

Slide 1

Administrivia

- Reminder: Homework 1 due Monday (now split into 1a (written problems) and 1b (programming problem)).
- First reading quiz posted (at last)! Also due Monday.

Slide 2

Minute Essay From Last Lecture

- Most people got the intended answer. Most common error was not realizing that it's not really possible to have all 100 processes in the "ready" state — *something* would be running on each CPU.

Interprocess Communication

Slide 3

- Processes almost always need to interact with other processes:
 - “Ordering constraints” — e.g., process B uses as input some data produced by process A.
 - Use of shared resources — files, shared memory locations, etc.
- Use of shared resources can lead to “race conditions” — output depends on details of interleaving.
- Processes must communicate to avoid race conditions and otherwise synchronize.
- “Classical IPC problems” — simplified versions of things real programs (both applications and O/S) need to do.

Mutual Exclusion Problem

Slide 4

- In many situations, we want only one process at a time to have access to some shared resource.
- Generic/abstract version: Multiple processes, each with a “critical region” (“critical section”):

```
while (true) {  
    do_cr();           // must be "finite"  
    do_non_cr();      // need not be "finite"  
}
```
- Goal is to add something to this code such that:
 1. No more than one process at a time can be “in its critical region”.
 2. No process not in its critical region can block another process.
 3. No process waits forever to enter its critical region.
 4. No assumptions are made about how many CPUs, their speeds.

Mutual Exclusion Problem, Continued

Slide 5

- We'll look at various solutions (some correct, some not):
 - Using only hardware features always present (some notion of shared variable).
 - Using optional hardware features.
 - Using “synchronization mechanisms” (abstractions that help solve this and other problems).
- Recall that a correct solution
 - Must work for more than one CPU.
 - Must work even in the face of unpredictable context switches — whatever we're doing, another process can pull the rug out from under us between “atomic operations” (machine instructions).

Sidebar: Atomic Operations

Slide 6

- “Atomic” operation — indivisible, executes without interference from other processes.
- Which of the following are atomic?
 - `x = 1;`
 - `x = x + 1;`
 - `++x;`
 - `if (x == 0) x = 1;`(Or does it depend? On what?)

Proposed Solution — Disable Interrupts

- Pseudocode for each process:

```
while (true) {  
    disable_interrupts();  
    do_cr();  
    enable_interrupts();  
    do_non_cr();  
}
```

- Does it work? reviewing the criteria ...

Slide 7

Disable Interrupts, Continued

- (1) okay – context switches take place only in response to interrupts, so yes *if one CPU*.
- (4) not okay — fails if more than one CPU (unless there is a way to disable interrupts on all CPUs).
- Also, user-level programs shouldn't be able to do this (though might be okay for O/S).

Slide 8

Slide 9

Proposed Solution — Simple Lock Variable

- Shared variables:

```
int lock = 0;
```

Pseudocode for each process:

```
while (true) {  
    while (lock != 0);  
    lock = 1;  
    do_cr();  
    lock = 0;  
    do_non_cr();  
}
```

- Does it work? reviewing the criteria ...

Slide 10

Simple Lock Variable, Continued

- Can easily fail (1).

Proposed Solution — Strict Alternation

- Shared variables:

```
int turn = 0;
```

Pseudocode for process p0:

```
while (true) {  
    while (turn != 0);  
    do_cr();  
    turn = 1;  
    do_non_cr();  
}
```

Pseudocode for process p1:

```
while (true) {  
    while (turn != 1);  
    do_cr();  
    turn = 0;  
    do_non_cr();  
}
```

Slide 11

- Does it work? reviewing the criteria ...

Strict Alternation, Continued

- (Yes, we're simplifying to only two processes.)
- (1) okay.
- (2) / (3) not okay, since non-critical region need not be finite.

Slide 12

Slide 13

Sidebar: Reasoning about Concurrent Algorithms

- For concurrent algorithms (such as various solutions proposed for mutual exclusion problem), testing is less helpful than for sequential algorithms. (Why?)
- May be helpful, then, to try to think through whether they work. How? Idea of “invariant” may be useful:
 - Loosely speaking — “something about the program that’s always true”. (If this reminds you of “loop invariants” in CSCI 1323 — good.)
 - Goal is to come up with an invariant that’s easy to verify by looking at the code and implies the property you want (here, “no more than one process in its critical region at a time”).
 - We will do this quite informally, but it can be done much more formally — mathematical “proof of correctness” of the algorithm.

