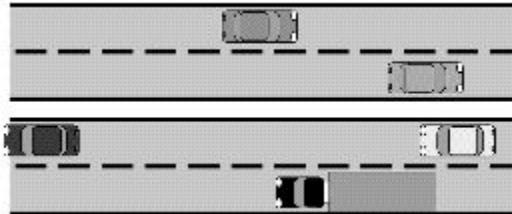


Concurrency in Java

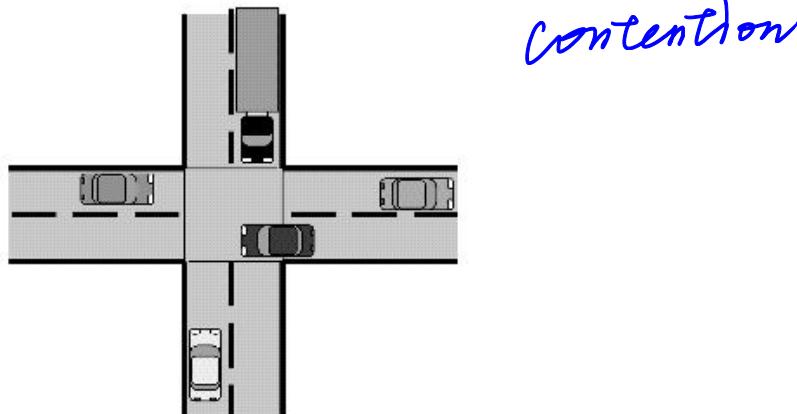
多线程开发

Concurrency in General

- Multiple events happening at the same time.



- The trouble (in real life and in software development) is the interaction between these events.



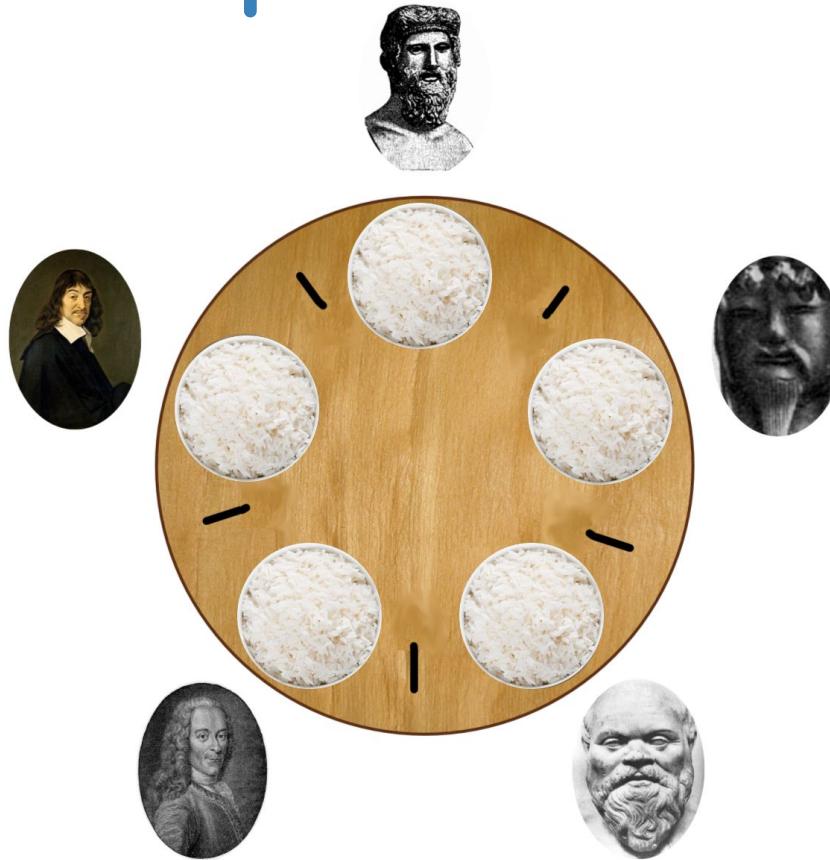
Lecture(s) Outline

- Review common concurrency examples
- Define concurrent principles / terms
- Rephrase common concurrency examples with the defined terms
- Introduction to concurrency specific to the Java programming language.
- Discuss Java programming language threading mechanisms

Dining Philosophers Problem

- Five silent philosophers sit around a table. In front of each is a bowl of rice. To each side of the philosopher is a single chopstick. Each philosopher must alternatively eat and think. A philosopher can pick up a chopstick on either side (right or left) but can only begin eating when holding two chopsticks. After eating the philosopher places both chopsticks back down. Assume the bowl of rice is bottomless (infinite rice within it).
- Problem - devise a scheme such that no philosopher starves.
 - Easy way to starve everyone is to have them all start by picking up the chopstick to their left.

