

Slide 1

Administrivia

- Reminder: Homework 5 due today. (Questions?)
- Quiz 5 a week from today.

Slide 2

Sharing Pages

- Shared pages can be useful, but can also present problems.
- Multiple processes running the same program is relatively easy (why?) but has one potential downside (what?)
- UNIX `fork` system call is — interesting in this context. POSIX definition says that child process's address space is basically a copy of the parent's address space. What's the easy-to-implement way to do this? What downside does that have in current systems? Is there a way to reduce its impact? And why duplicate in the first place?

Sharing Pages and `fork`

Slide 3

- Duplicating pages is easy but inefficient, especially if the child process is going to call `execve` or something similar right away. Some systems use “copy-on-write” to improve efficiency.
- Why did the people who designed UNIX require this duplication . . . Possibly because it makes some things easy (such as setting up parent/child pipes) and wasn't very costly when designed. Windows' system call for creating processes takes a different approach. Maybe that's better!

Sharing Pages, Continued

Slide 4

- One use for shared pages is multiple processes running the same program.
- What about sharing code at a level below whole programs (UNIX “shared libraries”, Windows DLLs)? Seems attractive; are there potential problems?

Shared Libraries

Slide 5

- One attraction is somewhat obvious — if code for library functions (e.g., `printf`) is statically linked into every program that uses it, programs need more memory — seems wasteful if processes can share one copy of code in memory.
- Another attraction is that library code can be updated independently of programs that use it. (Is there a downside to that?)
- How to make this happen . . . At link time, programs get “stub” versions of functions. References to real versions resolved at load time. Does this remind you of anything? and suggest a possible problem? how to fix?

Shared Libraries, Continued

Slide 6

- Downside of replacing shared libraries — may break applications that call their function. UNIX provides a way around this.
- Resolving references to shared code at load time — finer-grained version of “relocation problem”, no? and fixable by making sure library contains only “position-independent code”.

Slide 7

Memory-Mapped File I/O

- Worth mentioning here that some systems also provide a mechanism (e.g., via system calls) to allow reading/writing whole files into/from memory. If there's enough memory, this could improve performance.
- Example of how this works in Linux — `man` page for `mmap`.

Slide 8

Paging — One More Hardware Issue

- (Not discussed in class but worth reading about.)
- What if page to be replaced is waiting for I/O? probably trouble if we replace it anyway.
- One solution — allow pages to be “locked”.
- Another solution — do all I/O to o/s pages, then move to user pages.

Slide 9

Processing Memory References — Details Still To Fill In

- How to keep track of pages on disk.
- How to keep track of which page frames are free.
- How to “schedule I/O” (but that’s later).

Slide 10

Keeping Track of Pages on Disk

- To implement virtual memory, need space on disk to keep pages not in main memory. Reserve part of disk for this purpose (“swap space”); (conceptually) divide it into page-sized chunks. How to keep track of which pages are where?
- One approach — give each process a contiguous piece of swap space. Advantages/disadvantages?
- Another approach — assign chunks of swap space individually. Advantages/disadvantages?
- Either way — processes must know where “their” pages are (via page table and some other data structure), operating system must know where free slots are (in memory and in swap space).

Slide 11

One More Memory Management Strategy — Segmentation

- (Not discussed in class but worth reading about.)
- Idea — make program address “two-dimensional” / separate address space into logical parts. So a virtual address has two parts, a segment and an offset.
- To map virtual address to memory location, need “segment table”, like page table except each entry also requires a length/limit field. (So this is like a cross between contiguous-allocation schemes and paging.)

Slide 12

Segmentation, Continued

- Benefits?
 - Nice abstraction; nice way to share memory.
 - Flexible use of memory — can have many areas that grow/shrink as required, not just heap and stack — especially if we combine with paging.
- Drawbacks?
 - External fragmentation possible (can offset by also paging).
 - More complex.
 - “Paging” in/out more complex — issues similar to with contiguous-allocation.

Slide 13

Memory Management in Windows

- Apparently very complex, but basic idea is paging.
- Intraprocess memory management is in terms of code regions (some shared — DLLs), data regions, stack, and area for o/s. “Virtual Address Descriptor” for each contiguous group of pages tracks location on disk, etc.
- Memory-mapped files can make I/O faster and allow processes to (in effect) share memory.
- Demand-paged, with six (!) background threads that try to maintain a store of free page frames. Page replacement algorithm is based on idea of working set.

Slide 14

Memory Management in UNIX/Linux

- Very early UNIX used contiguous-allocation or segmentation with swapping. Later versions use paging. Linux uses multi-level page tables; details depend on architecture (e.g., three levels for Alpha, two for Pentium).
- Intraprocess memory management is in terms of text (code) segment, data segment, and stack segment. Linux reserves part of address space for o/s. For each contiguous group of pages, “vm_area_struct” tracks location on disk, etc.
- Memory-mapped files can make I/O faster and allow processes to (in effect) share memory.
- Demand-paged, with background process (“page daemon”) that tries to maintain a store of free page frames. Page replacement algorithms are mostly variants of clock algorithm.

Minute Essay

- Anything about memory management you'd like to hear more about / have clarified?

Slide 15