# last time

races: inconsistent results with multiple threads

atomic operations $\neq$ instructions

not constructing locks with atomic load/store

lock abstraction:
  acquire/lock: keep anyone else from acquiring/lock
  release/unlock: let someone else acquire/lock
  intended usage: lock before accessing shared thing, unlock after

(started) locks by disabling interrupts

# implementing locks: single core

intuition: context switch only happens on interrupt
    timer expiration, I/O, etc. causes OS to run

solution: disable them
    reenable on unlock

# implementing locks: single core

intuition: context switch only happens on interrupt
    timer expiration, I/O, etc. causes OS to run

solution: disable them
    reenable on unlock

x86 instructions:
    cli — disable interrupts
    sti — enable interrupts

# naive interrupt enable/disable (1)

```
Lock() {                    Unlock() {
    disable interrupts          enable interrupts
}                           }
```

# naive interrupt enable/disable (1)

```
Lock() {                        Unlock() {
    disable interrupts              enable interrupts
}                               }
```

problem: user can hang the system:

```
            Lock(some_lock);
            while (true) {}
```

# naive interrupt enable/disable (1)

```
Lock() {                        Unlock() {
    disable interrupts              enable interrupts
}                               }
```

problem: user can hang the system:

```
Lock(some_lock);
while (true) {}
```

problem: can't do I/O within lock

```
Lock(some_lock);
read from disk
    /* waits forever for (disabled) interrupt
       from disk IO finishing */
```

# naive interrupt enable/disable (2)

```
Lock() {                    Unlock() {
    disable interrupts          enable interrupts
}                           }
```

# naive interrupt enable/disable (2)

```
Lock() {                    Unlock() {
    disable interrupts          enable interrupts
}                           }
```

# naive interrupt enable/disable (2)

```
Lock() {                        Unlock() {
    disable interrupts              enable interrupts
}                               }
```

# naive interrupt enable/disable (2)

```
Lock() {                        Unlock() {
    disable interrupts              enable interrupts
}                               }
```

problem: nested locks

```
        Lock(milk_lock);
        if (no milk) {
            Lock(store_lock);
            buy milk
            Unlock(store_lock);
            /* interrupts enabled here?? */
        }
        Unlock(milk_lock);
```

# xv6 interrupt disabling (1)

```
...
acquire(struct spinlock *lk) {
  pushcli(); // disable interrupts to avoid deadlock
  ... /* this part basically just for multicore */
}
release(struct spinlock *lk)
{
  ... /* this part basically just for multicore */
  popcli();
}
```

# xv6 push/popcli

pushcli / popcli — need to be in pairs

pushcli — disable interrupts if not already

popcli — enable interrupts if corresponding pushcli disabled them
  don't enable them if they were already disabled

# a simple race

```
thread_A:                           thread_B:
    movl $1, x    /* x ← 1 */           movl $1, y    /* y ← 1 */
    movl y, %eax  /* return y */        movl x, %eax  /* return x */
    ret                                 ret

    x = y = 0;
    pthread_create(&A, NULL, thread_A, NULL);
    pthread_create(&B, NULL, thread_B, NULL);
    pthread_join(A, &A_result); pthread_join(B, &B_result);
    printf("A:%d B:%d\n", (int) A_result, (int) B_result);
```

# a simple race

```
thread_A:                           thread_B:
    movl $1, x    /* x ← 1 */           movl $1, y    /* y ← 1 */
    movl y, %eax  /* return y */        movl x, %eax  /* return x */
    ret                                 ret

    x = y = 0;
    pthread_create(&A, NULL, thread_A, NULL);
    pthread_create(&B, NULL, thread_B, NULL);
    pthread_join(A, &A_result); pthread_join(B, &B_result);
    printf("A:%d B:%d\n", (int) A_result, (int) B_result);
```

if loads/stores atomic, then possible results:

    A:1 B:1 — both moves into x and y, then both moves into eax execute

    A:0 B:1 — thread A executes before thread B

    A:1 B:0 — thread B executes before thread A

# a simple race: results

```
thread_A:                              thread_B:
    movl $1, x   /* x ← 1 */               movl $1, y   /* y ← 1 */
    movl y, %eax /* return y */             movl x, %eax /* return x */
    ret                                    ret

    x = y = 0;
    pthread_create(&A, NULL, thread_A, NULL);
    pthread_create(&B, NULL, thread_B, NULL);
    pthread_join(A, &A_result); pthread_join(B, &B_result);
    printf("A:%d B:%d\n", (int) A_result, (int) B_result);
```

my desktop, 100M trials:

| frequency | result | |
|-----------|--------|---|
| 99 823 739 | A:0 B:1 | ('A executes before B') |
| 171 161 | A:1 B:0 | ('B executes before A') |
| 4 706 | A:1 B:1 | ('execute moves into x+y first') |
| 394 | A:0 B:0 | ??? |

# a simple race: results

```
thread_A:                                    thread_B:
    movl $1, x    /* x ← 1 */                    movl $1, y    /* y ← 1 */
    movl y, %eax /* return y */                  movl x, %eax /* return x */
    ret                                          ret


    x = y = 0;
    pthread_create(&A, NULL, thread_A, NULL);
    pthread_create(&B, NULL, thread_B, NULL);
    pthread_join(A, &A_result); pthread_join(B, &B_result);
    printf("A:%d B:%d\n", (int) A_result, (int) B_result);
```

my desktop, 100M trials:

| frequency | result | |
|---|---|---|
| 99 823 739 | A:0 B:1 | ('A executes before B') |
| 171 161 | A:1 B:0 | ('B executes before A') |
| 4 706 | A:1 B:1 | ('execute moves into x+y first') |
| 394 | A:0 B:0 | ??? |

