

```

1 Handout for CS 372H
2 Class 8
3 9 February 2012

```

1. How can we implement lock, acquire(), and release()?

1a. Here is A BADLY BROKEN implementation:

```

9 struct Lock {
10     int locked;
11 }
12
13 void [BROKEN] acquire(Lock *lock) {
14     while (1) {
15         if (lock->locked == 0) { // C
16             lock->locked = 1;    // D
17             break;
18         }
19     }
20 }
21
22 void release (Lock *lock) {
23     lock->locked = 0;
24 }

```

What's the problem? Two acquire()'s on the same lock on different CPUs might both execute line C, and then both execute D. Then both will think they have acquired the lock. This is the same kind of race that we were trying to eliminate in insert(). But we have made a little progress: now we only need a way to prevent interleaving in one place (acquire()), not for many arbitrary complex sequences of code.

1b. Here's a way that is correct but only sometimes appropriate: Use an atomic instruction on the CPU. For example, on the x86, doing

```

37     "xchg addr, %eax"
38 does the following:

```

```

39 (i) freeze all CPUs' memory activity for address addr
40 (ii) temp = *addr
41 (iii) *addr = %eax
42 (iv) %eax = temp
43 (v) un-freeze memory activity

```

```

44
45 /* pseudocode */
46 int xchg_val(addr, value) {
47     %eax = value;
48     xchg (*addr), %eax
49 }

```

```

50
51 struct Lock {
52     int locked;
53 }

```

```

54
55 /* bare-bones version of acquire */
56 void acquire (Lock *lock) {
57     pushcli(); /* what does this do? */
58     while (1) {
59         if (xchg_val(&lock->locked, 1) == 0)
60             break;
61     }
62 }

```

```

63
64 /* optimization in acquire; call xchg_val() less frequently */
65 void acquire(Lock* lock) {
66     pushcli();
67     while (xchg_val(&lock->locked, 1) == 1) {
68         while (lock->locked) ;
69     }
70 }
71
72

```

```

73 void release(Lock *lock){
74     xchg_val(&lock->locked, 0);
75     popcli(); /* what does this do? */
76 }
77

```

The above is called a \*spinlock\* because acquire() spins.

The spinlock above is great for some things, not so great for others. The main problem is that it \*busy waits\*: it spins, chewing up CPU cycles. Sometimes this is what we want (e.g., if the cost of going to sleep is greater than the cost of spinning for a few cycles waiting for another thread or process to relinquish the spinlock). But sometimes this is not at all what we want (e.g., if the lock would be held for a while: in those cases, the CPU waiting for the lock would waste cycles spinning instead of running some other thread or process).

1c. Here's an object that does not involve busy waiting; it can work as the list\_lock mentioned above. Note: the "threads" here can be user-level threads, kernel threads, or threads-inside-kernel. The concept is the same in all cases.

```

95
96 struct Mutex {
97     bool is_held; /* true if mutex held */
98     thread_id owner; /* thread holding mutex, if locked */
99     thread_list waiters; /* queue of thread TCBS */
100     Lock wait_lock; /* as in 1b */
101 }

```

Now, instead of acquire(&list\_lock) and release(&list\_lock) as above, we'd write, mutex\_acquire(&list\_mutex) and mutex\_release(&list\_mutex). The implementation of the latter two would be something like this:

```

102
103 void mutex_acquire(Mutex *m) {
104
105     acquire(&m->wait_lock); /* we spin to acquire wait_lock */
106     while (m->is_held) { /* someone else has the mutex */
107         m->waiters.insert(current_thread)
108         release(&m->wait_lock);
109         schedule(); /* run a thread that is on the ready list */
110         acquire(&m->wait_lock); /* we spin again */
111     }
112     m->is_held = true; /* we now hold the mutex */
113     m->owner = self;
114     release(&m->wait_lock);
115 }

```

```

116
117 void mutex_release(Mutex *m) {
118
119     acquire(&m->wait_lock); /* we spin to acquire wait_lock */
120     m->is_held = false;
121     m->owner = 0;
122     wake_up_a_waiter(m->waiters); /* select and run a waiter */
123     release(&m->wait_lock);
124 }

```

[Please let me (MW) know if you see bugs in the above.]

