```
1   Handout for CS 372H
2   Class 12
3   24 February 2011
4
5   1. CAS / CMPXCHG
6
7      Useful operation: compare-and-swap, known as CAS. Says: "atomically
8      check whether a given memory cell contains a given value, and if it
9      does, then replace the contents of the memory cell with this other
10     value; in either case, return the original value in the memory
11     location".
12
13     On the X86, we implement CAS with the CMPXCHG instruction, but note
14     that this instruction is not atomic by default, so we need the LOCK
15     prefix.
16
17     Here's pseudocode:
18
19         int cmpxchg_val(int* addr, int oldval, int newval) {
20             LOCK: // remember, this is pseudocode
21             int was = *addr;
22             if (*addr == oldval)
23                 *addr = newval;
24             return was;
25         }
26
27     Here's inline assembly:
28
29         uint32_t cmpxchg_val(uint32_t* addr, uint32_t oldval, uint32_t newval) {
30             uint32_t was;
31             asm volatile("lock cmpxchg %3, %0"
32                             : "+m" (*addr), "=a" (was)
33                             : "a" (oldval), "r" (newval)
34                             : "cc");
35             return was;
36         }
37
38  2. MCS locks
39
40     Citation: Mellor-Crummey, J. M. and M. L. Scott.  Algorithms for
41     Scalable Synchronization on Shared-Memory Multiprocessors, ACM
42     Transactions on Computer Systems, Vol. 9, No.  1, February, 1991,
43     pp.21-65.
44
45     Each CPU has a qnode structure in *local* memory. Here, local can
46     mean local memory in NUMA machine or its own cache line that other
47     CPUs are not allowed to cache (i.e., the cache line is in exclusive
48     mode):
49
50     typedef struct qnode {
51         struct qnode* next;
52         bool someoneelse_locked;
53     } qnode;
54
55     typedef qnode* lock;  // a lock is a pointer to a qnode
56
57     --The lock itself is literally the *tail* of the list of CPUs holding
58     or waiting for the lock.
59
60     --While waiting, a CPU spins on its local "locked" flag. Here's the
61     code for acquire:
62
```

```
63         // lockp is a qnode**. I points to our local qnode.
64         void acquire(lock* lockp, qnode* I) {
65
66             I->next = NULL;
67             qnode* predecessor;
68
69             // next line makes lockp point to I (that is, it sets *lockp <-- I)
70             // and returns the old value of *lockp. Uses atomic operation
71             // XCHG. see l09 handout for implementation of xchg_val.
72
73             predecessor = xchg_val(lockp, I);    // "A"
74             if (predecessor != NULL) { // queue was non-empty
75                 I->someoneelse_locked = true;
76                 predecessor->next = I;           // "B"
77                 while (I->someoneelse_locked) ;    // spin
78             }
79             // we hold the lock!
80         }
81
82     What's going on?
83
84     --If the lock is unlocked, then *lockp == NULL.