# load/store reordering

load/stores atomic, but run *out of order*

recall?: out-of-order processors

processor optimization: sometimes execute instructions in non-program order
> hide delays from slow caches, variable computation rates, etc.
> documneted limits on when this is/is not allowed

track side-effects *within a thread* to make as if in-order
> but common choice: don't worry as much between cores/threads
> design decision: if programmer cares, they worry about it

want to avoid this *special instructions ensure strict ordering*

# why load/store reordering?

prior example: load of x executing before store of y

why do this? otherwise delay the load
    if x and y unrelated — no benefit to waiting

# compilers changes loads/stores too (1)

```
void Alice() {
    note_from_alice = 1;
    do {} while (note_from_bob);
    if (no_milk) {++milk;}
}
```

```
Alice:
  movl $1, note_from_alice    // note_from_alice ← 1
  movl note_from_bob, %eax    // eax ← note_from_bob
.L2:
  testl %eax, %eax
  jne .L2                     // while (eax == 0) repeat
  cmpl $0, no_milk            // if (no_milk != 0) ...
  ...
```

# compilers changes loads/stores too (1)

```
void Alice() {
    note_from_alice = 1;
    do {} while (note_from_bob);
    if (no_milk) {++milk;}
}
```

```
Alice:
  movl $1, note_from_alice   // note_from_alice ← 1
  movl note_from_bob, %eax   // eax ← note_from_bob
.L2:
  testl %eax, %eax
  jne .L2                    // while (eax == 0) repeat
  cmpl $0, no_milk           // if (no_milk != 0) ...
  ...
```

# compilers changes loads/stores too (2)

```
void Alice() {
    note_from_alice = 1;  // "Alice waiting" signal for Bob()
    do {} while (note_from_bob);
    if (no_milk) {++milk;}
    note_from_alice = 2;
}
```

---

```
Alice:
  // compiler optimization: don't set note_from_alice to 1,
  // (why? it will be set to 2 anyway)
  movl note_from_bob, %eax  // eax ← note_from_bob
.L2:
  testl %eax, %eax
  jne .L2                    // while (eax == 0) repeat
  ...
  movl $2, note_from_alice  // note_from_alice ← 2
```

# compilers changes loads/stores too (2)

```
void Alice() {
    note_from_alice = 1;  // "Alice waiting" signal for Bob()
    do {} while (note_from_bob);
    if (no_milk) {++milk;}
    note_from_alice = 2;
}
```

```
Alice:
  // compiler optimization: don't set note_from_alice to 1,
  // (why? it will be set to 2 anyway)
  movl note_from_bob, %eax  // eax ← note_from_bob
.L2:
  testl %eax, %eax
  jne .L2                   // while (eax == 0) repeat
  ...
  movl $2, note_from_alice  // note_from_alice ← 2
```

# compilers changes loads/stores too (2)

```
void Alice() {
    note_from_alice = 1;   // "Alice waiting" signal for Bob()
    do {} while (note_from_bob);
    if (no_milk) {++milk;}
    note_from_alice = 2;
}
```

---

```
Alice:
  // compiler optimization: don't set note_from_alice to 1,
  // (why? it will be set to 2 anyway)
  movl note_from_bob, %eax   // eax ← note_from_bob
.L2:
  testl %eax, %eax
  jne .L2                    // while (eax == 0) repeat
  ...
  movl $2, note_from_alice   // note_from_alice ← 2
```

# pthreads and reordering

many pthreads functions prevent reordering
    everything before function call actually happens before

includes preventing some optimizations
    e.g. keeping global variable in register for too long

pthread_mutex_lock/unlock, pthread_create, pthread_join, …
    basically: if pthreads is waiting for/starting something, no weird ordering

implementation of this: pthread functions use special instructions
    example: x86 `mfence` instruction

# mfence

x86 instruction `mfence`

make sure all loads/stores in progress finish

...and make sure no loads/stores were started early
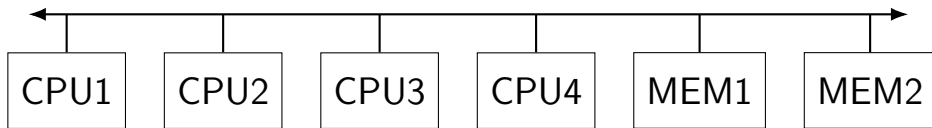
fairly expensive
  Intel 'Skylake': order 33 cycles + time waiting for pending stores/loads

# connecting CPUs and memory

multiple processors, common memory

how do processors communicate with memory?

