CS 33

Multithreaded Programming VII

CS33 Intro to Computer Systems

XXXVII-1 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

Implementing Mutexes

- Strategy
 - make the usual case (no waiting) very fast
 - can afford to take more time for the other case (waiting for the mutex)

CS33 Intro to Computer Systems

XXXVII-2 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

Futexes

- · Safe, efficient kernel conditional queueing in Linux
- All operations performed atomically
 - futex_wait(futex_t *futex, int val)
 - » if futex->val is equal to val, then sleep
 - » otherwise return
 - futex wake(futex t *futex)
 - » wake up one thread from futex's wait queue, if there are any waiting threads

CS33 Intro to Computer Systems

XXXVII-3 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

For details on futexes, avoid the Linux pages, but look man at http://people.redhat.com/drepper/futex.pdf, from which this material was obtained. Note that there's actually just one futex system call; whether it's a wait or a wakeup is specified by an argument.

Ancillary Functions • int atomic_inc(int *val) - add 1 to *val, return its original value • int atomic_dec(int *val) - subtract 1 from *val, return its original value • int CAS(int *ptr, int old, int new) { int tmp = *ptr; if (*ptr == old) *ptr = new; return tmp; }

These functions are available on most architectures, particularly on the x86. Note that their effect must be **atomic**: everything happens at once.

How can these instructions be made to be atomic? What's done is memory is accessed via special instructions that cause the memory controller to respond to a load then a store without anything happening in between. Thus, for the example of **atomic_inc**, **val** is loaded from memory, then incremented (in the processor), then stored back to memory. While this happens, no other load or stores may be done. If this were done for every instruction, memory access would slow down considerably, but doing it just occasionally has no severe effect.

Attempt 1 void lock(futex_t *futex) { int c; while ((c = atomic_inc(&futex->val)) != 0) futex_wait(futex, c+1); } void unlock(futex_t *futex) { futex->val = 0; futex_wake(futex); } csss Intro to Computer Systems XXXVII-5 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

If the futex's value is 0, it represents an unlocked mutex. If it's 1, it represents a locked mutex.

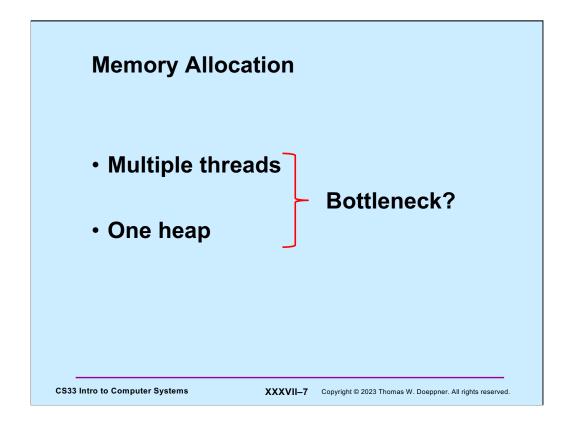
void lock(futex_t *futex) { int c; if ((c = CAS(&futex->val, 0, 1) != 0) do { if (c == 2 || (CAS(&futex->val, 1, 2) != 0)) futex_wait(futex, 2); while ((c = CAS(&futex->val, 0, 2)) != 0)) } void unlock(futex_t *futex) { if (atomic_dec(&futex->val) != 1) { futex_>val = 0; futex_wake(futex); }

In this version, if the futex's value is 0, it represents an unlocked mutex; if it's one it represents a locked mutex that has no threads are waiting for it; if it's greater than one it represents a locked mutex that might have threads waiting for it.

XXXVII-6 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

}

CS33 Intro to Computer Systems



In a naïve multithreaded implementation of malloc/free, there is one mutex protecting the heap, resulting in a bottleneck – a multithreaded program might be slowed down considerably since all threads that manipulate the heap must compete for the mutex.

Solution 1

- Divvy up the heap among the threads
 - each thread has its own heap
 - no mutexes required
 - no bottleneck
- How much heap does each thread get?

CS33 Intro to Computer Systems

XXXVII-8 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

Solution 2

- Multiple "arenas"
 - each with its own mutex
 - thread allocates from the first one it can find whose mutex was unlocked
 - » if none, then creates new one
 - deallocations go back to original arena

CS33 Intro to Computer Systems

XXXVII-9 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

Solution 3

- Global heap plus per-thread heaps
 - threads pull storage from global heap
 - freed storage goes to per-thread heap
 - » unless things are imbalanced
 - then thread moves storage back to global heap
 - mutexes on each heap
- What if one thread allocates and another frees storage?

