# CSE150
# Operating Systems
# Lecture 6

# Synchronization, Atomic operations, Locks

# Why allow cooperating threads?

- People cooperate; computers help/enhance people's lives, so computers must cooperate
  - By analogy, the non-reproducibility/non-determinism of people is a notable problem for "carefully laid plans"
- Advantage 1: Share resources
  - One computer, many users
  - One bank balance, many ATMs
    - » What if ATMs were only updated at night?
  - Embedded systems (robot control: coordinate arm & hand together)
- Advantage 2: Speedup
  - Overlap I/O and computation
  - Multiprocessors – chop up program into parallel pieces
- Advantage 3: Modularity
  - Chop large problem up into simpler pieces
    - » To compile, for instance, gcc calls cpp | cc1 | cc2 | as | ld
    - » Makes system easier to extend

# Interrupts and Two Functions

- The state of a thread is contained in the TCB
  - Registers, PC, stack pointer
  - States: New, Ready, Running, Waiting, or Terminated
- Interrupts: hardware mechanism for returning control to operating system
  - Used for important/high-priority events
  - Can force dispatcher to schedule a different thread (premptive multithreading)
- New Threads Created with `ThreadFork()`
  - Create initial TCB and stack to point at `ThreadRoot()`
  - `ThreadRoot()` calls thread code, then `ThreadFinish()`
  - `ThreadFinish()` wakes up waiting threads then prepares TCB/stack for distruction
- Threads can wait for other threads using `ThreadJoin()`
- Many scheduling options
  - Decision of which thread to run complex enough for complete lecture

# Today

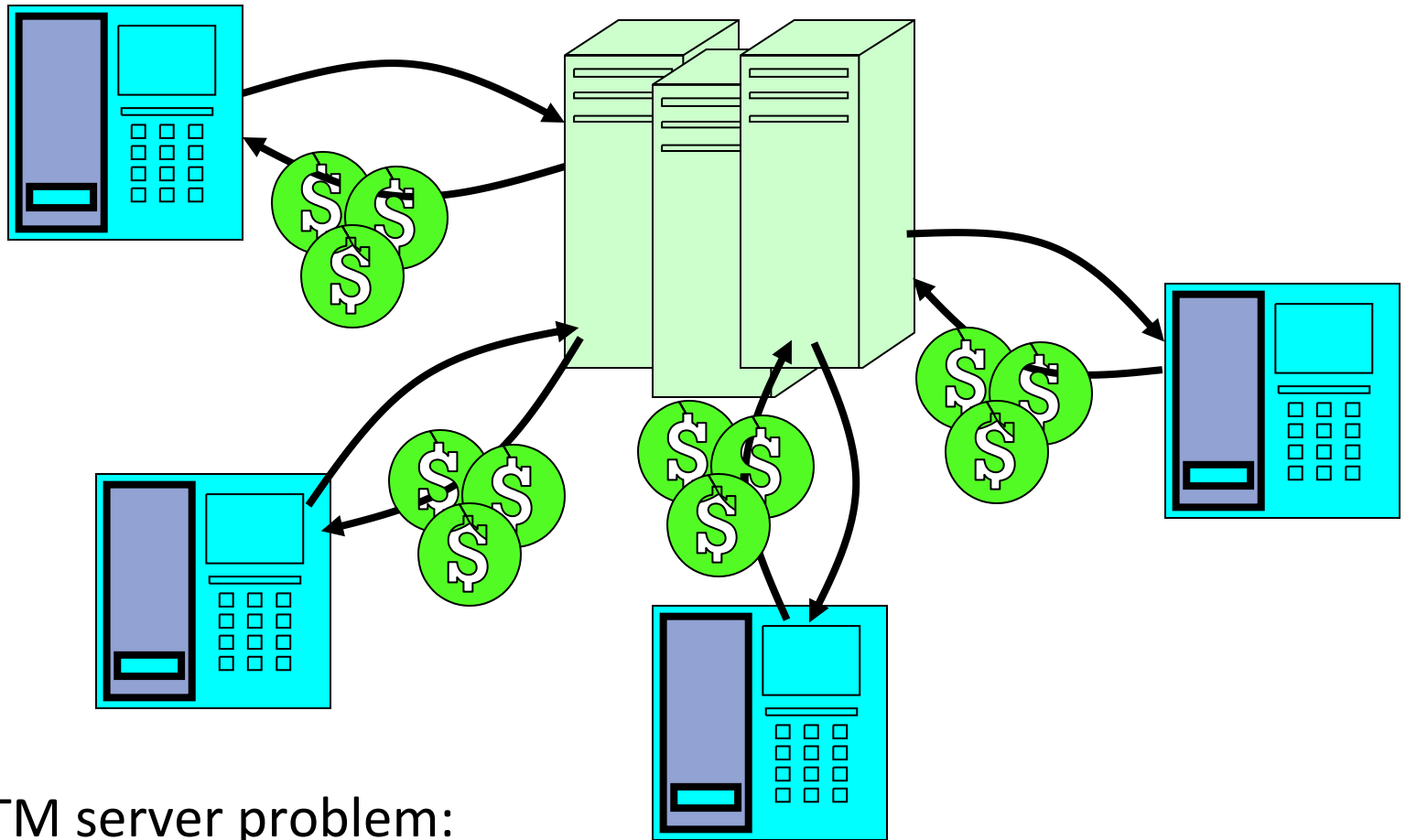- Concurrency examples and sharing
- Synchronization