NOTE: Unfortunately, insert() with these locks is correct only if there are some constraints on the order in which the CPU carries out memory reads and writes. For example, if insert() were executed so that the read at A appeared to another processor (and to memory) to be executed before the acquire(), then insert() would be incorrect even with locks.

How do we get the required guarantee? Answer: by ensuring that neither the programmer nor the processor reorders instructions with respect to the acquire().

## 145 2. Terminology

146  
147 To avoid confusion, we will use the following terminology in this  
148 course (you will hear other terminology elsewhere):

149 --A "lock" is an abstract object that provides mutual exclusion

150  
151 --A "spinlock" is a lock that works by busy waiting, as in 6b

152  
153 --A "mutex" is a lock that works by having a "waiting" queue and  
154 then protecting that waiting queue with atomic hardware  
155 instructions, as in 6c. The most natural way to "use the hardware"  
156 is with a spinlock, but there are others, such as turning off  
157 interrupts, which works if we're on a single CPU machine.  
158  
159  
160

## 161 3. Producer/consumer example [also known as bounded buffer]

162  
163 3a. Recall buggy implementation  
164  
165 /\*  
166 "buffer" stores BUFFER\_SIZE items  
167 "count" is number of used slots. a variable that lives in memory  
168 "out" is next empty buffer slot to fill (if any)  
169 "in" is oldest filled slot to consume (if any)  
170 \*/  
171  
172 void producer (void \*ignored) {  
173 for (;;) {  
174 /\* next line produces an item and puts it in nextProduced \*/  
175 nextProduced = means\_of\_production();  
176 while (count == BUFFER\_SIZE)  
177 ; // do nothing  
178 buffer [in] = nextProduced;  
179 in = (in + 1) % BUFFER\_SIZE;  
180 count++;  
181 }  
182 }  
183  
184 void consumer (void \*ignored) {  
185 for (;;) {  
186 while (count == 0)  
187 ; // do nothing  
188 nextConsumed = buffer[out];  
189 out = (out + 1) % BUFFER\_SIZE;  
190 count--;  
191 /\* next line abstractly consumes the item \*/  
192 consume\_item(nextConsumed);  
193 }  
194 }  
195

196 --Review: what's the problem?

197 --Answer: count++ and count-- might compile to, respectively:

```
198
199     reg1 <-- count      # load
200     reg1 <-- reg1 + 1  # increment register
201     count <-- reg1     # store
202
203     reg2 <-- count     # load
204     reg2 <-- reg2 - 1  # decrement register
205     count <-- reg2    # store
206
```

207 --Review: why not use instructions like "addl \$0x1, \_count"?

208 --Answer: not atomic if there are multiple CPUs.

209  
210 --Review: so why not use "LOCK addl \$0x1, \_count"?

211 --Answer: we could do that here, but LOCK won't save us every time

212  
213 --Review: so use general-purpose approach to protecting  
214 critical sections: locks (or mutexes).  
215  
216

Mar 19, 12 16:21

I08-handout-1.txt

Page 5/11

```

217
218 3b. Producer/consumer [bounded buffer] using mutexes
219
220     Mutex mutex;
221
222     void producer (void *ignored) {
223         for (;;) {
224             /* next line produces an item and puts it in nextProduced */
225             nextProduced = means_of_production();
226
227             acquire(&mutex);
228             while (count == BUFFER_SIZE) {
229                 release(&mutex);
230                 yield(); /* or schedule() */
231                 acquire(&mutex);
232             }
233
234             buffer [in] = nextProduced;
235             in = (in + 1) % BUFFER_SIZE;
236             count++;
237             release(&mutex);
238         }
239     }
240
241     void consumer (void *ignored) {
242         for (;;) {
243
244             acquire(&mutex);
245             while (count == 0) {
246                 release(&mutex);
247                 yield(); /* or schedule() */
248                 acquire(&mutex);
249             }
250
251             nextConsumed = buffer[out];
252             out = (out + 1) % BUFFER_SIZE;
253             count--;
254             release(&mutex);
255
256             /* next line abstractly consumes the item */
257             consume_item(nextConsumed);
258         }
259     }
260