# shared bus



tagged messages — everyone gets everything, filters

contention if multiple communicators
    some hardware enforces only one at a time

# shared buses and scaling

shared buses perform poorly with "too many" CPUs

so, there are other designs

we'll gloss over these for now

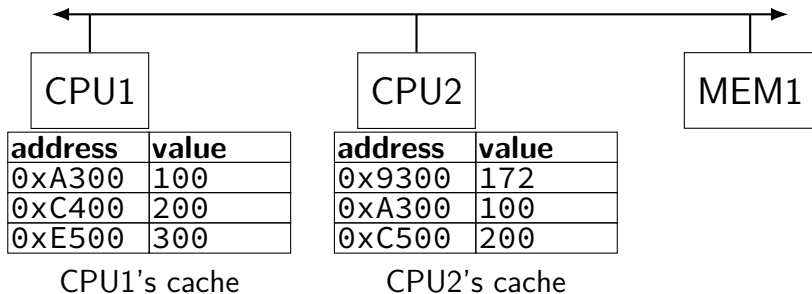# shared buses and caches

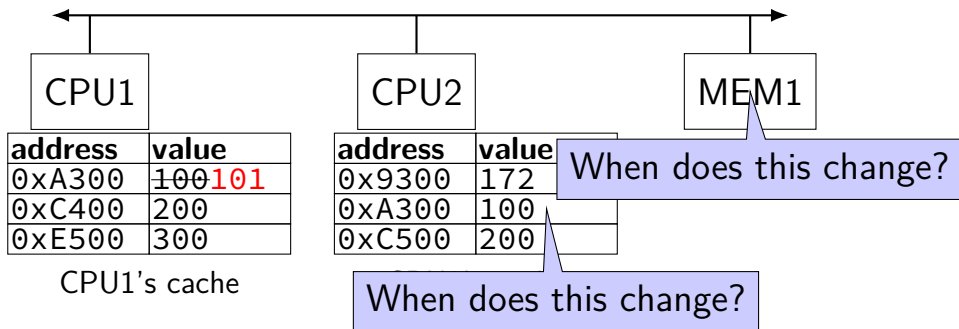remember caches?

memory is <span style="color:red">pretty slow</span>

each CPU wants to keep local copies of memory

what happens when multiple CPUs cache same memory?

# the cache coherency problem



| address | value |
| --- | --- |
| 0xA300 | 100 |
| 0xC400 | 200 |
| 0xE500 | 300 |

CPU1's cache

| address | value |
| --- | --- |
| 0x9300 | 172 |
| 0xA300 | 100 |
| 0xC500 | 200 |

CPU2's cache

# the cache coherency problem



| address | value |
|---------|-------|
| 0xA300  | ~~100~~101 |
| 0xC400  | 200   |
| 0xE500  | 300   |

CPU1's cache

| address | value |
|---------|-------|
| 0x9300  | 172   |
| 0xA300  | 100   |
| 0xC500  | 200   |

When does this change?

When does this change?

CPU1 writes 101 to 0xA300?

# "snooping" the bus

want to change a value other processors might have?

use bus to tell them "get rid of your copy"

want to start using value other processor might have reserved?

use bus to say "I'd like to use this value now"

# modifying cache blocks in parallel

cache coherency works on cache blocks

but typical memory access — less than cache block
> e.g. one 4-byte array element in 64-byte cache block

what if two processors modify different parts same cache block?
> 4-byte writes to 64-byte cache block

cache coherency — write instructions happen one at a time:
> processor 'locks' 64-byte cache block, fetching latest version
> processor updates 4 bytes of 64-byte cache block
> later, processor might give up cache block
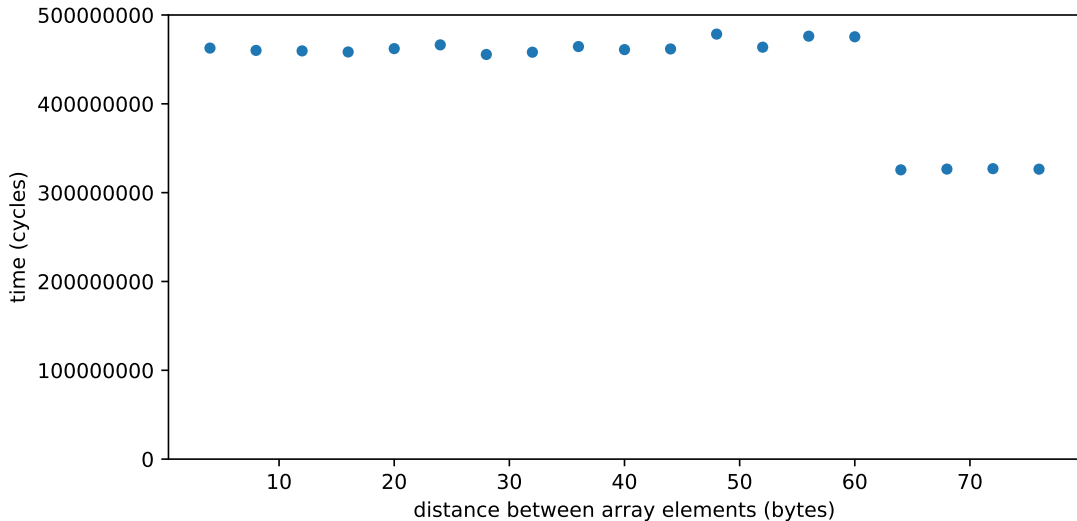
# modifying things in parallel (code)

```c
void *sum_up(void *raw_dest) {
    int *dest = (int *) raw_dest;
    for (int i = 0; i < 64 * 1024 * 1024; ++i) {
        *dest += data[i];
    }
}

__attribute__((aligned(4096)))
int array[1024];  /* aligned = address is mult. of 4096 */

void sum_twice(int distance) {
    pthread_t threads[2];
    pthread_create(&threads[0], NULL, sum_up, &array[0]);
    pthread_create(&threads[1], NULL, sum_up, &array[distance]);
    pthread_join(threads[0], NULL);
    pthread_join(threads[1], NULL);
}
```

# performance v. array element gap

(assuming `sum_up` compiled to not omit memory accesses)

# false sharing

synchronizing to access two independent things

two parts of same cache block

solution: separate them

# atomic read-modfiy-write

really hard to build locks for atomic load store
     and normal load/stores aren't even atomic...