# Announcements - Important time points

- Lab sessions start in **the week 02/10**.

- Project 1: 4 weeks

- Project 2: 4 weeks

- Three or Four homework (Feb. 13$^{th}$; Feb 27$^{th}$, March 10$^{th}$, April 9$^{th}$)

  - One week to finish each.

- In-class exam (March. 19$^{th}$)

- Final exam (11:30 AM - 2:30 PM May 13$^{th}$)

# Correctness for systems with concurrent threads

- **Independent Threads (ideal):**
  - No state shared with other threads
  - Deterministic $\Rightarrow$ Input state determines results
  - Reproducible $\Rightarrow$ Can recreate Starting Conditions, I/O
  - Scheduling order doesn't matter (if `switch()` works!!!)
- **Cooperating Threads:**
  - Shared State between multiple threads
  - Non-deterministic
  - Non-reproducible (performance??)
- Non-deterministic and Non-reproducible means that bugs can be intermittent

# ATM Bank Server



- ATM server problem:
  - Service a set of requests
  - Do so without corrupting database
  - Don't hand out too much money

# ATM bank server example

- Suppose we wanted to implement a server process to handle requests from an ATM network:

```
BankServer() {
    while (TRUE) {
        ReceiveRequest(&op, &acctId, &amount);
        ProcessRequest(op, acctId, amount);
    }
}
ProcessRequest(op, acctId, amount) {
    if (op == deposit) Deposit(acctId, amount);
    else if …
}
Deposit(acctId, amount) {
    acct = GetAccount(acctId); /* may use disk I/O */
    acct->balance += amount;
    StoreAccount(acct); /* Involves disk I/O */
}
```

- How could we speed this up?
  - More than one request being processed at once
  - Multiple threads (multi-proc, or overlap comp and I/O)

# Can Threads Help?

- One thread per request!

- Requests proceed to completion, blocking as required:

```
Deposit(acctId, amount) {
    acct = GetAccount(actId); /* May use disk I/O */
    acct->balance += amount;
    StoreAccount(acct);       /* Involves disk I/O */
}
```

- Unfortunately, shared state can get corrupted:

|            Thread 1             |            Thread 2             |
| ------------------------------- | ------------------------------ |
| load r1, acct->balance          |                                |
|                                 | load r1, acct->balance         |
|                                 | add r1, amount2                |
|                                 | store r1, acct->balance        |
| add r1, amount1                 |                                |
| store r1, acct->balance         |                                |