CS33 Intro to Computer Systems

XXXVII-10 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

Mutexes are required on per-thread heaps for the case when the freeing thread is different from the mallocing thread.

Malloc/Free Implementations

- ptmalloc
 - based on solution 2
 - in glibc (i.e., used by default)
- tcmalloc
 - based on solution 3
 - from Google
- · Which is best?

CS33 Intro to Computer Systems

XXXVII-11 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

Test Program const unsigned int N=64, nthreads=32, iters=10000000; int main() { void *tfunc(void *); pthread t thread[nthreads]; for (int i=0; i<nthreads; i++) {</pre> pthread create(&thread[i], 0, tfunc, (void *)i); pthread detach(thread[i]); pthread exit(0); void *tfunc(void *arg) { long i; for (i=0; i<iters; i++) {</pre> long *p = (long *) malloc(sizeof(long) *((i%N)+1)); free(p); return 0; **CS33 Intro to Computer Systems** XXXVII-12 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

In this test program, each thread does a sequence of mallocs and frees.

Not a Quiz

Which is fastest?

- a) glibc (i.e., standard Linux)
- b) Google

CS33 Intro to Computer Systems

XXXVII-13 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

Compiling It ...

```
% gcc -o ptalloc alloc.c -lpthread
% gcc -o tcalloc alloc.c -lpthread -ltcmalloc
```

CS33 Intro to Computer Systems

XXXVII-14 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

```
Running It (2014) ...
$ time ./ptalloc
          0m5.142s
real
          0m20.501s
user
          0m0.024s
sys
$ time ./tcalloc
          0m1.889s
real
          0m7.492s
user
          0m0.008s
SVS
CS33 Intro to Computer Systems
                                XXXVII-15 Copyright © 2023 Thomas W. Doeppner. All rights reserved.
```

The code was run on an Intel(R) Core(TM)2 Quad CPU Q6600 @ 2.40GHz (4 cores).

The rows labelled **user** show the sums of the amount of time each thread spent running in user mode. The rows labelled **sys** show the sums of the amount of time each thread spent running in kernel mode. The rows labelled **real** show the time that elapsed from when the command started to when it ended. It's less than the sum of the **user** and **sys** times because multiple cores were employed: for example, if two threads running simultaneously (on different cores) each used 1 second of user time, the total user time is 2 seconds, but the real time is one second.

Running It (2023) ... \$ time ./ptalloc real 0m0.666s user 0m5.815s sys 0m0.004s \$ time ./tcalloc real 0m0.496s user 0m4.197s sys 0m0.008s CS33 Intro to Computer Systems XXXVII-16 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

This was run on a 2023 CS department computer: AMD Ryzen 5 3600 @ $7.20 \mathrm{GHz}$ (6 cores). There were 4 times as many iterations as was done in 2014.

What's Going On (2014)? \$ strace -c -f ./ptalloc ... \$ time seconds usecs/call calls errors syscall 100.00 0.040002 13 3007 520 futex ... \$ strace -c -f ./tcalloc ... \$ time seconds usecs/call calls errors syscall ... 0.00 0.000000 0 59 13 futex ... CS33 Intro to Computer Systems XXXVII-17 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

strace is a system facility that supplies information about the system calls a process uses. The –c flag tells it to print the cumulative statistics after the process terminates. The –f flag tells it to include information on all threads and child processes.

Note that the times reported are the total times taken by all threads and don't account for concurrency: i.e., two threads might each take two seconds, totalling to 4 seconds, but the real time used is just two seconds. What's signficant are the counts: the number of calls and the number of errors. Thus its clear that ptalloc makes significantly more calls to futex than does totalloc. Errors indicates the number of times that futex_wait returned because its second argument (val) was not equal to futex->val.

What's Going On (2023)?

```
$ strace -c -f ./ptalloc
...

% time seconds usecs/call calls errors syscall
...

0.02 0.000001 0 5 1 futex
...

$ strace -c -f ./tcalloc
...

% time seconds usecs/call calls errors syscall
...

0.26 0.000006 0 23 3 futex
...

CS33 Intro to Computer Systems XXXVII-18 Copyright © 2023 Thomas W. Doeppner. All rights reserved.
```

Test Program 2, part 1

```
#define N 64
 #define npairs 16
 #define allocsPerIter 1024
 const long iters = 8*1024*1024/allocsPerIter;
 #define BufSize 10240
 typedef struct buffer {
   int *buf[BufSize];
   unsigned int nextin;
   unsigned int nextout;
   sem_t empty;
   sem_t occupied;
   pthread_t pthread;
   pthread t cthread;
 } buffer_t;
CS33 Intro to Computer Systems
                                  XXXVII-19 Copyright © 2023 Thomas W. Doeppner. All rights reserved.
```

This program creates pairs of threads: one thread allocates storage, the other deallocates storage. They communicate using producer-consumer communication.

