**Slide 1**

<span style="color:red">Administrivia</span>

- Homework 2 to be on the Web soon. I will send mail.

**Slide 2**

<span style="color:blue">Minute Essay From Last Lecture</span>

- Most people had some exposure to programming involving some kind of concurrency. Several had done things that sounded interesting and more ambitious than I might have thought!

- Several people commented that for the programming problem they spent a lot of time figuring out how to write a few lines of code. That was kind of the plan!

- Only one person mentioned the problem involving `strace`. I find it interesting how very many system calls . . .

## Mutual Exclusion Solutions So Far

- Solutions so far have some problems: inefficient, dependent on whether scheduler/etc. guarantees fairness.

- Also, they're very low-level, so might be hard to use for more complicated problems.

**Slide 3**

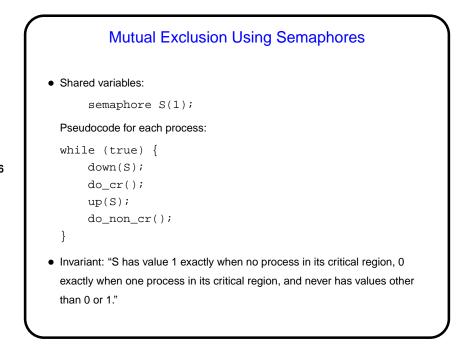- So, people have proposed various "synchronization mechanisms" . . .
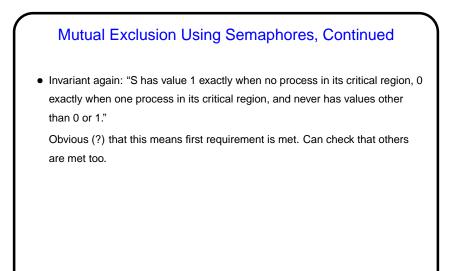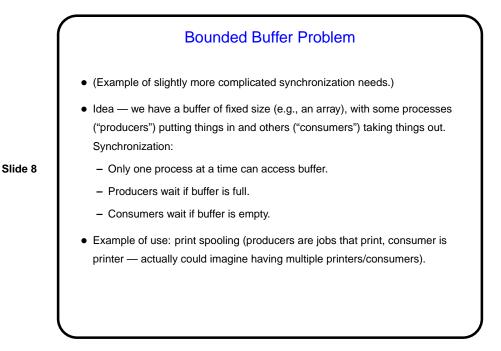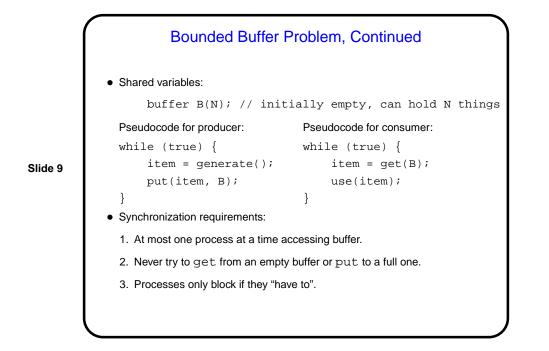
## Synchronization Mechanisms — Overview

- Synchronization using only shared variables seems to be tedious and inefficient.

- "Synchronization mechanisms" are more-abstract ways of coordinating what processes do. A key point is providing *something* that potentially makes a process wait.

**Slide 4**

## Semaphores

**Slide 5**

- History — 1965 paper by Dijkstra (possibly earlier work by Iverson, of APL/J fame).

- Idea — define semaphore ADT:
  - "Value" — non-negative integer.
  - Two operations, *both atomic*:
    * up (V) — add one to value.
    * down (P) — block until value is nonzero, then subtract one.

- Ignoring for now how to implement this — is it useful?

## Mutual Exclusion Using Semaphores

**Slide 6**

- Shared variables:

  ```
  semaphore S(1);
  ```

  Pseudocode for each process:

  ```
  while (true) {
      down(S);
      do_cr();
      up(S);
      do_non_cr();
  }
  ```

- Invariant: "S has value 1 exactly when no process in its critical region, 0 exactly when one process in its critical region, and never has values other than 0 or 1."

## Mutual Exclusion Using Semaphores, Continued

**Slide 7**

- Invariant again: "S has value 1 exactly when no process in its critical region, 0 exactly when one process in its critical region, and never has values other than 0 or 1."

  Obvious (?) that this means first requirement is met. Can check that others are met too.

## Bounded Buffer Problem

**Slide 8**

- (Example of slightly more complicated synchronization needs.)

- Idea — we have a buffer of fixed size (e.g., an array), with some processes ("producers") putting things in and others ("consumers") taking things out. Synchronization:
  - Only one process at a time can access buffer.
  - Producers wait if buffer is full.
  - Consumers wait if buffer is empty.

