## CSE 505: Concepts of Programming Languages

Dan Grossman Fall 2006 Lecture 14

Concurrency and Shared Memory

# Concurrency

- PL support for concurrency a huge topic
  - And increasingly important (first time in lecture in 505)
- We'll just do *explicit threads* plus:
  - Shared memory (*locks* and *transactions*)
  - Synchronous message passing (Concurrent ML)
- We'll skip

. . .

- Process calculi (foundational message-passing)
- Futures and asynchronous methods
- Data-parallel languages (Snyder)
- Mostly in ML syntax (inference rules where convenient)

### <u>Threads</u>

High-level: "Communicating sequential processes"

Low-level: "Multiple stacks plus communication"

From Caml's thread.mli:

type t (\* a thread handle; remember we're in module Thread \*)
val create : ('a->'b) -> 'a -> t (\* run new thread \*)
val self : unit -> t (\* what thread is executing this? \*)

The *code* for a thread is in a closure (with hidden fields) and Thread.create actually *spawns* the thread.

Most languages make the same distinction, e.g., Java:

- Create a Thread object (just the code and data)
- Call its run method to actually spawn the thread.

## Why use threads?

Any one of:

- 1. Performance (multiprocessor or mask I/O latency)
- 2. Isolation (separate errors or responsiveness)
- 3. Natural code structure (1 stack awkward)

It's not just performance.

Terminology sometimes used (but not universally known):

- Concurrency: interleaved pre-emptive scheduling
- Parallelism: multiple actually at the same time

## One possible formalism (no thread-ids)

- Program state is one heap and multiple expressions
- Any  $e_i$  might "take the next step" and potentially spawn a thread
- A value in the "thread-pool" is removable
- Nondeterministic with *interleaving granularity* determined by rules

Some example rules for  $H; e \rightarrow H'; e'; o$  (where  $o ::= \cdot | e$ ):

 $H; !l \rightarrow H; H(l); \cdot H; \mathsf{spawn}(v_1, v_2) \rightarrow H; 0; (v_1 \ v_2)$ 

$${H;e_1 
ightarrow H';e_1';o \over H;e_1e_2 
ightarrow H';e_1'e_2;o}$$

#### Formalism continued

The  $H; e \to H'; e'; o$  judgment is just a helper-judgment for  $H; T \to H'; T'$  where  $T ::= \cdot \mid e; T$ 

$$egin{aligned} H;e o H';e';\cdot\ H';e_1;\ldots;e;\ldots;e_n o H';e_1;\ldots;e';\ldots;e_n\ H;e o H';e';e''\ H';e_1;\ldots;e;\ldots;e_n o H';e_1;\ldots;e';\ldots;e_n;e'' \end{aligned}$$

 $H; e_1; \ldots; e_{i-1}; v; e_{i+1}; \ldots; e_n \rightarrow H; e_1; \ldots; e_{i-1}; e_{i+1}; \ldots; e_n$ 

Program termination:  $H; \cdot$ 

# Equivalence just changed

Expressions equivalent in a single-threaded world are not necessarily equivalent in a multithreaded context!

Example in Caml:

```
let x, y = ref 0, ref 0
let _ = create (fun () -> if (!y)=1 then x:=(!x)+1)
let _ = create (fun () -> if (!x)=1 then y:=(!y)+1) (* 1 *)
Can we replace line (1) with:
    create (fun () -> y:=(!y)+1; if (!x)<>1 then y:=(!y)-1)
```

For more compiler gotchas, see "Threads cannot be implemented as a library" by Hans-J. Boehm in PLDI2005

### <u>Communication</u>

If threads do nothing other threads needed to see, we are done

- Best to do as little communication as possible
- E.g., do not mutate shared data unnecessarily

One way to communicate: Shared memory

- One thread writes to a ref, another reads it
- Sounds nasty with pre-emptive scheduling
- Hence synchronization mechanisms
  - Taught in O/S for historical reasons!
  - Fundamentally about restricting interleavings

### <u>Join</u>

"Fork-join" parallelism a simple approach good for "farm out subcomputations then merge results"

```
(* suspend caller until/unless arg terminates *)
val join : t -> unit
```

Common pattern:

Apply the second argument to each element of the 'b array in parallel, then use third argument *after* they are done.

See lec14.ml for an (untested) implementation.

### Locks (a.k.a. mutexes)

```
(* mutex.mli *)
type t (* a mutex *)
val create : unit -> t
val lock : t -> unit (* may block *)
val unlock : t -> unit
```

Caml locks do not have two common features:

- Reentrancy (changes semantics of lock)
- Banning nonholder release (changes unlock semantics)

Also want condition variables (condition.mli), but skipping

# Using locks

Among infinite correct idioms using locks (and more incorrect ones), the most common:

- Determine what data must be "kept in sync"
- Always acquire a lock before accessing that data and release it afterwards
- Have a *partial order* on all locks and if a thread holds  $m_1$  it can acquire  $m_2$  only if  $m_1 < m_2$ .

See canonical "bank account" example in lec14.ml.

Coarser locking (more data with same lock) trades off parallelism with synchronization. (Related: Performance-bug of *false sharing*.)

