

# Concurrency and Mutual Exclusion

CSE 332 Spring 2021

**Instructor:** Hannah C. Tang

## Teaching Assistants:

Aayushi Modi	Khushi Chaudhari	Patrick Murphy
Aashna Sheth	Kris Wong	Richard Jiang
Frederick Huyan	Logan Milandin	Winston Jodjana
Hamsa Shankar	Nachiket Karmarkar	

- ❖ Define “non-determinism”
- ❖ Consider this sequential code from *L5: Recursive Algorithm Analysis*, which sums the elements of an array (!!)
  - Where in memory is `ARRLEN`? `mid`? `arr`'s contents?

```
int sum(int[] arr) {
    return help(arr, 0, arr.length);
}
int help(int[] arr, int lo, int hi) {
    if(lo == hi)    return 0;
    if(lo == hi-1) return arr[lo];
    int mid = (hi+lo)/2;
    return help(arr, lo, mid) + help(arr, mid, hi);
}

static final int ARRLEN = 4;
int main(int[] arr) {
    int[] arr = new int[ARRLEN];
    sum(arr);
    return 0;
}
```

# Announcements

# Lecture Outline

- ❖ **Farewell to Parallelism**
- ❖ Sharing Resources
- ❖ Concurrency: Managing Correct Access to Shared Resources
- ❖ Mutual Exclusion and Critical Sections

# Parallelism Recap (1 of 3)

- ❖ We studied two parallelism primitives
  - Map
  - Reduce
  
- ❖ We combined these primitives into complex parallel algorithms
  - sum: reduction
  - prefix: reduction + map
  - pack: map + prefix + map (or just prefix)
  - quicksort: parallelized recursive calls; partition using pack + pack
  - mergesort: parallelized recursive calls; merge using parallelized recursive calls

## Parallelism Recap (2 of 3)

- ❖ We studied one parallelism model in detail: ForkJoin
  - ... and noted where we “plugged in” our logic

```
class SumThread extends java.lang.Thread {
    // ... member fields and constructors elided ...
    public void run() { // override: implement "main"
        if (hi - lo < SEQUENTIAL_CUTOFF) {
            // Just do the calculation in this thread
            for (int i=lo; i < hi; i++)
                ans += arr[i];
        }
        else {
            // Create two new threads to calculate the left and right sums
            SumThread left = new SumThread(arr, lo, (hi+lo)/2);
            SumThread right = new SumThread(arr, (hi+lo)/2, hi);
            left.start();
            right.start();

            // Combine their results
            left.join();
            right.join();
            ans = left.ans + right.ans;
        }
    }
}
```

## Parallelism Recap (3 of 3)

- ❖ ... which will also help us understand other parallelism models and where we need to “plug in” our code

```
output = new array of size bitsum[n-1]
FORALL (i=0; i < input.length; i++){
    output[          ] =
}
```

```
class SumMapper extends org.apache.Hadoop.mapreduce.Mapper {
    // ... member fields and constructors elided ...
    public void map(Object mapkey, Integer mapval, Context context) {
        context.write(mapkey, new IntWritable(mapval));
    }
}

class SumReducer extends org.apache.Hadoop.mapreduce.Reducer {
    public void reduce(Object redkey, Iterable<IntWritable> redvals,
        Context context) {
        int sum = 0;
        for (IntWritable v : redvals) {
            sum += v.get();
        }
        context.write("result", new IntWritable(sum));
    }
}
```

# Lecture Outline

- ❖ Farewell to Parallelism
- ❖ **Sharing Resources**
- ❖ Concurrency: Managing Correct Access to Shared Resources
- ❖ Mutual Exclusion and Critical Sections



# Review: Parallelism and Sharing Resources

- ❖ We've studied **parallel algorithms** using the fork-join model and focused on reducing span via parallel tasks
- ❖ This model has a simple structure to avoid **race conditions**
  - Each thread had a part of memory the “only it accessed”
    - Example: each array sub-range accessed by only one thread
  - Result of forked executor not accessed until after `join()` called
  - Structure (mostly) ensures bad simultaneous access wouldn't occur

