# Java Concurrency In a Nutshell

COM 3563: Database Implementation

Avraham Leff

Yeshiva University

*avraham.leff@yu.edu*

COM3563: Fall 2020

# Today's Lecture: Overview

> - Note: we will likely not finish this material in class.
> - Because this material is not "core" *Database Implementation* material, am not allocating more than one lecture
> - For your (project implementation etc) sake: review (and think about) <u>all</u> the material

# Why A Lecture On "Java Concurrency"?

# How Is "Concurrency" Relevant To "Database Implementation"?

- ▸ As you probably know, this course is not *Parallel Programming* ☺
- ▸ Q: so, why a lecture on Java concurrency?
- ▸ A: because databases, like all "service providing" software, must perform well
  - ▸ Or: they won't have any customers ☺
- ▸ Throughput: a fundamental performance metric
  - ▸ *"How many client requests serviced per given time unit?"*
- ▸ Implication: your DBMS must support concurrent, multi-client access to its data
- ▸ Your implementation must therefore address a fundamental tension
  1. Must allow concurrent access ("the more clients, the better")
  2. Must ensure that concurrent access doesn't corrupt data-structures ("the fewer clients, the easier to implement")

# Why Java Concurrency?

- ▸ Having established that understanding concurrency issues & solutions <u>is</u> relevant to our course …
- ▸ Why: Java concurrency?
  - ▸ After all: many (most?) enterprise databases are written in C or C++
  - ▸ (See this discussion, for example)
- ▸ More importantly: the Java concurrency model is definitely different from C/C++
  - ▸ Java concurrency libraries are irrelevant to other languages ☹
- ▸ Frankly: was tempted to ignore the whole issue (as the textbook does ☺)
- ▸ But: you'll be implementing the course project in Java
  - ▸ To succeed, you must be competent in and comfortable with Java concurrency
  - ▸ Otherwise: your DBMS & my test cases will not get along ☹

You will have to teach yourself Java concurrency skills as you complete the course project

- This lecture will <u>not</u> teach you all or even much of what you need to know about Java concurrency
  - Consider taking *Parallel Programming* ☺
- Today's lecture is a guided tour, tailored to helping you succeed in the course project
  - Make you aware of the pitfalls and some solution techniques
  - Steer you away from useless or incorrect information on the Internet

# Broader Importance Of This Material

- ▸ <u>Any</u> server-side software must allow for (and <u>correctly manage</u>) concurrency
- ▸ Keep in mind
  - ▸ Java is a heavy-duty, industrial-strength language
  - ▸ Its approach to concurrency is one of its strengths
  - ▸ Being able to do concurrent programming is a key differentiator between amateur and professional programmers
- ▸ I want you to incorporate these skills into your software portfolio
- ▸ <u>Final point:</u>
  - ▸ Yes: the <u>solutions</u> we'll discuss today are specific to Java
  - ▸ But: the <u>issues</u> in (shared memory) concurrent programming that we're discussing are language independent

- ▸ Primary goal of using concurrency is maximizing use of CPU (s)
- ▸ But: the strategy for maximizing CPU usage has changed over time
- ▸ Single-core era: use non-blocking I/O as an alternative to "blocking" code
    - ▸ Also: prioritized background tasks
- ▸ Multi-core era: Coarse-grained, task-based concurrency
    - ▸ Largely about throughput: pushing more requests through a server
    - ▸ That's what we'll focus on for this course
- ▸ Many-core era: Fine-grained data parallelism
    - ▸ Largely about latency: "use more cores to get the answer faster"

Credits to Brian Goetz, Java Language Architect for these insights

- ▸ Hardware (previous slide) shapes the languages, library, and frameworks we write
- ▸ Java 1: supported threads, locks, condition queues
- ▸ Java 5: added thread pools, blocking queues, concurrent collections
- ▸ Java 7: added the fork-join library
- ▸ Java 8: added parallel streams

---

- ▸ You therefore want to focus on the Java 5 additions
- ▸ `java.util.concurrent` and sub-packages (`atomic` and `locks`)

- ▸ Since COM 1300, you've been "code-centric"
    1. You write code
    2. Computer executes that code "one line at a time"
- ▸ Your mind-set (for coding & debugging): *"I only have to worry about a single line of code at any given moment"*
    - ▸ Because the computer will execute exactly one line of code at any given moment

That mind-set is (mostly) incorrect ☺

- ▶ What was really going on:
    1. You write code
    2. You create a thread: an independent execution "engine"
    3. That thread executes your code "one line at a time"

- ▶ Key point:
    - ▶ You can create multiple threads, and then direct each thread to independently execute that code

- ▶ Until now:
    - ▶ By default: the `java` command creates only one thread to execute your program
    - ▶ That thread is labeled as `main`
    - ▶ Typically, one thread suffices to get your work done, but not any more ☺

