

COMP2221 Networks

David Head

University of Leeds

Lecture 10

Reminder of the last lectures

Two lectures ago we saw how to implement a simple Java server using `ServerSocket`:

- Called the **blocking** `accept()` method to wait for clients.
- Could only handle one client at a time.
- Examples: `DailyAdviceServer` and `KnockKnockServer`.

Last lecture we looked at **parallel computation** in general:

- Historical development of **parallel hardware**.
- The difference between:
 - **Shared** and **Distributed** memory.
 - **Parallelism** and **Concurrency**.
 - **Processes** and **Threads**.

Today's lecture

In today's lecture we will see how to implement a **multi-threaded server** in Java.

- Use Java's built-in Thread class.
- **Shared memory**, *i.e.* utilises multiple cores for one machine (or one node in a cluster), but **only** one machine/node.
- One thread per client.

In the next lecture we will see two alternatives to this **thread-per-client** model.

First, we will look at how to program in parallel in Java.

The Java Thread class

We have two basic options for threading in Java:

- ➊ Subclass the Thread class (in `java.lang`).
- ➋ Use the Runnable interface.

In both cases, you will put your code into:

```
public void run()
```

Request a new thread to execute the code in `run()` by calling:

```
public void start()
```

Other sometimes useful methods:

- `join()`, waits for the thread to finish.
- `sleep(long millisecs)`, causes the thread to sleep.

Scheduler

Threads do **not** start running the instant we create them.

- **When** a thread runs is managed by a component of the operating system called the **scheduler**.
- Allocates core time **dynamically**, depending on other applications, background tasks *etc.*

Since you cannot control when other applications or background tasks are active, a parallel program may generate **different results** each time it is executed.

- This **non-determinism** is usually undesirable.
- Can be corrected with **synchronisation** (*later this lecture*).

Method 1: Inheritance

Code on Minerva: HelloT.java

Subclass `java.lang.Thread`, overriding the `run()` method:

```
1 public class HelloT extends Thread {  
2  
3     public void run() {  
4         System.out.println( "Hello from a thread." );  
5     }  
6  
7     public static void main( String args[] ) {  
8         HelloT t = new HelloT();  
9         t.start();    // Schedule the execution.  
10    }  
11 }
```

When `t.run()` is actually called depends on the scheduler.

Running multiple threads

```
1 public class HelloT extends Thread {  
2     public void run() { ... }  
3  
4     public static void main( String args[] ) {  
5         HelloT t1 = new HelloT();  
6         t1.start(); // Start thread 1.  
7  
8         HelloT t2 = new HelloT();  
9         t2.start(); // Start thread 2.  
10  
11         // May have thread 1 and thread 2 running now,  
12         // in addition to the main thread.  
13     }  
14 }
```

Note that, as with the previous example, although calling `start()` starts the thread, the **scheduler** decides when `run()` is called.

Method 2: Interface

Code on Minerva: HelloI.java

Implement the `java.lang.Runnable` interface:

```
1 public class HelloI implements Runnable {  
2  
3     public void run() {  
4         System.out.println( "Hello from a thread." );  
5     }  
6  
7     public static void main( String args[] ) {  
8         HelloI h = new HelloI();  
9         Thread t = new Thread(h); // Initialise with h.  
10        t.start();  
11    }  
12 }
```

This time we **use** a thread, rather than **being** a thread.

Inheritance *versus* Interface

The **interface** approach is more general:

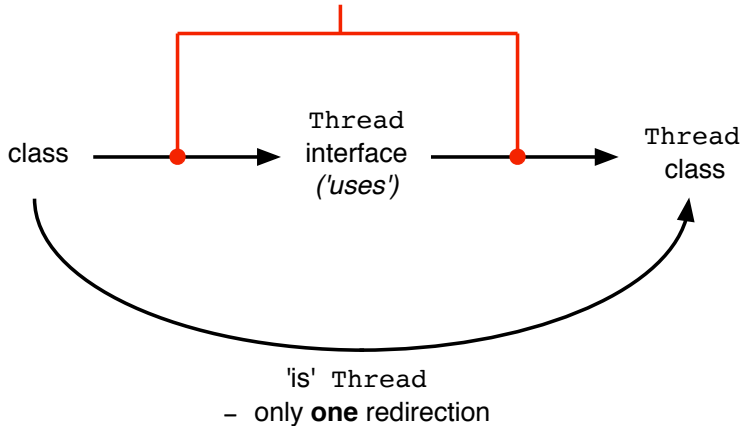
- Can inherit from a separate class.
- Can implement multiple interfaces, but cannot have multiple inheritance¹.
- Threading is a **property** of a class, not essential to its definition.