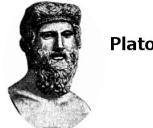
Dining Philosophers Problem (cont)



Dining Philosophers Problem (cont)

- Dijkstra's Solution (he also devised the problem!)
 - Number the chopsticks. The philosopher must always first pick up the lower of the two chopsticks around him, then the higher.

都拿身邊的小圓片



Plato



Descartes



Confucius



Voltaire



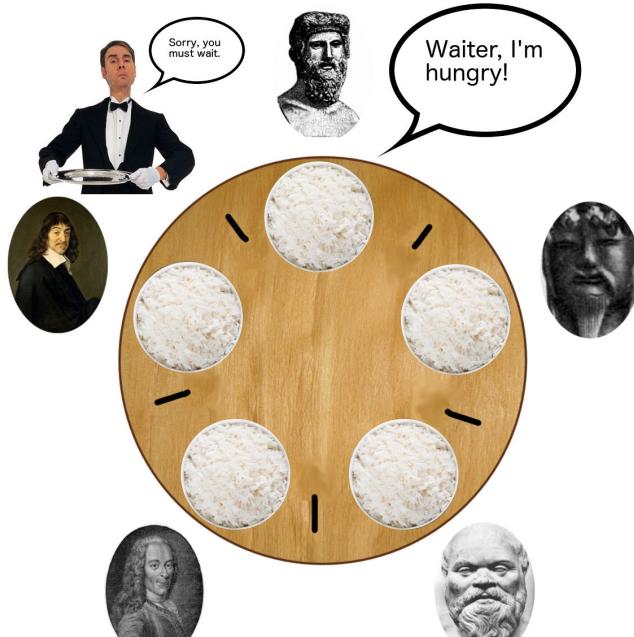
Socrates

不是取低
低者先取
為時序順
2人以上

Dining Philosophers Problem (cont)

- Arbitrator “Waiter” Solution
 - Philosophers must ask the waiter for permission to pick up a chopstick and the waiter only gives permission to one philosopher at a time.

