

Slide 1

Administrivia

- Reminder: Homework 6 due today.
- Reminder: Exam 2 next Monday. Review sheet on Web.
Right now I'm planning to have this one cover only material since Exam 1, but there *has* been less of it. How many would prefer to have some questions from material also covered earlier?
- Reminder: Quiz 6 next time. Likely topic is high-level questions about pipelining.
- Reminder: If you want (virtual?) attendance points for the days for which lectures were on video, send me a minute-essay e-mail.

Slide 2

Homework 6 Help — Tracing Operation of the Processor Circuit

- For the first problem, my intent was that you would trace through what the circuit in Figure 4.17 is actually doing rather than what you think it's supposed to be doing.
- So, you start with what you know — current saved value of the PC and what's at that address (in instruction memory) and contents of selected registers and data memory locations — and work from there. Taking the first few steps . . .
- Right away you can write down output of PC and input/output of instruction memory. For output of instruction memory, it's probably more helpful to write down the various fields of the instruction rather than the hex value of the whole instruction. You can do this the same way you did earlier (translating MIPS assembler language to machine language).
- (Continued on next slide.)

Tracing Operation of the Processor Circuit, Continued

Slide 3

- Now you can write down all the control signals, the inputs and output of the top left adder, and the register-number inputs to the register file. You can get the control signals from the table in Figure 4.18.
- Once you have those, you can write down outputs of the register file and start figuring out what the main ALU is doing. You can also determine whether the top right adder and the data memory will be used (based on control signals).
- Figuring out what the ALU does . . . You need to determine what operation it's doing (based on the `ALUop` control signal and the instruction function field, as shown in Figure 4.13). You also need to determine what the second operand is (contents of a register? sign-extended value from instruction?), again using control signals.
- "And so forth" . . .

Memory Hierarchy — Overview

Slide 4

- Significant overlap between Chapter 5 and material covered in operating-systems course (as I teach it anyway). In previous years pretty much all students went on to that course. Now not all do. Either way, not a bad idea to discuss briefly now.
- 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 5

- 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. (At any given time, a program isn’t going to be working with *all* of its data.)
- “Spatial locality” — memory locations close together in space likely to be referenced close together in time. (Examples include processing arrays sequentially, accessing local variables, which are apt to be located together in memory).

Memory Hierachy and Caching

Slide 6

- 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.
- *Note* 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 7

- Idea here is to interpose a “cache” (small but fast) between the processor and the memory, and use it to hold frequently-referenced data, and have this managed mostly by the hardware.
- 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 8

- 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(!).
- It turns out that a common way to do this also facilitates protecting one process’s data from other processes.

Cache Coherence — Executive-Level Summary

- Clearly(?) possible for cached data to be out of synch with data in memory. Probably of most concern if multiple processing elements, each with its own cache, share memory.
- Various schemes exist for ensuring that programs don't have to be aware of this complication. Details in the textbook.

Slide 9

Caches and Applications Programming

- 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. (Matrix-multiplication example.)
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).

Slide 10

Minute Essay

- None — quiz.

Slide 11