```
1   CS 202, Spring 2020
2   Handout 6 (Class 7)
3
4   Implementation of spinlocks and mutexes
5
6   1. Here is a BROKEN spinlock implementation:
7
8           struct Spinlock {
9             int locked;
10          }
11
12          void acquire(Spinlock *lock) {
13            while (1) {
14              if (lock->locked == 0) { // A
15                lock->locked = 1;      // B
16                break;
17              }
18            }
19          }
20
21          void release (Spinlock *lock) {
22            lock->locked = 0;
23          }
24
25          What's the problem? Two acquire()s on the same lock on different
26          CPUs might both execute line A, and then both execute B. Then
27          both will think they have acquired the lock. Both will proceed.
28          That doesn't provide mutual exclusion.
29
```

```
29
30  2. Correct spinlock implementation
31
32      Relies on atomic hardware instruction. For example, on the x86 (32-bit),
33          doing
34                  "xchg addr, %rax"
35          does the following:
36
37          (i)    freeze all CPUs' memory activity for address addr
38          (ii)   temp <-- *addr
39          (iii)  *addr <-- %rax
40          (iv)   %rax <-- temp
41          (v)    un-freeze memory activity
42
43      /* pseudocode */
44      int xchg_val(addr, value) {
45          %rax = value;
46          xchg (*addr), %rax
47      }
48
49      /* bare-bones version of acquire */
50      void acquire (Spinlock *lock) {
51        pushcli();     /* what does this do? */
52        while (1) {
53          if (xchg_val(&lock->locked, 1) == 0)
54            break;
55        }
56      }
57
58      void release(Spinlock *lock){
59          xchg_val(&lock->locked, 0);
60          popcli();      /* what does this do? */
61      }
62
63
64      /* optimization in acquire; call xchg_val() less frequently */
65      void acquire(Spinlock* lock) {
66          pushcli();
67          while (xchg_val(&lock->locked, 1) == 1) {
68              while (lock->locked) ;
69          }
70      }
71
72      The above is called a *spinlock* because acquire() spins. The
73      bare-bones version is called a "test-and-set (TAS) spinlock"; the
74      other is called a "test-and-test-and-set spinlock".
75
76      The spinlock above is great for some things, not so great for
77      others. The main problem is that it *busy waits*: it spins,
78      chewing up CPU cycles. Sometimes this is what we want (e.g., if
79      the cost of going to sleep is greater than the cost of spinning
80      for a few cycles waiting for another thread or process to
81      relinquish the spinlock). But sometimes this is not at all what we
82      want (e.g., if the lock would be held for a while: in those
83      cases, the CPU waiting for the lock would waste cycles spinning
84      instead of running some other thread or process).
85
86      NOTE: the spinlocks presented here can introduce performance issues
87      when there is a lot of contention. (This happens even if the
88      programmer is using spinlocks correctly.) The performance issues
89      result from cross-talk among CPUs (which undermines caching and
90      generates traffic on the memory bus). If we have time later, we will
91      study a remediation of this issue (search the Web for "MCS locks").
92
93      ANOTHER NOTE: In everyday application-level programming, spinlocks
94      will not be something you use (use mutexes instead). But you should
95      know what these are for technical literacy, and to see where the
96      mutual exclusion is truly enforced on modern hardware.
97
```

```
98   3. Mutex implementation
99
100     The intent of a mutex is to avoid busy waiting: if the lock is not
101     available, the locking thread is put to sleep, and tracked by a
102     queue in the mutex. The next page has an implementation.
103
```

```
1   #include <sys/queue.h>
2
3   typedef struct thread {
4       // ... Entries elided.
5       STAILQ_ENTRY(thread_t) qlink; // Tail queue entry.
6   } thread_t;
7
8   struct Mutex {
9       // Current owner, or 0 when mutex is not held.
10      thread_t *owner;
11
12      // List of threads waiting on mutex
13      STAILQ(thread_t) waiters;
14
15      // A lock protecting the internals of the mutex.
