

# CS 33

## Multithreaded Programming VII

# 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)

# 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

For details on futexes, avoid the Linux man pages, but look 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);
}
```

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

## Quiz 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);
}
```

Which of the following won't happen if the futex's value is zero and three threads call *lock* at the same time?

- a) one might return immediately, but at least two will call *futex\_wait*.
- b) even though *unlock* is called appropriately, one thread will never return from *futex\_wait*.
- c) threads might return from *futex\_wait* immediately, because the futex's value is not equal to *c+1*.


## Attempt 2

```
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.

## Memory Allocation

- Multiple threads
  - One heap
- 
- Bottleneck?**

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?**

## 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

## 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
  - mutex on only the global heap
- **What if one thread allocates and another frees storage?**

The latter case implies that there is a mutex on per-thread heaps, for use 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?**

## 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++) {
        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++) {
        long *p = (long *)malloc(sizeof(long)*((i%N)+1));
        free(p);
    }
    return 0;
}
```

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

## Quiz 2

**Which is fastest?**

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

## Compiling It ...

```
% gcc -o ptalloc alloc.cc -lpthread  
% gcc -o talloc alloc.cc -lpthread -ltcmalloc
```

## Running It (2014) ...

```
$ time ./ptalloc
real    0m5.142s
user    0m20.501s
sys     0m0.024s
$ time ./tcalloc
real    0m1.889s
user    0m7.492s
sys     0m0.008s
```

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 (2022) ...

```
$ time ./ptalloc
real    0m1.156s
user    0m3.456s
sys     0m0.004s
$ time ./tcalloc
real    0m0.876s
user    0m3.460s
sys     0m0.004s
```

This was run on a current CS department computer: Intel(R) Core(TM) i5-4690 CPU @ 3.50GHz (4 cores).

## 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

```
...
```

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

## What's Going On (2022)?

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

```
...
```

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

```
-----
```

```
...
```

31.23	0.019968	416	48	6	futex
-------	----------	-----	----	---	-------

```
...
```

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

```
...
```

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

```
-----
```

```
...
```

0.00	0.000000	0	42	3	futex
------	----------	---	----	---	-------

```
...
```

## 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;
```

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

## Test Program 2, part 2

```
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;
}
```

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

## 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.

## Test Program 2, part 4

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

## Running It (2014) ...

```
$ time ./ptalloc2
real    0m1.087s
user    0m3.744s
sys     0m0.204s
$ time ./tcalloc2
real    0m3.535s
user    0m11.361s
sys     0m2.112s
```

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



## Running It (2022) ...

```
$ time ./ptalloc2
real    0m0.367s
user    0m1.187s
sys     0m0.179s
$ time ./tccalloc2
real    0m0.426s
user    0m1.211s
sys     0m0.290s
```

This was run on a current CS department computer: Intel(R) Core(TM) i5-4690 CPU @ 3.50GHz (4 cores).

## What's Going On (2014)?

```
$ strace -c -f ./ptalloc2
```

```
...  
% time      seconds  usecs/call   calls   errors syscall  
-----  
94.96      2.347314      44      53653      14030 futex
```

```
...
```

```
$ strace -c -f ./tccalloc2
```

```
...  
% time      seconds  usecs/call   calls   errors syscall  
-----  
93.86      6.604632      36     185731      45222 futex
```

```
...
```

## What's Going On (2022)?

```
$ strace -c -f ./ptalloc2
```

```
...  
% time      seconds  usecs/call   calls   errors syscall  
-----  
 92.26    4.544802      66    68250    13340  futex
```

```
...
```

```
$ strace -c -f ./tccalloc2
```

```
...  
% time      seconds  usecs/call   calls   errors syscall  
-----  
 91.40    3.439416      52    65165    12182  futex
```

```
...
```

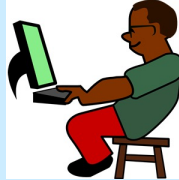
# You'll Soon Finish CS 33 ...

- You might
  - celebrate



- take another systems course

- » 320
- » 1380
- » 1660
- » 1670
- » 1680



- become a 33 TA



## 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, etc.
- **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!

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 12/16

**Happy Coding and Happy Holidays!**