- I've already motivated this shift in your programming paradigm: performance, performance, performance
- Contrast to (for example) a *Data Structures* project
  - Your focus: write <u>correct</u> code that meets requirements
  - Elegant code is a bonus
- *Introduction to Algorithms* added the idea of "pick the right algorithm"
  - To break the "Big-O" barrier
  - But fundamentally: <u>only one</u> thread of computation is executing your code
- <u>Now</u>: your DBMS code is a "service", available to multiple, <u>independent</u> clients
  - In other words: each of these "independent clients" is (in OS terms) a thread (or a process)
  - Because **OATSdb** is an embedded database, each of these "independent clients" is a Java thread (not a process)

- ▸ Theoretically: your DBMS <u>could</u> insist "I'll service only one client at a time"
  - ▸ This would definitely make it much easier to write "correct & elegant code" ☺
  - ▸ But: your throughput metric will be very disappointing
- ▸ Implication

  > - ▸ Your development <u>test-bed</u> will have to create and manage multiple clients (threads)
  > - ▸ Your DBMS implementation will have to ensure that these multiple threads don't corrupt your data-structures

- ▸ We say: *"your DBMS must support concurrency"*

- Stop thinking of your code as a "live" entity that can do only one thing at a given moment
  - Example: a *Stack* which can do <u>either</u> pop <u>or</u> push but <u>not both</u> at the same time
- Instead: your code is a "static" entity that is "infused with life" by being executed by a Thread
- If multiple threads execute your code, any number of your code points can be executed at the same time
  - Example: Stack.pop <u>and</u> Stack.push
- Murphy's Law says: the worst possible combination of concurrent executions <u>will</u> occur ☹

1. Explain why multi-threaded programs require a paradigm shift on your part
   - In other words: the "concurrency problem"
2. Present solutions to the "concurrency problem"
3. Turn from the "server-side" DBMS issues to "client-side" test-bed issues

Since COM 1300 you've become accustomed to rely on the notion of invariants to establish program correctness

*In computer programming, specifically object-oriented programming, a class invariant (or type invariant) is an invariant used for constraining objects of a class. Methods of the class should preserve the invariant. The class invariant constrains the state stored in the object.*

*Class invariants are established during construction and constantly maintained between calls to public methods. Code within functions may break invariants as long as the invariants are restored before a public function ends.*

*Source: Wikipedia*

- Once we enter the world of multi-threaded programming, we can no longer rely on invariants to establish program correctness
- In reality, methods are not executed atomically
- $Thread_1$, executing $method_1$, may have its execution interleaved with $Thread_2$, executing $method_2$
- Implication
  - $thread_1$ may see "intermediate" state affected by $method_2$ execution
  - $thread_2$ may see "intermediate" state affected by $method_1$ execution

# Bad Interleavings: Shared State Transitions Are Now Exposed

- ▸ The whole point of "object-oriented" programming is to encapsulate shared state
- ▸ Concurrent thread executions imply that (<u>unless we're careful</u>) we've broken encapsulation ☹
- ▸ Typical scenarios of different methods accessing shared state
  - ▸ One thread deposits money in a bank-account, other thread does a withdrawal
  - ▸ "Producer" thread enqueues a task in a queue, "consumer" thread dequeues a task
  - ▸ One thread does a *Hashtable.get*, other thread does a *Hashtable.put*
- ▸ Method abstraction conceals the fact that e.g., a single *Hashtable.get* is "really" multiple lower-level machine instructions
  - ▸ And thus can be interleaved with the execution of *Hashtable.put*
  - ▸ With disastrous results ☹

This code is absolutely correct in a single-threaded environment

```
1  class BankAccount {
2    private int balance = 0;
3    int getBalance() { return balance; }
4    void setBalance(int x) { balance = x; }
5    void withdraw(int amount) {
6      int b = getBalance();
7      if (amount > b) {
8        throw new WithdrawTooLargeException();
9      }
10
11     setBalance(b - amount);
12   }
13
14   // other operations like deposit, etc.
15 } // class BankAccount
```

# Multi-Threaded Environment: Interleaving

Scenario (*x* and *y* are *BankAccount* instances):

> - Thread $T_1$ calls x.withdraw(100)
> - Thread $T_2$ calls y.withdraw(100)

- If $T_1$ and $T_2$ executions overlap (to any extent) , we say the executions interleave
- Because of OS "time-slicing", can happen even with a single processor
- No problem if the threads are accessing different account instances
- But: if the threads access the same account, we have absolutely no idea what's going to happen

| Thread 1 | Thread 2 |
|---|---|
| `int b = getBalance();` |  |
|  | `int b = getBalance();`<br>`if(amount > b)`<br>`  throw new …;`<br>`setBalance(b - amount);` |
| `if(amount > b)`<br>`  throw new …;`<br>`setBalance(b - amount);` |  |

▸ Both withdrawals are successful, customer gets to "double-dip"

▸ This is a "lost-withdrawal" scenario, bank is very unhappy ☹

# Rearranging Code Doesn't Work

```java
1  void withdraw(int amount) {
2    if (amount > getBalance()) {
3      throw new WithdrawTooLargeException();
4    }
5    setBalance(getBalance() - amount);
6  }
```

