

Apr 21, 16 12:03

c23-handout.txt

Page 1/8

```

1 Handout for CS 480
2 Class 23
3 21 April 2016
4
5 1. Example to illustrate interleavings: say that thread A executes f()
6 and thread B executes g(). (Here, we are using the term "thread"
7 abstractly, to refer to execution contexts that share memory.)
8
9 a.
10
11     int x;
12
13     f() { x = 1; }
14
15     g() { x = 2; }
16
17 What are possible values of x after A has executed f() and B has
18 executed g()?
19
20 b.
21
22     int y = 12;
23
24     f() { x = y + 1; }
25     g() { y = y * 2; }
26
27 What are the possible values of x?
28
29 c.
30     int x = 0;
31     f() { x = x + 1; }
32     g() { x = x + 2; }
33
34 What are the possible values of x?
35
36 2. Linked list example
37
38 struct List_elem {
39     int data;
40     struct List_elem* next;
41 };
42
43 List_elem* head = 0;
44
45 insert(int data) {
46     List_elem* l = new List_elem;
47     l->data = data;
48     l->next = head;
49     head = l;
50 }
51
52 What happens if two threads execute insert() at once and we get the
53 following interleaving?
54
55 thread 1: l->next = head
56 thread 2: l->next = head
57 thread 2: head = l;
58 thread 1: head = l;
59

```

Apr 21, 16 12:03

c23-handout.txt

Page 2/8

```

59
60 3. Producer/consumer example:
61
62 /*
63 "buffer" stores BUFFER_SIZE items
64 "count" is number of used slots. a variable that lives in memory
65 "out" is next empty buffer slot to fill (if any)
66 "in" is oldest filled slot to consume (if any)
67 */
68
69 void producer (void *ignored) {
70     for (;;) {
71         /* next line produces an item and puts it in nextProduced */
72         nextProduced = means_of_production();
73         while (count == BUFFER_SIZE)
74             ; // do nothing
75         buffer [in] = nextProduced;
76         in = (in + 1) % BUFFER_SIZE;
77         count++;
78     }
79 }
80
81 void consumer (void *ignored) {
82     for (;;) {
83         while (count == 0)
84             ; // do nothing
85         nextConsumed = buffer[out];
86         out = (out + 1) % BUFFER_SIZE;
87         count--;
88         /* next line abstractly consumes the item */
89         consume_item(nextConsumed);
90     }
91 }
92
93 /*
94 what count++ probably compiles to:
95 reg1 <- count      # load
96 reg1 <- reg1 + 1   # increment register
97 count <- reg1      # store
98
99 what count-- could compile to:
100 reg2 <- count      # load
101 reg2 <- reg2 - 1   # decrement register
102 count <- reg2      # store
103 */
104
105 What happens if we get the following interleaving?
106
107 reg1 <- count
108 reg1 <- reg1 + 1
109 reg2 <- count
110 reg2 <- reg2 - 1
111 count <- reg1
112 count <- reg2
113

```

Apr 21, 16 12:03

c23-handout.txt

Page 3/8

```

114 4. Protecting the linked list.....
115
116     Lock list_lock;
117
118     insert(int data) {
119         List_elem* l = new List_elem;
120         l->data = data;
121
122         acquire(&list_lock);
123
124         l->next = head;           // A
125         head = l;                // B
126
127         release(&list_lock);
128     }
129
130 5. How can we implement list_lock, acquire(), and release()?
131
132     Here is A BADLY BROKEN implementation:
133
134     struct Lock {
135         int locked;
136     }
137
138     void [BROKEN] acquire(Lock *lock) {
139         while (1) {
140             if (lock->locked == 0) { // C
141                 lock->locked = 1;   // D
142                 break;
143             }
144         }
145     }
146
147     void release (Lock *lock) {
148         lock->locked = 0;
149     }

```

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.

Apr 21, 16 12:03

c23-handout.txt

Page 4/8

```

159 5a. Test-and-set spinlock
160
161     Relies on atomic instruction on the CPU. For example, on the x86,
162     doing
163         "xchg addr, %eax"
164     atomically swaps the contents of %eax with the contents of
165     (virtual) memory address addr. No other instructions can be
166     interleaved. One can think of xchg like this:
167
168     (i)   freeze all CPUs' memory activity for address addr
169     (ii)  temp = *addr
170     (iii) *addr = %eax
171     (iv)  %eax = temp
172     (v)   un-freeze memory activity
173
174     /* pseudocode */
175     int xchg_val(addr, value) {
176         %eax = value;
177         xchg (*addr), %eax
178     }
179
180     struct Lock {
181         int locked;
182     }
183
184     /* bare-bones version of acquire */
185     void acquire (Lock *lock) {
186         pushcli(); /* what does this do? */
187         while (1) {
188             if (xchg_val(&lock->locked, 1) == 0)
189                 break;
190         }
191     }
192
193     void release(Lock *lock){
194         xchg_val(&lock->locked, 0);
195         popcli(); /* what does this do? */
196     }
197
198 5b. Test-and-test-and-set lock
199
200     /* optimization in acquire; call xchg_val() less frequently */
201     void acquire(Lock* lock) {
202         pushcli();
203         while (xchg_val(&lock->locked, 1) == 1) {
204             while (lock->locked) ;
205         }
206     }
207
208
209