...so processors provide read/modify/write operations


one instruction that
*atomically*
reads *and* modifies *and* writes back a value

# x86 atomic exchange

`lock xchg (%ecx), %eax`

atomic exchange

$temp \leftarrow M[ECX]$

$M[ECX] \leftarrow EAX$

$EAX \leftarrow temp$

…without being interrupted by other processors, etc.

# test-and-set: using atomic exchange

one instruction that…

writes a fixed new value

and reads the old value

# test-and-set: using atomic exchange

one instruction that…

writes a fixed new value

and reads the old value

write: mark a locked as TAKEN (no matter what)

read: see if it was already TAKEN (if so, only us)

# implementing atomic exchange

make sure other processors don't have cache block

do read+modify+write operation

recall: Modified state = "I am the only one with a copy"

# x86-64 spinlock with xchg

lock variable in shared memory: `the_lock`

if 1: someone has the lock; if 0: lock is free to take

```
acquire:
    movl $1, %eax              // %eax ← 1
    lock xchg %eax, the_lock   // swap %eax and the_lock
                               // sets the_lock to 1 (taken)
                               // sets %eax to prior val. of th
    test %eax, %eax            // if the_lock wasn't 0 before:
    jne acquire                //   try again
    ret

release:
    mfence                     // for memory order reasons
    movl $0, the_lock          // then, set the_lock to 0 (not taken)
    ret
```

# x86-64 spinlock with xchg

lock variable in shared memory: `the_lock`

if 1: someone has the lock; if 0: lock is free to take

```
acquire:
    movl $1, %eax              // %eax ← 1
    lock xchg %eax, the_lock   // swap %eax and the_lock
                               // sets the_lock to 1 (taken)
                               // sets %eax to prior val. of th
    test %eax, %eax            // if    set lock variable to 1 (taken)
    jne acquire                //       read old value
    ret

release:
    mfence                     // for memory order reasons
    movl $0, the_lock          // then, set the_lock to 0 (not taken)
    ret
```

# x86-64 spinlock with xchg

lock variable in shared memory: `the_lock`

if 1: someone has the lock; if 0: lock is free to take

```
acquire:
    movl $1, %eax              // %eax ← 1
    lock xchg %eax, the_lock  // swap %eax and the_lock
                              // sets the_lock to 1 (taken)
                              // sets %eax to prior val. of th

    test %eax, %eax
    jne acquire
    ret

release:
    mfence                    // for memory order reasons
    movl $0, the_lock         // then, set the_lock to 0 (not taken)
    ret
```

if lock was already locked retry
"spin" until lock is released elsewhere

# x86-64 spinlock with xchg

lock variable in shared memory: `the_lock`

if 1: someone has the lock; if 0: lock is free to take

```
acquire:
    movl $1, %eax              // %eax ← 1
    lock xchg %eax, the_lock   // swap %eax and the_lock
                               // sets the_lock to 1 (taken)
                               // sets %eax to prior val. of th
    test %eax, %eax            release lock by setting it to 0 (not taken)
    jne acquire                allows looping acquire to finish
    ret

release:
    mfence                     // for memory order reasons
    movl $0, the_lock          // then, set the_lock to 0 (not taken)
    ret
```

# x86-64 spinlock with xchg

lock variable in shared memory: `the_lock`

if 1: someone has the lock; if 0: lock is free to take

```
acquire:
    movl $1, %eax            // %eax ← 1
    lock xchg %eax, the_lock // swap %eax and the_lock
                             // sets the_lock to 1 (taken)
                                                    of th
    test %eax, %eax
    jne acquire
    ret

release:
    mfence                   // for memory order reasons
    movl $0, the_lock        // then, set the_lock to 0 (not taken)
    ret
```

Intel's manual says:
no reordering of loads/stores across a `lock`
or `mfence` instruction

# exercise: spin wait

consider implementing 'waiting' functionality of pthread_join

thread calls ThreadFinish() when done

complete code below:

```
finished: .quad 0

ThreadFinish:
    _____
    ret
ThreadWaitForFinish:
    _____
    lock xchg %eax, finished
    cmp $0, %eax
    ____ ThreadWaitForFinish
    ret
```

A. mfence; mov $1, finished & C. mov $0, %eax & E. je

B. mov $1, finished; mfence & D. mov $1, %eax & F. jne

# spinlock problems

lock abstraction is not powerful enough
  lock/unlock operations don't handle "wait for event"
  common thing we want to do with threads
  solution: other synchronization abstractions

spinlocks waste CPU time more than needed
  want to run another thread instead of infinite loop
  solution: lock implementation integrated with scheduler

spinlocks can send a lot of messages on the shared bus
  more efficient atomic operations to implement locks