# Parallelism's Pitfall

- ❖ Fork-join model doesn't work well when:
  - Executors implementing the same algorithm access **overlapping memory**
  - Executors implementing different algorithms access **the same resources**
    - (rather than implementing the same algorithm)

# Parallelism: Non-overlapping Sharing

```
class SumTask extends RecursiveTask<Integer> {
    int lo; int hi; int[] arr;    // just the "input" arguments!

    protected Integer compute() { // override: implement "main"
        if(hi - lo < SEQUENTIAL_CUTOFF) {
            // Just do the calculation in this thread
            int ans = 0; // local variable instead of a member field
            for (int i=lo; i < hi; i++)
                ans += arr[i];
            return ans; // direct return of answer
        } else {
            // Create ONE new thread to calculate the left sum
            SumTask left = new SumTask(arr, lo, (hi+lo)/2);
            SumTask right = new SumTask(arr, (hi+lo)/2, hi);
            left.fork(); // create a thread and call its compute()
            int rightAns = right.compute(); // call compute() directly

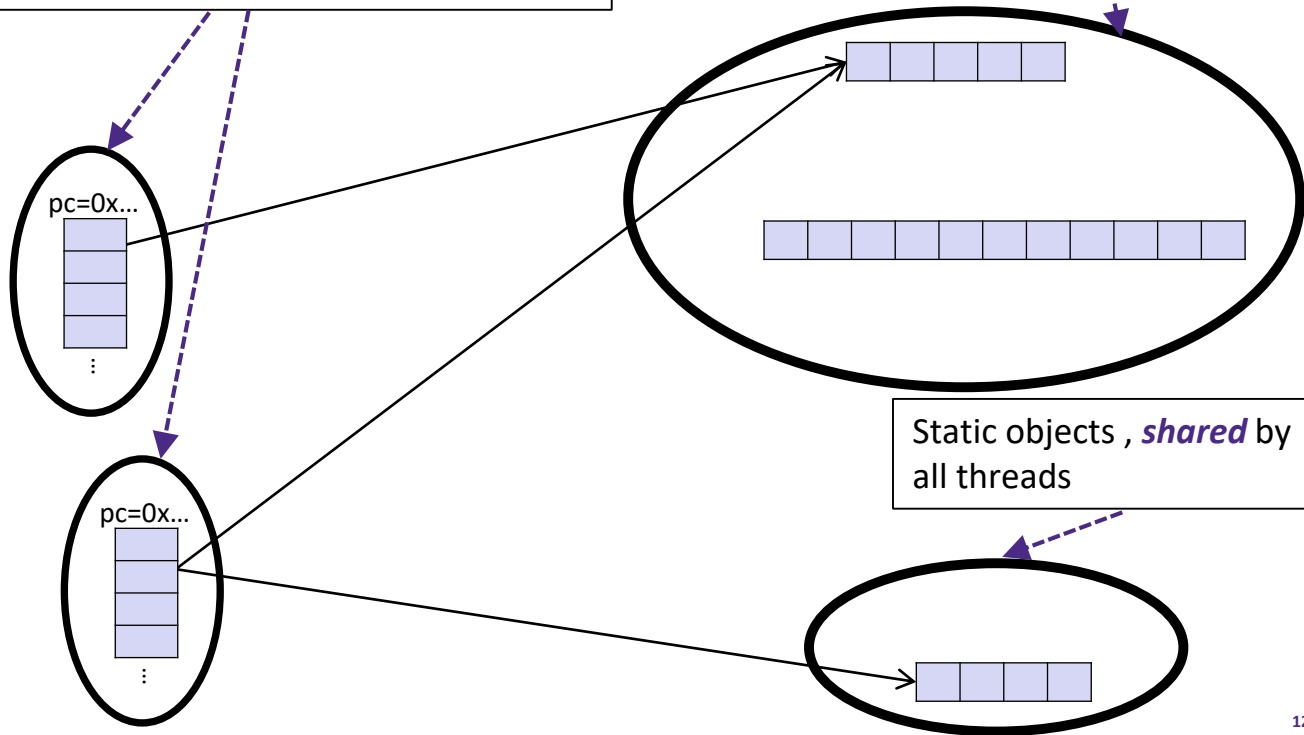
            // Combine results
            int leftAns = left.join();
            return leftAns + rightAns;
        }
    }
}
```