```

Apr 21, 16 12:03

c23-handout.txt

Page 5/8

```

210 6. Ticket locks
211
212 The spinlocks presented above have fairness issues on NUMA machines
213 (cores closer to the memory containing the 'locked' variable are
214 more likely to succeed in acquiring the lock).
215
216 Ticket locks address that issue.
217
218 They rely on an atomic primitive known as "fetch and increment."
219 On the x86, we implement fetch and increment with the XADD
220 instruction, but note that this instruction is not atomic by
221 default, so we need the LOCK prefix.
222
223 Here's pseudocode:
224
225     int fetch_and_increment (int* addr) {
226         LOCK: // remember, this is pseudocode
227         int was = *addr;
228         *addr = was + 1;
229         return was;
230     }
231
232 Here's inline assembly:
233
234     inline int fetch_and_increment(int *addr) {
235         int was = 1;
236         asm volatile("lock xaddl %1, %0"
237                     : "+m" (*addr), "=r" (was) // Output
238                     : "1" (was), "m" (*addr) // Input
239                     );
240         return was;
241     }
242
243     struct Lock {
244         int current_ticket;
245         int next_ticket;
246     }
247
248     void acquire (Lock *lock)
249     {
250         int t = fetch_and_increment (&lock->next_ticket);
251         while (t != lock->current_ticket) ;
252     }
253
254     void release (Lock *lock) {
255         lock->current_ticket++;
256     }
257

```

Apr 21, 16 12:03

c23-handout.txt

Page 6/8

```

258 7. MCS locks (a kind of queue lock)
259
260     Ticket locks are fair, as noted above, but they (and baseline
261     spinlocks) have performance issues when there is a lot of
262     contention. These issues fundamentally result from cross-talk among
263     CPUs (which undermines caching and generates traffic on the memory
264     bus). This phenomenon is investigated in depth in the "Scalable
265     Locks are Dangerous" paper.
266
267     The locks presented below address that issue. These are known as MCS
268     locks.
269
270     Citation: Mellor-Crummey, J. M. and M. L. Scott. Algorithms for
271     Scalable Synchronization on Shared-Memory Multiprocessors, ACM
272     Transactions on Computer Systems, Vol. 9, No. 1, February, 1991,
273     pp. 21-65.
274
275     A. CAS / CMPXCHG
276
277     Useful operation: compare-and-swap, known as CAS. Says: "atomically
278     check whether a given memory cell contains a given value, and if it
279     does, then replace the contents of the memory cell with this other
280     value; in either case, return the original value in the memory
281     location".
282
283     On the X86, we implement CAS with the CMPXCHG instruction, but note
284     that this instruction is not atomic by default, so we need the LOCK
285     prefix.
286
287     Here's pseudocode:
288
289     int cmpxchg_val(int* addr, int oldval, int newval) {
290         LOCK: // remember, this is pseudocode
291         int was = *addr;
292         if (*addr == oldval)
293             *addr = newval;
294         return was;
295     }
296
297     Here's inline assembly:
298
299     uint32_t cmpxchg_val(uint32_t* addr, uint32_t oldval, uint32_t newval) {
300         uint32_t was;
301         asm volatile("lock cmpxchg %3, %0"
302                     : "+m" (*addr), "=a" (was)
303                     : "a" (oldval), "r" (newval), "m" (*addr)
304                     : "cc");
305         return was;
306     }
307
308     B. The MCS lock
309
310     Each CPU has a qnode structure in *local* memory. Here, local can
311     mean local memory in NUMA machine or its own cache line that other
312     CPUs are not allowed to cache (i.e., the cache line is in exclusive
313     mode):
314
315     typedef struct qnode {
316         struct qnode* next;
317         bool someoneelse_locked;
318     } qnode;
319
320     typedef qnode* lock; // a lock is a pointer to a qnode
321
322     --The lock itself is literally the *tail* of the list of CPUs holding
323     or waiting for the lock.
324
325     --While waiting, a CPU spins on its local "locked" flag.
326

