

CS162 Operating Systems and Systems Programming Lecture 6

Semaphores, Conditional Variables, Deadlocks

February 7, 2011

Ion Stoica

<http://inst.eecs.berkeley.edu/~cs162>

Review: Definition of Monitor

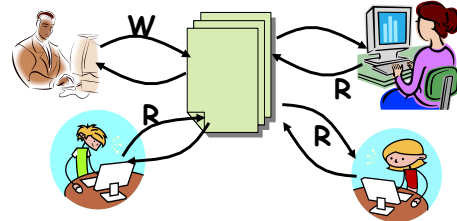
- Semaphores are confusing because dual purpose:
 - Both mutual exclusion and scheduling constraints
 - Cleaner idea: Use *locks* for mutual exclusion and *condition variables* for scheduling constraints
- **Monitor**: a lock and zero or more condition variables for managing concurrent access to shared data
- **Lock**: provides mutual exclusion to shared data:
 - Always acquire before accessing shared data structure
 - Always release after finishing with shared data
- **Condition Variable**: a queue of threads waiting for something *inside* a critical section
 - Key idea: allow sleeping inside critical section by atomically releasing lock at time we go to sleep
 - Contrast to semaphores: Can't wait inside critical section

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.2

Review: Readers/Writers Problem



- Motivation: Consider a shared database
 - Two classes of users:
 - » Readers – never modify database
 - » Writers – read and modify database
 - Is using a single lock on the whole database sufficient?
 - » Like to have many readers at the same time
 - » Only one writer at a time

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.3

Review: Code for a Reader

```
Reader() {  
    // First check self into system  
    lock.Acquire();  
    while ((AW + WW) > 0) { // Is it safe to read?  
        WR++; // No. Writers exist  
        okToRead.wait(&lock); // Sleep on cond var  
        WR--; // No longer waiting  
    }  
    AR++; // Now we are active!  
    lock.release();  
    // Perform actual read-only access  
    AccessDatabase(ReadOnly);  
    // Now, check out of system  
    lock.Acquire();  
    AR--; // No longer active  
    if (AR == 0 && WW > 0) // No other active readers  
        okToWrite.signal(); // Wake up one writer  
    lock.Release();  
}
```

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.4

Review: Code for a Writer

```
Writer() {
  // First check self into system
  lock.Acquire();
  while ((AW + AR) > 0) { // Is it safe to write?
    WW++; // No. Active users exist
    okToWrite.wait(&lock); // Sleep on cond var
    WW--; // No longer waiting
  }
  AW++; // Now we are active!
  lock.release();
  // Perform actual read/write access
  AccessDatabase(ReadWrite);
  // Now, check out of system
  lock.Acquire();
  AW--; // No longer active
  if (WW > 0) { // Give priority to writers
    okToWrite.signal(); // Wake up one writer
  } else if (WR > 0) { // Otherwise, wake reader
    okToRead.broadcast(); // Wake all readers
  }
  lock.Release();
}
2/7/11 Ion Stoica CS162 @UCB Spring 2011 Lec 1.5
```

Simulation of Readers/Writers solution

- Consider the following sequence of operators:
 - R1, R2, W1, R3 (AR = WR = AW = WW = 0)
 - On entry, each reader checks the following:

```
while ((AW + WW) > 0) { // Is it safe to read?
  WR++; // No. Writers exist
  okToRead.wait(&lock); // Sleep on cond var
  WR--; // No longer waiting
}
AR++; // Now we are active!
```
 - First, R1 comes along:
AR = 1, WR = 0, AW = 0, WW = 0
 - Next, R2 comes along:
AR = 2, WR = 0, AW = 0, WW = 0
 - Now, readers may take a while to access database
 - Situation: Locks released
 - Only AR is non-zero
- 2/7/11 Ion Stoica CS162 @UCB Spring 2011 Lec 1.6

Simulation(2)

- Next, W1 comes along:

```
while ((AW + AR) > 0) { // Is it safe to write?
  WW++; // No. Active users exist
  okToWrite.wait(&lock); // Sleep on cond var
  WW--; // No longer waiting
}
AW++;
```
 - Can't start because of readers, so go to sleep:
AR = 2, WR = 0, AW = 0, WW = 1
 - Finally, R3 comes along:
AR = 2, WR = 1, AW = 0, WW = 1
 - Now, say that R2 finishes before R1:
AR = 1, WR = 1, AW = 0, WW = 1
 - Finally, last of first two readers (R1) finishes and wakes up writer:

```
if (AR == 0 && WW > 0) // No other active readers
  okToWrite.signal(); // Wake up one writer
```
- 2/7/11 Ion Stoica CS162 @UCB Spring 2011 Lec 1.7