# spinlock problems

lock abstraction is not powerful enough
> lock/unlock operations don't handle "wait for event"
> common thing we want to do with threads
> solution: other synchronization abstractions

spinlocks waste CPU time more than needed
> want to run another thread instead of infinite loop
> solution: lock implementation integrated with scheduler

spinlocks can send a lot of messages on the shared bus
> more efficient atomic operations to implement locks

## mutexes: intelligent waiting

want: locks that wait better
     example: POSIX mutexes

instead of running infinite loop, give away CPU

lock = go to sleep, add self to list
     sleep = scheduler runs something else

unlock = wake up sleeping thread

# mutexes: intelligent waiting

want: locks that wait better
    example: POSIX mutexes

instead of running infinite loop, give away CPU

lock = go to sleep, add self to list
    sleep = scheduler runs something else

unlock = wake up sleeping thread

# better lock implementation idea

*shared* list of waiters

spinlock protects list of waiters from concurrent modification

lock = use spinlock to add self to list, then wait without spinlock

unlock = use spinlock to remove item from list

# better lock implementation idea

*shared* list of waiters

spinlock protects list of waiters from concurrent modification

lock = use spinlock to add self to list, then wait without spinlock

unlock = use spinlock to remove item from list

# one possible implementation

```
struct Mutex {
    SpinLock guard_spinlock;
    bool lock_taken = false;
    WaitQueue wait_queue;
};
```

# one possible implementation

```
struct Mutex {
    SpinLock guard_spinlock;
    bool lock_taken = false;
    WaitQueue wait_queue;
};
```

spinlock protecting `lock_taken` and `wait_queue`
only held for very short amount of time (compared to mutex itself)

# one possible implementation

```
struct Mutex {
    SpinLock guard_spinlock;
    bool lock_taken = false;
    WaitQueue wait_queue;
};
```

tracks whether any thread has locked and not unlocked

# one possible implementation

```
struct Mutex {
    SpinLock guard_spinlock;
    bool lock_taken = false;
    WaitQueue wait_queue;
};
```

list of threads that discovered lock is taken
and are waiting for it be free
these threads are not runnable

# one possible implementation

```
struct Mutex {
    SpinLock guard_spinlock;
    bool lock_taken = false;
    WaitQueue wait_queue;
};
```

```
LockMutex(Mutex *m) {
  LockSpinlock(&m->guard_spinlock);
  if (m->lock_taken) {
    put current thread on m->wait_queue
    mark current thread not runnable
    /* xv6: myproc()->state = SLEEPING; */
    UnlockSpinlock(&m->guard_spinlock);
    run scheduler
  } else {
    m->lock_taken = true;
    UnlockSpinlock(&m->guard_spinlock);
  }
}
```

```
UnlockMutex(Mutex *m) {
  LockSpinlock(&m->guard_spinlock);
  if (m->wait_queue not empty) {
    remove a thread from m->wait_queue
    make that thread runnable
    /* xv6: myproc()->state = RUNNABLE; */
  } else {
    m->lock_taken = false;
  }
  UnlockSpinlock(&m->guard_spinlock);
}
```

# one possible implementation

```
struct Mutex {
    SpinLock guard_spinlock;
    bool lock_taken = false;
    WaitQueue wait_queue;
};
```

instead of setting lock_taken to false
choose thread to hand-off lock to

```
LockMutex(Mutex *m) {
  LockSpinlock(&m->guard_spinlock);
  if (m->lock_taken) {
    put current thread on m->wait_queue
    mark current thread not runnable
    /* xv6: myproc()->state = SLEEPING; */
    UnlockSpinlock(&m->guard_spinlock);
    run scheduler
  } else {
    m->lock_taken = true;
    UnlockSpinlock(&m->guard_spinlock);
  }
}
```

```
UnlockMutex(Mutex *m) {
  LockSpinlock(&m->guard_spinlock);
  if (m->wait_queue not empty) {
    remove a thread from m->wait_queue
    make that thread runnable
    /* xv6: myproc()->state = RUNNABLE; */
  } else {
    m->lock_taken = false;
  }
  UnlockSpinlock(&m->guard_spinlock);
}
```

# one possible implementation

```
struct Mutex {
    SpinLock guard_spinlock;
    bool lock_taken = false;
    WaitQueue wait_queue;
};
```

subtle: what if UnlockMutex runs on another core between these lines?
this thread hasn't saved registers yet, so scheduler can't switch to it yet
one solution: flag indicating "hasn't run scheduler yet"

```
LockSpinlock(&m->guard_spinlock);          UnlockMutex(Mutex *m) {
if (m->lock_taken) {                          LockSpinlock(&m->guard_spinlock);
  put current thread on m->wait_queue         if (m->wait_queue not empty) {
  mark current thread not runnable              remove a thread from m->wait_queue
  /* xv6: myproc()->state = SLEEPING; */        make that thread runnable
  UnlockSpinlock(&m->guard_spinlock);           /* xv6: myproc()->state = RUNNABLE; */
  run scheduler                               } else {
} else {                                        m->lock_taken = false;
  m->lock_taken = true;                       }
  UnlockSpinlock(&m->guard_spinlock);         UnlockSpinlock(&m->guard_spinlock);
}                                           }
}
```