85
86     --If the lock is locked, and there are no waiters, then *lockp
87     points to the qnode of the owner
88
89     --If the lock is locked, and there are waiters, then *lockp points
90     to the qnode at the tail of the waiter list.
91
92  --Here's the code for release:
93
94     void release(lock* lockp, qnode* I) {
95         if (!I->next)   { // no known successor
96             if (cmpxchg_val(lockp, I, NULL) == I) {      // "C"
97                 // swap successful: lockp was pointing to I, so now
98                 // *lockp == NULL, and the lock is unlocked. we can
99                 // go home now.
100                return;
101            }
102            // if we get here, then there was a timing issue: we had
103            // no known successor when we first checked, but now we
104            // have a successor: some CPU executed the line "A"
105            // above. Wait for that CPU to execute line "B" above.
106            while (!I->next) ;
107        }
108        // handing the lock off to the next waiter is as simple as
109        // just setting that waiter's "someoneelse_locked" flag to false
110        I->next->someoneelse_locked = false;
111    }
112
113    What's going on?
114
115    --If I->next == NULL and *lockp == I, then no one else is
116    waiting for the lock. So we set *lockp == NULL.
117
118    --If I->next == NULL and *lockp != I, then another CPU is in
119    acquire (specifically, it executed its atomic operation, namely
120    line "A", before we executed ours, namely line "C"). So wait for
121    the other CPU to put the list in a sane state, and then drop
122    down to the next case:
123
124    --If I->next != NULL, then we know that there is a spinning
125    waiter (the oldest one). Hand it the lock by setting its flag to
126    false.
127
```

```
128  3. Some examples related to sequential consistency
129
130      [From S.V. Adve and K. Gharachorloo, IEEE Computer, December 1996,
131      66-76. http://rsim.cs.uiuc.edu/~sadve/Publications/computer96.pdf]
132
133      a. What might p2 return if run concurrently with p1?
134
135          int data = 0, ready = 0;
136
137          void p1 () {
138              data = 2000;
139              ready = 1;
140          }
141          int p2 () {
142              while (!ready) {}
143              return data;
144          }
145
146          [answer depends on the memory model given *by* the hardware *to*
147          the software. if the model is sequential consistency, then the
148          code does what you expect. but if not, then p2 can return 0.]
149
150      b. Can both "critical sections" run?
151
152          int flag1 = 0, flag2 = 0;
153
154          int main () {
155              tid id = thread_create (p1, NULL);
156              p2 (); thread_join (id);
157          }
158
159          void p1 (void *ignored) {
160              flag1 = 1;
161              if (!flag2) {
162                  critical_section_1 ();
163              }