# Getting it wrong

Races result from too little synchronization

- Data races: simultaneous read-write or write-write of same data
  - Lots of PL work in last 10 years on types and tools to prevent/detect.
  - Provided language has some guarantees, may not be a bug
     \* Canonical example: parallel search and "done" bits
- Higher-level races: much tougher to prevent in the language
  - Amount of correct nondeterminism inherently app-specific

Deadlock results from too much synchronization

- Cycle of threads waiting for someone else to do something
- Easy to detect dynamically with locks, but then what?

# The Evolution Problem

Write a new function that needs to update o1 and o2 together.

• What locks should you acquire? In what order?

There may be no answer that avoids races and deadlocks without breaking old code. (Need a stricter partial order.)

See xfer code in lec14.ml, which is yet another binary-method problem for OOP. Real example from Java:

```
synchronized append(StringBuffer sb) {
  int len = sb.length(); //synchronized call
  if(this.count+len > this.value.length) this.expand(...);
  sb.getChars(0,len,this.value,this.count); //synchronized call
```

```
}
```

Undocumented in 1.4; in 1.5 caller synchronizes on sb if necessary.

## Software Transactions

One of the hottest areas in CS research right now (me too).

```
Java: atomic { s }
```

```
Caml: atomic : (unit -> 'a) -> 'a
```

Execute the body/thunk as though there is no interleaving by other threads, while ensuring some scheduling fairness.

Most research on implementation (preserve parallelism unless there are true memory conflicts at run-time), but 505 not an implementation course.

### Transactions compose

Problems like append and xfer become trivial.

So does mixing coarse-grained and fine-grained operations (e.g., hashtable lookup and hashtable resize).

Transactions *are* great, but not a panacea:

- Application-level races can remain
- Application-level deadlock can remain
- Implementations generally try-and-abort, which is hard for "launch missiles" (e.g., I/O)
- Many software implementations provide a weaker and under-specified semantics (come ask me)
- Memory-model questions appear worse than with locks (ongoing research) ...

# Memory models

A *memory model* for a concurrent shared-memory language specifies "which write a read can see".

The gold standard is *sequential consistency* (Lamport): "the results of any execution is the same as if the operations of all the processors were executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by its program"

Under sequential consistency, this assert cannot fail:

## Relaxed memory models

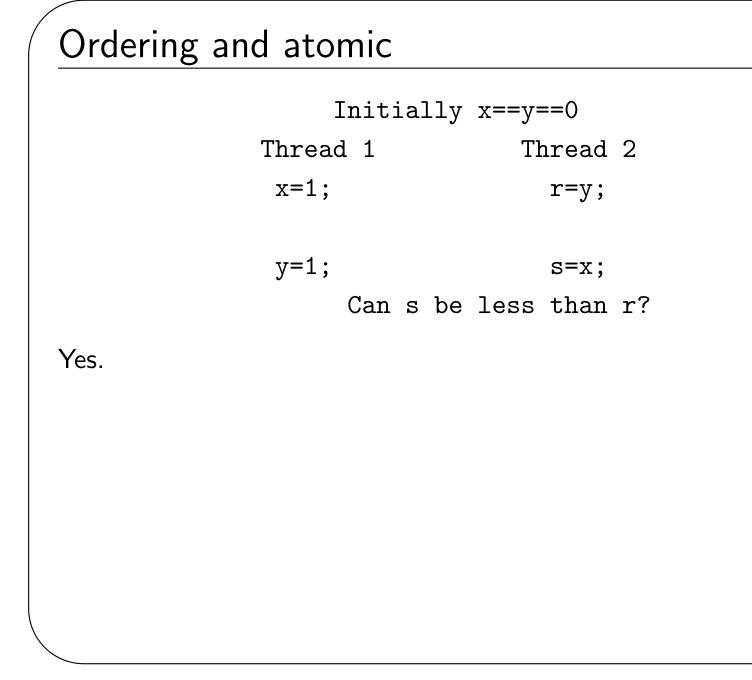
Modern imperative and OO languages do not promise sequential consistency (if they say anything at all)

- The hardware makes it prohibitively expensive
- Renders unsound almost every compiler optimization (e.g., common-subexpression elimination)

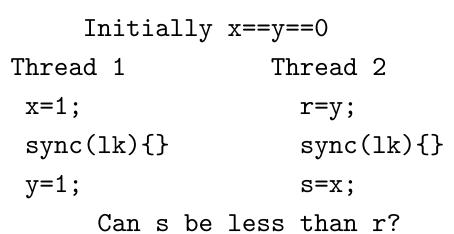
But (especially in a safe language) have to promise something

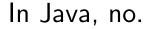
- When is code "correctly synchronized"?
- What can a compiler do in the presence of races? (E.g., cannot seg-fault Java)

The definitions are very complicated and programmers can usually ignore them, but do *not* assume sequential consistency.



### Ordering and atomic





## Ordering and atomic

Initially	x==y==0
Thread 1	Thread 2
x=1;	r=y;
atomic{}	atomic{}
y=1;	s=x;
Can s be	less than r?

Nobody has decided (in practice, yes)! (See my October 06 paper.)