

Topic 12

Concurrency

CH 13 Scott

skip: 13.1.2, 13.2.3 (except fork/join and Box 13.3) 13.3 except 13.3.5, 13.4.1, 13.4.2, 13.4.4, 13.4.5

Why?

Problem has a logically parallel structure
characters in a video game

Exploit machine capabilities
SPEED!

Physical distribution
the internet, some things arrive very slowly

Parallel and the web

see firefox developer tools / network (rightmost column)
imagine if all of the images were retrieved serially

problem: you started a bunch of threads, how do you know when they are done???

And how do you even start the threads in the first place???

Rust:

see `wait_rust/h[1,2,3].rs`

Go:

`wait_go/h3.go`

Java

`wait_java/H3.java`

Concurrent == 2 or more task may be active at any time (not that active does not imply actually doing anything. So you can have concurrent processing of threads on a single processor machine.

Parallel == more than one task can be running. So all parallel are concurrent. To get parallel need more than one processor.

Parallel Level

Instruction level — in the core recognizing that instructions take different route through memory (no need to anything at the application level.)

Vector — operations repeatedly on every element of a dataset (GPUs) The classic example of this form is weather forecasting ... and various graphics operations. GPUs are essentially vector parallel devices.

Thread — do different things on different data through (possibly) different processors

Thread == every concurrent activity. So how many threads are running on computer??

== “active entity that the programmer thinks of as running concurrently with other threads”

UNIX: “ps -e | wc” is a good start. But that undercounts as these are processes. Many processes are multi-threaded. So need to get information on a thread basis

```

-e == every
ps -eL | wc
-L == show threads

```

heavyweight process == own address space (a unix process?)

lightweight process == groups share address space

'lightweight processes were added to Unix to accommodate shared-memory multiprocessors"
 "One of the major characteristics of shared memory multiprocessors is that all processors have equally direct access to one large memory address space."

"In systems with no shared memory, each CPU must have its own copy of the operating system, and processes can only communicate through message passing."

===== Stopped Here on 12/5 =====

Communication and Synchronization

comm == "any method that allows one thread to obtain information from another thread"

shared memory — some or all memory is accessible to threads
 message passing — (pure) no shared memory "send" and receive

Java — shared and can do message passing (<https://www.coderscampus.com/java-multithreading-java-util-concurrent/>) but it really often ends up looking like shared memory

Go — shared and message passing

Rust — both BUT strong preference to message passing

A basic syncro problem is referred to as a "race" condition.

"race condition" Consider the code below

//Thread One	// Thread 2
do something that takes a while	do something that takes a while
set a=1	set a=2

Assume that a is a common object share by the threads. Then, these two threads are in a race in that if thread 1 finishes first then the value of a will be 1 for a while then 2. OTOH, if thread 2 finishes first ...

So in this case need a way of "synchronize" the threads in order to get a deterministic result. One of the principle things that a PL offers for parallel programming is synchronization tools.

All shared memory systems have to deal with handling of the shared resource. For example, counter as in counter_go/counter.go

Go **counter_go/counter.go** // discussed this on 12/5

same issue, big undercount

Java **counter_java/BCounter.java**

massive undercounting

Java & Go have at least 3 different approaches will discuss 2. Choice of which is stylistic and situational

ReentrantLock — much like a semaphore

Atomic — a datatype with its own lock.

atomic == inseparable. imagine everything happens at same time

see [counter_go/atomiccounter.go](#), // discussed this on 12/5

[counter_go/atomiccounterB.go](#),

Mutex locks

“mutex” == “MUTual EXclusion”

idea is that a thread locks a door. Until that thread unlocks the door, no other thread can get in. So whenever you lock, you need to be really sure to unlock.

[counter_go/lockcounter.go](#) // here make point of using defer in a function every time you unlock a lock so you guarantee unlock.

Synchronized — much like monitor

see [counter_java/ExecCounter.java](#)

Java: label a whole function as “synchronized” then only one thread can do anything with it. Vector and Hashtable are synchronized but this is really done by making every writing method synchronized. Is it sufficient to have each method simply labeled “synchronized”?

