

Slide 1

Administrivia

- Reminder: Quiz 5 Tuesday. Topics from Chapter 4.
- Reminder: Homework 4 due Tuesday. There will likely be one more homework after all, at least one programming problem, to be assigned by next Tuesday and due around the date of the final.

Slide 2

Memory Hierarchy — Overview

- Significant overlap between Chapter 5 and material covered in operating-systems course (as I teach it anyway). In previous years most students went on to that course. Now possibly not. Oh well!
- A key idea (borrowed from one writer of o/s textbooks): In a perfect world, we could have as much memory as we wanted, and it would be very fast and very cheap. In the real world, there are tradeoffs (e.g., fast versus cheap, fast versus large).

“Principle of Locality”

Slide 3

- Basic underlying idea — most applications exhibit locality with regard to memory.
- “Temporal locality” — memory locations referenced in the near past are likely to be referenced again in the near future.
- “Spatial locality” — memory locations close together in space likely to be referenced close together in time.

Memory Hierachy and Caching

Slide 4

- To exploit temporal locality, can use “caching” — keep copies of frequently-used data in faster but smaller memory. Can do this on multiple levels.
- To exploit spatial locality, can move data between levels in blocks.
- Terminology — cache hit, cache miss, cache block/line.
- *Notice* that while impact of caching on performance can be significant, it should not affect results (which is why it makes some sense to just ignore it initially).

Caches (Between Processor and RAM) — Executive-Level Summary

Slide 5

- To make these work, we need:
 - Some way to map memory address to cache location — can be simple (“direct map”) or not.
 - Some way to say, for each cache location, what memory address it's currently associated with, and whether the data is valid.
- Read “from memory” tries cache first, and then if not found there goes to RAM and updates cache.
- Write “to memory” is maybe more interesting — writes to cache but then must at some point write to RAM also — maybe right away (easier to get right but can be slow) or later.

Virtual Memory — Executive-Level Summary

Slide 6

- Basic idea here is to fake having more RAM than you really have, by keeping some data that would be in RAM on disk. In a sense, RAM is a cache for the “real” memory, on disk(!).
- Also provides a nice way to support multitasking — notion of “processes”, each with its own “address space”, with an operating system that maps this abstraction onto the hardware, by mapping “program addresses” (in a process's address space) to “physical addresses” (in RAM). *Lots of details here, but the basic idea is fairly simple. One big advantage is more control over what data each process can access.*

Cache Coherence — Executive-Level Summary

Slide 7

- Clearly(?) possible for cached data to be out of synch with data in memory.
- First step toward managing possible impact is to clearly define what we want, and then figure out how to either solve the possible problems or work around them.
- Clearly(?) much more difficult with multiple processing elements each of which has its own cache.

Caches and Applications Programming

Slide 8

- Mostly the memory hierarchy (including virtual memory) is managed transparently by a combination of hardware and (operating-system) software, so the first approximation presented in introductory courses (memory is essentially a really big array of bytes, with addresses as indices) is okay, especially if you just want right answers.
- However, effects on performance can be significant, so if you want right answers fast . . .

For single-threaded programs, key idea is to maximize locality (temporal and spatial). Rearranging order in which data is accessed can have a big effect.

For multi-threaded programs, also need to consider whether multiple threads need to share access to the same data (problem for correctness too!) or even nearby data (“false sharing” — no effect on correctness but can be slow).

Minute Essay

Slide 9

- A question about material from previous lecture(s): Many processors have a notion of two modes of operation, a privileged one for when they're doing operating-system stuff and an unprivileged one for regular applications. Attempts to do privileged things while in unprivileged mode generate exceptions. What if anything can you say about how this might help in making the whole system (hardware plus software) robust and secure? (Speculate!)
- How much of today's discussion was familiar?

Minute Essay Answer

Slide 10

- If regular applications execute in unprivileged mode, the hardware can enforce some restrictions on what they can do (e.g., only request I/O by going through the operating system). How do you get from unprivileged mode to privileged mode then? As part of exception processing — hardware transfers control to fixed location(s) and switches to privileged mode.