16      Spinlock splock;    // as in item 2 (prev page)
17  };
18
19  void mutex_acquire(struct Mutex *m) {
20
21      acquire(&m->splock);
22
23      // Check if the mutex is held, if not current thread gets mutex and returns
24      if (m->owner == 0) {
25          m->owner = id_of_this_thread;
26          release(&m->splock);
27      } else {
28          // Add thread to waiters.
29          STAILQ_INSERT_TAIL(&m->waiters, id_of_this_thread, qlink);
30
31          // Tell the scheduler to add current thread to the list
32          // of blocked threads. The scheduler needs to be careful
33          // when a corresponding sched_wakeup call is executed to
34          // make sure that it treats running threads correctly.
35          sched_mark_blocked(&id_of_this_thread);
36
37          // Unlock spinlock.
38          release(&m->splock);
39
40          // Stop executing until woken.
41          sched_swtch();
42
43          // When we get to this line, we are guaranteed to hold the mutex. This
44          // is because we can get here only if context-switched-TO, which itself
45          // can happen only if this thread is removed from the waiting queue,
46          // marked "unblocked", and set to be the owner (in mutex_release()
47          // below). However, we might actually have held the mutex at line 39 or
48          // 40 (if we were context-switched out after the spinlock release(),
49          // followed by being run as a result of another thread's release of the
50          // mutex). But if that happens, it just means that we are
51          // context-switched out an "extra" time before proceeding.
52      }
53  }
54
55  void mutex_release(struct Mutex *m) {
56      // Acquire the spinlock in order to make changes.
57      acquire(&m->splock);
58
59      // Assert that the current thread actually owns the mutex
60      assert(m->owner == id_of_this_thread);
61
62      // Check if anyone is waiting.
63      m->owner = STAILQ_GET_HEAD(&m->waiters);
64
65      // If so, wake them up.
66      if (m->owner) {
67          sched_wakeone(&m->owner);
68          STAILQ_REMOVE_HEAD(&m->waiters, qlink);
69      }
70
71      // Release the internal spinlock
72      release(&m->splock);
73  }
```

```
1   4. Simple deadlock example
2
3       T1:
4           acquire(mutexA);
5           acquire(mutexB);
6
7           // do some stuff
8
9           release(mutexB);
10          release(mutexA);
11
12      T2:
13          acquire(mutexB);
14          acquire(mutexA);
15
16          // do some stuff
17
18          release(mutexA);
19          release(mutexB);
20
```

```
21  5. More subtle deadlock example
22
23      Let M be a monitor (shared object with methods protected by mutex)
24      Let N be another monitor
25
26      class M {
27          private:
28              Mutex mutex_m;
29
30              // instance of monitor N
31              N another_monitor;
32
33              // Assumption: no other objects in the system hold a pointer
34              // to our "another_monitor"
35
36          public:
37              M();
38              ~M();
39              void methodA();
40              void methodB();
41      };
42
43      class N {
44          private:
45              Mutex mutex_n;
46              Cond cond_n;
47              int navailable;
48
49          public:
50              N();
51              ~N();
52              void* alloc(int nwanted);
53              void  free(void*);
54      }
55
56      int
57      N::alloc(int nwanted) {
58          acquire(&mutex_n);
59          while (navailable < nwanted) {
60              wait(&cond_n, &mutex_n);
61          }
62
63          // peel off the memory
64
65          navailable -= nwanted;
66          release(&mutex_n);
67      }
68
69      void
70      N::free(void* returning_mem) {
71
72          acquire(&mutex_n);
73
74          // put the memory back
75
76          navailable += returning_mem;
77
78          broadcast(&cond_n, &mutex_n);
79
80          release(&mutex_n);
81      }
82
```

```
83     void
84     M::methodA() {
85
86          acquire(&mutex_m);
87
88          void* new_mem = another_monitor.alloc(int nbytes);
89
90          // do a bunch of stuff using this nice
91          // chunk of memory n allocated for us
92
93          release(&mutex_m);
94     }
95
96     void
97     M::methodB() {
98
99          acquire(&mutex_m);
100
101         // do a bunch of stuff
102
103         another_monitor.free(some_pointer);
104
105         release(&mutex_m);