Inheritance is simpler, but our class **must** be a descendent of Thread.

- May be a *very slight* performance improvement relative to the interface approach (less 'indirection').

¹Some languages, e.g. C++, *do* allow multiple inheritance.

if using an interface, would need **two** redirections



General issues in thread programming

Two important multi-threading issues to be aware of are **data races** and **synchronisation**.

A **data race** is when two threads read and write the **same memory location**:

Thread 1

```
globalData = 1;  
...  
a = globalData;  
...  
Is a 1 or 2?
```

Thread 2

```
...  
globalData = 2;  
...
```

- The result of the **read** depends on whether or not the other thread has changed the data.
- Since we cannot control when threads run, the result may vary each time - **non-determinism**.

Synchronisation (1)

If one thread cannot continue until another one finishes, must **synchronise**. There are two common approaches:

- 1 Calling `t.join()` will pause the calling thread until its subthread `t` has finished — a **blocking** call (*below*).
- 2 Defining a **synchronised block** or **method** (*next slide*).

Consider the following code called by the main thread:

```
t = new Thread();  
t.start();  
... // Main thread and thread t run concurrently.  
t.join();  
... // Only main thread; t.run() has completed.
```

Synchronisation (2)

Can specify **blocks** of code that only one thread can enter at a time¹.

In Java, this is achieved using `synchronised(object)`:

```
synchronized( System.out ) {  
    System.out.print("Socket host:  " + so.getInetAddress());  
    System.out.print(" on port:  " + so.getPort());  
    System.out.println();  
}
```

Without `synchronised()`, the `print()` statements for each thread could become **interleaved**.

¹There is a single **lock** for the block with a thread must **acquire** before entering the block. It relinquishes the lock when leaving the block.

Synchronisation (3)

Java also allows you to synchronise **methods**. For example:

```
public synchronized void printStatus() {  
    System.out.print("Socket host:  " + so.getInetAddress());  
    System.out.print(" on port:    " + so.getPort());  
    System.out.println();  
}
```

This works the same way as the previous example, *i.e.* only one thread can enter the method at a time.

Note that use of `synchronized` can affect performance, sometimes severely.

Synchronisation (4)

These last two examples could be resolved by combining each `print()` statement into a single `println()`.

- `PrintStream`, `PrintWriter` internally synchronise.

`OutputStream` objects do **not** synchronise in general, so e.g. outputting to the same log file may require synchronisation.

For multi-threaded servers, must consider if multiple threads can write to the same global data, file *etc.*

- Data races can arise if **at least** one thread **writes**.
- Synchronisation may be required for e.g. outputting to a log file, accessing and updating a database, *etc.*

Client-server architecture

The software architecture we choose should be driven by the **application requirements**.

Relevant considerations are:

- How **many** clients to we expect (concurrently)?
- How **long** are clients connected for?
- What type of **protocol** are we implementing? For instance, is communication continual, or are there idle periods?

Examples:

- An ftp server might expect many, short connections.
- A chat client might expect fewer, longer connections.

Thread-per-client architecture

Code on Minerva: `KKMultiServer.java`, `KKClientHandler.java`,
`KnockKnockProtocol.java`

Uses a **single thread per connected client**.

- Creates the thread as the connection is made.
- Destroys the thread once the connection closes.

Assumes resources are available for many threads:

- Each thread requires some **local data**.
- Creating and destroying threads also costs **CPU cycles**.

At large enough **scale** (*i.e.* large number of threads), these costs outweigh the benefits

- Adding more threads makes the system less responsive.
- Typically limited to 100's to 1000's of unique threads on a typical desktop.

Multi-threaded KnockKnockServer

KnockKnockServer is replaced with KKMultiserver:

- Still a single thread (the 'main' thread).
- For each `accept()`, creates a new `KKClientHandler` object.