# Can Overlapped Sharing Happen?

Multiple threads, each with its own *independent* call stack and program counter

Heap for allocated objects, *shared* by all threads

Static objects, *shared* by all threads



# Overlapped Sharing (1 of 2)

- ❖ Threads are not just useful for parallelism
  - i.e., not always about implementing algorithms faster
  
- ❖ Threads are useful for:
  - Responsiveness
    - Respond to events in one thread while another is performing computation
  - Processor utilization (hide I/O latency)
    - If 1 thread “goes to disk,” process still has something else to do
  - Failure isolation
    - Prevent an exception in one task from stopping conceptually-parallel tasks

## Overlapped Sharing (2 of 2)

- ❖ What if we have multiple threads:
  - Processing different bank-account operations
    - What if 2 threads modify the same account at the same time?
  - Using a shared cache (e.g., hashtable) of recent files
    - What if 2 threads insert the same file at the same time?
  - Creating a pipeline (think assembly line) with a queue for handing work from one thread to next thread in sequence
    - What if enqueueer and dequeuer adjust a circular array queue at the same time?

# Sharing a Queue

- ❖ Imagine 2 threads
  - Running at the same time
  - Accessing a *shared linked-list-based queue*, initially empty

```
enqueue(x) {  
    if (back == null) {  
        back = new Node(x);  
        front = back;  
    }  
    else {  
        back.next = new Node(x);  
        back = back.next;  
    }  
}
```

# Overlapped Sharing Needs Concurrency

- ❖ **Concurrency**: *Correctly and efficiently* managing access to shared resources from multiple possibly-simultaneous clients
  - Requires *coordination*, particularly **synchronization**, to avoid incorrect simultaneous access
    - Make thread *block* (wait) until the resource is free
    - **join** is not what we want
    - Want other thread to be “done using *what we need*”, not “completely done executing”
- ❖ Correct concurrent applications are usually highly **non-deterministic**
  - How threads are scheduled affects order of operations
  - Non-repeatability complicates testing and debugging



# Attributes of Concurrent Programs

- ❖ In concurrent programs, it is common that:
  - Threads access the same resources in an *unpredictable order*
  - Threads access the same resources at (*approx.*) *the same time*
  - Correctness requires that simultaneous access be prevented
  - Simultaneous access is rare
    - Makes testing and debugging difficult
    - Rare != Impossible; need to be disciplined when designing / implementing
  
- ❖ In other words: concurrent programs are non-deterministic

# Lecture Outline

- ❖ Farewell to Parallelism
- ❖ Sharing Resources
- ❖ **Concurrency: Managing Correct Access to Shared Resources**
- ❖ Mutual Exclusion and Critical Sections

# Concurrency: Canonical Example

- ❖ In a single-threaded world, this code is correct!

```
class BankAccount {
    private int balance = 0;

    protected int getBalance()      { return balance; }
    protected void setBalance(int x) { balance = x; }

    public void withdraw(int amount) {
        int b = getBalance();
        if (amount > b)
            throw new WithdrawTooLargeException();
        setBalance(b - amount);
    }

    // ... other operations like deposit(), etc.
}
```

# Interleaving

- ❖ Suppose:
  - Thread **T1** calls **`x.withdraw(100)`**
  - Thread **T2** calls **`y.withdraw(100)`**
- ❖ If second call starts before first finishes, we say they **interleave**
  - e.g. T1 runs for 50 ms, pauses somewhere, T2 picks up for 50ms
  - Can happen with one processor; if **pre-empted** due to time-slicing
- ❖ If **`x`** and **`y`** refer to different accounts, no problem
  - “You cook in your kitchen while I cook in mine”
  - But if **`x`** and **`y`** alias, possible trouble...