106     }
107
108    QUESTION: What's the problem?
109
```

```
110  6. Locking brings a performance vs. complexity trade-off
111
112  /*
113   *      linux/mm/filemap.c
114   *
115   * Copyright (C) 1994-1999  Linus Torvalds
116   */
117
118  /*
119   * This file handles the generic file mmap semantics used by
120   * most "normal" filesystems (but you don't /have/ to use this:
121   * the NFS filesystem used to do this differently, for example)
122   */
123  #include <linux/export.h>
124  #include <linux/compiler.h>
125  #include <linux/dax.h>
126  #include <linux/fs.h>
127  #include <linux/sched/signal.h>
128  #include <linux/uaccess.h>
129  #include <linux/capability.h>
130  #include <linux/kernel_stat.h>
131  #include <linux/gfp.h>
132  #include <linux/mm.h>
133  #include <linux/swap.h>
134  #include <linux/mman.h>
135  #include <linux/pagemap.h>
136  #include <linux/file.h>
137  #include <linux/uio.h>
138  #include <linux/hash.h>
139  #include <linux/writeback.h>
140  #include <linux/backing-dev.h>
141  #include <linux/pagevec.h>
142  #include <linux/blkdev.h>
143  #include <linux/security.h>
144  #include <linux/cpuset.h>
145  #include <linux/hugetlb.h>
146  #include <linux/memcontrol.h>
147  #include <linux/cleancache.h>
148  #include <linux/shmem_fs.h>
149  #include <linux/rmap.h>
150  #include "internal.h"
151
152  #define CREATE_TRACE_POINTS
153  #include <trace/events/filemap.h>
154
155  /*
156   * FIXME: remove all knowledge of the buffer layer from the core VM
157   */
158  #include <linux/buffer_head.h> /* for try_to_free_buffers */
159
160  #include <asm/mman.h>
161
162  /*
163   * Shared mappings implemented 30.11.1994. It's not fully working yet,
164   * though.
165   *
166   * Shared mappings now work. 15.8.1995  Bruno.
167   *
168   * finished 'unifying' the page and buffer cache and SMP-threaded the
169   * page-cache, 21.05.1999, Ingo Molnar <mingo@redhat.com>
170   *
171   * SMP-threaded pagemap-LRU 1999, Andrea Arcangeli <andrea@suse.de>
172   */
173
174  /*
175   * Lock ordering:
176   *
177   *  ->i_mmap_rwsem              (truncate_pagecache)
178   *    ->private_lock            (__free_pte->__set_page_dirty_buffers)
179   *      ->swap_lock             (exclusive_swap_page, others)
180   *        ->i_pages lock
181   *
182   *  ->i_mutex
```

```
183    *    ->i_mmap_rwsem              (truncate->unmap_mapping_range)
184    *
185    *   ->mmap_sem
186    *     ->i_mmap_rwsem
187    *       ->page_table_lock or pte_lock   (various, mainly in memory.c)
188    *         ->i_pages lock         (arch-dependent flush_dcache_mmap_lock)
189    *
190    *   ->mmap_sem
191    *     ->lock_page                (access_process_vm)
192    *
193    *   ->i_mutex                    (generic_perform_write)
194    *     ->mmap_sem                 (fault_in_pages_readable->do_page_fault)
195    *
196    *   bdi->wb.list_lock
197    *     sb_lock                    (fs/fs-writeback.c)
198    *     ->i_pages lock             (__sync_single_inode)
199    *
200    *   ->i_mmap_rwsem
201    *     ->anon_vma.lock            (vma_adjust)
202    *
203    *   ->anon_vma.lock
204    *     ->page_table_lock or pte_lock     (anon_vma_prepare and various)
205    *
206    *   ->page_table_lock or pte_lock
207    *     ->swap_lock               (try_to_unmap_one)
208    *     ->private_lock            (try_to_unmap_one)
209    *     ->i_pages lock            (try_to_unmap_one)
210    *     ->zone_lru_lock(zone)     (follow_page->mark_page_accessed)
211    *     ->zone_lru_lock(zone)     (check_pte_range->isolate_lru_page)
212    *     ->private_lock           (page_remove_rmap->set_page_dirty)
