + (1 - r) * .0001 = .0002$$

where r is the page fault rate. Solving for r we get:

$$2.4r + 1.4r - r * .0001 = 3.7999r = .0001$$

$$r = .0001 / 3.7999 = .0000263$$

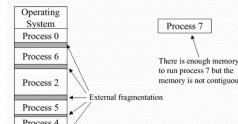
To compute the effective access time, we must know how much time is needed to service a page fault. A page fault causes the following sequence to occur:

- Trap to the operating system.
- Save the registers and process state.
- Determine that the interrupt was a page fault.
- Check that the page reference was legal, and determine the location of the page in secondary storage.
- Issue a read from the storage to a free frame:
 - Wait in a queue until the read request is serviced.
 - Wait for the device seek and/or latency time.
 - Begin the transfer of the page to a free frame.
- While waiting, allocate the CPU core to some other process.
- Receive an interrupt from the storage I/O subsystem (I/O completed).
- Save the registers and process state for the other process (if step 6 is executed).
- Determine that the interrupt was from the secondary storage device.
- Correct the page table and other tables to show that the desired page is now in memory.
- Wait for the CPU core to be allocated to this process again.
- Restore the registers, process state, and new page table, and then resume the interrupted instruction.

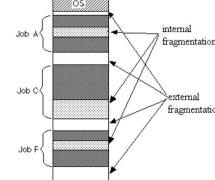
- Example:**
- Virtual memory: separation of user logical memory and physical memory
 - Only part of program needs to be in memory for execution → logical address space > physical address space
 - Allows address spaces to be shared by multiple processes → less swapping
 - Allows pages to be shared during fork, speeding process creation
 - Page fault results from the first time there is a reference to a specific page → traps the OS
 - Must decide to abort if the reference is invalid, or if the desired page is just not in memory yet
 - If the last page is empty frame, swap page into it, reestablishes to indicate page now in memory, set invalid bit, then repeat the process until the page is full
 - If an instruction accesses multiple pages near each other → less "pain" because of locality of reference
 - Demand Paging only brings a page into memory when it is needed → less I/O and memory needed
 - Lazy swapper → never swaps a page into memory unless page will be needed
 - Could result in a lot of page faults
 - Performance: $EAT = [(1-p) \text{memory access} + p(\text{page fault overhead} + \text{swap page out} + \text{swap page in} + \text{restart overhead})]$; Where Page Fault Rate $\theta = p^2$
 - If $p = 0$, then $EAT = 1$ → page fault is always a fault
 - Can minimize demand paging by loading entire process image to swap space at process load time
 - Pure Demand Paging starts with a page in memory
 - Copy-on-Write (CoW) allows both parent and child processes to initially share the same pages in memory
 - If either process modifies a shared page, only then is the page copied
 - Modify (dirty) bit can be used to reduce overhead of page transfers → only modified pages written to disk
 - When a page is replaced, write to disk if it has been dirty and swap in desired page
 - Pages can be replaced using different algorithms: FIFO, LRU (below)
 - Stack can be used to record the most recent page references (LRU is a "stack" algorithm)

- Second-chance algorithm uses a reference bit
 - If 1, decrement and leave in memory
 - If 0, replace next page
- Fixed page allocation: Proportional allocation → Allocate according to size of process
 - $s_i = \text{size of process } P_i, S = \sum s_i, m = \text{total number of frames}, a_i = \text{allocation for } P_i, a_i = (s_i / S)m$
- Global replacement: process selects a replacement frame from set of all frames
 - One process can take frame from another
 - Process execution time can vary greatly
 - Greater throughput
 - More consistent performance
 - Possible under-utilization of memory
- Local replacement: each process selects from only its own set of allocated frames
 - One process can take frame from another
 - Process execution time can vary greatly
 - Greater throughput
 - More consistent performance
 - Possible under-utilization of memory
- Page-fault rate is very high if a process does not have "enough" pages
- Memory-mapped file I/O allows file I/O to be treated as routine memory access by mapping a disk block to a page in memory
- I/O Interlock: Pages must sometimes be locked into memory

External fragmentation
Total memory space is enough to satisfy a request or to reside a process in it, but it is not contiguous so it can not be used.