# Activity: What is the Balance at the End?

- ❖ Two threads both `withdraw()` from the same account:
  - Assume initial balance == 150

```
class BankAccount {
    private int balance = 0;

    protected int getBalance() { return balance; }
    protected void setBalance(int x) { balance = x; }

    public void withdraw(int amount) {
        int b = getBalance();
        if (amount > b)
            throw new WithdrawTooLargeException();
        setBalance(b - amount);
    }

    // ... other operations, etc.
}
```

Thread A

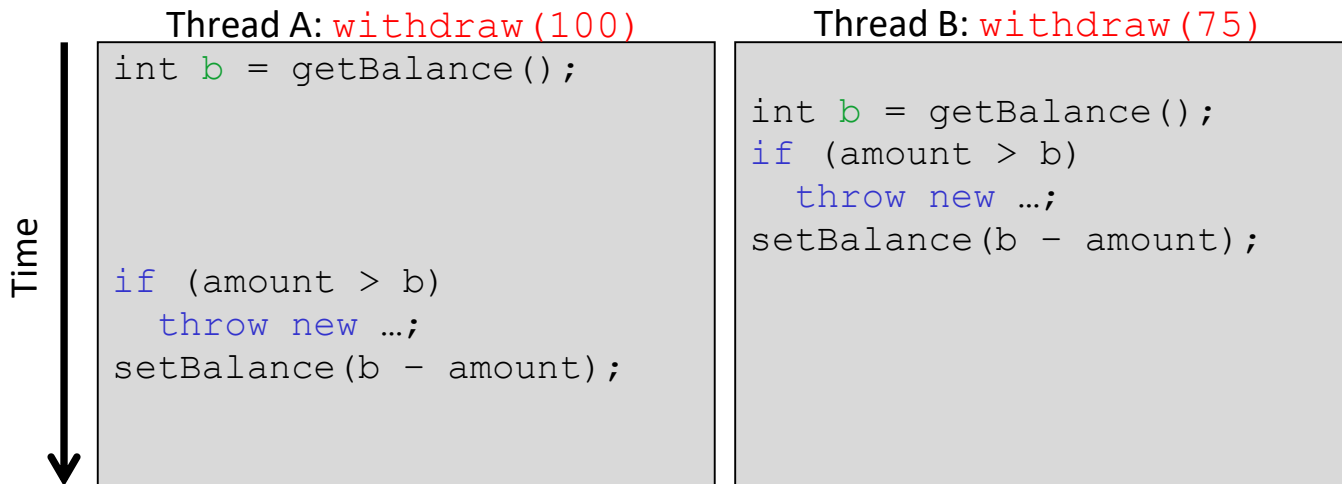
```
x.withdraw(100);
```

Thread B

```
x.withdraw(75);
```

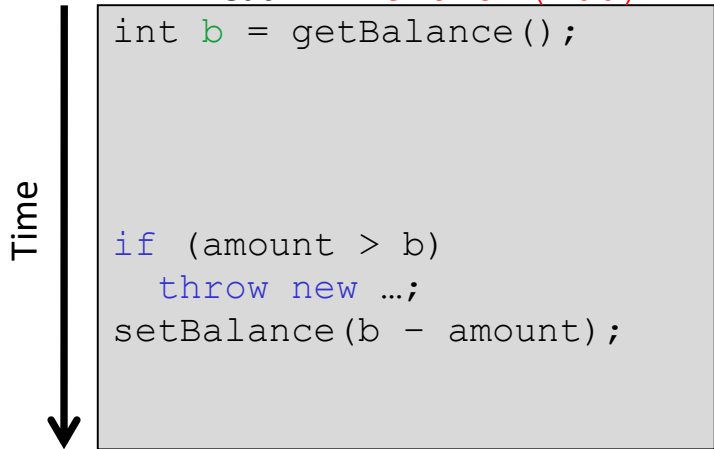
# A Bad Interleaving

- ❖ Interleaved `withdraw()` calls on the same account
  - Assume initial balance == 150
  - This *should* cause a `WithdrawTooLarge` exception (but doesn't)



- ❖ Find two more bad interleavings for `withdraw()`

Thread A: `withdraw(100)`

A vertical arrow on the left side of the diagram points downwards and is labeled "Time". To the right of the arrow are two rectangular boxes representing the execution of Thread A and Thread B. Thread A's box contains the code for a `withdraw(100)` operation. Thread B's box contains the code for a `withdraw(75)` operation. The boxes are positioned such that they represent a specific interleaving of the two threads' execution over time.