Slide 14

Sidebar of Sidebar: Reasoning About Loops

- (I probably won’t have time to through these slides in much detail in class but will leave them here for anyone interested.)
- Usually want to prove two things: (1) the loop eventually terminates, and (2) it establishes some desired postcondition.
- Proving that it terminates: Define a *metric* that you know decreases by some minimum amount with every trip through the loop, and when it goes below some threshold value, the loop ends.
- Proving that it establishes the postcondition: Use a *loop invariant*.
- (I say “prove” here, since this can be done very rigorously, but in practical situations an informal version is good enough.)

Slide 15

Reasoning About Loops, Continued

- What's a loop invariant? in the context of reasoning about programs, it's a *predicate* (boolean expression using program variables) that
 - is true before the loop starts, and
 - if true before a trip through the loop, with the loop condition true, is also true after the trip through the loop.

If you can prove that a particular predicate is a loop invariant, then after the loop exits, you know it's still true, and the loop condition is not. With a well-chosen invariant, this is enough to prove useful things.

- (Might be worth noting that compiler writers have a different definition — some computation that can be moved outside the loop.)

Slide 16

Reasoning About Loops, Simple Example

- Loop to compute sum of elements of array `a` of size `n`:

```
i = 0; sum = 0;
while (i != n) {
    sum = sum + a[i];
    i = i + 1;
}
```

At end, `sum` is sum of elements of `a`.

- Does this work? well, you probably believe it does, but you could prove it using the invariant:

`sum` is the sum of `a[0]` through `a[i-1]`

Reasoning About Loops, Example

Slide 17

- Euclid's algorithm for computing greatest common divisor of nonnegative integers a and b :

```

i = a; j = b;
while (j != 0) {
    q = i / j; r = i % j;
    i = j; j = r;
}

```

At end, $i = \text{gcd}(a, b)$.

- Does this work? work through some examples and gain some confidence — or prove using invariant:

$$\text{gcd}(i, j) = \text{gcd}(a, b)$$

and the math fact $\text{gcd}(n, 0) = n$

Strict Alternation, Revisited

Slide 18

- Shared variables:

```
int turn = 0;
```

Pseudocode for process p_0 :

```

while (true) {
    while (turn != 0);
    do_cr();
    turn = 1;
    do_non_cr();
}

```

Pseudocode for process p_1 :

```

while (true) {
    while (turn != 1);
    do_cr();
    turn = 0;
    do_non_cr();
}

```

- Proposed invariant: "If p_n is in its critical region, turn has value n , and turn is either 0 or 1" (interpreting "in its critical region" as "from just after the `while` to the line after `do_cr()`").

Slide 19

Strict Alternation, Continued

- Proposed invariant again: "If p_n is in its critical region, turn has value n , and turn is either 0 or 1".
- How would this help? would mean that if p_0 and p_1 are both in their critical regions, turn has two different values — impossible. So the first requirement would be met. (Still have to think about the other three.)
- Is it an invariant? check whether true initially and remains true even when one process changes something it mentions. Fairly obvious that it's initially true, so check ...

Slide 20

Strict Alternation, Continued

- Proposed invariant: "If p_n is in its critical region, turn has value n , and turn is either 0 or 1". True initially. When could it become false?
- When either process enters its critical region. But this happens for p_n only when turn is n , so invariant stays true (okay).
- When either process leaves its critical region. Also okay.
- When either process changes turn . Only happens after process leaves its critical region. So also okay.

Actual Solution(s)

- There are low-level solutions using only shared variables, highly portable but unsatisfactory in other ways.
- There are other low-level solutions that need hardware support, arguably better though still not great.
- And then there are "synchronization mechanisms".
- To be continued!

Slide 21

Minute Essay

- Did you learn about loop invariants in CSCI 1323 (Discrete Structures)?
- Tell me about your exposure to concurrent programming (multi-threading, message passing, etc. — anything involving multiple threads of control).

Slide 22