```

Mar 19, 12 16:21

I08-handout-1.txt

Page

```

261
262 3c. Producer/consumer [bounded buffer] using mutexes and condition
263 variables
264
265     Mutex mutex;
266     Cond nonempty;
267     Cond nonfull;
268
269     void producer (void *ignored) {
270         for (;;) {
271             /* next line produces an item and puts it in nextProduced */
272             nextProduced = means_of_production();
273
274             acquire(&mutex);
275             while (count == BUFFER_SIZE)
276                 cond_wait(&nonfull, &mutex);
277
278             buffer [in] = nextProduced;
279             in = (in + 1) % BUFFER_SIZE;
280             count++;
281             cond_signal(&nonempty, &mutex);
282             release(&mutex);
283         }
284     }
285
286     void consumer (void *ignored) {
287         for (;;) {
288
289             acquire(&mutex);
290             while (count == 0)
291                 cond_wait(&nonempty, &mutex);
292
293             nextConsumed = buffer[out];
294             out = (out + 1) % BUFFER_SIZE;
295             count--;
296             cond_signal(&nonfull, &mutex);
297             release(&mutex);
298
299             /* next line abstractly consumes the item */
300             consume_item(nextConsumed);
301         }
302     }
303
304
305     Question: why does cond_wait need to both release the mutex and
306     sleep? Why not:
307
308         while (count == BUFFER_SIZE) {
309             release(&mutex);
310             cond_wait(&nonfull);
311             acquire(&mutex);
312         }
313

```

```

314 3d. Producer/consumer [bounded buffer] with semaphores
315
316 Semaphore mutex(1); /* mutex initialized to 1 */
317 Semaphore empty(BUFFER_SIZE); /* start with BUFFER_SIZE empty slots */
318 Semaphore full(0); /* 0 full slots */
319
320 void producer (void *ignored) {
321     for (;;) {
322         /* next line produces an item and puts it in nextProduced */
323         nextProduced = means_of_production();
324
325         /*
326          * next line diminishes the count of empty slots and
327          * waits if there are no empty slots
328          */
329         sem_down(&empty);
330         sem_down(&mutex); /* get exclusive access */
331
332         buffer [in] = nextProduced;
333         in = (in + 1) % BUFFER_SIZE;
334
335         sem_up(&mutex);
336         sem_up(&full); /* we just increased the # of full slots */
337     }
338 }
339
340 void consumer (void *ignored) {
341     for (;;) {
342         /*
343          * next line diminishes the count of full slots and
344          * waits if there are no full slots
345          */
346         sem_down(&full);
347         sem_down(&mutex);
348
349         nextConsumed = buffer[out];
350         out = (out + 1) % BUFFER_SIZE;
351
352         sem_up(&mutex);
353         sem_up(&empty); /* one further empty slot */
354
355         /* next line abstractly consumes the item */
356         consume_item(nextConsumed);
357     }
358 }
359
360 Semaphores *can* (not always) lead to elegant solutions (notice
361 that the code above is fewer lines than 3c) but they are much
362 harder to use.
363
364 The fundamental issue is that semaphores make implicit (counts,
365 conditions, etc.) what is probably best left explicit. Moreover,
366 they *also* implement mutual exclusion.
367
368 For this reason, you should not use semaphores. This example is
369 here mainly for completeness and so you know what a semaphore
370 is. But do not code with them. Solutions that use semaphores in
371 this course will receive no credit.
372
373

```

```

374 4. Example of a monitor: MyBuffer
375
376 // This is pseudocode that is inspired by C++.
377 // Don't take it literally.
378
379 class MyBuffer {
380     public:
381         MyBuffer();
382         ~MyBuffer();
383         void Enqueue(Item);
384         Item = Dequeue();
385     private:
386         int count;
387         int in;
388         int out;
389         Item buffer[BUFFER_SIZE];
390         Mutex* mutex;
391         Cond* nonempty;
392         Cond* nonfull;
393 }
394
395 void
396 MyBuffer::MyBuffer()
397 {
398     in = out = count = 0;
399     mutex = new Mutex;
400     nonempty = new Cond;
401     nonfull = new Cond;
402 }
403
404 void
405 MyBuffer::Enqueue(Item item)
406 {
407     mutex.acquire();
408     while (count == BUFFER_SIZE)
409         cond_wait(&nonfull, &mutex);
410
411     buffer[in] = item;
412     in = (in + 1) % BUFFER_SIZE;
413     ++count;
414     cond_signal(&nonempty, &mutex);
415     mutex.release();
416 }
417
418 Item
419 MyBuffer::Dequeue()
420 {
421     mutex.acquire();
422     while (count == 0)
423         cond_wait(&nonempty, &mutex);
424
425     Item ret = buffer[out];
426     out = (out + 1) % BUFFER_SIZE;
427     --count;
428     cond_signal(&nonfull, &mutex);
429     mutex.release();
430     return ret;
431 }
432