- ▸ Rearranging or repeating code does not solve interleaving problems
- ▸ Here: only narrows the exposure, but other thread can still change the balance between second call to *getBalance* and reducing by *amount*
  - ▸ May not even do that ☺... since compiler may optimize back to original version
- ▸ May even change to "negative balance" scenario
  - ▸ Second call to *getBalance* occurs after other thread has already reduced balance
- ▸ Again: you simply can't assume that this won't happen

> A race condition occurs when the computation result depends on scheduling: specifically, how hardware & software interact to interleave the concurrent execution of threads

- ▸ I'm only introducing the term "race condition" because it's so prevalent
- ▸ Unfortunately, it's an overloaded term
  - ▸ It refers to "bad" interleavings (the problem we've been discussing so far)
  - ▸ <u>And</u> it refers to data races: a simultaneous read/write or write/write (by two threads) of the same memory location
- ▸ The common denominator of these two types of race conditions is that they're concurrency bugs
  - ▸ If only one thread, then no problem manifests

The good news: we'll use the same solution approach to both "bad interleavings" and "data races"

▸ (I'm therefore going to postpone discussion of "data race" problem until we've discussed this solution)
▸ Conceptual difference between "data race" and "bad interleaving" is that a "data race" condition is always an error!
▸ In contrast: only you (the programmer) can determine whether an interleaving is "bad"
▸ Depends on how you've defined your invariants (the "code specification")
  ▸ "Bad interleaving" exists *iff* your code exposes an intermediate state that violates an invariant

- ▸ Think of it this way: our concurrency problems are caused by allowing independent threads to run wild in our code
- ▸ Solution idea: impose restrictions on thread executions
  - ▸ Specifically: annotate our code to state *"only one thread can execute this block of code at a time"*
- ▸ Example: allow at most one thread to withdraw from account at a time
  - ▸ Note: this requires that we also exclude other concurrent operations such as *deposit*
- ▸ This technique is called mutual exclusion or the creation of critical sections

- ▸ Unfortunately: mutual exclusion cannot be enforced by the compiler ☹
  - ▸ The semantics are at a higher-level than the compiler understand
  - ▸ How <u>can</u> the compiler possibly understand what interleavings are OK and which ones are dangerous?
- ▸ We do need, and can get, support from Java language primitives
  - ▸ The language must be able to enforce mutual exclusion

# One Approach: Locks or Mutex Variables

- ▸ Some languages and libraries use the concept of lock or mutex variable to provide mutual exclusion
- ▸ Conceptually: two methods
  - ▸ `acquire`
  - ▸ `release`
- ▸ A critical section (such as accessing an account's `balance` variable) is …
- ▸ …Preceded by a call to `acquire` (*thread*$_2$ will be forced to block if *thread*$_1$ currently holds the lock)
- ▸ …Followed by a call to `release` (*thread*$_1$ will release the lock after it modifies `balance`, allowing *thread*$_2$ to acquire the lock and proceed)
- ▸ Locks require special hardware and os support to implement
  - ▸ We're requiring that *"check if lock is held and if not held acquire the lock"* be an atomic operation

# Lock Variables: Getting It Wrong

- ▸ <u>Wrong</u>: if you use <u>different</u> locks for `withdraw` and `deposit`
  - ▸ Mutual exclusion works only when the shared resource being protected (`balance` instance variable) is protected by the <u>same lock</u>
- ▸ <u>Bad Performance</u>: too coarse a lock granularity
  - ▸ Example: <u>don't use</u> one lock for all bank accounts ☺
- ▸ <u>Wrong</u>: if *thread*$_1$ "forgets" to release a lock, then it will block other threads forever
- ▸ Example:
  1. Thread acquires lock
  2. Thread invokes `withdraw`
  3. Withdrawal attempt triggers *WithdrawTooLargeException*
  4. The lock will never be released
  5. All threads are prevented from using this code forever

# Locks: Reentrantcy Problem

- ▸ Assume you've used locks correctly: `withdraw` and `deposit` use the same lock
- ▸ Now you implement `getBalance` and `setBalance`
- ▸ Two scenarios
    1. The new code uses a different lock than the "withdraw" lock
        - ▸ Incorrect: you've enabled a race condition between `setBalance` and `withdraw`
    2. The new code uses the same lock as the "withdraw lock"
        - 2.1 `withdraw` acquires the lock
        - 2.2 `withdraw` (internally) invokes `getBalance`
        - 2.3 Thread blocks forever as it tries to acquire the `getBalance` which it already has ☹
    3. One solution: allow locks to be reentrant

- ▸ A reentrant lock (sometimes called a *recursive lock*)
  - ▸ "Remembers" which thread has acquired it
  - ▸ Maintains a count variable
- ▸ When locks goes from "not held" to "acquired" state …
  - ▸ count is set to 1
- ▸ If acquiring thread (re)invokes "acquire" on this lock …
  - ▸ Thread does not block
  - ▸ Instead: count is incremented
- ▸ When thread invokes "release" on this lock …
  - ▸ count is decremented
  - ▸ If count now reaches 0, lock enters "not held" state

