

Slide 1

Administrivia

- Reading Quiz 2 posted; due at end of the week.
- Homework 1 posted; due next Monday. Separated into two parts: 1a is written problems; 1b is programming problem.

Slide 2

Limited Direct Execution — Recap

- First an unrelated bit of usage/terminology: I'll use O/S to mean operating system henceforth!
- "Direct execution" means as opposed to, say, emulation.
- "Limited" in that we want to do *something* to ensure that the O/S can "defend itself" and also protect processes from each other.
Note that this is only even possible if hardware provides some support. May explain why early O/S's for PCs didn't always provide this kind of safety (as in my old war story!).

Slide 3

Limited Direct Execution — Review/Clarification

- Textbook figure 6.2 very helpful, but I'm skeptical of some details — seems to indicate that O/S always returns to caller after system call, but clearly that can't be true if sometimes it terminates the process!
- "Trap table" is a name I had not encountered before, and I wonder about the name. Could it be specific to x86?
- More broadly: System calls (whatever the name is — trap, MIPS `syscall`, etc.) are a type of interrupt. There are other kinds of interrupts. Part of the interaction between the O/S and the hardware is the address(es) of handlers for various kinds of interrupts — possibly only one, or possibly different ones for different interrupts. Details might vary among architectures, but in general, textbook is (as far as I know) right that the O/S has to set this up at boot time — if nothing else, put its code at the hardware-specified fixed address.

Slide 4

Limited Direct Execution — Review/Clarification, Continued

- Textbook figure 6.3 also helpful, though initially I was somewhat skeptical about details being applicable to all O/S's and architectures.
However, on reflection it makes sense to talk about two separate save/restore operations:
If the scheduler says "keep running the interrupted process" then there's only a need to restore any machine state the interrupt handler messed up. If it says "switch processes" then more may be needed.
- In any case, in my usage (and I think this is standard), "context switch" usually refers to what happens when the O/S switches from one process to another.

Slide 5

“What About Concurrency?” Indeed

- Last section in the chapter on limited direct execution is well-titled. Very complicated topic, and not a bad idea to address it later and not as part of a discussion of “process management”, as some textbooks do.
- For now worth mentioning that indeed the question of what happens in the case of nested(?) interrupts is interesting. Two points:
 - They complicate things enormously (like recursion and how it can overflow the stack).
 - Could fix that by disabling further interrupts until we finish one. But that might mean missing one.
So what we do while interrupts disabled should be short!
- Also, as I understand things, some things that might generate interrupts will keep trying until they get an “acknowledged” response.

Slide 6

Concurrency — One More Thing

- A key problem — how to ensure that a sequence of actions happens, or appears to happen, as one “atomic” thing — i.e., without interference from anything else.
- This is what the textbook was getting at in talking about “atomically” — but I found their explanation unclear and possibly misleading.
- From an application programmer’s point of view, not guaranteeing atomicity of a sequence of operations can lead to race conditions, which can be solved via “locks”.
- Actually implementing locks is not as easy as it might sound! (Later.)

Scheduling — Overview

Slide 7

- Textbook likes to distinguish between “mechanisms” and “policies”. (And it *is* usually a good thing to separate them.)
- Previous chapters in this part have been about mechanisms for virtualizing the CPU. Some (which process to run next) require decision-making — i.e.(?), involve policies.
- (Textbook says the question of scheduling has a history before computers. Interesting!)
- Many policies possible; many have been tried over the years.

Scheduling — Simple View

Slide 8

- Scheduling algorithms (textbook calls them disciplines) usually based on “jobs” (units of work — name goes back to batch systems, where users submitted “jobs” to system operator, and there was no notion of interactive users).
- Textbook lays out some simplifying assumptions. I say we can start with slightly less restrictive ones:
 - Each job arrives at some predefined time.
 - Each job runs for some fixed predefined amount of time.
 - Once started, a job runs to completion (so, no switching back and forth among processes).
 - Jobs use only the CPU (i.e., no I/O).

Scheduling Algorithms — Metrics

- Useful to be able to compare different algorithms/policies/disciplines.
- A simple one is based on “turnaround time” (completion time minus arrival time) for jobs. Average turnaround time could be a simple metric for comparing.
- (To be continued.)

Slide 9

Minute Essay

- Are you finding the reading quizzes useful or interesting?

Slide 10