```

Apr 21, 16 12:03

c23-handout.txt

Page 7/8

```

326
327     Here's the code for acquire:
328
329     // lockp is a qnode**. I points to our local qnode.
330     void acquire(lock* lockp, qnode* I) {
331
332         I->next = NULL;
333         qnode* predecessor;
334
335         // next line makes lockp point to I (that is, it sets *lockp <- I)
336         // and returns the old value of *lockp. Uses atomic operation
337         // XCHG. see earlier in handout (or earlier handouts)
338         // for implementation of xchg_val.
339
340         predecessor = xchg_val(lockp, I); // "A"
341         if (predecessor != NULL) { // queue was non-empty
342             I->someoneelse_locked = true;
343             predecessor->next = I; // "B"
344             while (I->someoneelse_locked); // spin
345         }
346         // we hold the lock!
347     }
348
349     What's going on?
350
351     --If the lock is unlocked, then *lockp == NULL.
352
353     --If the lock is locked, and there are no waiters, then *lockp
354     points to the qnode of the owner
355
356     --If the lock is locked, and there are waiters, then *lockp points
357     to the qnode at the tail of the waiter list.
358
359     --Here's the code for release:
360
361     void release(lock* lockp, qnode* I) {
362         if (!I->next) { // no known successor
363             if (cmpxchg_val(lockp, I, NULL) == I) { // "C"
364                 // swap successful: lockp was pointing to I, so now
365                 // *lockp == NULL, and the lock is unlocked. we can
366                 // go home now.
367                 return;
368             }
369             // if we get here, then there was a timing issue: we had
370             // no known successor when we first checked, but now we
371             // have a successor: some CPU executed the line "A"
372             // above. Wait for that CPU to execute line "B" above.
373             while (!I->next);
374         }
375         // handing the lock off to the next waiter is as simple as
376         // just setting that waiter's "someoneelse_locked" flag to false
377         I->next->someoneelse_locked = false;
378     }
379
380     What's going on?
381
382     --If I->next == NULL and *lockp == I, then no one else is
383     waiting for the lock. So we set *lockp == NULL.
384
385     --If I->next == NULL and *lockp != I, then another CPU is in
386     acquire (specifically, it executed its atomic operation, namely
387     line "A", before we executed ours, namely line "C"). So wait for
388     the other CPU to put the list in a sane state, and then drop
389     down to the next case:
390
391     --If I->next != NULL, then we know that there is a spinning
392     waiter (the oldest one). Hand it the lock by setting its flag to
393     false.
394

```

Apr 21, 16 12:03

c23-handout.txt

Page 8/8

```

395     8. Mutexes
396
397     Motivation: all of the aforementioned locks were called spinlocks
398     because acquire() spins. A mutex avoids busy waiting. Usually, in
399     user space code, you want to be using mutexes, not spinlocks.
400
401     Spinlocks are good for some things, not so great for others. The
402     main problem is that it *busy waits*: it spins, chewing up CPU
403     cycles. Sometimes this is what we want (e.g., if the cost of going
404     to sleep is greater than the cost of spinning for a few cycles
405     waiting for another thread or process to relinquish the spinlock).
406     But sometimes this is not at all what we want (e.g., if the lock
407     would be held for a while: in those cases, the CPU waiting for the
408     lock would waste cycles spinning instead of running some other
409     thread or process).
410
411     With a mutex, if the lock is not available, the locking thread is
412     put to sleep, and tracked by a queue in the mutex.
413
414     struct Mutex {
415         bool is_held; // true if mutex held */
416         thread_id owner; // thread holding mutex, if locked */
417         thread_list waiters; // queue of thread TCBs */
418         Lock wait_lock; // a spinlock, as above */
419     }
420
421     The implementation of mutex_acquire() and mutex_release() would
422     be something like:
423
424     void mutex_acquire(Mutex *m) {
425
426         acquire(&m->wait_lock); // we spin to acquire wait_lock */
427
428         while (m->is_held); // someone else has the mutex */
429
430         m->waiters.insert(current_thread);
431         release(&m->wait_lock);
432
433         /*
434         * NOTE! Right here, mutex_release() could execute. To
435         * avoid "losing the wakeup", we check whether we are
436         * on the scheduler's ready list. If we are, we
437         * shouldn't yield().
438
439         yield_if_we_are_not_ready();
440
441         acquire(&m->wait_lock); // we spin again */
442         m->waiters.remove(current_thread);
443
444         }
445
446         m->is_held = true; // we now hold the mutex */
447         m->owner = self;
448
449         release(&m->wait_lock);
450     }
451
452     void mutex_release(Mutex *m) {
453
454         acquire(&m->wait_lock); // we spin to acquire wait_lock */
455
456         m->is_held = false;
457         m->owner = 0;
458
459         /* tell scheduler to run a waiter */
460         place_a_waiter_on_ready_list(m->waiters);
461
462         release(&m->wait_lock);
463
464     }
465

```