- ▸ I've included this discussion of the lock approach because
  - ▸ Many languages and libraries use it
  - ▸ Because this is the approach used in the JDK 5 *java.util.concurrent.locks.ReentrantLock* class(es) ☺
  - ▸ Because it will help you understand the Java language solution: object monitors
- ▸ My opinion only: the Java object monitor approach is superior for "vanilla" mutual exclusion
  - ▸ You can't beat built-in language support
  - ▸ Specifically "automatic `acquire` and `release`"
- ▸ Use the JDK 5 facilities for more complicated situations
  - ▸ And read the Javadoc carefully!

# Java Object Monitors: Introduction

- ▸ In Java, every object instance and every class is associated with a monitor
  - ▸ A monitor is implicitly associated with a lock
  - ▸ Java manages the magic transparently
- ▸ When a thread acquires the monitor lock, all other threads are blocked from acquiring that lock
- ▸ How does a thread acquire the monitor lock?
  - ▸ Not by accessing the object's <u>code</u>!
  - ▸ By (literally) synchronizing on the object <u>instance</u>

```
1    synchronized (o) {
2      // code goes here
3    }
```

- ▸ We say that the thread has synchronized on Object o
- ▸ JVM will block this thread if some other thread currently has "Object o lock"

```java
1    synchronized (o) {
2      // code goes here
3    }
```

▸ Thread has the lock if it makes it past beginning of synchronized block
▸ JVM will automatically release that lock when thread executions moves past the closing curly brace of the synchronized block
  ▸ Even if it exits via an Exception ☺
▸ Note (I can't stress this too much): complete decoupling of the thread …
  ▸ A "live thing", something which **executes code**
▸ From the lock …
  ▸ Which is associated with an object (or class) instance
  ▸ Static code and associated mutable state

# "Synchronization" Is Just a Lock

- ▸ I introduced Java synchronized keyword <u>after</u> discussing reentrant lock approach because
  - ▸ That's essentially what a Java monitor is: a reentrant lock associated with an object or class instance
  - ▸ For which the language and JVM support provide automatic "acquire" and "release" semantics
- ▸ But conceptually the same
- ▸ With all the perils thereto ☺
  - ▸ You are responsible for ensuring that the lock is used consistently to protect the relevant critical section throughout your code

## Synchronize On Private Object?

You <u>can</u> synchronize on explicit lock variable

```
1   class BankAccount {
2     private int balance = 0;
3     private Object lk = new Object();
4     // etc
5
6     void withdraw(int amount) {
7       synchronized(lk) {
8         // modify balance
9       }
10    }
11  }
```

But …
- ▸ Makes it harder for other code to synchronize on this BankAccount's operations
- ▸ And …why bother? Just synchronize on this
  - ▸ The BankAccount instance itself

```
1   class BankAccount {
2     private int balance = 0;
3     int getBalance() { synchronized (this) { return balance; }
4         }
5   }
```

```
1    class BankAccount {
2      private int balance = 0;
3      synchronized int getBalance() { return balance; }
```

▸ Syntactic sugar ...

▸ Putting synchronized before a method declaration is equivalent to surrounding the method body with a synchronized(this)

▸ Looks nicer ☺

▸ May cost you if the method consists of multiple, "long-duration" critical sections involving different shared resources

  ▸ Prevents other threads from accessing all of those resource

## "Data Races"

- ▸ In addition to solving the "bad interleaving" problem, mutual exclusion solves the "data race" problem
- ▸ <u>Key point</u>: sensitize yourself to the existence of such problems

*What is a Data Race?*

*The Thread Analyzer detects data-races that occur during the execution of a multi-threaded process. A data race occurs when:*

1. *two or more threads in a single process access the same memory location concurrently, and*
2. *at least one of the accesses is for writing, and*
3. *the threads are not using any exclusive locks to control their accesses to that memory.*

*When these three conditions hold, the order of accesses is non-deterministic, and the computation may give different results from run to run depending on that order. Some data-races may be benign (for example, when the memory access is used for a busy-wait), but many data-races are bugs in the program.*

*Sun Studio 12: Thread Analyzer User's Guide*

# Data Race: This Code Is Incorrect

```
1   class Stack<E> {
2     private E[] array = (E[]) new Object[SIZE];
3     int index = -1;
4     boolean isEmpty() { // unsynchronized: wrong?
5       return index==-1;
6     }
7
8     synchronized void push(E val) {
9       array[++index] = val;
10    }
11
12    synchronized E pop() {
13      return array[index--];
14    }
15
16    E peek() { // unsynchronized: wrong!
17      return array[index];
18    }
19  }
```

▸ Perhaps *isEmpty* and *peek* can "get away" with dropping synchronized?
  ▸ After all: looks like the state changes written or read by these methods are atomic