# one possible implementation

```
struct Mutex {
    SpinLock guard_spinlock;
    bool lock_taken = false;
    WaitQueue wait_queue;
};
```

```
LockMutex(Mutex *m) {
  LockSpinlock(&m->guard_spinlock);
  if (m->lock_taken) {
    put current thread on m->wait_queue
    mark current thread not runnable
    /* xv6: myproc()->state = SLEEPING; */
    UnlockSpinlock(&m->guard_spinlock);
    run scheduler
  } else {
    m->lock_taken = true;
    UnlockSpinlock(&m->guard_spinlock);
  }
}
```

```
UnlockMutex(Mutex *m) {
  LockSpinlock(&m->guard_spinlock);
  if (m->wait_queue not empty) {
    remove a thread from m->wait_queue
    make that thread runnable
    /* xv6: myproc()->state = RUNNABLE; */
  } else {
    m->lock_taken = false;
  }
  UnlockSpinlock(&m->guard_spinlock);
}
```

# mutex and scheduler subtly

| core 0 (thread A) | core 1 (thread B) | core 2 |
|---|---|---|
| start LockMutex | | |
| acquire spinlock | | |
| discover lock taken | | |
| enqueue thread A | | |
| thread A set not runnable | | |
| release spinlock | start UnlockMutex | |
| | dequeue thread A | |
| | thread A set runnable | |
| | | run scheduler |
| | | scheduler switches to A |
| | | …with old verison of registers |
| thread A runs scheduler | | … |
| …finally saving registers | | … |

xv6 soln.: hold scheduler lock until thread A saves registers

Linux soln.: track that/check if thread A is still on core 0

# mutex and scheduler subtly

| core 0 (thread A) | core 1 (thread B) | core 2 |
|---|---|---|
| start LockMutex | | |
| acquire spinlock | | |
| discover lock taken | | |
| enqueue thread A | | |
| thread A set not runnable | | |
| release spinlock | start UnlockMutex | |
| | dequeue thread A | |
| | thread A set runnable | |
| | | run scheduler |
| | | scheduler switches to A |
| | | …with old verison of registers |
| thread A runs scheduler | | … |
| …finally saving registers | | … |

xv6 soln.: hold scheduler lock until thread A saves registers

Linux soln.: track that/check if thread A is still on core 0

# mutex efficiency

'normal' mutex **uncontended** case:
>     lock: acquire + release spinlock, see lock is free
>     unlock: acquire + release spinlock, see queue is empty

not much slower than spinlock

# recall: pthread mutex

```
#include <pthread.h>

pthread_mutex_t some_lock;
pthread_mutex_init(&some_lock, NULL);
// or: pthread_mutex_t some_lock = PTHREAD_MUTEX_INITIALIZER;
...
pthread_mutex_lock(&some_lock);
...
pthread_mutex_unlock(&some_lock);
pthread_mutex_destroy(&some_lock);
```

# pthread mutexes: addt'l features

mutex attributes (`pthread_mutexattr_t`) allow:
    (reference: `man pthread.h`)


error-checking mutexes
    locking mutex twice in same thread?
    unlocking already unlocked mutex?
    …

mutexes shared between processes
    otherwise: must be only threads of same process
    (unanswered question: where to store mutex?)

…

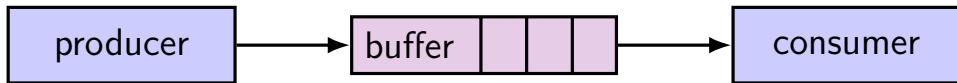# POSIX mutex restrictions

pthread_mutex rule: <span style="color:red">unlock from same thread you lock in</span>

implementation I gave before — not a problem

…but there other ways to implement mutexes
    e.g. might involve comparing with "holding" thread ID

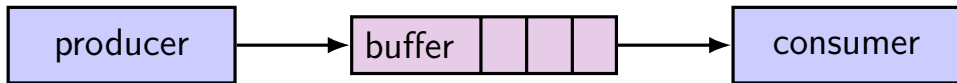# example: producer/consumer



shared buffer (queue) of fixed size

one or more producers inserts into queue

one or more consumers removes from queue

# example: producer/consumer



shared buffer (queue) of fixed size
> one or more producers inserts into queue
> one or more consumers removes from queue

producer(s) and consumer(s) don't work in lockstep
> (might need to wait for each other to catch up)

# example: producer/consumer

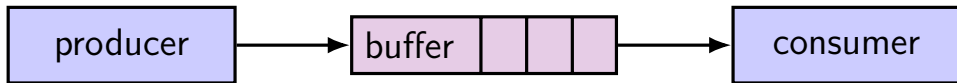

shared buffer (queue) of fixed size
    one or more producers inserts into queue
    one or more consumers removes from queue

producer(s) and consumer(s) don't work in lockstep
    (might need to wait for each other to catch up)

example: C compiler
    preprocessor $\rightarrow$ compiler $\rightarrow$ assembler $\rightarrow$ linker

# monitors/condition variables

locks for mutual exclusion

condition variables for waiting for event
    operations: wait (for event); signal/broadcast (that event happened)

related data structures

monitor = lock + 0 or more condition variables + shared data
    Java: every object is a monitor (has instance variables, built-in lock, cond. var)
    pthreads: build your own: provides you locks + condition variables