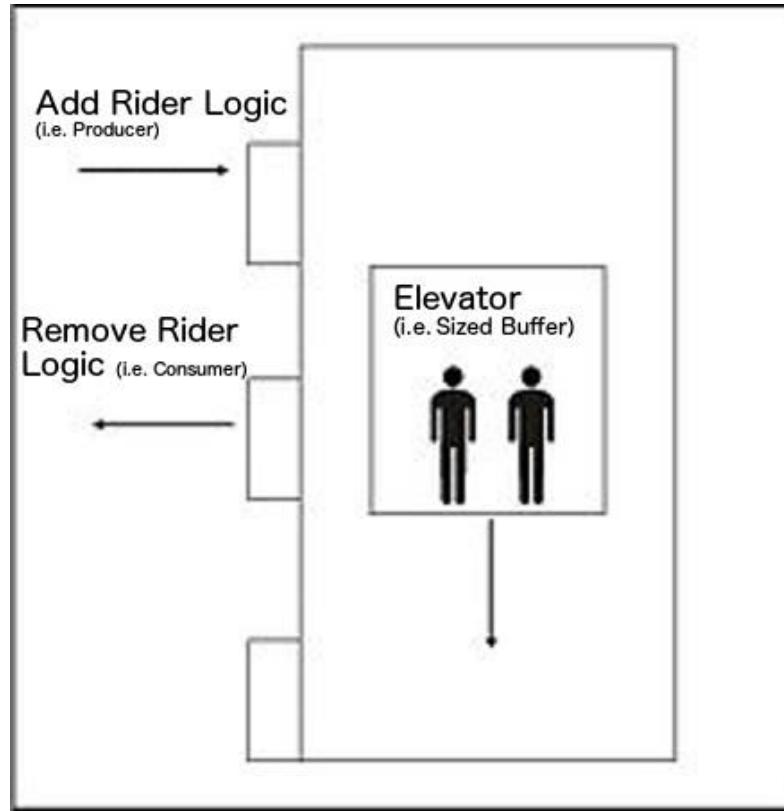
Waiter
for 5 philosophers
5 resources



Producer / Consumer Problem

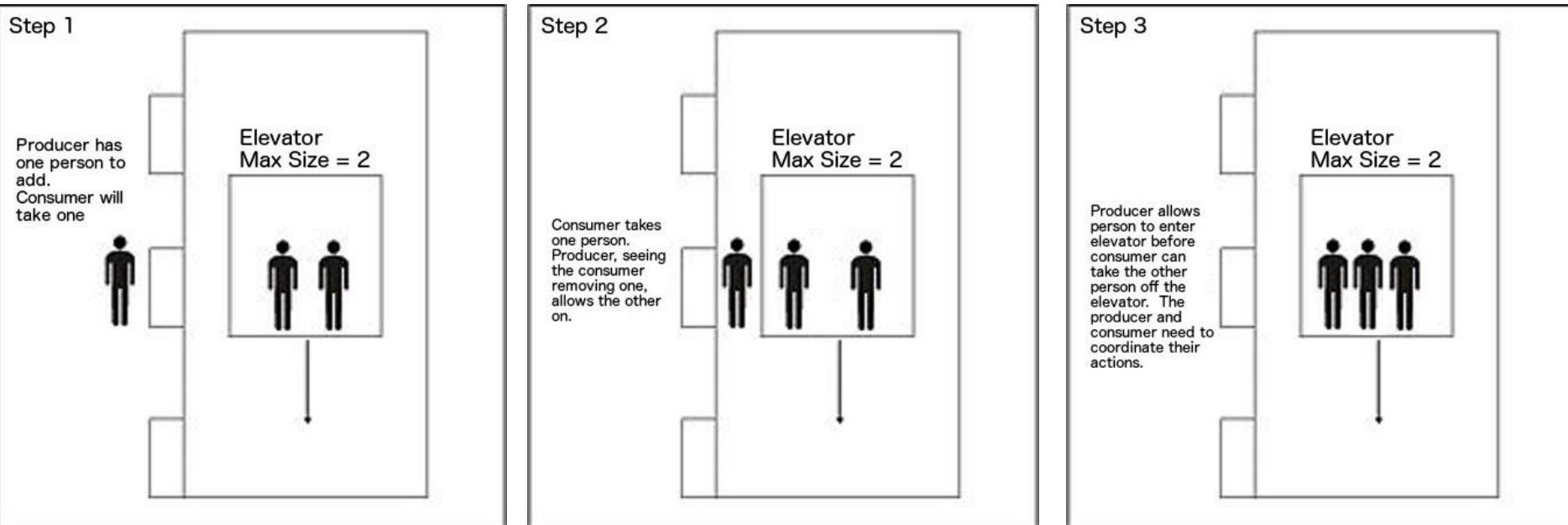
- There's a producer and a consumer sharing a fixed sized buffer. The job of the producer is to create data to put into the buffer. The job of the consumer is to remove data from the buffer.
 - The problem is to ensure the producer does not put data into a full buffer and the consumer does not try to remove data from an empty buffer.
- A real life example of this might be an elevator. The elevator itself is the buffer as it's fixed size (there is a maximum amount of people, or weight, allowed before it breaks). The producer is the logic responsible for allowing people onto the elevator. The consumer is the logic responsible for allowing people off.

Producer / Consumer Problem (cont)



Producer / Consumer Problem (cont)

- Problem can arise for the elevator if the same floor has a producer and a consumer and they do not coordinate their movements.