Internal fragmentation
Memory block assigned to process is bigger. Some portion of memory is left unused as it can not be used by another process.



9.7 Consider the two-dimensional array A:

int A[10][10] = new int[10][10];

where A[0][0] = 1 at location 200 in a paged memory system with pages of size 200. A small process that manipulates the matrix resides in page 0 (location 0 to 199). Thus, every instruction fetch will be from page 0.

For three page frames, how many page faults are generated by the following code? Assume that all initializations loops using LRU replacement and assuming that page frame 0 contains the process and the other two are initially empty?

a. for (int i = 0; i < 100; i++)
 for (int j = 0; j < 100; j++)
 A[i][j] = 0;

b. for (int i = 0; i < 100; i++)
 for (int j = 0; j < 100; j++)
 A[i][j] = 1;

Answer:

a. 5000

b. 50

9.12 Consider a demand-paged computer system where the degree of multiprogramming is currently fixed at four. The system was recently measured to determine utilization of CPU and the paging disk. The results are one entry per row in the table below. For each case, what is happening? Can the degree of multiprogramming be increased to increase the CPU utilization? Is the paging helping?

- a. CPU utilization 13 percent; disk utilization 97 percent
b. CPU utilization 87 percent; disk utilization 3 percent
c. CPU utilization 13 percent; disk utilization 3 percent

Answer:

a. Threading is occurring.

b. CPU utilization is sufficiently high to leave things alone, and increase degree of multiprogramming.

c. Increase the degree of multiprogramming.

9.13 We have an operating system for a machine that uses base and limit registers, but we have modified the machine to provide a page table. Can the page tables be set up to simulate base and limit registers? How can this be done, or why can't they be?

Answer:

The page table can be set up to simulate base and limit registers provided that the memory is allocated in fixed-size segments. In this way, the base of a segment can be entered into the page table and the valid/invalid bit used to indicate that portion of the segment as resident in the memory. There will be some problem with internal fragmentation.

But what happens if the process tries to access a page that was not brought into memory? Access to a page marked invalid causes a page fault. The paging hardware, in translating the address through the page table, will notice that the invalid bit is set, causing a trap to the operating system. This trap is the result of the operating system's failure to bring the desired page into memory. The procedure for handling this page fault is straightforward (Figure 10.5):

- We check an internal table (usually kept with the process control block) for this process to determine whether the reference was a valid or an invalid memory access.
- If the reference was invalid, we terminate the process. If it was valid but we have not yet brought in that page, we now page it in.
- We find a free frame (by taking one from the free-frame list, for example).
- We schedule a secondary storage operation to read the desired page into the newly allocated frame.
- When the storage read is complete, we modify the internal table kept with the process and the page table to indicate that the page is now in memory.
- We restart the instruction that was interrupted by the trap. The process can now access the page as though it had always been in memory.

LRU = LEAST RECENTLY USED

When a page frame is full, then replace the least recently used request with the current page request.

| PAGE REQUESTS | | | | | | | | | | | | | | | | | | |
|---------------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 7 | 2 | 3 | 1 | 2 | 5 | 3 | 4 | 6 | 7 | 1 | 0 | 5 | 4 | 6 | 2 | 3 | 0 | 1 |
| 7 | 1 | 1 | 1 | 1 | 3 | 3 | 3 | 7 | 7 | 1 | 1 | 5 | 5 | 2 | 2 | 1 | | |
| 2 | 2 | 2 | 2 | 2 | 2 | 2 | 4 | 4 | 4 | 1 | 1 | 4 | 4 | 3 | 3 | | | |
| 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 6 | 6 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| F | F | F | F | V | F | F | F | V | F | F | F | F | F | F | F | F | F | F |

3 page frames.

