

Feb 26, 10 8:50

l12-handout.txt

Page 1/7

```

1 Handout for CS 372H
2 Class 12
3 25 February 2010

```

## 1. CAS / CMPXCHG

Useful operation: compare-and-swap, known as CAS. Says: "atomically check whether a given memory cell contains a given value, and if it does, then replace the contents of the memory cell with this other value; in either case, return the original value in the memory location".

On the X86, we implement CAS with the CMPXCHG instruction, but note that this instruction is not atomic by default, so we need the LOCK prefix.

Here's pseudocode:

```

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     }

```

Here's inline assembly:

```

28     uint32_t cmpxchg_val(uint32_t* addr, uint32_t oldval, uint32_t newval) {
29         uint32_t was;
30         asm volatile("lock cmpxchg %3, %0"
31                     : "+m" (*addr), "=a" (was)
32                     : "a" (oldval), "r" (newval)
33                     : "cc");
34         return was;
35     }

```

## 2. MCS locks

Citation: Mellor-Crummey, J. M. and M. L. Scott. Algorithms for Scalable Synchronization on Shared-Memory Multiprocessors, ACM Transactions on Computer Systems, Vol. 9, No. 1, February, 1991, pp.21-65.

Each CPU has a qnode structure in \*local\* memory. Here, local can mean local memory in NUMA machine or its own cache line that other CPUs are not allowed to cache (i.e., the cache line is in exclusive mode):

```

50     typedef struct qnode {
51         struct qnode* next;
52         bool someoneelse_locked;
53     } qnode;

```

```

54     typedef qnode* lock; // a lock is a pointer to a qnode

```

--The lock itself is literally the tail of the list of CPUs holding or waiting for the lock.

--While waiting, a CPU spins on its local "locked" flag. Here's the code for acquire:

Feb 26, 10 8:50

l12-handout.txt

Page 2/7

```

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 109 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!

```

What's going on?

--If the lock is unlocked, then \*lockp == NULL.

--If the lock is locked, and there are no waiters, then \*lockp points to the qnode of the owner

--If the lock is locked, and there are waiters, then \*lockp points to the qnode at the tail of the waiter list.

--Here's the code for release:

```

93     void release(lock* lockp, qnode* I) {
94         if (!I->next) { // no known successor
95             if (cmpxchg_val(lockp, I, NULL) == I) { // "C"
96                 // swap successful: lockp was pointing to I, so now
97                 // *lockp == NULL, and the lock is unlocked. we can
98                 // go home now.
99                 return;
100             }
101             // if we get here, then there was a timing issue: we had
102             // no known successor when we first checked, but now we
103             // have a successor: some CPU executed the line "A"
104             // above. Wait for that CPU to execute line "B" above.
105             while (!I->next) ;
106         }
107         // handing the lock off to the next waiter is as simple as
108         // just setting that waiter's "someoneelse_locked" flag to false
109         I->next->someoneelse_locked = false;
110     }

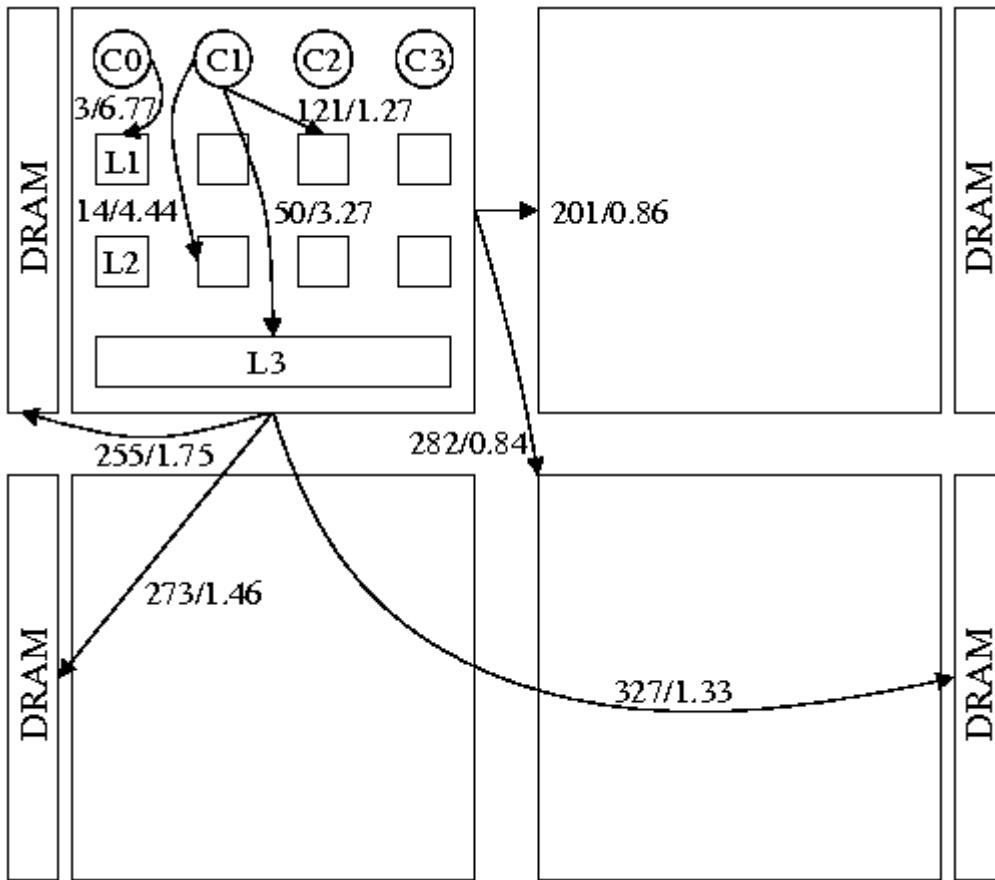
```