Rust — ownership model makes shared memory hard ... Had to add things to language .. we will not discuss. You can get it, but you really have to try

Message Passing

There is no shared memory!!! so no sync problems

Advantages? naturally applies to distributed systems

Disadvantages?? speed

Need to add discussion of message passing and its implementation in Rust.

mpsc == “multiple produce — single consumer” That is, many threads are allowed to send messages, but only one thread is allowed to read them

See the sequence of main functions in [messages_rust/src/main.rs](#)

main_1_1 send a message

main_1_10 send at least 10 messages

main_10_10 send at least 10 messages from 10 threads to main

main_2d receive some number of messages from threads, then tell each thread to go away.

synch mechanisms

busy-wait /spinning

while (1)

sleep(100)

check if there is something to do

blocking — idea give control to another thread and leave a note somewhere saying you are blocked; In the future, another thread sees the note and unblocks
Alternate: in one thread: in a section of code, put up a block. Run code.

Unblock.

in Java definitely want to unblock in a try/finally. In Go, put unblock in a defer.

Parallelism and Event-driven programming

Instead of parallel programming could achieve some of the same effect through a “dispatch loop” And this is basically what you have to do when using Javascript (and effectively when doing event driven programming).

Dispatch loop and Javascript

JS is single threaded. But it can send stuff off to things outside the language, and do other things while those dispatched items are being processed. “in effect the dispatch loop turns the program inside out, making the management of tasks explicit and the control flow within tasks implicit”

For instance, consider task “get a web page then ask user if it relevant, if yes, save local copy”. in JS, need to write this as:

event Handler for web page loaded — ask if relevant
event Handler for user input about relevance — if yes save, if no download next web page
event Handler for save complete — download next web page

Main

get list of web pages to ask about
start web page download

events_java/GTReader.java

Here we have our own event and handler Idea is that we want to read from the internet. But we want to keep doing stuff. So do the read in a separate thread. When complete, tell the main thread that reading is complete. Done using busy/wait syncing

this example is way to complex

fork == start a new thread
join == merge two threads

Problem with shared memory — synchronization. What happens when two threads read/write to shared memory at the same time?

As in the above programs

More generally, shared-memory systems can be subject to all of these issues:

- *Data race*, when two or more threads attempt to access and write the same data.
 - (examples in Go and java above) Question: how can I reduce/eliminate atomicity violations in the counter programs
- *Deadlock*: when two or more threads attempt to access a shared resource protected with locks (usually in an attempt to fix a data race!)
- see **ThreadLock_java/TL.java**
- *Livelock*: similar to deadlock, except one thread is not actually locked — it executes continuously trying to acquire a shared resource.

- *Starvation*: in which a process is unable to complete in the allotted time because a resource is given to higher priority processes. (Can also occur when there are too many threads)
 - On the counter.go program
 - 100000000 * 10 threads: 38.23s user 0.47s system 707% cpu 5.469 total
 - 100 * 100000000 threads: 91.51s user 12.24s system 484% cpu 21.413 total
 - 10 * 1000000000 threads: 704.96s user 126.58s system 385% cpu 3:35.83 total
- *Blocking suspension*: when one thread waits an unacceptably long time for a resource. Similar to starvation, but it eventually succeeds.
- *Order violation*: when the desired order of operations between multiple threads is violated.
- *Atomicity violation*: in which two code blocks in a thread are interleaved with code blocks in another thread, such that the result of the computation is inconsistent. In other words, a block of code that a developer intended or expected to be atomic is not executed atomically.

Handling the starvation problem: limit threads

First issue after simply starting threads. How do you limit the number of threads running?

Why do you want to limit?

Java **counter_java/EPCounter.java**

effectively the same as CCounter, but stylistically better. Idea is to use

java system that limits the number of active threads. Question: why is this a good idea?

Can accomplish the same things in Go/Rust

Two ways of breaking down a problem:

task parallel == break down problem into a set of "tasks". Run these in parallel.

Problem: can only get as much parallelism as you can define tasks

data Parallel == break problem into smaller pieces by way of data. Then do very similar operations on each piece. (Similar to vector level parallel from above)

Results from tasks

When tasks are running async, you do not get their return value. Why?

So you need a way of capturing the the results from the task.

Shared memory handling

need a method of making sure that shared memory is only used by one thread at a time

Lock

when you set to a point in the code, check if "lock" has been set. If yes, wait

for unlock (bathroom)

Monitor

set a block of code/memory that only one thread at a time can use

Java "synchronized" is essentially a monitor

Monitors probably use locks under the covers, the point is that users

(programmers) do not need to explicitly lock and unlock. (an autolocking door on a bathroom?)

Semaphores

basically a counter which you cannot get by unless it is 0 for you

each thread is responsible for incrementing/decrementing the semaphore

Locks are semaphores with a max count of 1. But semaphores can allow more than one. (Tables in a diner.)

For real world problem using shared memory
see **timeJSON.go** getting info about all public lab machines.

Checkpointing and long parallel runs.