164          }
165
166          void p2 (void *ignored) {
167              flag2 = 1;
168              if (!flag1) {
169                  critical_section_2 ();
170              }
171          }
172
173          [answer again depends on the memory model. if there's no
174          sequential consistency, both "critical sections" can run.]
175
176      c. If a processor can read its own writes early, then both functions
177      below can return 2:
178
179          /*
180           * keyword "register" tells compiler to place the variable in a
181           * register, not on the stack.
182           */
183
184          int flag1 = 0, flag2 = 0;
185
186          int p1 (void *ignored)      int p2 (void *ignored)
187          {                           {
188            register int f, g;          register int f, g;
189            flag1 = 1;                  flag2 = 1;
190            f = flag1;                  f = flag2;
191            g = flag2;                  g = flag1;
192            return 2*f + g;             return 2*f + g;
193          }                           }
194
195      The point of these examples: one processor's writes may not
196      show up in program order to another processor. Which means: if
197      you're using synchronization primitives, don't access shared data
198      outside of the mutex. If you're implementing synchronization
199      primitives, read the processor and compiler manuals carefully.
200
```

```
201  4. Performance v complexity trade-off with locks
202
203  /*
204   *      linux/mm/filemap.c
205   *
206   * Copyright (C) 1994-1999  Linus Torvalds
207   */
208
209  /*
210   * This file handles the generic file mmap semantics used by
211   * most "normal" filesystems (but you don't /have/ to use this:
212   * the NFS filesystem used to do this differently, for example)
213   */
214  #include <linux/config.h>
215  #include <linux/module.h>
216  #include <linux/slab.h>
217  #include <linux/compiler.h>
218  #include <linux/fs.h>
219  #include <linux/aio.h>
220  #include <linux/capability.h>
221  #include <linux/kernel_stat.h>
222  #include <linux/mm.h>
223  #include <linux/swap.h>
224  #include <linux/mman.h>
225  #include <linux/pagemap.h>
226  #include <linux/file.h>
227  #include <linux/uio.h>
228  #include <linux/hash.h>
229  #include <linux/writeback.h>
230  #include <linux/pagevec.h>
231  #include <linux/blkdev.h>
232  #include <linux/security.h>
233  #include <linux/syscalls.h>
234  #include "filemap.h"
235  /*
236   * FIXME: remove all knowledge of the buffer layer from the core VM
237   */
238  #include <linux/buffer_head.h> /* for generic_osync_inode */
239
240  #include <asm/uaccess.h>
241  #include <asm/mman.h>
242
243  static ssize_t
244  generic_file_direct_IO(int rw, struct kiocb *iocb, const struct iovec *iov,
245          loff_t offset, unsigned long nr_segs);
246
247  /*
248   * Shared mappings implemented 30.11.1994. It's not fully working yet,
249   * though.
250   *
251   * Shared mappings now work. 15.8.1995  Bruno.
