```
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
```
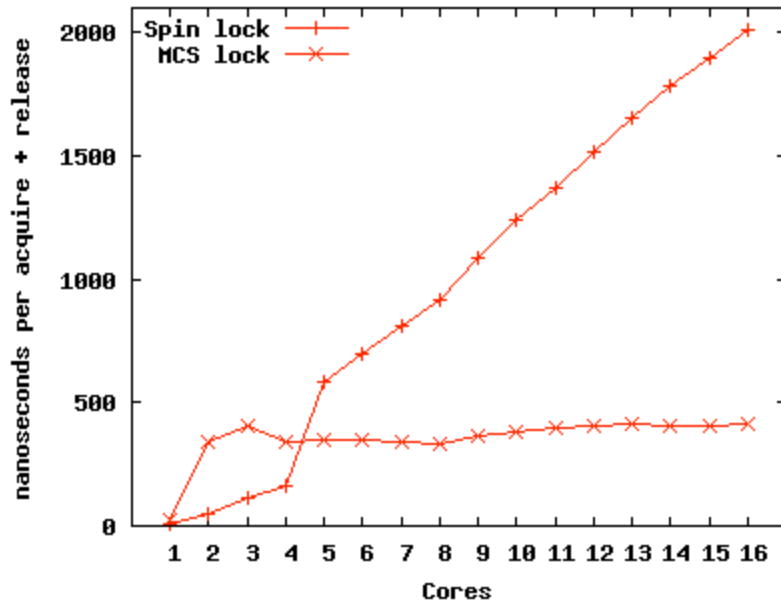
```
60          // lockp is a qnode**. I points to our local qnode.
61          void acquire(lock* lockp, qnode* I) {
62
63              I->next = NULL;
64              qnode* predecessor;
65
66              // next line makes lockp point to I (that is, it sets *lockp <-- I)
67              // and returns 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
75                  predecessor->next = I;          // "B"
76                  while (I->someoneelse_locked) ;     // spin
77              }
78              // we hold the lock!
79          }
80
81      What's going on?
82
83      --If the lock is unlocked, then *lockp == NULL.
84
85      --If the lock is locked, and there are no waiters, then *lockp
86      points to the qnode of the owner
87
88      --If the lock is locked, and there are waiters, then *lockp points
89      to the qnode at the tail of the waiter list.
90
91  --Here's the code for release:
92
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
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.
```

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.]

```
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)
```

```
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]
```

```
1   1. 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
22
23
24
25
26  2. More subtle deadlock example
27
28
29      Let M be a monitor (shared object with methods protected by mutex)
30      Let N be another monitor
31
32      class M {
33          private:
34              Mutex mutex_m;
35
36              // instance of monitor N
37              N another_monitor;
38
39              // Assumption: no other objects in the system hold a pointer
40              // to our "another_monitor"
41
42          public:
43              M();
44              ~M();
45              void methodA();
46              void methodB();
47      };
48
49      class N {
50          private:
51              Mutex mutex_n;
52              Cond cond_n;
53              int navailable;
54
55          public:
56              N();
57              ~N();
58              void* alloc(int nwanted);
59              void  free(void*);
60      }
61
62      int
63      N::alloc(int nwanted) {
64          acquire(&mutex_n);
65          while (navailable < nwanted) {
66              wait(&cond_n, &mutex_n);
67          }
68
69          // peel off the memory
70
71          navailable -= nwanted;
72          release(&mutex_n);
73      }
```

```
74
75      void
76      N::free(void* returning_mem) {
77
78          acquire(&mutex_n);
79
80          // put the memory back
81
82          navailable += returning_mem;
83
84          broadcast(&cond_n, &mutex_n);
85
86          release(&mutex_n);
87      }
88
89      void
90      M::methodA() {
91
92          acquire(&mutex_m);
93
94          void* new_mem = another_monitor.alloc(int nbytes);
95
96          // do a bunch of stuff using this nice
97          // chunk of memory n allocated for us
98
99          release(&mutex_m);
100     }
101
102     void
103     M::methodB() {
104
105         acquire(&mutex_m);
106
107         // do a bunch of stuff
108
109         another_monitor.free(some_pointer);
110
111         release(&mutex_m);
112     }
113
114     QUESTION: What's the problem?
115
116
```