FIFO = FIRST IN, FIRST OUT

When a page frame is full, then replace the first request with the current page request.

| PAGE REQUESTS | | | | | | | | | | | | | | | | | | |
|---------------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 7 | 1 | 2 | 3 | 1 | 2 | 5 | 3 | 4 | 6 | 7 | 1 | 0 | 5 | 4 | 6 | 2 | 3 | 0 |
| 7 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 2 | 2 | 2 | 2 | 5 | 5 | 5 | 7 | 7 | 7 | 1 | 1 | 5 | 5 | 2 | 2 | 1 | | |
| 3 | 3 | 3 | 3 | 3 | 3 | 4 | 4 | 4 | 1 | 1 | 1 | 4 | 4 | 3 | 3 | | | |
| F | F | F | V | F | V | F | F | V | V | F | F | F | F | F | F | F | F | F |

3 page frames.

Optimal Replacement

When a page frame is full, then replace page we don't use for the longest time w/ current page request.

| PAGE REQUESTS | | | | | | | | | | | | | | | | | | |
|---------------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 7 | 1 | 2 | 3 | 1 | 2 | 5 | 3 | 4 | 6 | 7 | 1 | 0 | 5 | 4 | 6 | 2 | 3 | 0 |
| 7 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 2 | 2 | 2 | 2 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 5 | 4 | 6 | 7 | 7 | 1 | 1 | 1 |
| 3 | 3 | 3 | 3 | 3 | 3 | 4 | 4 | 4 | 1 | 1 | 1 | 4 | 4 | 3 | 3 | 3 | 3 | 3 |
| F | F | F | V | F | V | F | F | V | V | F | F | F | F | F | F | F | F | F |

3 page frames.

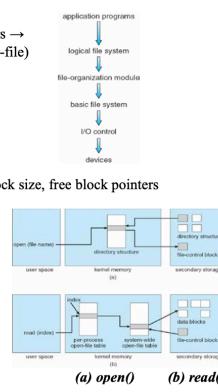
NS = 200 NANOSEC. = 0.0002 MILISPC.
MS = 200 NANOSEC. = 0.2 MICROSEC. US ↗

Chapter 14 - File System Implementation

Ch.12 – File System Implementation

- File system resides on secondary storage – disks; file system is organized into layers →
- File control block: storage structure consisting of information about a file (exist per-file)
- Device driver: controls the physical device; manage I/O devices
- File organization module: understands files, logical addresses, and physical blocks
 - Translates logical block number to physical block number
 - Manages free space, disk allocation
- Logical file system: manages metadata information – maintains file control blocks
- Boot control block: contains info needed by system to boot OS from volume
- Volume control block: contains volume details; ex: total # blocks, # free blocks, block size, free block pointers
- Root partition: contains OS; mounted at boot time
- For all partitions, system is consistency checked at mount time
 - Check metadata for correctness – only allow mount to occur if so
- Virtual file systems** provide object-oriented way of implementing file systems
- Directories can be implemented as **Linear Lists** or **Hash Tables**
 - Linear list of file names with pointer to data blocks – simple but slow
 - Hash table – linear list with hash data structure – decreased search time
 - Good if entries are fixed size
 - Collisions can occur in hash tables when two file names hash to same location
- Contiguous allocation**: each file occupies set of contiguous blocks
 - Simple, best performance in most cases; problem – finding space for file, external fragmentation
 - Extent based file systems are modified contiguous allocation schemes – extent is allocated for file allocation
- Linked Allocation**: each file is a linked list of blocks – no external fragmentation
 - Locating a block can take many I/Os and disk seeks
- Indexed Allocation**: each file has its own index block(s) of pointers to its data blocks
 - Need index table; can be random access; dynamic access without external fragmentation but has overhead
- Best methods: linked good for sequential, not random; contiguous good for sequential and random
- File system maintains **free-space list** to track available blocks/clusters
- Bit vector or bit map (n blocks): block number calculation → (#bits/word)*(# 0-value words)+(offset for 1st bit)
- Example:


```
block size = 4KB = 212 bytes
disk size = 240 bytes (1 terabyte)
n = 240/212 = 228 bits (or 256 MB)
if clusters of 4 blocks > 64MB of memory
```
- Space maps (used in ZFS) divide device space into **metaslab** units and manages metaslabs
 - Each metaslab has associated space map
- Buffer cache – separate section of main memory for frequently used blocks
- Synchronous writes sometimes requested by apps or needed by OS – no buffering
- Asynchronous writes are more common, bufferable, faster
- Free-behind and read-ahead techniques to optimize sequential access
- Page cache caches pages rather than disk blocks using virtual memory techniques and addresses
 - Memory mapped I/O uses page cache while routine I/O through file system uses buffer (disk) cache
- Unified buffer cache: uses same page cache to cache both memory-mapped pages and ordinary file system I/O to avoid double caching



