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.]

```
1   Handout for CS 372H
2   Class 18
3   21 March 2013
4
5   1. Recall implementation of acquire() and release() in spinlocks
6   context:
7
8      [this item is fully review.]
9
10     It uses an atomic instruction on the CPU. For example, on the
11     x86, doing
12            "xchg addr, %eax"
13     does the following:
14
15     (i)   freeze all CPUs' memory activity for address addr
16     (ii)  temp = *addr
17     (iii) *addr = %eax
18     (iv)  %eax = temp
19     (v)   un-freeze memory activity
20
21     /* pseudocode */
22     int xchg_val(addr, value) {
23         %eax = value;
24         xchg (*addr), %eax
25     }
26
27     struct Lock {
28       int locked;
29     }
30
31     /* bare-bones version of acquire */
32     void acquire (Lock *lock) {
33       pushcli();    /* what does this do? */
34       while (1) {
35         if (xchg_val(&lock->locked, 1) == 0)
36           break;
37       }
38     }
39
40     /* optimization in acquire; call xchg_val() less frequently */
41     void acquire(Lock* lock) {
42         pushcli();
43         while (xchg_val(&lock->locked, 1) == 1) {
44             while (lock->locked) ;
45         }
46     }
47
48     void release(Lock *lock){
49         xchg_val(&lock->locked, 0);
50         popcli();    /* what does this do? */
51     }
52
53     The above is called a *spinlock* because acquire() spins.
54
```

```
55  2. Here's an alternative.....
56
57     Instead of using the XCHG instruction, it uses CMPXCHG.
58
59  A. CAS / CMPXCHG
60
61     Useful operation: compare-and-swap, known as CAS. Says: "atomically
62     check whether a given memory cell contains a given value, and if it
63     does, then replace the contents of the memory cell with this other
64     value; in either case, return the original value in the memory
65     location".
66
67     On the X86, we implement CAS with the CMPXCHG instruction, but note
68     that this instruction is not atomic by default, so we need the LOCK
69     prefix.
70
71     Here's pseudocode:
72
73         int cmpxchg_val(int* addr, int oldval, int newval) {
74             LOCK: // remember, this is pseudocode
75             int was = *addr;
76             if (*addr == oldval)
77                 *addr = newval;
78             return was;
79         }
80
81     Here's inline assembly:
82
83         uint32_t cmpxchg_val(uint32_t* addr, uint32_t oldval, uint32_t newval) {
84             uint32_t was;
85             asm volatile("lock cmpxchg %3, %0"
86                               : "+m" (*addr), "=a" (was)
87                               : "a" (oldval), "r" (newval)
88                               : "cc");
89             return was;
90         }
91
92  B. MCS locks
93
94     Citation: Mellor-Crummey, J. M. and M. L. Scott.  Algorithms for
95     Scalable Synchronization on Shared-Memory Multiprocessors, ACM
96     Transactions on Computer Systems, Vol. 9, No.  1, February, 1991,
97     pp.21-65.
98
99     Each CPU has a qnode structure in *local* memory. Here, local can
100    mean local memory in NUMA machine or its own cache line that other
101    CPUs are not allowed to cache (i.e., the cache line is in exclusive
102    mode):
103
104    typedef struct qnode {
105        struct qnode* next;
106        bool someoneelse_locked;
107    } qnode;
108
109    typedef qnode* lock;  // a lock is a pointer to a qnode
110
111    --The lock itself is literally the *tail* of the list of CPUs holding
112    or waiting for the lock.
113
114    --While waiting, a CPU spins on its local "locked" flag. Here's the
115    code for acquire:
116
```
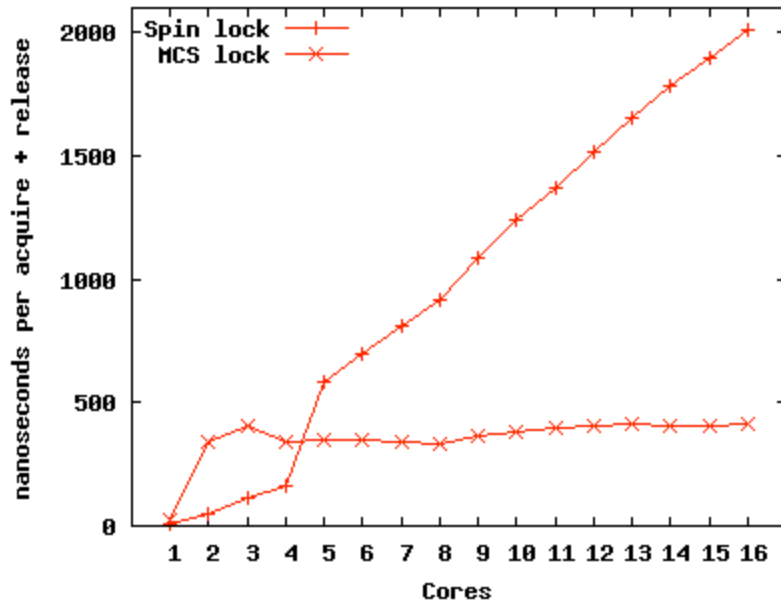
```
117        // lockp is a qnode**. I points to our local qnode.
118        void acquire(lock* lockp, qnode* I) {
119
120            I->next = NULL;
121            qnode* predecessor;
122
123            // next line makes lockp point to I (that is, it sets *lockp <-- I)
124            // and returns the old value of *lockp. Uses atomic operation
125            // XCHG. see earlier in handout (or earlier handouts)
126            // for implementation of xchg_val.
127
128            predecessor = xchg_val(lockp, I);     // "A"
129            if (predecessor != NULL) { // queue was non-empty
130                I->someoneelse_locked = true;
131
132                predecessor->next = I;          // "B"
133                while (I->someoneelse_locked) ;     // spin
134            }
135            // we hold the lock!
136        }
137
138        What's going on?
139
140        --If the lock is unlocked, then *lockp == NULL.
141
142        --If the lock is locked, and there are no waiters, then *lockp
143        points to the qnode of the owner
144
145        --If the lock is locked, and there are waiters, then *lockp points
146        to the qnode at the tail of the waiter list.
147
148    --Here's the code for release:
149
150        void release(lock* lockp, qnode* I) {
151            if (!I->next)   { // no known successor
152                if (cmpxchg_val(lockp, I, NULL) == I) {      // "C"
153                    // swap successful: lockp was pointing to I, so now
154                    // *lockp == NULL, and the lock is unlocked. we can
155                    // go home now.
156                    return;
157                }
158                // if we get here, then there was a timing issue: we had
159                // no known successor when we first checked, but now we
160                // have a successor: some CPU executed the line "A"
161                // above. Wait for that CPU to execute line "B" above.
162                while (!I->next) ;
163            }
164
165            // handing the lock off to the next waiter is as simple as
166            // just setting that waiter's "someoneelse_locked" flag to false
167            I->next->someoneelse_locked = false;
168        }
169
170        What's going on?
171
172        --If I->next == NULL and *lockp == I, then no one else is
173        waiting for the lock. So we set *lockp == NULL.
174
175        --If I->next == NULL and *lockp != I, then another CPU is in
176        acquire (specifically, it executed its atomic operation, namely
177        line "A", before we executed ours, namely line "C"). So wait for
178        the other CPU to put the list in a sane state, and then drop
179        down to the next case:
180
181        --If I->next != NULL, then we know that there is a spinning
182        waiter (the oldest one). Hand it the lock by setting its flag to
183        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.]