

Slide 1

### Administrivia

- (I guess there's not any!)

Slide 2

### Minute Essay From Last Lecture

- Alleged joke (from some random Usenet person):  
A man's P should exceed his V else what's a sema for?  
Do you understand this?  
(P is down, V is up — if not more P's than V's, no point in having a semaphore?)
- Anything else unclear?

## Semaphores — Recap

Slide 3

- Idea — define ADT that will be easier to use for interprocess communication/synchronization, maybe we can implement without (as much) busy-waiting.
- Definition as 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.
- How does this relate to operating systems?
  - Process abstraction (and its use within the o/s) means we have to solve “synchronization problems”.
  - Solution should somehow be part of o/s.

## Implementing Semaphores

Slide 4

- 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.

## 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 5

## Implementing Semaphores, Continued

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

```

down() {
    bool zero;
    enter_cr();
    zero = (value == 0);
    if (!zero)
        value -= 1;
    else
        enqueue(current_process, queue);
    leave_cr();
    if (zero)
        block(); // mark current process blocked
}

up() {
    process p = null;
    enter_cr();
    if (empty(queue))
        value += 1;
    else
        p = dequeue(queue);
    leave_cr();
    if (p != null)
        unblock(p); // mark p runnable
}

```

- `enter_cr()`, `leave_cr()` mostly like before; see p. 113.

Slide 6

Slide 7

## Monitors

- History — Hoare (1975) and Brinch Hansen (1975).
- Idea — combine synchronization and object-oriented paradigm.
- A monitor consists of
  - Data for a shared object (and initial values).
  - Procedures — only one at a time can run (e.g., whole procedure is a critical region).
- “Condition variable” ADT allows us to wait for specified conditions (e.g., buffer not empty):
  - Value — queue of suspended processes.
  - Operations:
    - \* Wait — suspend execution (and release mutual exclusion).
    - \* Signal — *if* there are processes suspended, allow *one* to continue. (if not, signal is “lost”).

Slide 8

## Bounded Buffer Problem, Revisited

- Define a `bounded_buffer` monitor with a `queue` and `insert` and `remove` procedures.

- Shared variables:

```
bounded_buffer B(N);
```

Pseudocode for producers:

```
while (true) {
    item = generate();
    B.insert(item);
}
```

Pseudocode for consumers:

```
while (true) {
    B.remove(item);
    use(item);
}
```

## Bounded-Buffer Monitor

- Data:

```
buffer B(N); // N is a constant, buffer initially empty
int count = 0;
condition full;
condition empty;
```

```
insert(item itm) {
    while (count == N)
        wait(full);
    put(itm, B);
    count += 1;
    signal(empty);
}

remove(item &itm) {
    while (count == 0)
        wait(empty);
    itm = get(B);
    signal(full);
}
```

Slide 9

## Implementing Monitors

- Requires compiler support, so more difficult to implement than (e.g.) semaphores.
- Java's methods for thread synchronization are based on monitors:
  - Data for monitor is instance variables (data for class).
  - Procedures for monitor are `synchronized` methods/blocks — mutual exclusion provided by implicit object lock.
  - `wait`, `notify`, `notifyAll` methods.
  - No condition variables, but above methods provide more or less equivalent functionality.

Slide 10

## Message Passing

Slide 11

- Previous synchronization mechanisms all involve shared variables, okay in some circumstances but not very feasible in others (e.g., multiple-processor system without shared memory).
- Idea of message passing — each process has a unique ID; two basic operations:
  - Send — specify destination ID, data to send (message).
  - Receive — specify source ID, buffer to hold received data. Usually some way to let source ID be “any”.

## Message Passing, Continued

Slide 12

- Exact specifications can vary, but typical assumptions include:
  - Sending a message never blocks a process (more difficult to implement but easier to work with).
  - Receiving a message blocks a process until there is a message to receive.
  - All messages sent are eventually available to receive (can be non-trivial to implement).
  - Messages from process A to process B arrive in the order in which they were sent.

### Implementing Message Passing

- On a machine with no physically shared memory (e.g., multicomputer), must send messages across interconnection network.
- On a machine with physically shared memory, can either copy (from address space to address space) or somehow be clever.  
(Why would you want to do this? programming model is in some ways simpler, doesn't require memory shared among processes.)
- Examples next time . . .

Slide 13

### Minute Essay

- Tell me one thing you've learned from the textbook.

Slide 14