# Problem is at the lowest level

- Most of the time, threads are working on separate data, so scheduling doesn't matter:

|          Thread A          |          Thread B          |
|:--------------------------:|:--------------------------:|
|           x = 1;           |           y = 2;           |

- However, What about (Initially, y = 12):

|          Thread A          |          Thread B          |
|:--------------------------:|:--------------------------:|
|           x = 1;           |           y = 2;           |
|          x = y+1;          |          y = y*2;          |

- What are the possible values of x?

|          Thread A          |          Thread B          |
|:--------------------------:|:--------------------------:|
|           x = 1;           |                            |
|          x = y+1;          |                            |
|                            |           y = 2;           |
|                            |          y = y*2           |

x=13

Preemption can occur at any time!

# Problem is at the lowest level

- Most of the time, threads are working on separate data, so scheduling doesn't matter:

| Thread A | Thread B |
|----------|----------|
| x = 1;   | y = 2;   |

- However, What about (Initially, y = 12):

| Thread A | Thread B |
|----------|----------|
| x = 1;   | y = 2;   |
| x = y+1; | y = y*2; |

  – What are the possible values of x?

| Thread A | Thread B |
|----------|----------|
|          | y = 2;   |
|          | y = y*2; |
| x = 1;   |          |
| x = y+1; |          |

x=5

Preemption can occur at any time!

# Problem is at the lowest level

- Most of the time, threads are working on separate data, so scheduling doesn't matter:

| Thread A | Thread B |
|----------|----------|
| x = 1; | y = 2; |

- However, What about (Initially, y = 12):

| Thread A | Thread B |
|----------|----------|
| x = 1; | y = 2; |
| x = y+1; | y = y*2; |

  - What are the possible values of x?

| Thread A | Thread B |
|----------|----------|
|  | y = 2; |
| x = 1; |  |
| x = y+1; |  |
|  | y= y*2; |

x=3

Preemption can occur at any time!

# Atomic Operations

- To understand a concurrent program, we need to know what the underlying indivisible operations are!

- Atomic Operation: an operation that always runs to completion or not at all
  - It is *indivisible:* it cannot be stopped in the middle and state cannot be modified by someone else in the middle
  - Fundamental building block – if no atomic operations, then have no way for threads to work together

- On most machines, memory references and assignments (i.e. loads and stores) of words are atomic

- Many instructions are not atomic
  - Double-precision floating point store often not atomic
  - VAX and IBM 360 had an instruction to copy a whole array

# Concurrency Challenges

- Multiple computations (threads) executing in parallel to

  - share resources, and/or

  - share data

- Examples:

  - Sharing CPU for 10ms vs. 1min

  - Sharing a database at the row vs. table granularity

- Fine grain sharing:

  ⇑  increase concurrency → better performance

  ⇓  more complex

- Coarse grain sharing:

  ⇑  Simpler to implement

  ⇓  Lower performance

# Definitions

- Synchronization: using atomic operations to ensure cooperation between threads
  - For now, only loads and stores are atomic
  - We will show that it is hard to build anything useful with only atomic reads and writes

- Critical Section: piece of code that only one thread can execute at once.

- Mutual Exclusion: ensuring that only one thread executes a critical section
  - One thread *excludes* the other while doing its task
  - Critical section and mutual exclusion are two ways of describing the same thing.
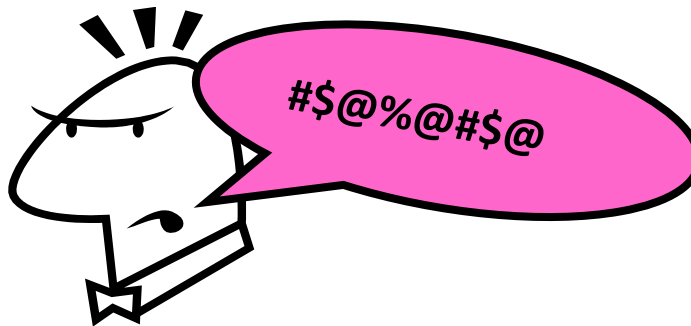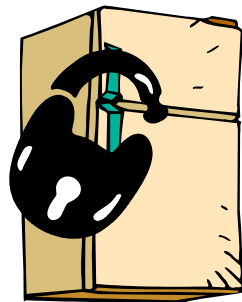
# Motivation: "Too much milk"