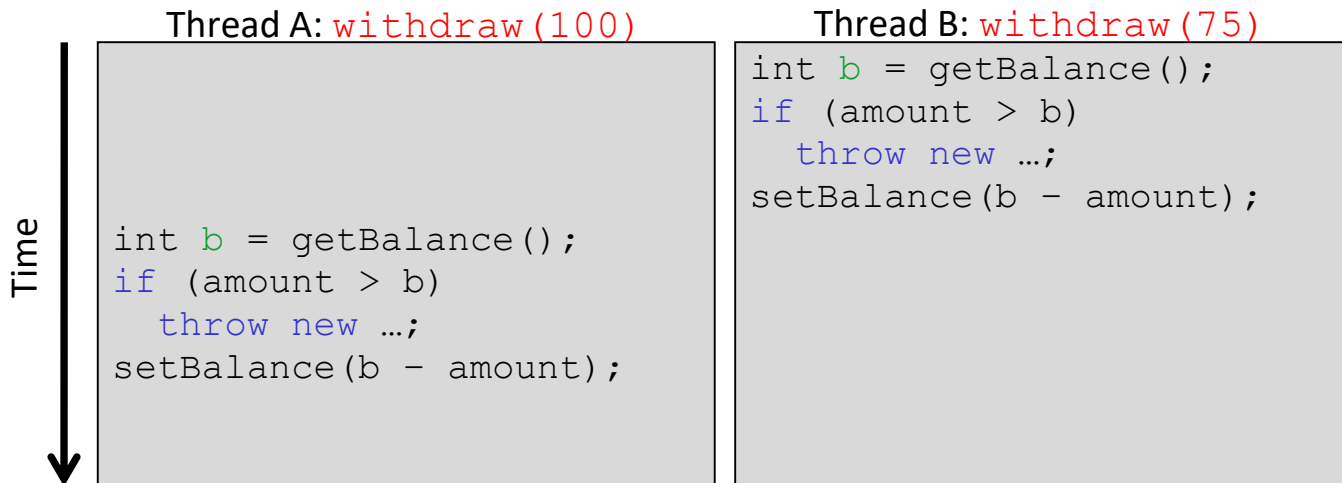
```
int b = getBalance();  
  
if (amount > b)  
    throw new ...;  
setBalance(b - amount);
```

Thread B: `withdraw(75)`

```
int b = getBalance();  
if (amount > b)  
    throw new ...;  
setBalance(b - amount);
```

# A Good Interleaving is Also Possible

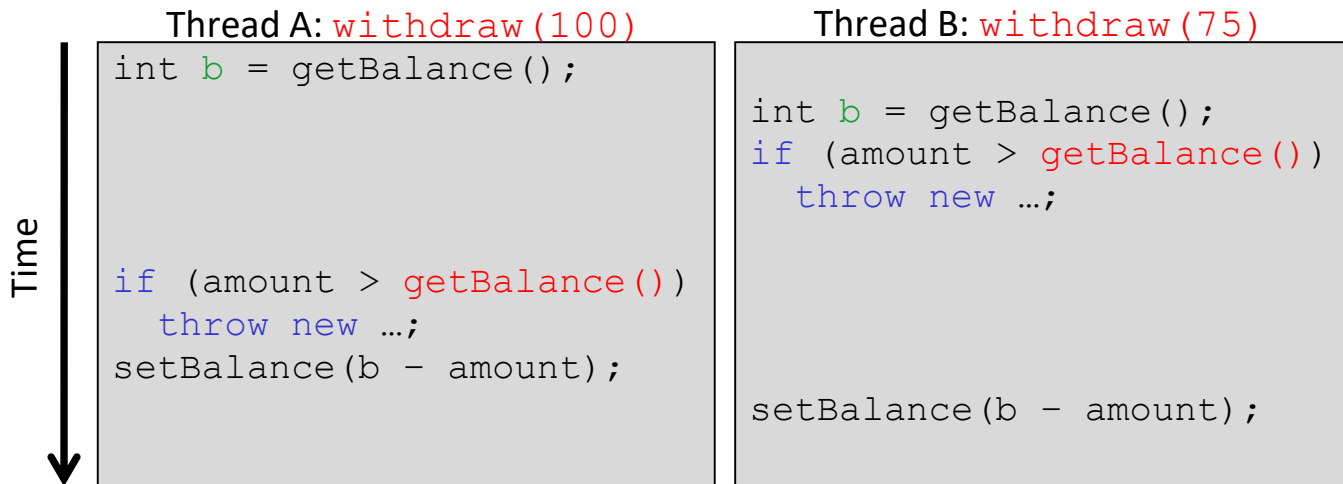
- ❖ Interleaved `withdraw()` calls on the same account
  - Assume initial balance == 150
  - This *does* cause a `WithdrawTooLarge` exception