# Data Race Problems: You Must Synchronize

- ▸ Even "tiny steps" may require multiple steps in the compiled code implementation
  - ▸ `array[++index] = val` (probably) takes at least two steps
- ▸ Useful to think of a "data race" as a "simultaneous access error"
- ▸ Compared to "bad interleavings", a "data race" is a "lower-level" error
  - ▸ Absolutely, fundamentally, a problem
  - ▸ You simply cannot reason about the behavior of code that has data races
  - ▸ From the C++ reference: "If a data race occurs, the behavior of the program is undefined"
- ▸ Not going into how data races actually cause problems
  - ▸ That's a discussion for a compiler and/or computer architecture course ☺
- ▸ See the `java.util.concurrent.atomic` classes if you want "atomic" operations on individual variables
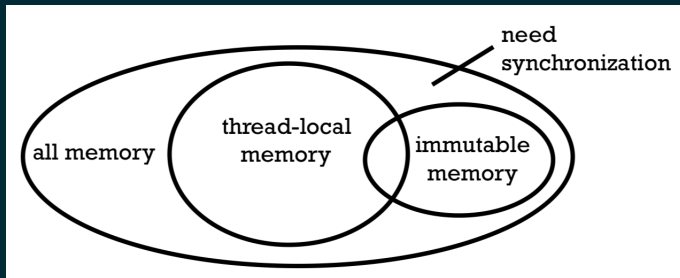
- ▸ Students' reaction to discussion so far is often *"no problem, I'll wrap all my methods with `synchronized`"* ☺
- ▸ One issue: may hurt performance because other threads can't make progress while this thread is executing a given method
  - ▸ They can execute neither <u>this</u> method <u>nor any other</u> method of this object
  - ▸ We'll elaborate on this issue later
- ▸ More fundamental problem: deadlock
  1. *Thread*$_1$ acquires the lock but blocks before releasing (e.g., "data structure is full")
  2. *Thread*$_2$, waiting to empty the data structure, can't acquire the lock
  3. Both threads block forever
- ▸ Classic example: a "producer/consumer" class
- ▸ Lesson: don't blindly "synchronize" everything

Concurrency Programming Guidelines

- Previous material is about making you aware of the multi-threading paradigm shift
- Why perfectly correct sequential code can break as soon as it's deployed in a multi-threaded environment
- The need for using mutual exclusion techniques to keep concurrent threads from stepping on one another
- Now: some guidelines as you begin writing your own concurrent code

Every memory location in your program should have (at least) one of the following properties

- ▸ Be thread-local: only one thread ever accesses it
- ▸ Be immutable: (after being initialized), memory is only read, never written
- ▸ Be synchronized: locks are used to ensure there are no race conditions

## Thread-Local Memory

- ▶ Simplest way to avoid data races is for only one thread to have access to memory ☺
- ▶ So: whenever possible, do not share resources among threads
- ▶ <u>Definitely</u> encourage you to read the `java.lang.ThreadLocal` Javadoc!
- ▶ If multiple threads need to access a resource, see if you can give each a copy of that resource
  - ▶ Example: `java.util.Random`
  - ▶ Only consider allowing threads to share a mutable resource if no other choice
- ▶ Note: this is why you don't need to synchronize on a method's local variables …
- ▶ The runtime's call-stack is (effectively) thread-local!

- ▸ Whenever possible, do not update objects
  - ▸ Instead: create new objects
- ▸ This is one of the key tenets of **functional programming**
  - ▸ Always (even in sequential environment) a good idea to avoid side-effects
  - ▸ An invaluable idea in a concurrent environment!
- ▸ Key point: concurrent read operations can never cause data races
- ▸ Programming tip: we're biased to "reusing" (mutating) object instances
  - ▸ We justify that bias by talking about "efficiency" ☺
  - ▸ Try to correct for that bias

▸ After you've minimized the amount of state shared & mutable state, the burden of providing mutual exclusion falls on you

▸ Some guidelines for these scenarios ...

▸ Guideline #0: <u>No data races</u>
  ▸ Never allow two threads to read/write or write/write the same location at the same time

▸ (Remember: avoiding "data races" will not necessarily fix "bad interleavings")

▸ Guideline #1: For each piece of state that requires mutual exclusion, associate a lock that must be acquired before thread can read or write that location
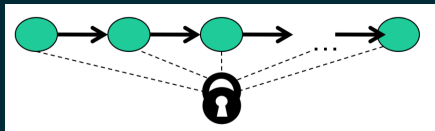  ▸ We say the lock guards the location

- Perfectly fine to use the same lock to guard multiple pieces of state
  - In Java, the guard is often the object containing the location
  - `this` inside an object's methods
  - But you can guard a larger data-structure with its own lock to provide more granular mutual exclusion
- Remember to document the guard for each location
- You – the programmer – get to decide how to partition the control of your "shared-and-mutable" state lock