What's going on?

--If I->next == NULL and \*lockp == I, then no one else is waiting for the lock. So we set \*lockp == NULL.

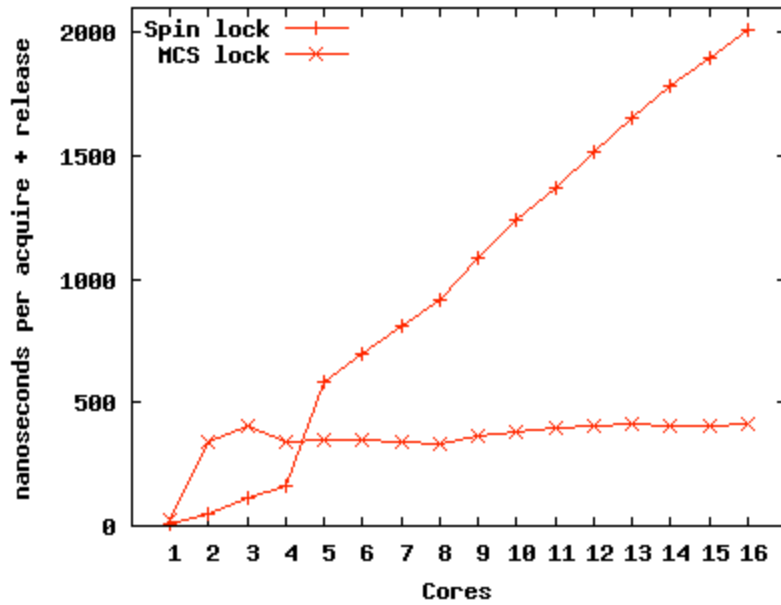
--If I->next != NULL and \*lockp != I, then another CPU is in acquire (specifically, it executed its atomic operation, namely line "A", before we executed ours, namely line "C"). So wait for the other CPU to put the list in a sane state, and then drop down to the next case:

--If I->next != NULL, then we know that there is a spinning waiter (the oldest one). Hand it the lock by setting its flag to false.



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.]



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 26, 10 8:50

l12-handout.txt

Page 3/7

```

128 3. Simple deadlock example
129
130 T1:
131     acquire(mutexA);
132     acquire(mutexB);
133
134     // do some stuff
135
136     release(mutexB);
137     release(mutexA);
138
139 T2:
140     acquire(mutexB);
141     acquire(mutexA);
142
143     // do some stuff
144
145     release(mutexA);
146     release(mutexB);
147

```

Feb 26, 10 8:50

l12-handout.txt

Page 4/7

```

148 4. More subtle deadlock example
149
150 Let M be a monitor (shared object with methods protected by mutex)
151 Let N be another monitor
152
153 class M {
154     private:
155         Mutex mutex_m;
156
157         // instance of monitor N
158         N another_monitor;
159
160         // Assumption: no other objects in the system hold a pointer
161         // to our "another_monitor"
162
163     public:
164         M();
165         ~M();
166         void methodA();
167         void methodB();
168 };
169
170 class N {
171     private:
172         Mutex mutex_n;
173         Cond cond_n;
174         int navailable;
175
176     public:
177         N();
178         ~N();
179         void* alloc(int nwanted);
180         void free(void*);
181 }
182
183 int
184 N::alloc(int nwanted) {
185     acquire(&mutex_n);
186     while (navailable < nwanted) {
187         wait(&cond_n, &mutex_n);
188     }
189
190     // peel off the memory
191
192     navailable -= nwanted;
193     release(&mutex_n);
194 }
195
196 void
197 N::free(void* returning_mem) {
198     acquire(&mutex_n);
199
200     // put the memory back
201
202     navailable += returning_mem;
203
204     broadcast(&cond_n, &mutex_n);
205
206     release(&mutex_n);
207 }
208
209 void
210 M::methodA() {
211
212     acquire(&mutex_m);
213
214     void* new_mem = another_monitor.alloc(int nbytes);
215
216     // do a bunch of stuff using this nice
217     // chunk of memory n allocated for us
218
219     release(&mutex_m);
220