252   *
253   * finished 'unifying' the page and buffer cache and SMP-threaded the
254   * page-cache, 21.05.1999, Ingo Molnar <mingo@redhat.com>
255   *
256   * SMP-threaded pagemap-LRU 1999, Andrea Arcangeli <andrea@suse.de>
257   */
258
259  /*
260   * Lock ordering:
261   *
262   *  ->i_mmap_lock              (vmtruncate)
263   *    ->private_lock           (__free_pte->__set_page_dirty_buffers)
264   *      ->swap_lock            (exclusive_swap_page, others)
265   *        ->mapping->tree_lock
266   *
267   *  ->i_mutex
268   *    ->i_mmap_lock            (truncate->unmap_mapping_range)
269   *
270   *  ->mmap_sem
271   *    ->i_mmap_lock
272   *      ->page_table_lock or pte_lock   (various, mainly in memory.c)
273   *        ->mapping->tree_lock (arch-dependent flush_dcache_mmap_lock)
```

```
274  *
275  *  −>mmap_sem
276  *    −>lock_page               (access_process_vm)
277  *
278  *  −>mmap_sem
279  *    −>i_mutex                 (msync)
280  *
281  *  −>i_mutex
282  *    −>i_alloc_sem             (various)
283  *
284  *  −>inode_lock
285  *    −>sb_lock                 (fs/fs−writeback.c)
286  *    −>mapping−>tree_lock      (__sync_single_inode)
287  *
288  *  −>i_mmap_lock
289  *    −>anon_vma.lock           (vma_adjust)
290  *
291  *  −>anon_vma.lock
292  *    −>page_table_lock or pte_lock    (anon_vma_prepare and various)
293  *
294  *  −>page_table_lock or pte_lock
295  *    −>swap_lock               (try_to_unmap_one)
296  *    −>private_lock            (try_to_unmap_one)
297  *    −>tree_lock               (try_to_unmap_one)
298  *    −>zone.lru_lock           (follow_page−>mark_page_accessed)
299  *    −>zone.lru_lock           (check_pte_range−>isolate_lru_page)
300  *    −>private_lock            (page_remove_rmap−>set_page_dirty)
301  *    −>tree_lock               (page_remove_rmap−>set_page_dirty)
302  *    −>inode_lock              (page_remove_rmap−>set_page_dirty)
303  *    −>inode_lock              (zap_pte_range−>set_page_dirty)
304  *    −>private_lock            (zap_pte_range−>__set_page_dirty_buffers)
305  *
306  *  −>task−>proc_lock
307  *    −>dcache_lock             (proc_pid_lookup)
308  */
309
310  /*
311   * Remove a page from the page cache and free it. Caller has to make
312   * sure the page is locked and that nobody else uses it − or that usage
313   * is safe.  The caller must hold a write_lock on the mapping's tree_lock.
314   */
315  void __remove_from_page_cache(struct page *page)
316  {
317          struct address_space *mapping = page−>mapping;
318
319  .............
320
321  [point of this item on the handout: fine−grained locking leads to complexity]
```
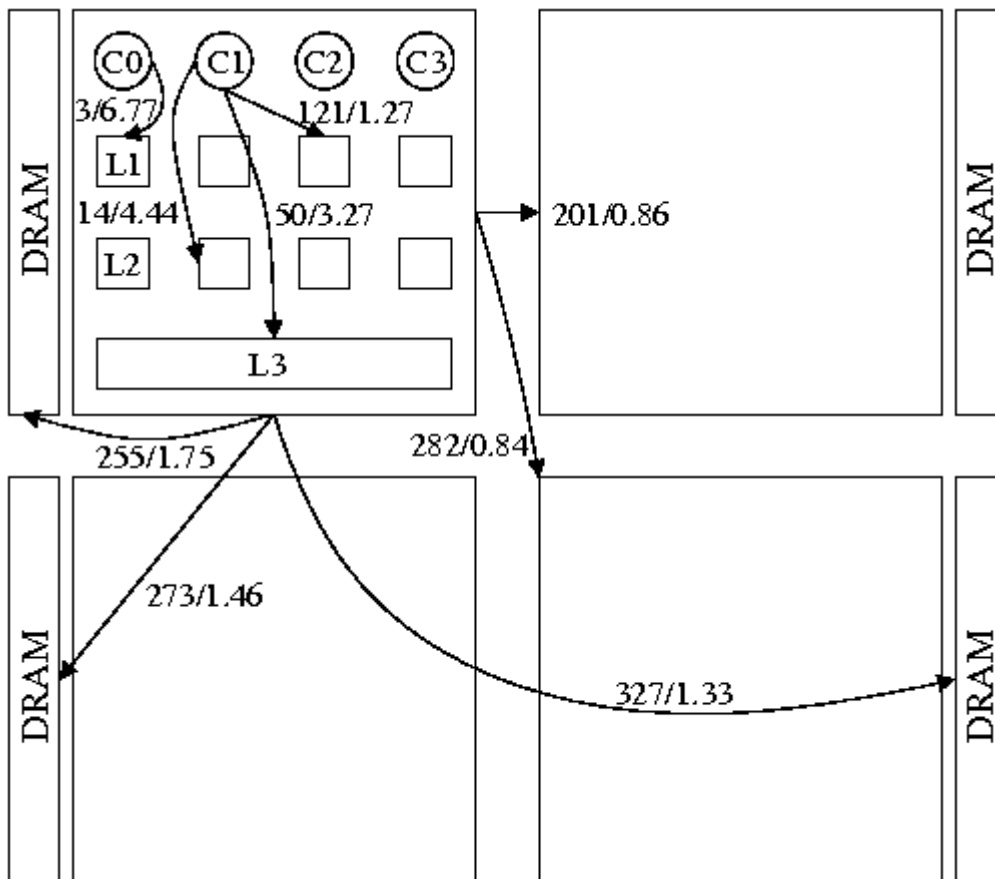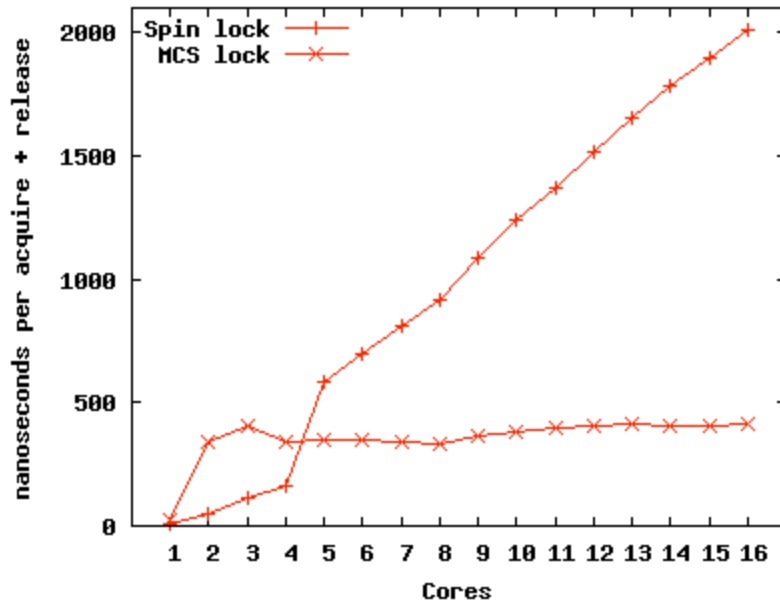
C0  C1  C2  C3

3/6.77

L1

121/1.27

14/4.44    50/3.27

L2

L3

DRAM

201/0.86

DRAM

255/1.75    282/0.84

273/1.46

DRAM    327/1.33    DRAM

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.

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.]