`KKClientHandler` is derived from `Thread`.

- Passed the client `Socket` in its constructor.
- Creates one instance of the protocol **per client**.
 - *i.e.* `KnockKnockProtocol`, unchanged from the earlier version.

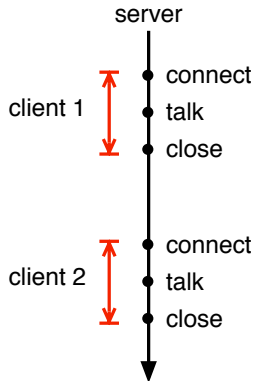
The client `KnockKnockClient` is unchanged.

- Will not consider multi-threaded clients, although they are sometimes useful, e.g. loading all images from a web page concurrently.

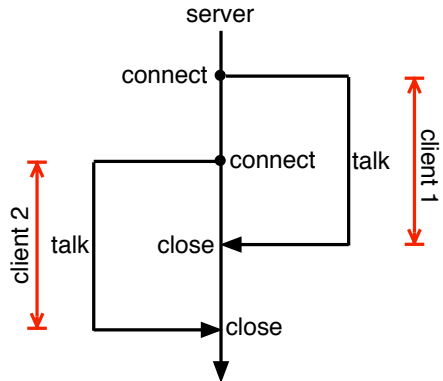
KKMultiServer

```
1 public class KKMultiServer {
2     public static void main(String[] args) throws IOException
3     {
4         boolean listening = true; // Always true here.
5
6         try {
7             serverSocket = new ServerSocket(2323);
8         } catch (IOException e) {
9             System.err.println("Could not listen on port: 2323.");
10            System.exit(-1);
11        }
12
13        // Each accept() creates a new thread.
14        while( listening )
15            new KKClientHandler(serverSocket.accept()).start();
16
17        serverSocket.close();
18    }
19 }
```

Single thread:

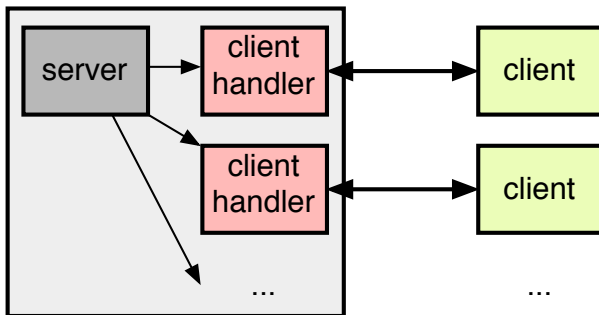


Multi-thread:



KKClientHandler

```
1 public class KKClientHandler extends Thread {
2     public void run() {
3         try {
4             PrintWriter out = new PrintWriter(
5                 socket.getOutputStream(), true);
6             BufferedReader in = new BufferedReader(
7                 new InputStreamReader(
8                     socket.getInputStream()));
9             String inputLine, outputLine;
10            KnockKnockProtocol kkp = new KnockKnockProtocol();
11            outputLine = kkp.processInput(null);
12            out.println(outputLine);
13
14            while ((inputLine = in.readLine()) != null) {
15                // As per the basic version
16            }
17            out.close();
18            in.close();
19            socket.close();
20        } catch (IOException e) { ... }
21    }
```



Thread-per-client pros and cons

Pros:

- Simple to understand and implement.
- Improvement over the non-threaded version we saw earlier.

Cons:

- Resources can potentially grow **without limit**.
 - Could try to estimate a maximum number of clients.
- A thread is created and destroyed for **each client**.
 - Although a much smaller overhead than for spawning processes¹, it can be significant.
- Better to create threads **once** and re-use for multiple clients.

¹Multiple processes **were** dynamically created/destroyed in early servers; see Harold, *Java Network Programming*, 4th ed., chapter 3.

Today we have looked at **multi-threading** in Java, and our first implementation of a **multi-threaded server**:

- The **thread-per-client** model where each client has a **client handler** on its own thread.
- Risks unlimited growth of resources.

Next time we will see two more models using **thread pools** that remove the risk of unlimited growth of resources.

Code for the thread-per-client model is on Minerva.