Concurrency Terminology

单核协调执行 (线程交替)

- Race condition
- Deadlock
- Resource Starvation
 - Livelock
- Atomicity
- Mutual Exclusion
- Locks / Semaphores 信号量
 - Binary Semaphore
 - Mutexes 互斥锁
- Monitors

区
域

parallelism 并行

多核同时执行

Race Condition

- When output is dependent upon the uncontrollable sequence of events.
 - Becomes a bug if the sequence of events happens in a way the programmer did not intend.
 - Means that a program can behave “correctly” even though it suffers from a race condition.
- The producer-consumer elevator problem suffers from a race condition. The output is dependent upon the order of the persons entering and leaving the elevator. In some cases it works (no bug), in some cases it fails (bug) but in all cases it suffers from a race condition.
- Example!

线程顺序, 由哪个线程先执行决定了结果
in 什么次序, 哪一个线程先执行决定了结果
上一个线程
下一个线程
在哪个位置

Deadlock

- Two (or more) resources are waiting for each other to finish and thus neither ever do finish.
 - We saw this in the Dining Philosophers. If we make a scheme where all philosophers pick up left chopstick and then wait to pick up right chopstick before putting down the left.
- Example!

Resource Starvation

犯(罪)做(恶)
perpetrate

- When a process/thread is perpetually denied resources.
 - An example of this is livelock. This is different than deadlock in that the program isn't halted it however cannot move forward because it is continuously denied a resource.
 - In the Dining Philosophers problem change the deadlock scheme such that after all picking up the left chopstick the philosophers wait ten minutes, but down the left chopstick and then all pick up the right chopstick, wait ten minutes, etc.

Example!

Atomicity!

Atomicity prevent race condition

一个 thread -> 完成操作是整体. 需要完全流动到另一个 thread, 才会让从那不流动到另一个 thread

- An operation (or set of operations) appears to occur to components of a system instantaneously.
○ From **Java Concurrency in Practice** (2.2.3, page 22): “Operations A and B are atomic with respect to each other if, from the perspective of a thread executing A, when another thread executes B, either all of B has executed or none of it has.”
- The race condition in the producer-consumer-elevator problem and in the example were an indirect result of the assumption that a set of operations were atomic.
 - In the case of the elevator, that the count of the persons in the elevator and the act of removing the person was atomic. **整体完成**
 - In the case of the example, the retrieval and addition of the unique counter.
- Example!
一行语句不能保证 atomic

Atomicity (cont)

- Operations on long and double are not atomic (unless you use the volatile keyword- more on that later)
 - <https://docs.oracle.com/javase/specs/jls/se7/html/jls-17.html#jls-17.7>

使用 long
long double
修改内部
atomic v

long 及 float 在 32 位机上，它是被拆成高 32 位 和 低 32 位
进行 逐字节修改。故其 两条语句都不是 atomic，可能会造成 race condition。
当两条语句间另一个 thread 进入时，将导致 race condition.

实现资源使用互斥和实现原子性 from atomicity

Mutual Exclusion

互斥

某 block 同时只能有一个 thread

在使用这个 block 互斥时，interact with it

比方说：在电梯上

two threads

通过 "two"

in block 使用

"电梯" 资源，

by two people

in thread 不同

通过 "two"

in block 使用

"电梯" 资源。

- For concurrent programs, mutual exclusion ensures that no two concurrent processes/threads are in a critical section at the same time.
 - A critical section is any portion of code which accesses a shared resource which cannot be accessed concurrently.
- E.g., in the elevator example, there needs to be a mutual exclusion guarantee around the consumer decrementing the count of persons within the elevator and the act of removing a person from the elevator.
- E.g., in the 'UniqueNumberGenerator' example, there needs to be mutual exclusion around returning a unique value and incrementing to the next value.

先 acquire 资源

再执行 For 和 no mutual exclusion 语句。

无锁

Interface

Locks & Semaphores

⇒ Interrupted Exception 沒有鎖
就沒有死線程!

- Locks - a synchronization mechanism to limit access to a shared resource; a way of implementing mutual exclusion.
 - Semaphores - a particular type of Lock.
 - Records how many units of a resource are available in conjunction with providing a safe (prevents race conditions) means of adjusting the record and potentially waiting for resources to become available.
 - A semaphore with an arbitrary resource count is called a counting semaphore.
 - A semaphore with value either 0 or 1 is called a binary semaphore.
- 用於隨意 thread 使用 resource 用於計數器 資源的前綴用在 resource 上
- encapsulated semaphore
- new Semaphore c1
- too unique number generator, 這樣用 binary semaphore
- dead resource
- 要先 release, 才能再 release, 才能死 lock
- Example!