Simulation(3)

- When writer wakes up, get:
AR = 0, WR = 1, AW = 1, WW = 0
 - Then, when writer finishes:

```
if (WW > 0) { // Give priority to writers
  okToWrite.signal(); // Wake up one writer
} else if (WR > 0) { // Otherwise, wake reader
  okToRead.broadcast(); // Wake all readers
}
```

 - Writer wakes up reader, so get:
AR = 1, WR = 0, AW = 0, WW = 0
 - When reader completes, we are finished
- 2/7/11 Ion Stoica CS162 @UCB Spring 2011 Lec 1.8

Questions

- Can readers starve? Consider Reader() entry code:

```
while ((AW + WW) > 0) { // Is it safe to read?
    WR++; // No. Writers exist
    okToRead.wait(&lock); // Sleep on cond var
    WR--; // No longer waiting
}
AR++; // Now we are active!
```

- What if we erase the condition check in Reader exit?

```
AR--; // No longer active
if (AR == 0 && WW > 0) // No other active readers
    okToWrite.signal(); // Wake up one writer
```

- Further, what if we turn the signal() into broadcast()

```
AR--; // No longer active
okToWrite.broadcast(); // Wake up one writer
```

- Finally, what if we use only one condition variable (call it "okToContinue") instead of two separate ones?
 - Both readers and writers sleep on this variable
 - Must use broadcast() instead of signal()

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.9

Can we construct Monitors from Semaphores?

- Locking aspect is easy: Just use a mutex
- Can we implement condition variables this way?

```
Wait() { semaphore.P(); }
Signal() { semaphore.V(); }
```

- Does this work better?

```
Wait(Lock lock) {
    lock.Release();
    semaphore.P();
    lock.Acquire();
}
Signal() { semaphore.V(); }
```

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.10

Construction of Monitors from Semaphores (con't)

- Problem with previous try:
 - P and V are commutative – result is the same no matter what order they occur
 - Condition variables are NOT commutative
- Does this fix the problem?

```
Wait(Lock lock) {
    lock.Release();
    semaphore.P();
    lock.Acquire();
}
Signal() {
    if semaphore queue is not empty
        semaphore.V();
}
```

- Not legal to look at contents of semaphore queue
- There is a race condition – signaler can slip in after lock release and before waiter executes semaphore.P()
- It is actually possible to do this correctly
 - Complex solution for Hoare scheduling in book
 - Can you come up with simpler Mesa-scheduled solution?

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.11

Monitor Conclusion

- Monitors represent the logic of the program
 - Wait if necessary
 - Signal when change something so any waiting threads can proceed
- Basic structure of monitor-based program:

```
lock
while (need to wait) {
    condvar.wait();
}
unlock
do something so no need to wait

lock
condvar.signal();
unlock
```

Check and/or update state variables
Wait if necessary

Check and/or update state variables

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.12

C-Language Support for Synchronization

- C language: Pretty straightforward synchronization
 - Just make sure you know *all* the code paths out of a critical section

```
int Rtn() {
    lock.acquire();
    ...
    if (exception) {
        lock.release();
        return errReturnCode;
    }
    ...
    lock.release();
    return OK;
}
```

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.13

C++ Language Support for Synchronization

- Languages with exceptions like C++
 - Languages that support exceptions are problematic (easy to make a non-local exit without releasing lock)
 - Consider:

```
void Rtn() {
    lock.acquire();
    ...
    DoFoo();
    ...
    lock.release();
}
void DoFoo() {
    ...
    if (exception) throw errException;
    ...
}
```

- Notice that an exception in DoFoo() will exit without releasing the lock

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.14

C++ Language Support for Synchronization (con't)

- Must catch all exceptions in critical sections
 - Catch exceptions, release lock, and re-throw exception:

```
void Rtn() {
    lock.acquire();
    try {
        ...
        DoFoo();
        ...
    } catch (...) { // catch exception
        lock.release(); // release lock
        throw; // re-throw the exception
    }
    lock.release();
}
void DoFoo() {
    ...
    if (exception) throw errException;
    ...
}
```

- Even Better: `auto_ptr<T>` facility. See C++ Spec.
 - » Can deallocate/free lock regardless of exit method

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.15

Java Language Support for Synchronization

- Java has explicit support for threads and thread synchronization
- Bank Account example:

```
class Account {
    private int balance;
    // object constructor
    public Account (int initialBalance) {
        balance = initialBalance;
    }
    public synchronized int getBalance() {
        return balance;
    }
    public synchronized void deposit(int amount) {
        balance += amount;
    }
}
```

- Every object has an associated lock which gets automatically acquired and released on entry and exit from a *synchronized* method.