Virtual File Systems (1)

- File systems have general and storage method-dependent parts
- Virtual file system** is specific to the OS
 - file system-generic operations
 - works with inodes (FCB), files, directories, superblocks (partitions)
 - the stuff that we have discussed
- Physical file system** is specific to how secondary storage will be used to manage data
 - converts the objects above into blocks

Virtual File Systems (2)

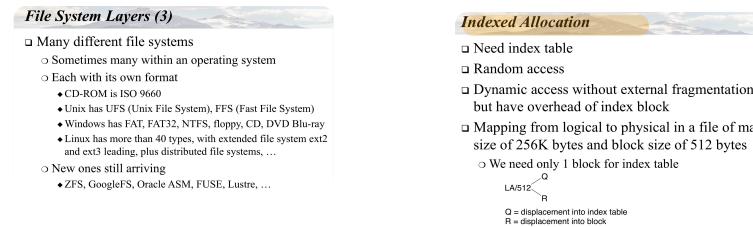
- Virtual File Systems (VFS)** on Unix provide an object-oriented way of implementing file systems
- VFS allows the same system call interface (the API) to be used for different types of file systems
 - Separates file-system generic operations from implementation details
 - Implementation can be one of many file system types, or network file system
 - implements **vnodes** which hold inodes or network file details
 - Then dispatches operation to appropriate file system implementation routines

File System Layers (1)

- Device drivers manage I/O devices at the **I/O control layer**
 - High-level commands like:
 - read drive1, cylinder 72, track 2, sector 10, into memory location 1060
 - Translate into low-level hardware specific commands to hardware controller
- Basic file system given command like “retrieve block 123” translates to device driver
- Also manages memory buffers and caches (allocation, freeing, replacement)
- Buffers hold data in transit
- Caches hold frequently used data
- File organization module** understands files, logical address, and physical blocks
 - Translates logical block # to physical block #
 - Manages free space, disk allocation

File System Layers (3)

- Many different file systems
 - Sometimes many within an operating system
 - Each with its own format
 - CD-ROM is ISO 9660
 - Unix has UFS (Unix File System), FFS (Fast File System)
 - Windows has FAT, FAT32, NTFS, floppy, CD, DVD Blu-ray
 - Linux has more than 40 types, with extended file system ext2 and ext3 leading, plus distributed systems, ...
 - New ones still arriving
 - ZFS, GoogleFS, Oracle ASM, FUSE, Lustre, ...



OSC 14.3: Consider a system that supports the strategies of contiguous, linked, and indexed allocation. What criteria should be used in deciding which strategy is best utilized for a particular file?

Answer: Contiguous allocation requires that each file occupy a set of contiguous blocks on disk. It facilitates file access and can be higher performing, but can cause problems with finding space for new files. Short files that are long-lived and do not fluctuate in size over time are good candidates for contiguous allocation.

Linked allocation can address contiguous allocation problems by maintaining a linked list of disk blocks that make up a file. It is easy to extend the file by adding more blocks to the list. It is effective for files of

medium size that might grow and shrink, but is really for sequential access only because positioning within a file requires traversing the links to find the right block.

Indexed allocation uses an index directory to find blocks for purposes of rapid positioning. It is better to use this strategy for larger files that are supporting non-sequential access.