```

Mar 19, 12 16:21

I08-handout-1.txt

Page 9/11

```

433 int main(int, char**)
434 {
435     MyBuffer buf;
436     int dummy;
437     tid1 = thread_create(producer, &buf);
438     tid2 = thread_create(consumer, &buf);
439     thread_join(tid1);
440
441     // never reach this point
442     return -1;
443 }
444
445 void producer(void* buf)
446 {
447     MyBuffer* sharedbuf = reinterpret_cast<MyBuffer*>(buf);
448     for (;;) {
449         /* next line produces an item and puts it in nextProduced */
450         Item nextProduced = means_of_production();
451         sharedbuf->Enqueue(nextProduced);
452     }
453 }
454
455 void consumer(void* buf)
456 {
457     MyBuffer* sharedbuf = reinterpret_cast<MyBuffer*>(buf);
458     for (;;) {
459         Item nextConsumed = sharedbuf->Dequeue();
460
461         /* next line abstractly consumes the item */
462         consume_item(nextConsumed);
463     }
464 }

```

465 Key point: \*Threads\* (the producer and consumer) are separate from  
466 \*shared object\* (MyBuffer). The synchronization happens in the  
467 shared object.

#### 469 5. Readers/writers

```

470 state variables:
471
472     AR = 0; // # active readers
473     AW = 0; // # active writers
474     WR = 0; // # waiting readers
475     WW = 0; // # waiting writers
476
477     Condition okToRead = NIL;
478     Condition okToWrite = NIL;
479     Mutex mutex = FREE;
480
481
482 Database::read() {
483     startRead(); // first, check self into the system
484     Access Data
485     doneRead(); // check self out of system
486 }
487
488 Database::startRead() {
489     acquire(&mutex);
490     while((AW + WW) > 0){
491         WR++;
492         wait(&okToRead, &mutex);
493         WR--;
494     }
495     AR++;
496     release(&mutex);
497 }
498
499 Database::doneRead() {
500     acquire(&mutex);
501     AR--;
502     if (AR == 0 && WW > 0) { // if no other readers still
503         signal(&okToWrite, &mutex); // active, wake up writer
504     }

```

Mar 19, 12 16:21

I08-handout-1.txt

Page 10

```

505     release(&mutex);
506 }
507
508 Database::write(){ // symmetrical
509     startWrite(); // check in
510     Access Data
511     doneWrite(); // check out
512 }
513
514 Database::startWrite() {
515     acquire(&mutex);
516     while ((AW + AR) > 0) { // check if safe to write.
517         // if any readers or writers, wait
518         WW++;
519         wait(&okToWrite, &mutex);
520         WW--;
521     }
522     AW++;
523     release(&mutex);
524 }
525
526 Database::doneWrite() {
527     acquire(&mutex);
528     AW--;
529     if (WW > 0) {
530         signal(&okToWrite, &mutex); // give priority to writers
531     } else if (WR > 0) {
532         broadcast(&okToRead, &mutex);
533     }
534     release(&mutex);
535 }
536
537 NOTE: what is the starvation problem here?
538

```

```
539 6. Shared locks
540
541     struct sharedlock {
542         int i;
543         Mutex mutex;
544         Cond c;
545     };
546
547     void AcquireExclusive (sharedlock *sl) {
548         acquire(&sl->mutex);
549         while (sl->i) {
550             wait (&sl->c, &sl->mutex);
551         }
552         sl->i = -1;
553         release(&sl->mutex);
554     }
555
556     void AcquireShared (sharedlock *sl) {
557         acquire(&sl->mutex);
558         while (sl->i < 0) {
559             wait (&sl->c, &sl->mutex);
560         }
561         sl->i++;
562         release(&sl->mutex);
563     }
564
565     void ReleaseShared (sharedlock *sl) {
566         acquire(&sl->mutex);
567         if (--sl->i)
568             signal (&sl->c, &sl->mutex);
569         release(&sl->mutex);
570     }
571
572     void ReleaseExclusive (sharedlock *sl) {
573         acquire(&sl->mutex);
574         sl->i = 0;
575         broadcast (&sl->c, &sl->mutex);
576         release(&sl->mutex);
577     }
578
579     QUESTIONS:
580     A. There is a starvation problem here. What is it? (Readers can keep
581        writers out if there is a steady stream of readers.)
582     B. How could you use these shared locks to write a cleaner version
583        of the code in item 5., above? (Though note that the starvation
584        properties would be different.)
585
586
```