- Example of use: print spooling (producers are jobs that print, consumer is printer — actually could imagine having multiple printers/consumers).

## Bounded Buffer Problem, Continued

**Slide 9**

- Shared variables:

  ```
  buffer B(N); // initially empty, can hold N things
  ```

  Pseudocode for producer:          Pseudocode for consumer:
  ```
  while (true) {                    while (true) {
      item = generate();                item = get(B);
      put(item, B);                     use(item);
  }                                 }
  ```

- Synchronization requirements:

  1. At most one process at a time accessing buffer.

  2. Never try to get from an empty buffer or put to a full one.

  3. Processes only block if they "have to".

## Bounded Buffer Problem, Continued

**Slide 10**

- We already know how to guarantee one-at-a-time access. Can we extend that?

- Three situations where we want a process to wait:

  - Only one get/put at a time.

  - If B is empty, consumers wait.

  - If B is full, producers wait.

## Bounded Buffer Problem, Continued

- What about three semaphores?
  - One to guarantee one-at-a-time access.
  - One to make producers wait if B is full — so, it should be zero if B is full — "number of empty slots"?

**Slide 11**
  - One to make consumers wait if B is empty — so, it should be zero if B is empty — "number of slots in use"?
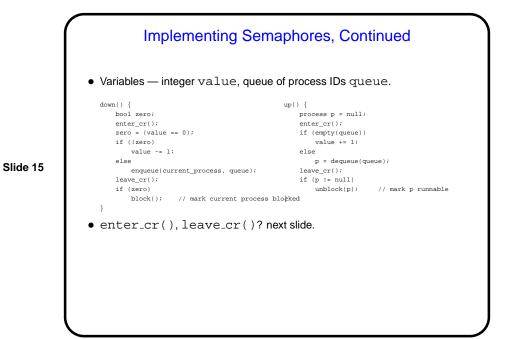
## Bounded Buffer Problem — Solution

- Shared variables:

```
    buffer B(N); // empty, capacity N
    semaphore mutex(1);
    semaphore empty(N);
    semaphore full(0);
```

**Slide 12**

Pseudocode for producer:            Pseudocode for consumer:

```
while (true) {                      while (true) {
    item = generate();                  down(full);
    down(empty);                        down(mutex);
    down(mutex);                        item = get(B);
    put(item, B);                       up(mutex);
    up(mutex);                          up(empty);
    up(full);                           use(item);
}                                   }
```

### Implementing Semaphores

- We want to define:
    - Data structure to represent a semaphore.
    - Functions up and down.
- up and down should work the way we said, and we'd like to do as little busy-waiting as possible.

**Slide 13**

### Implementing Semaphores, Continued

- Idea — represent semaphore as integer plus queue of waiting processes (represented as, e.g., process IDs).
- Then how should this work . . .

**Slide 14**

**Slide 15**

## Implementing Semaphores, Continued

- Variables — integer `value`, queue of process IDs `queue`.

```
down() {                                          up() {
    bool zero;                                        process p = null;
    enter_cr();                                       enter_cr();
    zero = (value == 0);                              if (empty(queue))
    if (!zero)                                            value += 1;
        value -= 1;                                   else
    else                                                  p = dequeue(queue);
        enqueue(current_process, queue);          leave_cr();
    leave_cr();                                       if (p != null)
    if (zero)                                             unblock(p);     // mark p runnable
        block();    // mark current process blocked
}
```

- `enter_cr()`, `leave_cr()`? next slide.

**Slide 16**

## Implementing Semaphores, Continued

- Revised functions to enter, leave critical region:

```
enter_cr:
    TSL registerX, lockVar
    compare registerX with 0
    if equal, jump to ok
    invoke scheduler # thread yields to another thread
    jump to enter_cr
ok:
    return


leave_cr:
    store 0 in lock
    return
```

**Slide 17**

## Sidebar: Shared Memory and Synchronization

- Solutions that rely on variables shared among processes assume that assigning a value to a variable actually changes its value in memory (RAM), more or less right away. Fine as a first approximation, but reality may be more complicated, because of various tricks used to deal with relative slowness of accessing memory:

  Optimizing compilers may keep variables' values in registers, only reading/writing memory when necessary to preserve semantics.

  Hardware may include cache, logically between CPU and memory, such that memory read/write goes to cache rather than RAM. Different CPUs' caches may not be in synch.

**Slide 18**

## Sidebar: Shared Memory and Synchronization, Continued

- So, actual implementations need notion of "memory fence" — point at which all apparent reads/writes have actually been done. Some languages provide standard ways to do this; others (e.g., C!) don't. C's `volatile` ("may be changed by something outside this code") helps some but may not be enough.

**Slide 19**

## Minute Essay

- None — quiz.