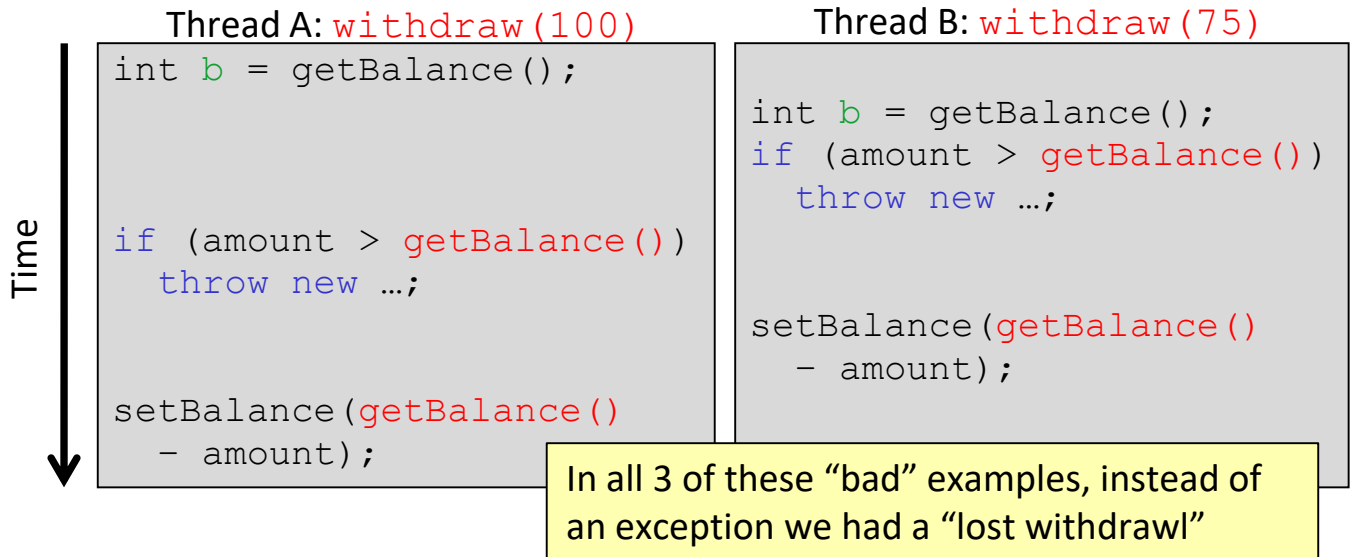
# A Bad Fix: Another Bad Interleaving

- ❖ Interleaved `withdraw()` calls on the same account
  - Assume initial balance == 150
  - This *should* cause a `WithdrawTooLarge` exception (but doesn't)



# ANOTHER Bad Fix: Another Bad Interleaving

- ❖ Interleaved `withdraw()` calls on the same account
  - Assume initial balance == 150
  - This *should* cause a `WithdrawTooLarge` exception (but doesn't)



# Incorrect “Fixes”

- ❖ It is tempting *and almost always wrong* to try fixing a bad interleaving by rearranging or repeating operations, such as:

```
public void withdraw(int amount) {  
    if (amount > getBalance())  
        throw new WithdrawTooLargeException();  
  
    // Maybe the balance was changed  
    setBalance(getBalance() - amount);  
}
```

- ❖ This fixes nothing!
  - Potentially narrows the problem by one statement
  - And that’s not even guaranteed!
    - The compiler could optimize it into the old version, because you didn’t indicate a need to synchronize