**13. [12 points]** Consider the function `doublecheck_alloc()` below, which is intended to be invoked from multiple threads on a multiprocessor machine. Its purpose is to avoid a mutex acquisition in the common case that `ptr` is already initialized. The requirements for this function are:

- (i) `doublecheck_alloc()` must call `alloc_foo()` no more than once over the whole execution.
- (ii) A caller of `doublecheck_alloc()` must, after the function returns, observe `ptr` as non-zero.

The machine does **not** offer sequential consistency. Thus, a processor is not guaranteed to see the memory operations of another processor in program order. However, each of `mutex_acquire()` and `mutex_release()` is implemented correctly; in particular, each of them internally contains a memory barrier (`mfence` on the x86). Recall that `mfence` ensures that all memory operations before the `mfence` barrier appear to all processors to have executed before all memory operations after the `mfence` barrier.

On the other hand, the compiler **preserves** program order (it does not reorder instructions).

```

struct foo {
    int abc;
    int def;
};
static int ready = 0;
static mutex_t mutex;
static struct foo* ptr = 0;

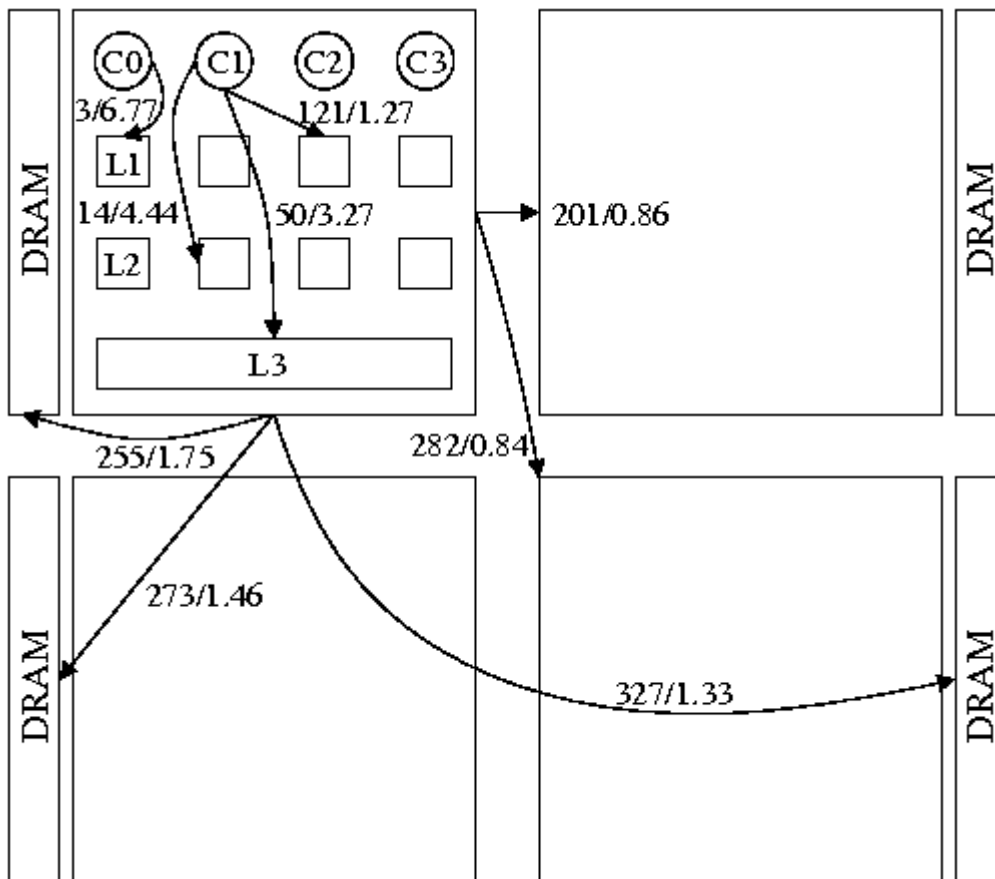
void
doublecheck_alloc()
{
    if (!ready) { /* <-- accesses shared variable w/out holding mutex */
        mutex_acquire(&mutex);
        if (!ready) {
            ptr = alloc_foo(); /* <-- sets ptr to be non-zero */
            ready = 1;
        }
        mutex_release(&mutex);
    }
    return;
}