# monitor idea

a monitor

| lock |
| shared data |
| condvar 1 |
| condvar 2 |
| … |
| operation1(…) |
| operation2(…) |

# monitor idea

a monitor

| |
|---|
| lock |
| shared data |
| condvar 1 |
| condvar 2 |
| … |
| operation1(…) |
| operation2(…) |

lock must be acquired
before accessing
any part of monitor's stuff

# monitor idea

a monitor

| | |
|---|---|
| lock | |
| shared data | |
| condvar 1 | |
| condvar 2 | |
| … | |

| |
|---|
| operation1(…) |
| operation2(…) |

threads waiting for lock

# monitor idea



a monitor

| | |
|---|---|
| lock | threads waiting for lock |
| shared data | |
| condvar 1 | |
| condvar 2 | threads waiting for |
| … | condition to be true |
| operation1(…) | about shared data |
| operation2(…) | |

# condvar operations

condvar operations:
Wait(cv, lock) — unlock lock, add current thread to cv queue
…and reacquire lock before returning
Broadcast(cv) — remove all from condvar queue
Signal(cv) — remove one from condvar queue

a monitor

| lock |
| shared data |
| condvar 1 |
| condvar 2 |
| … |
| operation1(…) |
| operation2(…) |

threads waiting for lock

threads waiting for
condition to be true
about shared data

# condvar operations

condvar operations:
**Wait(cv, lock)** — unlock lock, add current thread to cv queue
…and reacquire lock before returning
Broadcast(cv) — remove all from condvar queue
Signal(cv) — remove one from condvar queue

# condvar operations
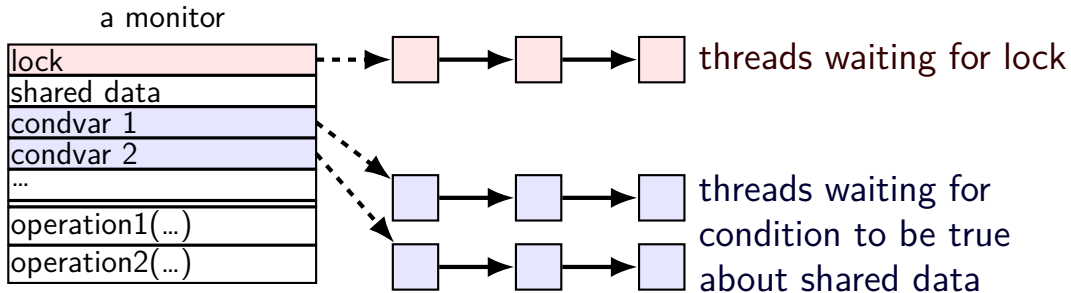
condvar operations:

Wait(cv, lock) — unlock lock, add current thread to cv queue

…and reacquire lock before returning

Broadcast(cv) — remove all from condvar queue

Signal(cv) — remove one from condvar queue



unlock lock — allow thread from queue to go

a monitor

| lock |
| shared data |
| condvar 1 |
| condvar 2 |
| ... |
| operation1(...) |
| operation2(...) |

threads waiting for lock

threads waiting for condition to be true about shared data
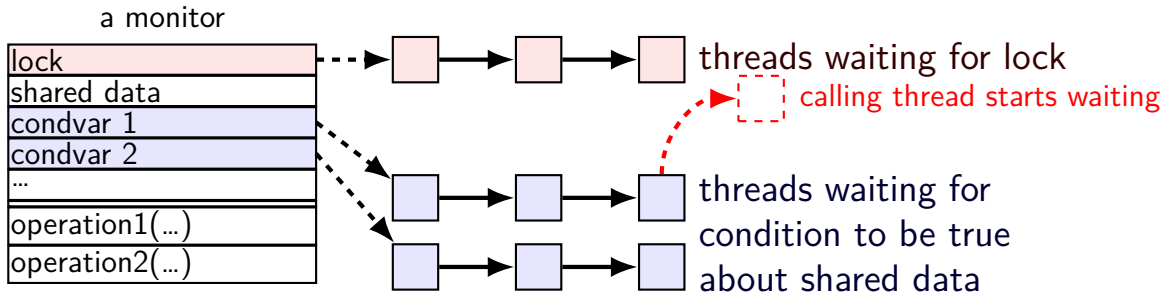
# condvar operations

condvar operations:
Wait(cv, lock) — unlock lock, add current thread to cv queue
…and reacquire lock before returning
**Broadcast(cv)** — remove all from condvar queue
Signal(cv) — remove one from condvar queue

a monitor



threads waiting for lock

all threads removed from cv queue
to start waiting for lock

threads waiting for
condition to be true
about shared data

lock
shared data
condvar 1
condvar 2
…
operation1(…)
operation2(…)