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.16

Java Language Support for Synchronization (con't)

- Java also has *synchronized* statements:

```
synchronized (object) {  
    ...  
}
```

- Since every Java object has an associated lock, this type of statement acquires and releases the object's lock on entry and exit of the body

- Works properly even with exceptions:

```
synchronized (object) {  
    ...  
    DoFoo();  
    ...  
}  
void DoFoo() {  
    throw errException;  
}
```

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.17

Java Language Support for Synchronization (2/2)

- In addition to a lock, every object has a **single** condition variable associated with it

- How to wait inside a synchronization method or block:

```
» void wait(long timeout); // Wait for timeout  
» void wait(long timeout, int nanoseconds); //variant  
» void wait();
```

- How to signal in a synchronized method or block:

```
» void notify(); // wakes up oldest waiter  
» void notifyAll(); // like broadcast, wakes everyone
```

- Condition variables can wait for a bounded length of time. This is useful for handling exception cases:

```
t1 = time.now();  
while (!ATMRequest()) {  
    wait (CHECKPERIOD);  
    t2 = time.now();  
    if (t2 - t1 > LONG_TIME) checkMachine();  
}
```

- Not all Java VMs equivalent!

» Different scheduling policies, not necessarily preemptive!

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.18

Summary: Semaphores and Cond. Variables

- Semaphores: Like integers with restricted interface

- Two operations:

» **P()**: Wait if zero; decrement when becomes non-zero
» **V()**: Increment and wake a sleeping task (if exists)
» Can initialize value to any non-negative value

- Use separate semaphore for each constraint

- Monitors: A lock plus one or more condition variables

- Always acquire lock before accessing shared data

- Use condition variables to wait inside critical section

» Three Operations: **Wait()**, **Signal()**, and **Broadcast()**

- Language support for synchronization:

- Java provides **synchronized** keyword and one condition-variable per object (with **wait()** and **notify()**)

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.19

5min Break

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.20

Tips for Programming in a Project Team



"You just have to get your synchronization right!"

- Big projects require more than one person (or long, long, long time)
 - Big OS: thousands of person-years!
- It's very hard to make software project teams work correctly
 - Doesn't seem to be as true of big construction projects
 - » Empire state building finished in **one** year: staging iron production thousands of miles away
 - » Or the Hoover dam: built towns to hold workers

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.21

Big Projects

- What is a big project?
 - Time/work estimation is hard
 - Programmers are eternal optimistics (it will only take two days!)
 - » This is why we bug you about starting the project early



- Can a project be efficiently partitioned?
 - Partitionable task decreases in time as you add people
 - But, if you require communication:
 - » Time reaches a minimum bound
 - » With complex interactions, time increases!
 - Mythical person-month problem:
 - » You estimate how long a project will take
 - » Starts to fall behind, so you add more people
 - » Project takes even more time!



2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.22

Techniques for Partitioning Tasks

- Functional
 - Person A implements threads, Person B implements semaphores, Person C implements locks...
 - Problem: Lots of communication across APIs
 - » If B changes the API, A may need to make changes
 - » Story: Large airline company spent \$200 million on a new scheduling and booking system. Two teams "working together." After two years, went to merge software. Failed! Interfaces had changed (documented, but no one noticed). Result: would cost another \$200 million to fix.
- Task
 - Person A designs, Person B writes code, Person C tests
 - May be difficult to find right balance, but can focus on each person's strengths (Theory vs systems hacker)
 - Since Debugging is hard, Microsoft has *two* testers for *each* programmer
- Most CS162 project teams are functional, but people have had success with task-based divisions

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.23

Communication

- More people mean more communication
 - Changes have to be propagated to more people
 - Think about person writing code for most fundamental component of system: everyone depends on them!
- Miscommunication is common
 - "Index starts at 0? I thought you said 1!"
- Who makes decisions?
 - Individual decisions are fast but trouble
 - Group decisions take time
 - Centralized decisions require a big picture view (someone who can be the "system architect")
- Often designating someone as the system architect can be a good thing
 - Better not be clueless
 - Better have good people skills
 - Better let other people do work