# Flexibility: Change The Locking Protocol Dynamically

- ▸ We already have three techniques to avoid problems
  1. Thread-local
  2. Immutable
  3. Consistent locking
- ▸ Your code can dynamically switch between one technique and another
  - ▸ As long as all threads agree on the protocol being used at any given time
- ▸ Example: data-structure initialization may use locking technique because multiple threads are used for performance reasons
  - ▸ Then kill all threads but "main": now protected by "thread-local" technique
  - ▸ Or: once initialized, data-structure becomes immutable

# "Coarse-Grained" Versus "Fine-Grained" Locking



- ▸ Fewer locks, or "more objects per lock"
  - ▸ Example: One lock for entire data structure (e.g., array)
  - ▸ Example: One lock for all bank accounts

- ▸ More locks, or "fewer objects per lock"
  - ▸ Example: One lock per data element (e.g., array index)
  - ▸ Example: One lock per bank account

Not a dichotomy: this distinction is a continuum!

# "Coarse-Grained" Versus "Fine-Grained": Tradeoffs

- ▸ Coarse-grained advantages
  - ▸ Simpler to implement
  - ▸ Faster/easier to implement operations that access multiple locations
    - ▸ Because all guarded by the same lock
- ▸ Fine-grained advantages
  - ▸ Increase possibilities for concurrent access
  - ▸ This can improve performance in situations where coarse-grained locking leads to unnecessary blocking
- ▸ Guideline #2: Start with coarse-grained (simpler) and move to fine-grained (performance) only if contention on the coarser locks becomes an issue
  - ▸ Use a profiler before trusting your intuition!
  - ▸ Can guarantee that (at least initially) you'll introduce lots of bugs ☺

# "Critical-Section" Granularity

- ▸ We've been discussing lock granularity: how many objects per lock
- ▸ Orthogonal issue: critical-section granularity
  - ▸ How much "work" in the code that must be done while holding locks
- ▸ Critical sections are too long?
  - ▸ Performance loss because other threads are blocked
- ▸ Critical sections are too short?
  - ▸ Bugs! You've written code where other threads will be able to see intermediate state
- ▸ Guideline #3: Make an effort to not do expensive computation or I/O in critical sections
  - ▸ And take even more care not to introduce race conditions ☺

# Understanding The "Critical-Section Granularity" Issue

- ▸ Tempting (and often a good idea) to simply wrap a chunk of code with a single lock
  - ▸ Example: code only contains operations like assignment statements and "cheap" method
    - ▸ Code runs so fast that it's not worth splitting into multiple critical sections
- ▸ Assumption fails when critical section will contain "expensive" computation
- ▸ Scenario: you need to replace a Hashtable value by performing some computation on its old value
  - ▸ We want other threads to see either the old value or the new value
  - ▸ We also need the old value in order to perform the expensive code required to compute the new value

# This Critical-Section Is <u>Too Large</u>

- ▸ Assume that `lock` guards the Hashtable
- ▸ The code below locks the entire table during the expensive computation

```
1   synchronized(lock) {
2     v1 = table.get(k);
3     v2 = expensive(v1);
4     table.put(k, v2);
5   }
```

# This Critical-Section Is Too Short

```
1        synchronized(lock) {
2          v1 = table.get(k);
3        }
4
5        v2 = expensive(v1);
6
7        synchronized(lock) {
8          table.put(k, v2);
9        }
```

- ▸ The above code is buggy!
- ▸ Scenario: *thread$_2$* updates to map $k \rightarrow v3$ while *thread$_1$* was executing expensive(v1)
- ▸ If *thread$_2$* does the update after *thread$_1$*, then table will end up with $k \rightarrow v3$ (not v2)
- ▸ If *thread$_2$* does the update before *thread$_1$*, then table will end up with $k \rightarrow$ expensive(v3) (not v2)
- ▸ Desired semantics: $k \rightarrow v2$

# This Critical-Section Is Just Right

```
1   boolean loop_done = false;
2   while (!loop_done) {
3     synchronized(lock) {
4       v1 = table.get(k);
5     }
6
7     v2 = expensive(v1);
8
9     synchronized(lock) {
10      // If true, we already have correctly
11      // computed ''v2''
12      //
13      // If false, we know that some other
14      // thread has updated the map: throw out
15      // pre-computed value and start over
16      if (table.get(k) == v1) {
17        loop_done = true;
18        table.put(k, v2);
19      }
20    } // while-loop
```

Note: this approach assumes that it doesn't matter if some arbitrary sequence of Hashtable operations occurred between the two critical section. All we care about is that we change mapping from $k \to v1$ to $k \to v2$

- ▸ Operations on your data-structures should be atomic
  - ▸ No other thread can see the "partial execution" of that operation
- ▸ Guideline #4: Think about atomicity first and locks second
  - ▸ "What operations need to be atomic?"
  - ▸ Make critical sections just long enough to preserve atomicity
  - ▸ Then design the locking protocol to implement the critical sections correctly