OSC 14.4: One problem with contiguous allocation is that the user must preallocate enough space for each file. If the file grows to be larger than the space allocated for it, special actions must be taken. One solution to this problem is to define a file structure consisting of an initial contiguous area of a specified size. If this area is filled, the operating system automatically defines an overflow area that is linked to the initial contiguous area. If the overflow area is filled, another overflow area is allocated. Compare this implementation of a file with the standard contiguous and linked implementations.

Answer: The solution is certainly an improvement over the standard contiguous allocation since it solves the problem mentioned, allowing the file to grow in size. It has similarities to the linked allocation, but does not necessarily require a particular blocking factor to be used.

OSC 14.11: Some file systems allow disk storage to be allocated at different levels of granularity. For instance, a file system could allocate 4KB of disk space as a single 4KB block or as eight 512B blocks. How could we take advantage of this flexibility to improve performance? What modifications would have to be made to the free-space management scheme in order to support this feature?

Answer: There are 2 general performance benefits that come from variable granularity in storage allocation. First, there is the benefit of smaller granularity in storage efficiency, due to less potential internal fragmentation. Second, larger granularity can benefit larger files, especially for sequential access, because larger blocks of data can be moved on/off the disk at a time, without having to reposition to another block. Essentially, the free-space management scheme will have to be modified cognizant of the block size. It could have a fixed partitioning of the disk for blocks of different sizes. Then allocation could be done for each granularity using only the blocks of that size. The free-space management then would be the same as before, except for each level of granularity.

Another strategy is to allocate everything at a small granularity, but allow contiguous smaller granularity blocks to be dynamically combined into a larger granularity block. This would make the free-space management more complicated as it would keep track of all of the free small granularity blocks, but also store information to indicate if a block's neighbor was also free. Here, the allocator might use some optimization to represent free blocks of different granularities and to decide which blocks to allocate to retain larger blocks in the free list.

14.9 Summary

- Most file systems reside on **secondary storage**, which is designed to hold a large amount of data permanently. The most common secondary-storage medium is the **disk**, but the use of **NVM devices** is increasing.
- Storage devices are segmented into **partitions** to control media use and to allow multiple, possibly varying, file systems on a single device. These file systems are mounted onto a logical file system architecture to make them available for use.
- File systems are often implemented in a layered or modular structure. The lower levels deal with the physical properties of storage devices and communicating with them. Upper levels deal with symbolic file names and logical properties of files.
- The various files within a file system can be allocated space on the storage device in three ways: through **contiguous**, **linked**, or **indexed allocation**. Contiguous allocation can suffer from external fragmentation. Direct access is very inefficient with linked allocation. Indexed allocation may require substantial overhead for its index block. These algorithms can be optimized in many ways. Contiguous space can be enlarged through extents to increase flexibility and to decrease external fragmentation. Indexed allocation can be done in clusters of multiple blocks to increase throughput and to reduce the number of index entries needed. Indexing in large clusters is similar to contiguous allocation with extents.
- Free-space allocation methods also influence the efficiency of disk-space use, the performance of the file system, and the reliability of secondary storage. The methods used include bit vectors and linked lists. Optimizations include grouping, counting, and the FAT, which places the linked list in one contiguous area.
- Directory-management routines must consider efficiency, performance, and reliability. A hash table is a commonly used method, as it is fast and efficient. Unfortunately, damage to the table or a system crash can result in inconsistency between the directory information and the disk's contents.
- A consistency checker can be used to repair damaged file-system structures. Operating-system backup tools allow data to be copied to magnetic tape or other storage devices, enabling the user to recover from data loss or even entire device loss due to hardware failure, operating system bug, or user error.
- Due to the fundamental role that file systems play in system operation, their performance and reliability are crucial. Techniques such as log structures and caching help improve performance, while log structures and RAID improve reliability. The WAFL file system is an example of optimization of performance to match a specific I/O load.