Int main() { long i; buffer_t b[npairs]; for (i=0; i<npairs; i++) { b[i].nextin = 0; b[i].nextout = 0; sem_init(&b[i].empty, 0, BufSize/allocsPerIter); sem_init(&b[i].occupied, 0, 0); pthread_create(&b[i].pthread, 0, prod, &b[i]); pthread_create(&b[i].cthread, 0, cons, &b[i]); } for (i=0; i<npairs; i++) { pthread_join(b[i].pthread, 0); pthread_join(b[i].cthread, 0); } return 0; }</pre>

XXXVII-20 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

The main function creates **npairs** (16) of communicating pairs of threads.

CS33 Intro to Computer Systems

Test Program 2, part 3 void *prod(void *arg) { long i, j; buffer_t *b = (buffer_t *) arg; for (i = 0; i<iters; i++) { sem_wait(&b->empty); for (j = 0; j<allocsPerIter; j++) { b->buf[b->nextin] = malloc(sizeof(int)*((j%N)+1)); if (++b->nextin >= BufSize) b->nextin = 0; } sem_post(&b->occupied); } return 0; }

To reduce the number of calls to **sem_wait** and **sem_post**, at each iteration the thread calls malloc **allocsPerIter** (1024) times.

XXXVII-21 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

CS33 Intro to Computer Systems

Test Program 2, part 4

```
void *cons(void *arg) {
 long i, j;
 buffer_t *b = (buffer_t *) arg;
 for (i = 0; i<iters; i++) {</pre>
   sem wait(&b->occupied);
   for (j = 0; j<allocsPerIter; j++) {</pre>
     free(b->buf[b->nextout]);
     if (++b->nextout >= BufSize)
       b->nextout = 0;
   sem_post(&b->empty);
  return 0;
```

CS33 Intro to Computer Systems

XXXVII-22 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

The code was run on a SunLab machine (an Intel(R) Core(TM)2 Quad CPU Q6600 @ $2.40 \mathrm{GHz}$).

This was run on a 2023 CS department computer: AMD Ryzen 5 3600 @ 7.20GHz (6 cores).

What's Going On (2014)?

```
$ strace -c -f ./ptalloc2
% time seconds usecs/call calls errors syscall
93.04 8.246196 117 70173 20775 futex
$ strace -c -f ./tcalloc2
% time seconds usecs/call calls errors syscall
99.92 47.796676 153 311012 7244 futex
```

CS33 Intro to Computer Systems XXXVII—25 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

What's Going On (2023)?

```
$ strace -c -f ./ptalloc2
% time seconds usecs/call calls errors syscall
_____ ______
98.48 42.883196 79 539757 179723 futex
$ strace -c -f ./tcalloc2
% time seconds usecs/call calls errors syscall
99.99 346.746205 146 2372684 44547 futex
```

CS33 Intro to Computer Systems XXXVII—26 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

You'll Soon Finish CS 33 ...

- You might
 - celebrate



- take another systems course
 - » 320
 - » 1380
 - » 1660
 - » 1670
 - » 1680





- become a 33 TA

CS33 Intro to Computer Systems

XXXVII-27 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

Systems Courses Next Semester

- CS 320 (Intro to Software Engineering)
 - you've mastered low-level systems programming
 - now do things at a higher level
 - learn software-engineering techniques using Java, XML,
- CS 1380 (Distributed Systems)
 - you now know how things work on one computer
 - what if you've got lots of computers?
 - some may have crashed, others may have been taken over by your worst (and smartest) enemy
- CS 1660/1620/2660 (Computer Systems Security)
 - liked buffer?
 - you'll really like 1660
- CS 1670/1690/2670 (Operating Systems)
 - still mystified about what the OS does?
 - write your own!

CS33 Intro to Computer Systems

XXXVII-28 Copyright © 2023 Thomas W. Doeppner. All rights reserved.

2660 is for graduate students only and combines 1660 and 1620.

2670 is for graduate students only and combines 1670 and 1690.

The End

Well, not quite ... Database is due on 12/15

Happy Coding and Happy Holidays!

CS33 Intro to Computer Systems

XXXVII-29 Copyright © 2023 Thomas W. Doeppner. All rights reserved.