2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.24

Coordination



- More people \Rightarrow no one can make all meetings!
 - They miss decisions and associated discussion
 - Example from earlier class: one person missed meetings and did something group had rejected
 - Why do we limit groups to 5 people?
 - » You would never be able to schedule meetings otherwise
 - Why do we require 4 people minimum?
 - » You need to experience groups to get ready for real world
- People have different work styles
 - Some people work in the morning, some at night
 - How do you decide when to meet or work together?
- What about project slippage?
 - It will happen, guaranteed!
 - Ex: everyone busy but not talking. One person way behind. No one knew until very end – too late!
- Hard to add people to existing group
 - Members have already figured out how to work together

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.25

How to Make it Work?

- People are human. Get over it.
 - People will make mistakes, miss meetings, miss deadlines, etc. You need to live with it and adapt
 - It is better to anticipate problems than clean up afterwards.
- Document, document, document
 - Why Document?
 - » Expose decisions and communicate to others
 - » Easier to spot mistakes early
 - » Easier to estimate progress
 - What to document?
 - » Everything (but don't overwhelm people or no one will read)
 - Standardize!
 - » One programming format: variable naming conventions, tab indents, etc.
 - » Comments (Requires, effects, modifies)—javadoc?



2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.26

Suggested Documents for You to Maintain

- Project objectives: goals, constraints, and priorities
- Specifications: the manual plus performance specs
 - This should be the first document generated and the last one finished
- Meeting notes
 - Document all decisions
 - You can often cut & paste for the design documents
- Schedule: What is your anticipated timing?
 - This document is critical!
- Organizational Chart
 - Who is responsible for what task?



2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.27

Test Continuously

- Integration tests all the time, not at 11pm on due date!
 - Write dummy stubs with simple functionality
 - » Let's people test continuously, but more work
 - Schedule periodic integration tests
 - » Get everyone in the same room, check out code, build, and test.
 - » Don't wait until it is too late!
- Testing types:
 - Unit tests: check each module in isolation (use JUnit?)
 - Daemons: subject code to exceptional cases
 - Random testing: Subject code to random timing changes
- Test early, test later, test again
 - Tendency is to test once and forget; what if something changes in some other part of the code?



2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.28

Resource Contention and Deadlock

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.29

Resources

- Resources – passive entities needed by threads to do their work
 - CPU time, disk space, memory
- Two types of resources:
 - Preemptable – can take it away
 - » CPU, Embedded security chip
 - Non-preemptable – must leave it with the thread
 - » Disk space, printer, chunk of virtual address space
 - » Critical section
- Resources may require exclusive access or may be sharable
 - Read-only files are typically sharable
 - Printers are not sharable during time of printing
- One of the major tasks of an operating system is to manage resources



2/7/11

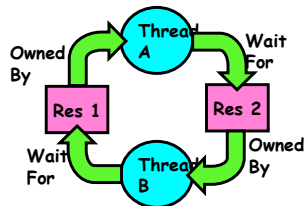
Ion Stoica CS162 ©UCB Spring 2011

Lec 1.30

Starvation vs Deadlock



- Starvation vs. Deadlock
 - Starvation: thread waits indefinitely
 - » Example, low-priority thread waiting for resources constantly in use by high-priority threads
 - Deadlock: circular waiting for resources
 - » Thread A owns Res 1 and is waiting for Res 2
 - » Thread B owns Res 2 and is waiting for Res 1



- Deadlock ⇒ Starvation but not vice versa
 - » Starvation can end (but doesn't have to)
 - » Deadlock can't end without external intervention

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.31

Conditions for Deadlock

- Deadlock not always deterministic – Example 2 mutexes:

Thread A	Thread B	Deadlock
x.P();	y.P();	A: x.P();
y.P();	x.P();	A: y.P();
y.V();	x.V();	B: y.P();
x.V();	y.V();	B: x.P();
		...

 - Deadlock won't always happen with this code
 - » Have to have exactly the right timing ("wrong" timing)
 - » So you release a piece of software, and you tested it, and there it is, controlling a nuclear power plant...
- Deadlocks occur with multiple resources
 - Means you can't decompose the problem
 - Can't solve deadlock for each resource independently
- Example: System with 2 disk drives and two threads
 - Each thread needs 2 disk drives to function
 - Each thread gets one disk and waits for another one

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.32

Bridge Crossing Example

The diagram shows a bridge crossing a road. On the left side of the road, a car is driving towards the bridge. On the right side, a car and a truck are driving away from the bridge. The bridge itself is a narrow strip across the road, and the cars are positioned on it.