Homework 9

<https://github.com/NYU-CS9053/Spring-2016/homework/week9>

```
while(!Thread.isInterrupted()) {
    try {
        Thread.sleep(1000) // 进程运行待机期间，如果别的东西想 interrupt 它
    } catch(InterruptedException e) {
        if(Thread.interrupted()) // 在线程结束前，如果别人调用了 interrupt()
            throw new InterruptedException();
    }
}
```

如果线程没有实际 interrupt / 被标记
被置过，while 检查
会认为 thread 已
interrupt

try acquire() 只问当前线程用
于拒绝 false;
acquire() 问，若无可用则挂
拒绝等待 → throw InterruptedException;
并行的 →

(blocking) (checked)
thread 在挂起时用资源时，

不能正常响应中断 (interrupt) 操作。
故若在此段中打断 thread 会有 InterruptedException

Monitor

带有条件判断的 resource lock

→ 即使有 resource 可用，
也要满足需求，才得到 resource

- A synchronization mechanism that allows threads to have mutual exclusion and the ability to wait (block) for a condition to become true (e.g., there's room within the elevator). Monitors have a mechanism to signal waiting threads that a condition may have been met.
 - Monitors are implemented with a mutex lock and condition variables. → 被 encapsulated in
- Java loosely based its means of synchronization of objects based on this concept. However, fields are not all private (exposing state), methods are not forced to be synchronized (mutual exclusion) and the intrinsic lock (mutex) is exposed (it's the object itself).
- Example! ^{从房}

电梯中往里加入人进程序，
只有在电梯不满时，才把资源给他

Monitors in Java

- Since Java 1.5, there is an explicit Lock/Condition object to give semantics of a monitor.
- General idea is to have one Lock object for all threads (per logical grouping of code needing mutual exclusion).

```
1 lock.lock(); 相当于 acquire
2 try {
3     // do something needing mutual exclusion
4 } finally {
5     lock.unlock();
6 }
```

Monitors in Java (cont)

- If the code within the mutual exclusion block should only be executed based on the state of other values use a Condition.

```
1 Lock lock = new ReentrantLock(); // 1个semaphore 加锁某个 thread  
2 Condition listNotEmpty = lock.newCondition(); // 拥有资源，再遇到  
3  
4 // elsewhere in code  
5  
6 lock.lock(); // 资源  
7 try {  
8     while (list.isEmpty()) {  
9         listNotEmpty.await();  
10    } // remove, fullCondition.signalAll()  
11    // do something knowing now the list is not empty  
12 } finally {  
13     lock.unlock();  
14 }
```

↑ fair!

↑ signalAll() =
↓ signal()
↑ fair!

↑ To implement mutual exclusion is to

↑ it's not thread safe
↑ if listNotEmpty is not thread safe
↑ if release(listNotEmpty) is not resource

↑ release lock and block thread.
↑ 通过线程使其
↑ block not full (isEmpty)
↑ in signalAll() 信号
↑ to awake, 并且线程否
↑ while in list (while not full)

Threads in Java

- How?
 - By default there is one main Thread.
 - To date, all of your code written in Java has been run on this thread. 逐行执行
 - You can create additional Thread objects and start them.
 - Key elements are the Runnable object and the start and join methods.
 - Never call run or stop
 - Careful about interruptions -> always handle them, rethrow or terminate. If not handling them, call Thread.currentThread().interrupt()
- Example!

```
Runnable code = new Runnable() {  
    override public void run() {  
        // TODO  
    }  
}  
  
Thread thread = new Thread(code);  
thread.start();
```

join(): main() blocks until thread terminates

Visibility

- What happens when this program runs?
 - A) Prints 100000
 - B) Prints 0
 - C) Prints nothing
 - D) Any of the above

↑ no synchronization

concurrency ↑ Deadlock 原因

↑ no thread
↑ share memory
↑ copy value
↑ ready
↑ number
↑ race condition of update
↑ map

```
1 public class Invisible {  
2     private static boolean ready; ready false  
3     private static int number; main thread  
4  
5     public static void main(String[] args) {  
6         Thread thread = new Thread(new Runnable() {  
7             @Override public void run() {  
8                 while (!ready) {  
9                     Thread.yield(); thread yielding  
10                } is main waiting  
11                System.out.printf("%d%n", number);  
12            }  
13        });  
14        thread.start();  
15        number = 100000; number 100000  
16        ready = true;  
17    }  
18}  
19 }  
20 }  
21 }
```