- Great thing about OS's – analogy between problems in OS and problems in real life
  - Help you understand real life problems better
  - But, computers are much stupider than people
- Example: People need to coordinate:

| Time | Person A | Person B |
|------|----------|----------|
| 3:00 | Look in Fridge. Out of milk | |
| 3:05 | Leave for store | |
| 3:10 | Arrive at store | Look in Fridge. Out of milk |
| 3:15 | Buy milk | Leave for store |
| 3:20 | Arrive home, put milk away | Arrive at store |
| 3:25 | | Buy milk |
| 3:30 | | Arrive home, put milk away |

# Lock

- Lock: prevents someone from doing something
  - Lock before entering critical section and before accessing shared data
  - Unlock when leaving, after accessing shared data
  - Wait if locked
    - » Important idea: all synchronization involves waiting
- For example: fix the milk problem by putting a key on the refrigerator
  - Lock it and take key if you are going to go buy milk
  - Fixes too much (coarse granularity): roommate angry if only wants OJ

  #$@%@#$@
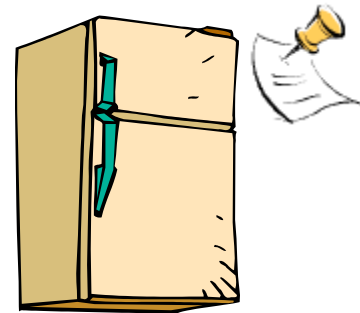
  - Of Course – We don't know how to make a lock yet

# Too Much Milk: Correctness Properties

- Need to be careful about correctness of concurrent programs, since non-deterministic
  - Always write down **desired** behavior first
  - Impulse is to start coding first, then when it doesn't work, pull hair out
  - Instead, think first, then code
- What are the correctness properties for the "Too much milk" problem?
  - Never more than one person buys
  - Someone buys if needed
- Restrict ourselves to use only atomic load and store operations as building blocks

# Too Much Milk: Solution #1

- Use a note to avoid buying too much milk:
  - Leave a note before buying (kind of "lock")
  - Remove note after buying (kind of "unlock")
  - Don't buy if note (wait)

- Suppose a computer tries this (remember, only memory read/write are atomic):

```
if (noMilk) {
    if (noNote) {
        leave Note;
        buy milk;
        remove note;
    }
}
```

- Result?

# Too Much Milk: Solution #1

- Still too much milk <span style="color:red">but only occasionally!</span>

```
     Thread A                    Thread B
  if (noMilk)
    if (noNote) {
                          if (noMilk)
                            if (noNote) {

      leave Note;
      buy milk;
      remove note;
    }
  }
                              leave Note;
                              buy milk;
                              …
```

- Thread can get context switched after checking milk and note but before leaving note!

- Solution makes problem worse since fails <span style="color:red">intermittently</span>
  - Makes it really hard to debug…
  - Must work despite what the thread dispatcher does!

# Too Much Milk: Solution #1½

- Clearly the Note is not quite blocking enough
  - Let's try to fix this by placing note first
- Another try at previous solution:

```
leave Note;
if (noMilk) {
    if (noNote) {
        leave Note;
        buy milk;
    }
}
remove note;
```

- What happens here?
  - Well, with human, probably nothing bad
  - With computer: no one ever buys milk

# Too Much Milk Solution #2

- How about labeled notes?
  - Now we can leave note before checking

- Algorithm looks like this:

<table>
<tr><th>Thread A</th><th>Thread B</th></tr>
</table>

```
leave note A;               leave note B;
if (noNote B) {             if (noNote A) {
    if (noMilk) {               if (noMilk) {
        buy Milk;                   buy Milk;
    }                           }
}                           }
remove note A;              remove note B;
```

- Does this work?