- Each segment of road can be viewed as a resource
 - Car must own the segment under them
 - Must acquire segment that they are moving into
- For bridge: must acquire both halves
 - Traffic only in one direction at a time
 - Problem occurs when two cars in opposite directions on bridge: each acquires one segment and needs next
- If a deadlock occurs, it can be resolved if one car backs up (preempt resources and rollback)
 - Several cars may have to be backed up
- Starvation is possible
 - East-going traffic really fast \Rightarrow no one goes west

2/7/11 Ion Stoica CS162 ©UCB Spring 2011 Lec 1.33

Train Example (Wormhole-Routed Network)

- Circular dependency (Deadlock!)
 - Each train wants to turn right
 - Blocked by other trains
 - Similar problem to multiprocessor networks
- Fix? Imagine grid extends in all four directions
 - Force ordering of channels (tracks)
 - » Protocol: Always go east-west first, then north-south
 - Called “dimension ordering” (X then Y)

The diagram shows a grid of tracks. A train is shown on a track that is blocked by a sign that says "Disallowed By Rule". The sign is pink and has a diagonal line through it. The train is blue and is positioned on a track that is blocked by another train on a perpendicular track.

2/7/11 Ion Stoica CS162 ©UCB Spring 2011 Lec 1.34

Train Example (Wormhole-Routed Network)

- Circular dependency (Deadlock!)
 - Each train wants to turn right
 - Blocked by other trains
 - Similar problem to multiprocessor networks
- Fix? Imagine grid extends in all four directions
 - Force ordering of channels (tracks)
 - » Protocol: Always go east-west first, then north-south
 - Called “dimension ordering” (X then Y)

The diagram shows a grid of tracks. A red 'X' is drawn on the grid, indicating a deadlock state. The 'X' is formed by two intersecting tracks, one horizontal and one vertical. The horizontal track is labeled 'A' and the vertical track is labeled 'B'.

2/7/11 Ion Stoica CS162 ©UCB Spring 2011 Lec 1.35

Dining Philosopher Problem

The diagram shows a round dining table with five chairs around it. The word "RICE" is written in the center of the table. Five philosophers are seated around the table, each with a chopstick. The philosophers are depicted as cartoon characters.

- Five chopsticks/Five philosopher (really cheap restaurant)
 - Free-for all: Philosopher will grab any one they can
 - Need two chopsticks to eat
- What if all grab at same time?
 - Deadlock!
- How to fix deadlock?
 - Make one of them give up a chopstick (Hah!)
 - Eventually everyone will get chance to eat
- How to prevent deadlock?
 - Never let philosopher take last chopstick if no hungry philosopher has two chopsticks afterwards

2/7/11 Ion Stoica CS162 ©UCB Spring 2011 Lec 1.36

Four requirements for Deadlock

- **Mutual exclusion**
 - Only one thread at a time can use a resource.
- **Hold and wait**
 - Thread holding at least one resource is waiting to acquire additional resources held by other threads
- **No preemption**
 - Resources are released only voluntarily by the thread holding the resource, after thread is finished with it
- **Circular wait**
 - There exists a set $\{T_1, \dots, T_n\}$ of waiting threads
 - » T_1 is waiting for a resource that is held by T_2
 - » T_2 is waiting for a resource that is held by T_3
 - » ...
 - » T_n is waiting for a resource that is held by T_1

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.37

Summary: Deadlock

- Suggestions for dealing with Project Partners
 - Start Early, Meet Often
 - Develop Good Organizational Plan, Document Everything, Use the right tools, Develop Comprehensive Testing Plan
 - (Oh, and add 2 years to every deadline!)
- Starvation vs. Deadlock
 - Starvation: thread waits indefinitely
 - Deadlock: circular waiting for resources
- Four conditions for deadlocks
 - **Mutual exclusion**
 - » Only one thread at a time can use a resource
 - **Hold and wait**
 - » Thread holding at least one resource is waiting to acquire additional resources held by other threads
 - **No preemption**
 - » Resources are released only voluntarily by the threads
 - **Circular wait**
 - » \exists set $\{T_1, \dots, T_n\}$ of threads with a cyclic waiting pattern

2/7/11

Ion Stoica CS162 ©UCB Spring 2011

Lec 1.38