# Administrivia

- Homework 2 on the Web. Due next Friday.

- Reminder: Quiz 2 Monday. Topics will come from the parts of chapter 2 we've talked about through today.

**Slide 1**

# Minute Essay From Last Lecture

- See slides from last time for solution. Several people got it right. Correlation with those who are keeping up with reading?

- One person mentioned that without locking the results might be unpredictable, another that SWAP wouldn't be atomic. Um, hardware designers' problem?

- Some people proposed solutions in which SWAP was — something other than what it is (exchange values of register and memory location).

**Slide 2**

## Semaphores — Review/Recap

- Abstract data type with values that are non-negative integers, two operations (up/down).

- Can be used to solve mutual-exclusion problem fairly nicely. How about other problems?

**Slide 3**

## Bounded Buffer Problem

- (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).

**Slide 4**

## Bounded Buffer Problem, Continued

**Slide 5**

- 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 6**

- 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 7**

  – 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 8**

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 9**

## 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 10**

## 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 11**

## 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 12**

## Another Synchronization Mechanism — 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.

**Slide 13**

- "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"). Some choices about whether signalling process continues, or signalled process awakens right away.

## Bounded Buffer Problem, Revisited

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

- Shared variables:

  `bounded_buffer B(N);`

**Slide 14**

Pseudocode for producers:                Pseudocode for consumers:

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

**Slide 15**

### Bounded-Buffer Monitor

- Data:

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

- Procedures:

```
insert(item itm) {          remove(item &itm) {
    if (count == N)             if (count == 0)
        wait(full);                wait(empty);
    put(itm, B);               itm = get(B);
    count += 1;                count -= 1;
    signal(empty);             signal(full);
}                           }
```

- Does this work? (Yes.)

**Slide 16**

### Implementing Monitors

- Requires compiler support, so more difficult to implement than (e.g.) semaphores.

- Java's methods for thread synchronization are based on monitors . . .

### Java's Adaptation of the Monitor Idea

- 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.

**Slide 17**

- No condition variables, but above methods provide more or less equivalent functionality.

  *Note* that the language specs for Java allow spurious wake-ups. So "best practice" is to `wait()` in a loop, re-checking the desired condition. The textbook's bounded-buffer code doesn't do this (?!).

### Minute Essay

- 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? (Remember that P is "down" and V is "up".)

**Slide 18**

**Slide 19**

### Minute Essay Answer

- It's a pun. The idea is roughly that if you never have a situation in which you've attempted more "down" operations than "up" operations, you didn't need a semaphore. (Or that's what I think it means. The author might have another idea!)