Contiguous Allocation

- Contiguous allocation is where each file occupies set of contiguous blocks
- Advantages:**
 - Simple – need to remember only starting location to access any block, plus length (number of blocks)
 - Good performance when reading successive blocks on disk
 - Best performance in most cases
- Disadvantages:**
 - File size has to be known a priori
 - need to find space for the file of that size
 - External fragmentation

Linked Allocation – (1)

- Linked allocation allocates a linked list of blocks to a file
 - Keep a pointer to the first block of a file
 - Each block contains pointer to next block
 - File ends at nil pointer
 - Free space management system called when new block needed
- Advantages:**
 - No external fragmentation
 - No compaction
- Disadvantages:**
 - Reliability can be a problem
 - Locating a block can take many I/Os and disk seeks

Chapter 15 File System Internals

- **tmpfs**—a “temporary” file system that is created in volatile main memory and has its contents erased if the system reboots or crashes
- **objfs**—a “virtual” file system (essentially an interface to the kernel that looks like a file system) that gives debuggers access to kernel symbols
- **cfts**—a virtual file system that maintains “contract” information to manage which processes start when the system boots and must continue to run during operation
- **lofs**—a “loop back” file system that allows one file system to be accessed in place of another one
- **procfs**—a virtual file system that presents information on all processes as a file system
- **ufs, zfs**—general-purpose file systems

The file systems of computers, then, can be extensive. Even within a file system, it is useful to segregate files into groups and manage and act on those groups. This organization involves the use of directories (see Section 14.3).

Page cache: cache of virtual memory pages including memory-mapped files

Buffer cache: stores blocks from disk, now unified with the page cache (why?)

Page scanner: schedules dirty pages to be written to disk

Dentry cache: remembers mappings from directory entry to inode (for what purpose?)

Inode cache: stores inodes in hash table for fast lookup

An inode contains:

Caches:

- File attributes (time of creation, permissions, ...)
- 10 direct pointers (logical disk block ids)
- 1 one-level indirect pointer (points to a disk block which in turn contains pointers)
- 1 two-level indirect pointer (points to a disk block of pointers to disk blocks of pointers)
- 1 three-level indirect pointer (points to a disk block of pointers to disk blocks of pointers to disk blocks)

Finding Files

- Now we know how to retrieve the blocks of a file once we know:
 - The FAT entry for DOS
 - The inode of the file in UNIX

□ But how do we find these in the first place?

- The directory where this file resides should contain this information

NFS Protocol

- Provides a set of remote procedure calls for remote file operations. The procedures support the following operations:
 - Searching for a file within a directory
 - Reading a set of directory entries
 - Manipulating links and directories
 - Accessing file attributes
 - Reading and writing files
- NFS servers are stateless; each request has to provide a full set of arguments (NFS V4 is different, stateful)
- Modified data must be committed to the server’s disk before results are returned to the client (lose advantages of caching)
- NFS protocol does not provide concurrency-control mechanisms

Three Major Layers of NFS Architecture

- UNIX file-system interface (based on the open, read, write, and close calls, and file descriptors)
- Virtual File System (VFS) layer – distinguishes local files from remote ones, and local files are further distinguished according to their file-system types
 - The VFS activates file-system-specific operations to handle local requests according to their file-system types
 - Calls the NFS protocol procedures for remote requests
- NFS service layer – bottom layer of the architecture
 - Implements the NFS protocol

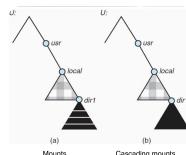
Sun Network File System (NFS) (1)

- Implementation and a specification of a software system for accessing remote files across LANs (or WANs)
- The implementation is part of the Solaris and SunOS operating systems running on Sun workstations using an unreliable datagram protocol (UDP/IP) and Ethernet

Sun Network File System (NFS) (2)

- Interconnected workstations viewed as a set of independent machines with independent file systems, which allows sharing among these file systems in a transparent manner
 - A remote directory is mounted over a local file system directory
 - mounted directory looks like an integral subgraph of the local file system, replacing the subtree descending from the local directory
 - Specification of the remote directory for the mount operation is nontransparent; the host name of the remote directory has to be provided
 - files in the remote directory can then be accessed in a transparent manner
 - Subject to access-rights accreditation, potentially any file system (or directory within a file system), can be mounted remotely on top of any local directory

MOUNTS:



NFS Mount Protocol

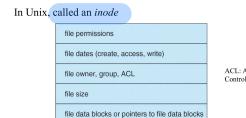
- Establishes initial logical connection between server and client
- Mount operation includes name of remote directory to be mounted and name of server machine storing it
 - Mount request is mapped to a response RPC or forwarded to mount daemon on server machine
- Export list – specifies local file systems that server exports for mounting, along with names of machines that are permitted to mount them
- Following a mount request that conforms to its export list, the server returns a file handle—a key for further accesses
- File handle – a file-system identifier, and an inode number to identify the mounted directory within the exported file system
- The mount operation changes only the user’s view and does not affect the server side

Just as a file must be opened before it can be used, a file system must be mounted before it can be available to processes on the system. More specifically, the directory structure may be built out of multiple file-system-containing volumes, which must be mounted to make them available within the file system name space.

The mount procedure is straightforward. The operating system is given the name of the device and the **mount point**—the location within the file structure where the file system is to be attached. Some operating systems require that a file-system type be provided, while others inspect the structures of the device

Performance (1)

- Best method depends on file access type
 - Contiguous great for sequential and random
 - Linked good for sequential, not random
- Declare access type at creation
 - Select either contiguous or linked
- Indexed more complex
 - Single block access could require 2 index block reads then data block read
 - Clustering can help improve throughput, reduce CPU overhead



- A race condition occurs when processes have concurrent access to shared data and the final result depends on the particular order in which concurrent accesses occur. Race conditions can result in corrupted values of shared data.
- A critical section is a section of code where shared data may be manipulated and a possible race condition may occur. The critical-section problem is to design a protocol whereby processes can synchronize their activity to cooperatively share data.
- A solution to the critical-section problem must satisfy the following three requirements: (1) mutual exclusion, (2) progress, and (3) bounded waiting. Mutual exclusion ensures that only one process at a time is active in its critical section. Progress ensures that programs will cooperatively determine what process will next enter its critical section. Bounded waiting limits how much time a program will wait before it can enter its critical section.
- Software solutions to the critical-section problem, such as Petersen’s solution, do not work well on modern computer architectures.
- Hardware support for the critical-section problem includes memory barriers; hardware instructions, such as the compare-and-swap instruction; and atomic variables.
- A mutex lock provides mutual exclusion by requiring that a process acquire a lock before entering a critical section and release the lock on exiting the critical section.
- Semaphores, like mutex locks, can be used to provide mutual exclusion. However, whereas a mutex lock has a binary value that indicates if the lock is available or not, a semaphore has an integer value and can therefore be used to solve a variety of synchronization problems.
- A monitor is an abstract data type that provides a high-level form of process synchronization. A monitor uses condition variables that allow processes to wait for certain conditions to become true and to signal one another when conditions have been set to true.
- Solutions to the critical-section problem may suffer from liveness problems, including deadlock.
- The various tools that can be used to solve the critical-section problem as well as to synchronize the activity of processes can be evaluated under varying levels of contention. Some tools work better under certain contention loads than others.

Flock – File System Consistency Check

- Blocks:
 - for every block keep 2 counters:
 - a) # occurrences in files
 - b) # occurrences in free list.
 - For every inode, increment all the (a)s for the blocks that the file covers.
 - For the free list, increment (b) for all blocks in the free list.
 - Ideally (a) + (b) = 1 for every block.
 - However:
 - If (a) + (b) < 0, missing block, add to free list.
 - If (a) = (b) = 1, remove the block from free list
 - If (b) > 1, remove duplicates from free list
 - If (a) > 1, make copies of this block, and insert into each of the other files.