```

Feb 26, 10 8:50

l12-handout.txt

Page 5/7

```

221     }
222
223     void
224     M::methodB() {
225
226         acquire(&mutex_m);
227
228         // do a bunch of stuff
229
230         another_monitor.free(some_pointer);
231
232         release(&mutex_m);
233     }
234
235     QUESTION: What's the problem?
236

```

Feb 26, 10 8:50

l12-handout.txt

Page 6/7

```

237 5. Performance v complexity trade-off with locks
238
239 /*
240  *      linux/mm/filemap.c
241  *
242  * Copyright (C) 1994-1999 Linus Torvalds
243  */
244
245 /*
246  * This file handles the generic file mmap semantics used by
247  * most "normal" filesystems (but you don't /have/ to use this:
248  * the NFS filesystem used to do this differently, for example)
249  */
250 #include <linux/config.h>
251 #include <linux/module.h>
252 #include <linux/slab.h>
253 #include <linux/compiler.h>
254 #include <linux/fs.h>
255 #include <linux/aio.h>
256 #include <linux/capability.h>
257 #include <linux/kernel_stat.h>
258 #include <linux/mm.h>
259 #include <linux/swap.h>
260 #include <linux/mman.h>
261 #include <linux/pagemap.h>
262 #include <linux/file.h>
263 #include <linux/uio.h>
264 #include <linux/hash.h>
265 #include <linux/writeback.h>
266 #include <linux/pagevec.h>
267 #include <linux/blkdev.h>
268 #include <linux/security.h>
269 #include <linux/syscalls.h>
270 #include "filemap.h"
271 /*
272  * FIXME: remove all knowledge of the buffer layer from the core VM
273  */
274 #include <linux/buffer_head.h> /* for generic_osync_inode */
275
276 #include <asm/uaccess.h>
277 #include <asm/mman.h>
278
279 static ssize_t
280 generic_file_direct_IO(int rw, struct kiocb *iocb, const struct iovec *iov,
281                       loff_t offset, unsigned long nr_segs);
282
283 /*
284  * Shared mappings implemented 30.11.1994. It's not fully working yet,
285  * though.
286  *
287  * Shared mappings now work. 15.8.1995 Bruno.
288  *
289  * finished 'unifying' the page and buffer cache and SMP-threaded the
290  * page-cache, 21.05.1999, Ingo Molnar <mingo@redhat.com>
291  *
292  * SMP-threaded pagemap-LRU 1999, Andrea Arcangeli <andrea@suse.de>
293  */
294
295 /*
296  * Lock ordering:
297  *
298  * ->i_mmap_lock                (vmtruncate)
299  * ->private_lock              (__free_pte->__set_page_dirty_buffers)
300  * ->swap_lock                  (exclusive_swap_page, others)
301  * ->mapping->tree_lock
302  *
303  * ->i_mutex
304  * ->i_mmap_lock                (truncate->unmap_mapping_range)
305  *
306  * ->mmap_sem
307  * ->i_mmap_lock
308  * ->page_table_lock or pte_lock (various, mainly in memory.c)
309  * ->mapping->tree_lock (arch-dependent flush_dcache_mmap_lock)

```

Feb 26, 10 8:50

l12-handout.txt

Page 7/7

```

310 *
311 * ->mmmap_sem
312 * ->lock_page (access_process_vm)
313 *
314 * ->mmmap_sem
315 * ->i_mutex (msync)
316 *
317 * ->i_mutex
318 * ->i_alloc_sem (various)
319 *
320 * ->inode_lock
321 * ->sb_lock (fs/fs-writeback.c)
322 * ->mapping->tree_lock (__sync_single_inode)
323 *
324 * ->i_mmap_lock
325 * ->anon_vma.lock (vma_adjust)
326 *
327 * ->anon_vma.lock
328 * ->page_table_lock or pte_lock (anon_vma_prepare and various)
329 *
330 * ->page_table_lock or pte_lock
331 * ->swap_lock (try_to_unmap_one)
332 * ->private_lock (try_to_unmap_one)
333 * ->tree_lock (try_to_unmap_one)
334 * ->zone.lru_lock (follow_page->mark_page_accessed)
335 * ->zone.lru_lock (check_pte_range->isolate_lru_page)
336 * ->private_lock (page_remove_rmap->set_page_dirty)
337 * ->tree_lock (page_remove_rmap->set_page_dirty)
338 * ->inode_lock (page_remove_rmap->set_page_dirty)
339 * ->inode_lock (zap_pte_range->set_page_dirty)
340 * ->private_lock (zap_pte_range->__set_page_dirty_buffers)
341 *
342 * ->task->proc_lock
343 * ->dcache_lock (proc_pid_lookup)
344 */
345
346 /*
347 * Remove a page from the page cache and free it. Caller has to make
348 * sure the page is locked and that nobody else uses it - or that usage
349 * is safe. The caller must hold a write_lock on the mapping's tree_lock.
350 */
351 void __remove_from_page_cache(struct page *page)
352 {
353     struct address_space *mapping = page->mapping;
354
355     .....
356 [point of this item on the handout: fine-grained locking leads to complexity]

```