```

The above code certainly violates our coding standards, but this problem is about whether it violates requirements (i) and (ii), above. The questions are given on the next page.

Name:

UT EID:



The AMD 16-core system topology. Memory access latency is in cycles and listed before the backslash. Memory bandwidth is in bytes per cycle and listed after the backslash. The measurements reflect the latency and bandwidth achieved by a core issuing load instructions. The measurements for accessing the L1 or L2 caches of a different core on the same chip are the same. The measurements for accessing any cache on a different chip are the same. Each cache line is 64 bytes, L1 caches are 64 Kbytes 8-way set associative, L2 caches are 512 Kbytes 16-way set associative, and L3 caches are 2 Mbytes 32-way set associative.

[Reprinted with permission from S. Boyd-Wickizer et al. Corey: An Operating System for Many Cores. Proceedings of Usenix Symposium on Operating Systems Design and Implementation (OSDI), December 2008.]



```

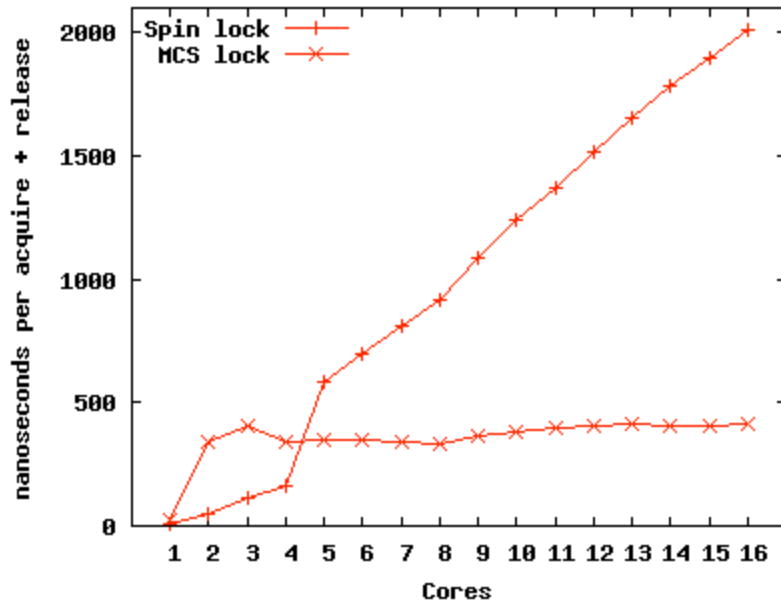
1
2 A. CAS / CMPXCHG
3
4 Useful operation: compare-and-swap, known as CAS. Says: "atomically
5 check whether a given memory cell contains a given value, and if it
6 does, then replace the contents of the memory cell with this other
7 value; in either case, return the original value in the memory
8 location".
9
10 On the X86, we implement CAS with the CMPXCHG instruction, but note
11 that this instruction is not atomic by default, so we need the LOCK
12 prefix.
13
14 Here's pseudocode:
15
16     int cmpxchg_val(int* addr, int oldval, int newval) {
17         LOCK: // remember, this is pseudocode
18         int was = *addr;
19         if (*addr == oldval)
20             *addr = newval;
21         return was;
22     }
23
24 Here's inline assembly:
25
26     uint32_t cmpxchg_val(uint32_t* addr, uint32_t oldval, uint32_t newval) {
27         uint32_t was;
28         asm volatile("lock cmpxchg %3, %0"
29                     : "+m" (*addr), "=a" (was)
30                     : "a" (oldval), "r" (newval)
31                     : "cc");
32         return was;
33     }
34
35 B. MCS locks
36
37 Citation: Mellor-Crummey, J. M. and M. L. Scott. Algorithms for
38 Scalable Synchronization on Shared-Memory Multiprocessors, ACM
39 Transactions on Computer Systems, Vol. 9, No. 1, February, 1991,
40 pp.21-65.
41
42 Each CPU has a qnode structure in *local* memory. Here, local can
43 mean local memory in NUMA machine or its own cache line that other
44 CPUs are not allowed to cache (i.e., the cache line is in exclusive
45 mode):
46
47     typedef struct qnode {
48         struct qnode* next;
49         bool someoneelse_locked;
50     } qnode;
51
52     typedef qnode* lock; // a lock is a pointer to a qnode
53
54 --The lock itself is literally the *tail* of the list of CPUs holding
55 or waiting for the lock.
56
57 --While waiting, a CPU spins on its local "locked" flag. Here's the
58 code for acquire:
59