# Don't "Roll Your Own" (I)

- ▸ You need to understand today's material because you'll be writing new code
- ▸ But: make sure you're aware of the JDK 5 `jdk.util.concurrent` concurrent collection APIs
  - ▸ These are high-performance concurrent implementations of standard collection interfaces such as *List*, *Queue*, and *Map*
  - ▸ To provide high concurrency, these implementations manage their own synchronization internally
  - ▸ Implication: you can't exclude concurrent activity from a concurrent collection
  - ▸ Locking a concurrent collection only slows your code down, accomplishes nothing!
- ▸ Take the time to understand what these classes do and how they do it before coding your own data-structure
  - ▸ See if you can re-purpose the concurrent collection classes even if it doesn't immediately look like it addresses your problem

# Don't "Roll Your Own" (II)

```java
// Code by Bloch, Item #81: "Simulate the
// behavior of String.intern:"
private static final ConcurrentMap<String, String> map =
    new ConcurrentHashMap<>();

public static String intern(String s) {
  // ConcurrentHashMap is optimized for
  // retrieval operations, such as
  // get. Therefore, it is worth invoking
  // get initially and calling putIfAbsent
  // only if get indicates that it is
  // necessary:
  String result = map.get(s);
  if (result == null) {
    result = map.putIfAbsent(s, s);
    if (result == null)
      result = s;
  }
  return result;
}
```

- ▸ (FWIW, my own performance tests confirm Bloch's statements below)
- ▸ "On my machine, the intern method above is over six times faster than String.intern"
- ▸ "Concurrent collections make synchronized collections largely obsolete"
- ▸ "For example, use ConcurrentHashMap in preference to Collections.synchronizedMap"
- ▸ "Simply replacing synchronized maps with concurrent maps can dramatically increase the performance of concurrent applications"

## Blocking Operations

- ▸ In JDK 5, some of the collection interfaces were extended with blocking operations
  - ▸ That is: if a method can't be performed at the time that a thread invoked it …
  - ▸ The thread will (automatically) block until the internal state of the collection permits the method to continue
- ▸ Example: `BlockingQueue` extends `Queue` and adds several methods, including `take`
  - ▸ `take` removes and returns the head element from the queue, waiting if the queue is empty
- ▸ Blocking queues can be used to implement work queues (also known as "producer-consumer" queues)
  - ▸ One or more producer threads enqueue work items
  - ▸ One or more consumer threads dequeue and process items as they become available
  - ▸ Threads block as long as queue state isn't appropriate for desired semantics

# Synchronizers

- ▸ Synchronizers are objects that enable threads to wait for one another
  - ▸ Use synchronizers when you need to coordinate decoupled activities
- ▸ Typically, the *CountDownLatch* and *Semaphore* classes suffice
  - ▸ Sometimes you need the more powerful (and complicated) *CyclicBarrier* and *Exchanger* classes
  - ▸ Or even the *Phaser*

# CountDownLatch (I)

- ▸ Countdown latches are "single-use" barriers
  - ▸ They allow one or more threads to wait for one or more other threads to do something else
  - ▸ CountDownLatch constructor takes an `int`: specifies the number of times the `countDown` method must be invoked on the latch before all waiting threads are allowed to proceed
- ▸ Scenario (see next slide for code)
  - ▸ You want to "time" code execution for a given "concurrency level"
  - ▸ That implies that you only start the clock when all worker threads are "ready to go"
  - ▸ Implementation: one latch is initialized to the given concurrency level, each worker decrements the count when ready
  - ▸ That latch prevents timer from starting until all are ready
  - ▸ Another latch ensures that the timer measurement is only taken when all workers have completed

# CountDownLatch (II)

```java
// Code from Goetz and Bloch
public class ConcurrentTimer {
  private ConcurrentTimer() { } // Noninstantiable

  public static long time(Executor executor, int concurrency,
                          Runnable action) throws InterruptedException {
    CountDownLatch ready = new CountDownLatch(concurrency);
    CountDownLatch start = new CountDownLatch(1);
    CountDownLatch done  = new CountDownLatch(concurrency);

    for (int i = 0; i < concurrency; i++) {
      executor.execute(() -> {
          ready.countDown(); // Tell timer we're ready
          try {
            start.await(); // Wait till peers are ready
            action.run();
          } catch (InterruptedException e) {
            Thread.currentThread().interrupt();
          } finally {
            done.countDown();  // Tell timer we're done
          }
      });
    }

    ready.await();     // Wait for all workers to be ready
    long startNanos = System.nanoTime();
    start.countDown(); // And they're off!
    done.await();      // Wait for all workers to finish
    return System.nanoTime() - startNanos;
  }
}
```

- Old-style discussions of how to create multiple, independent, tasks focus on *"should I extend `Thread` or should I implement `Runnable`?"*
- That's really "old-style": almost always, refrain from working directly with threads!
- When you work directly with threads, a Thread serves as <u>both</u> a "unit of work" (or "task") <u>and</u> the mechanism for executing the unit of work
- JDK 5 provides the *Executor* framework, which separates the task concept from the execution mechanism
- Two kinds of tasks: *Runnable* and its close cousin, *Callable*
  - A *Callable* is like *Runnable*, except that it returns a value and can throw <u>arbitrary</u> exceptions