# Too Much Milk Solution #2
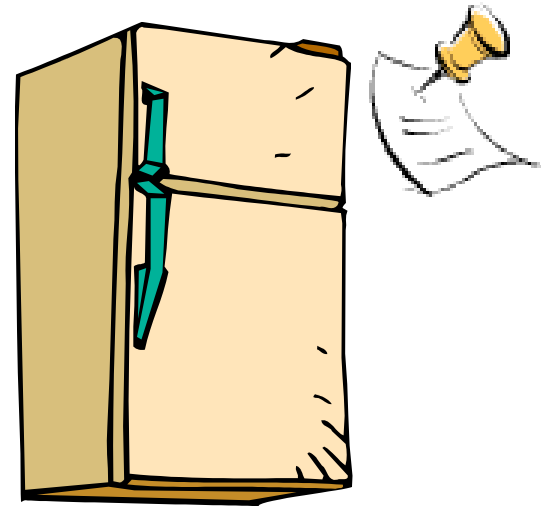
- Possible for neither thread to buy milk!

| Thread A | Thread B |
|---|---|
| `leave note A;` | |
| | `leave note B;` |
| | `if (noNote A) {` |
| | `    if (noMilk) {` |
| | `        buy Milk;` |
| | `    }` |
| | `}` |
| `if (noNote B) {` | |
| `    if (noMilk) {` | |
| `        buy Milk;` | |
| `        …` | |
| | `remove note B;` |

- Really insidious:
  - Unlikely that this would happen, but will at worse possible time

# Too Much Milk Solution #2: problem!



- *I'm* not getting milk, *You're* getting milk
- This kind of lockup is called "starvation!"

# Too Much Milk Solution #3

- Here is a possible two-note solution:

<table>
<tr><td>Thread A</td><td>Thread B</td></tr>
</table>

```
      Thread A                    Thread B
leave note A;               leave note B;
while (note B) { //X        if (noNote A) { //Y
    do nothing;                 if (noMilk) {
}                                   buy milk;
if (noMilk) {                   }
    buy milk;               }
}                           remove note B;
remove note A;
```

- Does this work? Yes. Both can guarantee that:
  - It is safe to buy, or
  - Other will buy, ok to quit
- At X:
  - if no note B, safe for A to buy,
  - otherwise wait to find out what will happen
- At Y:
  - if no note A, safe for B to buy
  - Otherwise, A is either buying or waiting for B to quit

# Solution #3 discussion

- Our solution protects a single "Critical-Section" piece of code for each thread:

```
if (noMilk) {
    buy milk;
}
```

- Solution #3 works, but it's really unsatisfactory
  - Really complex – even for this simple an example
    - » Hard to convince yourself that this really works
  - A's code is different from B's – what if lots of threads?
    - » Code would have to be slightly different for each thread
  - While A is waiting, it is consuming CPU time
    - » This is called "busy-waiting"
- There's a better way

# Too Much Milk: Solution #4

- Suppose we have some sort of implementation of a lock (more in a moment).
    - `Lock.Acquire()` – wait until lock is free, then grab
    - `Lock.Release()` – Unlock, waking up anyone waiting
    - These must be atomic operations – if two threads are waiting for the lock and both see it's free, only one succeeds to grab the lock

- Then, our milk problem is easy:

    ```
    milklock.Acquire();
    if (noMilk)
        buy milk;
    milklock.Release();
    ```

- Once again, section of code between `Acquire()` and `Release()` called a "Critical Section"

# Summary

- Concurrent threads are a very useful abstraction
  - Allow transparent overlapping of computation and I/O
  - Allow use of parallel processing when available
- Concurrent threads introduce problems when accessing shared data
  - Programs must be insensitive to arbitrary interleaving
  - Without careful design, shared variables can become completely inconsistent
- Important concept: Atomic Operations
  - An operation that runs to completion or not at all
  - These are the primitives on which to construct various synchronization primitives
- How to protect a critical section with only atomic load and store $\Rightarrow$ pretty complex!

# Next

- Reading: Chapter 6.