```

```

60 // lockp is a qnode**. I points to our local qnode.
61 void acquire(lock* lockp, qnode* I) {
62
63     I->next = NULL;
64     qnode* predecessor;
65
66     // next line makes lockp point to I (that is, it sets *lockp <-- I)
67     // and returns the old value of *lockp. Uses atomic operation
68     // XCHG. see earlier in handout (or earlier handouts)
69     // for implementation of xchg_val.
70
71     predecessor = xchg_val(lockp, I); // "A"
72     if (predecessor != NULL) { // queue was non-empty
73         I->someoneelse_locked = true;
74         predecessor->next = I; // "B"
75         while (I->someoneelse_locked); // spin
76     }
77     // we hold the lock!
78 }
79
80 What's going on?
81
82 --If the lock is unlocked, then *lockp == NULL.
83
84 --If the lock is locked, and there are no waiters, then *lockp
85 points to the qnode of the owner
86
87 --If the lock is locked, and there are waiters, then *lockp points
88 to the qnode at the tail of the waiter list.
89
90 --Here's the code for release:
91
92 void release(lock* lockp, qnode* I) {
93     if (!I->next) { // no known successor
94         if (cmpxchg_val(lockp, I, NULL) == I) { // "C"
95             // swap successful: lockp was pointing to I, so now
96             // *lockp == NULL, and the lock is unlocked. we can
97             // go home now.
98             return;
99         }
100        // if we get here, then there was a timing issue: we had
101        // no known successor when we first checked, but now we
102        // have a successor: some CPU executed the line "A"
103        // above. Wait for that CPU to execute line "B" above.
104        while (!I->next);
105    }
106    // handing the lock off to the next waiter is as simple as
107    // just setting that waiter's "someoneelse_locked" flag to false
108    I->next->someoneelse_locked = false;
109 }
110
111 What's going on?
112
113 --If I->next == NULL and *lockp == I, then no one else is
114 waiting for the lock. So we set *lockp == NULL.
115
116 --If I->next == NULL and *lockp != I, then another CPU is in
117 acquire (specifically, it executed its atomic operation, namely
118 line "A", before we executed ours, namely line "C"). So wait for
119 the other CPU to put the list in a sane state, and then drop
120 down to the next case:
121
122 --If I->next != NULL, then we know that there is a spinning
123 waiter (the oldest one). Hand it the lock by setting its flag to
124 false.

```



Time required to acquire and release a lock on a 16-core AMD machine when varying number of cores contend for the lock. The two lines show Linux kernel spin locks and MCS locks (on Corey). A spin lock with one core takes about 11 nanoseconds; an MCS lock about 26 nanoseconds.

[Reprinted with permission from S. Boyd-Wickizer et al. Corey: An Operating System for Many Cores. Proceedings of Symposium on Operating Systems Design and Implementation (OSDI), December 2008.]