Modified example from Goetz's Java Concurrency in Practice (Listing 3.1)

Visibility (cont)

- Writes to a shared variable on one thread can be seen by other threads.
 - With multiple threads, writing to a shared variable does not guarantee that other threads will immediately see the update (or ever see the update).
 - To ensure visibility synchronization techniques must be used.

The volatile Keyword

只用于对reference的 primitive value

double (atomic in)

不推荐, 反复使用

ex. lock 简单, lock 极端

by thread in
线程共享, 但

不是线程
可见性
保证

- Provides a weaker form of synchronization than using locking mechanisms.
- Marking a variable with volatile alerts the JVM that the variable will be shared amongst threads. Its value will not be cached (e.g. in registers); thus a read of a volatile variable will always be the most recent value written by a thread.
- CAREFUL! This does not make your variable/program thread-safe.
 - Good use of volatile is simply to ensure the value of a variable is seen by all threads.
 - E.g., the semantics of volatile are not strong enough to make counter++ atomic (and thus thread safe).
 - Use the Atomic variants of Boolean, Integer, Long.

只读
值
increment

Example Using volatile

```
1  public class Visibility implements Runnable {  
2  
3      private volatile boolean sleeping;  
4  
5      @Override public void run() {  
6          while (sleeping) {  
7              sleepLonger();  
8          }  
9          process();  
10     }  
11  
12     // other methods...  
13  
14 }
```

Publication & Escape

allow access to variable of other thread

- Is this program thread-safe?

Not

```
1 public interface Escape {  
2  
3     void process(Date date);  
4  
5 }
```

```
1 public class Publication {  
2  
3     public static void main(String[] args) {  
4         final Publication publication = new Publication(new Date());  
5         final Escape escape = EscapeFactory.create();  
6         Thread thread = new Thread(new Runnable() {  
7             @Override public void run() {  
8                 escape.process(publication.getDate());  
9             }  
10        });  
11        thread.start();  
12        escape.process(publication.getDate());  
13    }  
14    reference is immutable mutable (object本身)  
15    private final Date date;  
16    Date 不能改变 导致不同时  
17    public Publication(Date date) {  
18        this.date = date;  
19    }  
20    Publication: thread 可以  
21    access to date 同时  
22    not thread safe 不线程安全  
23    public Date getDate() {  
24        return date;  
25    }  
26}
```

Publication & Escape (cont)

- What does this program do incorrectly?

```
1 public interface Registry {  
2  
3     static interface Listener {  
4         void process();  
5     }  
6  
7     void register(Listener listener);  
8 }  
9 }
```

```
1 public class Worker {  
2  
3     public Worker(Registry registry) {  
4         registry.register(new Registry.Listener() {  
5             @Override public void process() {  
6                 handleCallback();  
7             }  
8         });  
9     }  
10    public void handleCallback() {  
11        // TODO - do something  
12    }  
13 }  
14 }
```

Don't do this
inside a constructor!

To object *@Override* *to worker object* *anonymous inner class*
constructor *handleCallback()* *public* *方法*

Immutability

不改就沒有

- “Immutable objects are always thread-safe” - Goetz’s Java Concurrency in Practice, 3.4
- Immutable definition (see Goetz 3.4); object is immutable if
 - All its fields are final
 - Its state cannot be modified after construction *不是 reference object*
 - It is properly constructed (does not escape during construction) *不會*
- For this (and all the other reasons prior to concurrency we mentioned in lectures) prefer immutable objects. Only use mutable objects if there’s a compelling reason not to *修改
做不到*

Thread Confinement

```
private final ThreadLocal<SimpleDateFormat>  
    formatter = new ThreadLocal<SimpleDateFormat>(){  
        @Override protected SimpleDateFormat initialValue(){  
            return new SimpleDateFormat("MM/dd/yyyy");  
        }  
    };
```

- Confining a resource/value to a thread means not sharing it across threads. Not sharing guarantees thread safety.
 - Stack confinement - seen in RaceConditionWithLock. Local variables are always thread safe as they're confined to the stack and consequently to the thread running the stack.
 - ThreadLocal - a more formal means of allowing the programmer to maintain thread confinement by associating a per-thread object with a value holder object.
- ThreadLocal example!