# Lecture Outline

- ❖ Farewell to Parallelism
- ❖ Sharing Resources
- ❖ Concurrency: Managing Correct Access to Shared Resources
- ❖ **Mutual Exclusion and Critical Sections**

# The Correct Fix: Mutual Exclusion

- ❖ Want at most one thread at a time to withdraw from account A
  - Exclude other simultaneous operations on A (e.g., deposit)
- ❖ More generally, we want **mutual exclusion**:
  - One thread using a resource means another thread must wait
- ❖ The area of code needing mutual exclusion is a **critical section**
- ❖ Programmer (you!) must identify and protect critical sections:
  - Compiler doesn't know which interleavings are allowed/disallowed
  - But you still need system-level primitives to do it!

# Why Do We Need System-level Primitives?

- ❖ Why can't we implement our own mutual-exclusion protocol?
  - Can we coordinate it ourselves using a boolean variable "**busy**"?
  - Possible under certain assumptions, but won't work in real languages

```
class BankAccount {
    private int balance = 0;
    private boolean busy = false;

    public void withdraw(int amount) {
        while (busy) { /* "spin-wait" */ }
        busy = true;
        int b = getBalance();
        if (amount > b)
            throw new WithdrawTooLargeException();
        setBalance(b - amount);
        busy = false;
    }
    // deposit() would spin on same boolean
}
```



# What We Actually Need: Lock ADT

- ❖ All ways out of this conundrum require system-level support
- ❖ One solution: **Mutual-Exclusion Locks** (aka **Mutex**, or just **Lock**)
  - For now, still discussing concepts; `Lock` is not a Java class
- ❖ We will define **Lock** as an ADT with operations:
  - **new**: make a new lock, initially “not held”
  - **acquire**: blocks current thread if this lock is “held”
    - Once “not held”, makes lock “held”
    - Checking & setting the “held” boolean is a single uninterruptible operation
    - Fixes problem we saw before!!
  - **release**: makes this lock “not held”
    - If  $\geq 1$  threads are blocked on it, another thread – but only one! – can now acquire



# Why a System-level Lock Works

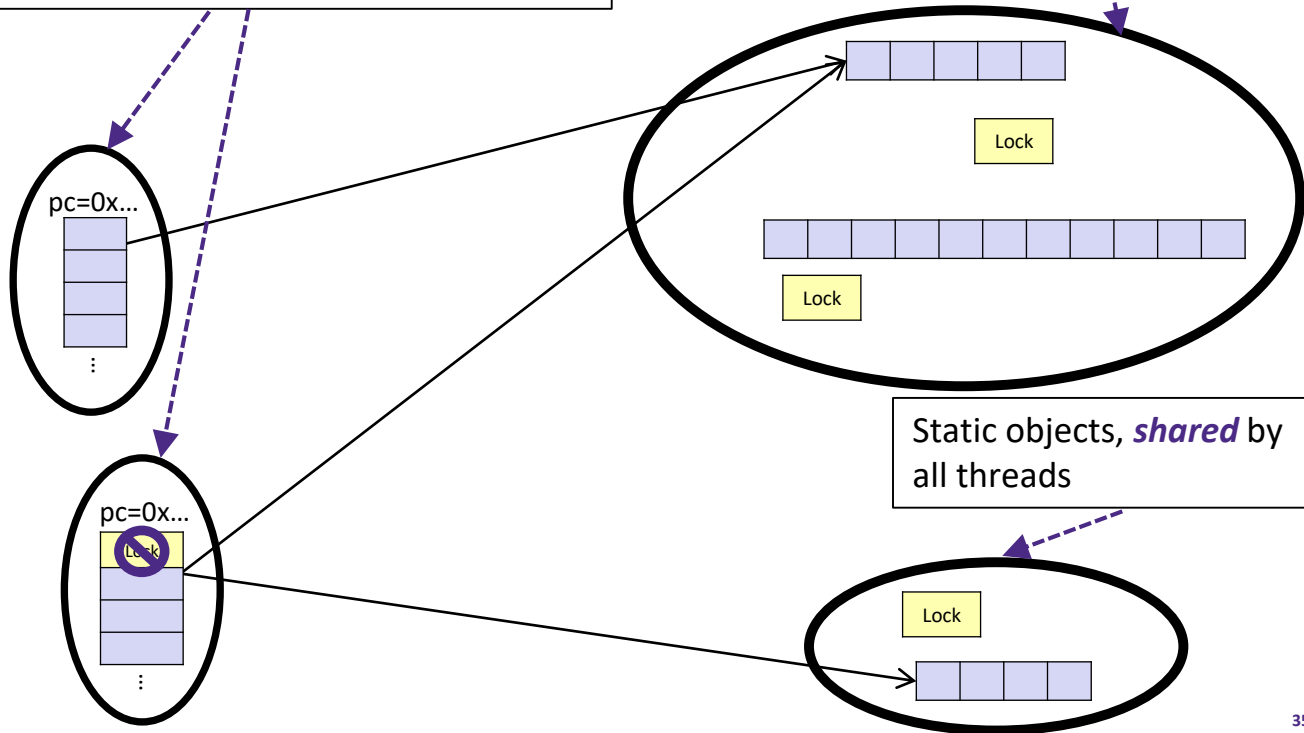
- ❖ Lock must ensure that, given simultaneous acquires/releases, “the correct thing” will happen
  - E.g.: if we have two acquires: one will “win” and one will block
  