213    *     ->i_pages lock           (page_remove_rmap->set_page_dirty)
214    *     bdi.wb->list_lock        (page_remove_rmap->set_page_dirty)
215    *     ->inode->i_lock          (page_remove_rmap->set_page_dirty)
216    *     ->memcg->move_lock       (page_remove_rmap->lock_page_memcg)
217    *     bdi.wb->list_lock        (zap_pte_range->set_page_dirty)
218    *     ->inode->i_lock          (zap_pte_range->set_page_dirty)
219    *     ->private_lock           (zap_pte_range->__set_page_dirty_buffers)
220    *
221    * ->i_mmap_rwsem
222    *   ->tasklist_lock            (memory_failure, collect_procs_ao)
223    */
224
225  static int page_cache_tree_insert(struct address_space *mapping,
226                                     struct page *page, void **shadowp)
227  {
228          struct radix_tree_node *node;
229          .....
230
231  [the point is: fine-grained locking leads to complexity.]
232
```

```
233  7. Cautionary tale
234
235  Consider the code below:
236
237      struct foo {
238          int abc;
239          int def;
240      };
241      static int ready = 0;
242      static mutex_t mutex;
243      static struct foo* ptr = 0;
244
245      void
246      doublecheck_alloc()
247      {
248          if (!ready) { /* <-- accesses shared variable w/out holding mutex */
249
250              mutex_acquire(&mutex);
251              if (!ready) {
252                  ptr = alloc_foo(); /* <-- sets ptr to be non-zero */
253                  ready = 1;
254              }
255
256              mutex_release(&mutex);
257
258          }
259          return;
260      }
261
262  This is an example of the so-called "double-checked locking pattern."
263  The programmer's intent is to avoid a mutex acquistion in the common
264  case that 'ptr' is already initialized.  So the programmer checks a flag
265  called 'ready' before deciding whether to acquire the mutex and
266  initialize 'ptr'. The intended use of doublecheck_alloc() is something
267  like this:
268
269      void f() {
270          doublecheck_alloc();
271          ptr->abc = 5;
272      }
273
274      void g() {
275          doublecheck_alloc();
276          ptr->def = 6;
277      }
278
279  We assume here that mutex_acquire() and mutex_release() are implemented
280  correctly (each contains memory barriers internally, etc.). Furthermore,
281  we assume that the compiler does not reorder instructions.
282
283  NEVERTHELESS, on multi-CPU machines that do not offer sequential
284  consistency, doublecheck_alloc() is broken. What is the bug?
285
286  -----------------------
287
288  Unfortunately, double-checked initialization (or double-checked locking
289  as it's sometimes known) is a common coding pattern. Even some
290  references on threads suggest it! Still, it's broken.
291
292  While you can fix it (in C) by adding another barrier (exercise:
293  where?), this is not recommended, as the code is tricky to reason about.
294  One of the points of this example is to show you why it's so important
295  to protect global data with a mutex, even if "all" one is doing is
296  reading memory, and even if the shortcut looks harmless.
297
```

```
298  Finally, here are some references on this topic:
299
300      --http://www.aristeia.com/Papers/DDJ_Jul_Aug_2004_revised.pdf
301       explores issues with this pattern in C++
302
303      --The "Double-Checked Locking is Broken" Declaration:
304      http://www.cs.umd.edu/~pugh/java/memoryModel/DoubleCheckedLocking.html
305
306      --C++11 provides a way to implement the pattern correctly and
307      portably (again, using memory barriers):
308      https://preshing.com/20130930/double-checked-locking-is-fixed-in-cpp11/
```