实际为一种线程变量类型，每一个使用该变量的其他进程，都直接访问该变量（本身共享）
而是自己即时 copy 一份

Synchronization Mechanisms

- Revisit Lock/Condition
 - Lock in Java - ReentrantLock
 - Lock with Condition
 - Lock without Condition
 - edu.nyu.cs9053.concurrency.chapter14.WithCondition
 - edu.nyu.cs9053.concurrency.chapter14.WithoutCondition
 - **IMPORTANT - Must own lock to await on condition and must own lock to signal a condition**
 - edu.nyu.cs9053.concurrency.chapter14.UnownedLockConditionInvoker

Java's Built-in Synchronization

- Java loosely based its means of synchronization of objects based on this concept. However, fields are not all private (exposing state), methods are not forced to be synchronized (mutual exclusion) and the intrinsic lock (mutex) is exposed (it's the object itself).
- Every Object has an intrinsic lock (and one condition). *不是只有一個 condition*
 - Not the same as Lock/Condition (these came in Java 1.5)
 - Every synchronized can be rewritten as Lock/Condition but not every Lock/Condition can be rewritten using synchronized (which is why Java 1.5 added Lock/Condition)

- Example!

```
public void print() throws InterruptedException {
    synchronized (this) {
        while (value < 100) {
            wait();
        }
        System.out.printf("%d\n", value);
    }
}
```

這個資源

→ synchronized (this) {
→ while (value < 100) {
→ wait();
→ }
→ System.out.printf("%d\n", value);
→ }

```
public void process() {
    synchronized (this) {
        value += 1;
        notifyAll();
    }
}
```

Java synchronized (cont)

- Can you synchronize as a block or around any method.
 - If around method, using the intrinsic lock of the Object
 - So can you synchronize a static method?
 - What's the Object for static methods? *整个类!*
 - If used within a block the programmer must specify which Object to use as the lock.
 - Can use this for the current Object
 - Can use Class.class for the Object's class (to lock all Objects of a given class).
 - Example!

```
private static void foo() {  
    synchronized (synchronized.class) {  
        }  
        }  
    }
```

Thread-safe Collections/Utilities

- Look at classes within package `java.util.concurrent`
- AtomicInteger, AtomicBoolean, AtomicLong, AtomicReference
 - Provides atomic interactions (get and set, etc) *一步完成操作*
- ConcurrentMap (interface) and ConcurrentHashMap (implementation)
 - Leverages “striped locking” - see <http://www.ibm.com/developerworks/library/j-jtp07233/>
- BlockingQueue - allows you to easily implement Producer/Consumer problems *让 thread 挑队，允许消费者向队列插入得到资源*
- CountDownLatch
 - Allows work to be done until x values reach the gate *线程若干个线程继续 main, 等待 join*
- CyclicBarrier
 - Similar to a CountDownLatch but can be reused *t.join() 代理 t thread 语句
才能继续 main*

Executors

- Java 1.5 provided a framework for interacting with Threads which abstracts logic of starting/stopping and executing work.
 - Executor - interface of something which can process work
 - ScheduledExecutor - extends Executor but allows the programmer to control over when and how frequent work is done
- ExecutorService - similar to a "thread pool" owns a number of Executor objects. Get an instance of this and give it work.
- ScheduledExecutorService - same as an ExecutorService but also can schedule jobs (ScheduledExecutor).
- Executors - utility class which creates ExecutorService and ScheduledExecutorService objects.
- Example!

Homework 10

<https://github.com/NYU-CS9053/Spring-2016/homework/week10>

```
if (random.nextBoolean()) {    定時BF  
    scheduledExecutorService.schedule(new Runnable() {  
        @Override public void run() {  
            // TODO  
            }  
        }, random.nextInt(1000), TimeUnit.MILLISECONDS);  
    } else {  
    executorService.submit(printer);  
    }    行動項加進程
```