Feb 09, 12 15:22

I08-handout-3.txt

Page 1/2

```

1 Performance v complexity trade-off with locks
2
3 /*
4 *      linux/mm/filemap.c
5 *
6 * Copyright (C) 1994-1999 Linus Torvalds
7 */
8
9 /*
10 * This file handles the generic file mmap semantics used by
11 * most "normal" filesystems (but you don't /have/ to use this:
12 * the NFS filesystem used to do this differently, for example)
13 */
14 #include <linux/config.h>
15 #include <linux/module.h>
16 #include <linux/slab.h>
17 #include <linux/compiler.h>
18 #include <linux/fs.h>
19 #include <linux/aio.h>
20 #include <linux/capability.h>
21 #include <linux/kernel_stat.h>
22 #include <linux/mm.h>
23 #include <linux/swap.h>
24 #include <linux/mman.h>
25 #include <linux/pagemap.h>
26 #include <linux/file.h>
27 #include <linux/uio.h>
28 #include <linux/hash.h>
29 #include <linux/writeback.h>
30 #include <linux/pagevec.h>
31 #include <linux/blkdev.h>
32 #include <linux/security.h>
33 #include <linux/syscalls.h>
34 #include "filemap.h"
35 /*
36 * FIXME: remove all knowledge of the buffer layer from the core VM
37 */
38 #include <linux/buffer_head.h> /* for generic_osync_inode */
39
40 #include <asm/uaccess.h>
41 #include <asm/mman.h>
42
43 static ssize_t
44 generic_file_direct_IO(int rw, struct kiocb *iocb, const struct iovec *iov,
45                       loff_t offset, unsigned long nr_segs);
46
47 /*
48 * Shared mappings implemented 30.11.1994. It's not fully working yet,
49 * though.
50 *
51 * Shared mappings now work. 15.8.1995 Bruno.
52 *
53 * finished 'unifying' the page and buffer cache and SMP-threaded the
54 * page-cache, 21.05.1999, Ingo Molnar <mingo@redhat.com>
55 *
56 * SMP-threaded pagemap-LRU 1999, Andrea Arcangeli <andrea@suse.de>
57 */
58
59 /*
60 * Lock ordering:
61 *
62 * ->i_mmap_lock          (vmtruncate)
63 * ->private_lock        (__free_pte->__set_page_dirty_buffers)
64 * ->swap_lock           (exclusive_swap_page, others)
65 * ->mapping->tree_lock
66 *
67 * ->i_mutex
68 * ->i_mmap_lock          (truncate->unmap_mapping_range)
69 *
70 * ->mmap_sem
71 * ->i_mmap_lock
72 * ->page_table_lock or pte_lock (various, mainly in memory.c)
73 * ->mapping->tree_lock (arch-dependent flush_dcache_mmap_lock)

```

Feb 09, 12 15:22

I08-handout-3.txt

Page 2/2

```

74 *
75 * ->mmap_sem
76 * ->lock_page          (access_process_vm)
77 *
78 * ->mmap_sem
79 * ->i_mutex             (msync)
80 *
81 * ->i_mutex
82 * ->i_alloc_sem        (various)
83 *
84 * ->inode_lock
85 * ->sb_lock             (fs/fs-writeback.c)
86 * ->mapping->tree_lock (__sync_single_inode)
87 *
88 * ->i_mmap_lock
89 * ->anon_vma.lock       (vma_adjust)
90 *
91 * ->anon_vma.lock
92 * ->page_table_lock or pte_lock (anon_vma_prepare and various)
93 *
94 * ->page_table_lock or pte_lock
95 * ->swap_lock           (try_to_unmap_one)
96 * ->private_lock        (try_to_unmap_one)
97 * ->tree_lock           (try_to_unmap_one)
98 * ->zone.lru_lock        (follow_page->mark_page_accessed)
99 * ->zone.lru_lock        (check_pte_range->isolate_lru_page)
100 * ->private_lock        (page_remove_rmap->set_page_dirty)
101 * ->tree_lock           (page_remove_rmap->set_page_dirty)
102 * ->inode_lock          (page_remove_rmap->set_page_dirty)
103 * ->inode_lock          (zap_pte_range->set_page_dirty)
104 * ->private_lock        (zap_pte_range->__set_page_dirty_buffers)
105 *
106 * ->task->proc_lock
107 * ->dcache_lock         (proc_pid_lookup)
108 */
109
110 /*
111 * Remove a page from the page cache and free it. Caller has to make
112 * sure the page is locked and that nobody else uses it - or that usage
113 * is safe. The caller must hold a write_lock on the mapping's tree_lock.
114 */
115 void __remove_from_page_cache(struct page *page)
116 {
117     struct address_space *mapping = page->mapping;
118
119     .....
120
121 [point of this item on the handout: fine-grained locking leads to complexity]

```