- ❖ How can this be implemented?
  - The key is that the “check if held; if not, make held” operation must happen “all-at-once”. It cannot be interrupted!
  - Thus, requires and uses hardware and O/S support
    - See computer-architecture or operating-systems course
  - In CSE 332, we’ll assume a lock is a primitive and just use it

# Locks Must Be Accessible By Multiple Threads!

Multiple threads, each with its own *independent* call stack and program counter

Heap for allocated objects, *shared* by all threads

Static objects, *shared* by all threads



# Almost-Correct Pseudocode

```
class BankAccount {
    private int balance = 0;
    private Lock lk = new Lock();

    public void withdraw(int amount) {
        lk.acquire(); // may block
        int b = getBalance();
        if (amount > b)
            throw new WithdrawTooLargeException();
        setBalance(b - amount);
        lk.release();
    }

    // deposit() would also acquire/release lk
}
```

Note: 'Lock' is not an actual Java class

1. Where is the critical section?
2. How many locks do we need?
  - a) One lock per BankAccount object?
  - b) Two locks per BankAccount object?
    - i.e., one for withdraw() and one for deposit()
  - c) One lock for the entire Bank
    - Bank contains multiple BankAccount instances
3. There is a bug in withdraw(), can you find it?
4. Do we need locks for:
  - a) getBalance?
  - b) setBalance?

# Some Common Locking Mistakes (1 of 2)

- ❖ A lock is very primitive; up to you to use correctly
- ❖ **Incorrect**: different locks for **withdraw** and **deposit**
  - Mutual exclusion works only when sharing same lock
  - **balance** field is the shared resource being protected
- ❖ **Poor performance**: same lock for entire Bank
  - No simultaneous operations on *different* accounts

## Some Common Locking Mistakes (2 of 2)

- ❖ **Bug**: forgot to release a lock when exiting early
  - Can block other threads forever if there's an exception

```
if (amount > b) {  
    lk.release(); // hard to remember!  
    throw new WithdrawTooLargeException();  
}
```

Remembering to release() before every exit is challenging!

- ❖ What about **getBalance** and **setBalance**?
  - Assume now that they are **public** (which may be reasonable)
  - If they **do not acquire the same lock**, then **setBalance** and **withdraw** could interleave badly and produce a wrong result
  - If they **do acquire the same lock**, then **withdraw** would block forever because it tries to acquire a lock it already has!

# Summary

- ❖ Threads are useful beyond just fork-join-style parallelism
  - But general use-cases require **concurrency** to ensure correctness when dealing with overlapped sharing
- ❖ Overlapped sharing introduces **non-determinism** because the system controls the scheduling of threads
  - Therefore, the system must also provide **locks** to ensure **mutual exclusion** in **critical sections** of code
  - Mutual exclusion is the technique we employ to prevent **bad interleavings**