Sun Network File System (NFS) (3)

- NFS is designed to operate in a heterogeneous environment of different machines, operating systems, and network architectures; the NFS specifications are independent of these media
- This independence is achieved through the use of RPC primitives built on top of an External Data Representation (XDR) protocol used between two implementation-independent interfaces
- The NFS specification distinguishes between the services provided by a mount mechanism and the actual remote-file-access services

- Volumes containing file systems can be mounted into the computer’s file-system space.
- Depending on the operating system, the file-system space is seamless (mounted file systems integrated into the directory structure) or distinct (each mounted file system having its own designation).
- At least one file system must be bootable for the system to be able to start—that is, it must contain an operating system. The boot loader is run first; it is a simple program that is able to find the kernel in the file system, load it, and start its execution. Systems can contain multiple bootable partitions, letting the administrator choose which to run at boot time.
- Most systems are multi-user and thus must provide a method for file sharing and file protection. Frequently, files and directories include metadata, such as owner, user, and group access permissions.
- Mass storage partitions are used either for raw block I/O or for file systems. Each file system resides in a volume, which can be composed of one partition or multiple partitions working together via a volume manager.
- To simplify implementation of multiple file systems, an operating system can use a layered approach, with a virtual file-system interface making access to possibly dissimilar file systems seamless.
- Remote file systems can be implemented simply by using a program such as **ftp** or the web servers and clients in the World Wide Web, or with more functionality via a client–server model. Mount requests and user IDs must be authenticated to prevent unapproved access.
- Client–server facilities do not natively share information, but a distributed information system such as DNS can be used to allow such sharing, providing a unified user name space, password management, and system identification. For example, Microsoft CIFS uses active directory, which employs a version of the Kerberos network authentication protocol to provide a full set of naming and authentication services among the computers in a network.
- Once file sharing is possible, a consistency semantics model must be chosen and implemented to moderate multiple concurrent access to the same file. Semantics models include UNIX, session, and immutable-shared-files semantics.
- **NFS** is an example of a remote file system, providing clients with seamless access to directories, files, and even entire file systems. A full-featured remote file system includes a communication protocol with remote operations and path-name translation.
- General-purpose operating systems provide many file-system types, from special-purpose through general.

- Classic problems of process synchronization include the bounded-buffer, readers–writers, and dining-philosophers problems. Solutions to these problems can be developed using the tools presented in Chapter 6, including mutex locks, semaphores, monitors, and condition variables.
- Windows uses dispatcher objects as well as events to implement process synchronization tools.
- Linux uses a variety of approaches to protect against race conditions, including atomic variables, spinlocks, and mutex locks.
- The POSIX API provides mutex locks, semaphores, and conditioned variables. POSIX provides two forms of semaphores: named and unnamed. Several unrelated processes can easily access the same named semaphore by simply referring to its name. Unnamed semaphores cannot be shared as easily, and require placing the semaphore in a region of shared memory.
- Java has a rich library and API for synchronization. Available tools include monitors (which are provided at the language level) as well as reentrant locks, semaphores, and condition variables (which are supported by the API).
- Alternative approaches to solving the critical-section problem include transactional memory, OpenMP, and functional languages. Functional languages are particularly intriguing, as they offer a different programming paradigm from procedural languages. Unlike procedural languages, functional languages do not maintain state and therefore are generally immune from race conditions and critical sections.

Log Structured File Systems

- Log structured (or journaling) file systems record each metadata update to the file system as a **transaction**
- All transactions are written to a log
 - A transaction is considered committed once it is written to the log (sequentially)
 - Sometimes to a separate device or section of disk
 - However, the file system may not yet be updated
- The transactions in the log are asynchronously written to the file system structures
 - When the file system structures are modified, the transaction is removed from the log
- If the file system crashes, all remaining transactions in the log must still be performed
- Faster recovery from crash, removes chance of inconsistency of metadata

NFS Remote Operations

- Nearly one-to-one correspondence between regular UNIX system calls and the NFS protocol RPCs (except opening and closing files)
- NFS adheres to the remote-service paradigm, but employs buffering and caching techniques for the sake of performance
- File-blocks cache – when a file is opened, the kernel checks with the remote server whether to fetch or revalidate the cached attributes
 - Cached file blocks are used only if the corresponding cached attributes are up to date
- File-attribute cache – the attribute cache is updated whenever new attributes arrive from the server
- Clients do not free delayed-write blocks until the server confirms that the data have been written to disk