# Executor Service

- The general mechanism for executing tasks is the JDK 5 (interface-based) executor service
- Think in terms of <u>tasks</u> (not Threads!), let an executor service execute those tasks for you
  - You'll gain the flexibility of scheduling different execution policies, concurrency levels, and other tunable knobs
  - Read the Javadoc!
- Some examples:
  - You can wait for any or all of a collection of tasks to complete (see `invokeAny` and `invokeAll`)
  - You can wait for the executor service to terminate (see `awaitTermination`)
  - You can retrieve the results of tasks one by one as they complete (using an *ExecutorCompletionService*)
  - You can schedule tasks to run at a particular time or to run periodically (using a *ScheduledThreadPoolExecutor*)

# Executor Example

```
1   // Code from Java Concurrency In Practice (Goetz)
2   public class PrimeGenerator implements Runnable {
3     private static ExecutorService exec = Executors.newCachedThreadPool();
4     private final List<BigInteger> primes = new ArrayList<BigInteger>();
5     private volatile boolean cancelled;
6
7     public void run() {
8       BigInteger p = BigInteger.ONE;
9       while (!cancelled) {
10        p = p.nextProbablePrime();
11        synchronized (this) {
12          primes.add(p);
13        }
14      }
15    }
16    public void cancel() {
17      cancelled = true;
18    }
19    public synchronized List<BigInteger> get() {
20      return new ArrayList<BigInteger>(primes);
21    }
22    static List<BigInteger> aSecondOfPrimes() throws InterruptedException {
23      PrimeGenerator generator = new PrimeGenerator();
24      exec.execute(generator);
25      try {
26        SECONDS.sleep(1);
27      } finally {
28        generator.cancel();
29      }
30      return generator.get();
31    }
32  }
```

# Some Closing Thoughts

# "Concurrent" Programming

- <u>Concurrency</u>: the problem of correctly and efficiently managing concurrent access to shared resources from multiple clients
- Requires coordination between threads
  - Typically done by requiring other threads to block until "first thread is no longer accessing the shared resource"
  - Compare to the "fork/join" framework from *Introduction to Algorithms*
    - Concurrency coordination is very different from "fork/join" coordination
    - *join* semantics: *thread₁* blocks until *thread₂* has completely finished its execution
    - Concurrency coordination: *thread₁* blocks until *thread₂* has finished some "arbitrary" section of code execution
- Major challenge: correct concurrent programs are non-deterministic
  - Non-repeatability → much more complicated to test & debug than single-threaded environments ☺

# Parallelism *Versus* Concurrency (I)

- ‣ We've previously discussed (in *Introduction to Algorithms*) using multiple threads to achieve parallelism
  - ‣ Defined as "how to break up a computation into parts that can be executed in parallel"
  - ‣ Remember the Fork/Join framework?
- ‣ These parallel algorithms all had a very simple structure that avoided the race conditions discussed in today's lecture
- ‣ Each thread is assigned its "own section of memory"
  - ‣ Example: algorithm partitions array into sub-arrays, one sub-array per thread
- ‣ Yes: all threads have access to the same shared memory ("entire array")
  - ‣ But algorithm ensures that threads agree as to what's "mine" and what's "theirs"
- ‣ On `fork`, algorithm completely assigns ownership of sub-array to "forkee"
- ‣ Main thread only resumes ownership when it does a `join` on the "forkee" thread

- ▸ Unfortunately for us, the approach we used for fork/join parallelism won't work when
  - ▸ Memory is accessed by threads in overlapping or unpredictable ways
  - ▸ Threads are doing un-coordinated tasks but still need concurrent access to the same resources
  - ▸ Such scenarios raise the mutual-exclusion issues discussed in today's lecture
- ▸ Thought experiment: can you "reduce" implementation of a data-structure (let alone a DBMS) to fork/join parallelism?
- ▸ In a nutshell:
  - ▸ Parallelism is about making algorithms run faster
  - ▸ Concurrency will not make "an algorithm" run faster

# Concurrency: Not Just About Throughput

- We've previously explained that concurrency is about achieving good throughput
  - Alternatively: enabling as many client threads to use a service (or data-structure) at the same time (while preserving "correctness"
- But note that concurrency also improves <u>another</u> performance metric: latency
- No time to pursue this idea, but here are some examples
  - Respond to GUI events in one thread while another thread performs an expensive computation
  - Web-application services another request while another thread does I/O for previous request
  - Failure in $task_1$ does not bring down $task_2$

Why A Lecture On "Java Concurrency"?

Transitioning To Multi-Threaded Code

Race Conditions: The Problem(s)

Synchronization: The Solution

Concurrency Programming Guidelines

Creating & Managing Multiple Clients

Some Closing Thoughts

- I've already strongly recommended Java Concurrency in Practice by *Brian Goetz et al* Amazon
- This lecture overlaps considerably with Chapters 2 & 3
- Am not assigning "homework", just a course project ☺