# condvar operations
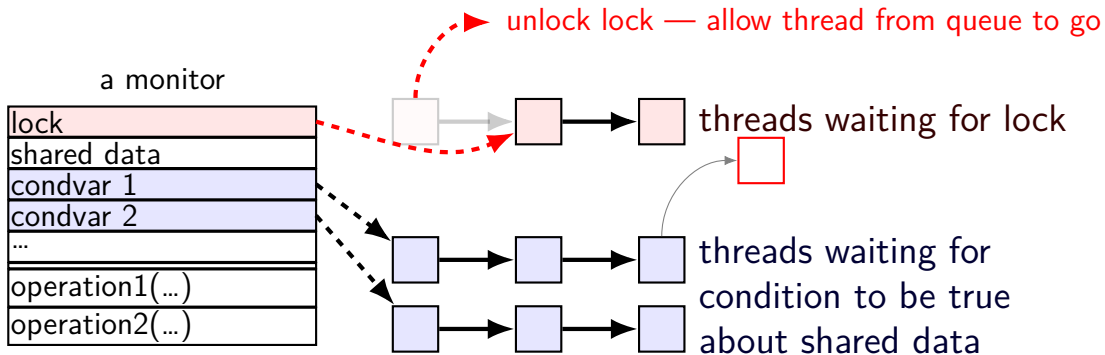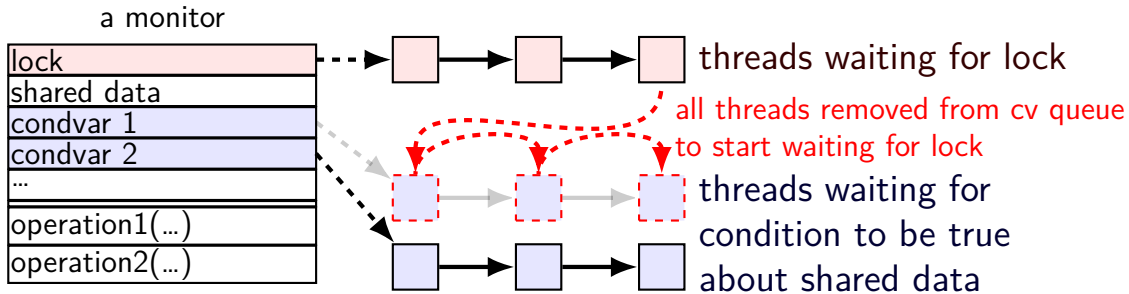
condvar operations:

Wait(cv, lock) — unlock lock, add current thread to cv queue

...and reacquire lock before returning

Broadcast(cv) — remove all from condvar queue

**Signal(cv)** — remove one from condvar queue



a monitor

threads waiting for lock

any one thread removed from cv queue
to start waiting for lock

threads waiting for
condition to be true
about shared data

# pthread cv usage

```
// MISSING: init calls, etc.
pthread_mutex_t lock;
bool finished;   // data, only accessed with after acquiring lock
pthread_cond_t finished_cv;  // to wait for 'finished' to be true

void WaitForFinished() {
  pthread_mutex_lock(&lock);
  while (!finished) {
    pthread_cond_wait(&finished_cv, &lock);
  }
  pthread_mutex_unlock(&lock);
}

void Finish() {
  pthread_mutex_lock(&lock);
  finished = true;
  pthread_cond_broadcast(&finished_cv);
  pthread_mutex_unlock(&lock);
}
```

# pthread cv usage

```
// MISSING: init calls, etc.
pthread_mutex_t lock;
bool finished;    // data, only accessed with after acquiring lock
pthread_cond_t finished_cv;  // to wait for 'finished' to be true

void WaitForFinished() {
  pthread_mutex_lock(&lock);
  while (!finished) {
    pthread_cond_wait(&finished_cv, &lock);
  }
  pthread_mutex_unlock(&lock);
}

void Finish() {
  pthread_mutex_lock(&lock);
  finished = true;
  pthread_cond_broadcast(&finished_cv);
  pthread_mutex_unlock(&lock);
}
```

acquire lock before
reading or writing `finished`

# pthread cv usage

```
// MISSING: init calls, etc.
pthread_mutex_t lock;
bool finished;   // data, only accessed with after acquiring lock
pthread_cond_t finished_cv;  // to wait for 'finished' to be true

void WaitForFinished() {
  pthread_mutex_lock(&lock);
  while (!finished) {
    pthread_cond_wait(&finished
  }
  pthread_mutex_unlock(&lock);
}

void Finish() {
  pthread_mutex_lock(&lock);
  finished = true;
  pthread_cond_broadcast(&finished_cv);
  pthread_mutex_unlock(&lock);
}
```

check whether we need to wait at all
(why a loop? we'll explain later)

# pthread cv usage

```
// MISSING: init calls, etc.
pthread_mutex_t lock;
bool finished;   // data, only accessed with after acquiring lock
pthread_cond_t finished_cv;  // to wait for 'finished' to be true

void WaitForFinished() {
  pthread_mutex_lock(&lock);
  while (!finished) {
    pthread_cond_wait(&finished_cv, &lock);
  }
  pthread_mutex_unlock(&lock);
}

void Finish() {
  pthread_mutex_lock(&lock);
  finished = true;
  pthread_cond_broadcast(&finished_cv);
  pthread_mutex_unlock(&lock);
}
```

know we need to wait
(finished can't change while we have lock)
so wait, releasing lock…

# pthread cv usage

```
// MISSING: init calls, etc.
pthread_mutex_t lock;
bool finished;   // data, only accessed with after acquiring lock
pthread_cond_t finished_cv;  // to wait for 'finished' to be true

void WaitForFinished() {
  pthread_mutex_lock(&lock);
  while (!finished) {
    pthread_cond_wait(&finished_cv, &lock);
  }
  pthread_mutex_unlock(&lock);
}

void Finish() {
  pthread_mutex_lock(&lock);
  finished = true;
  pthread_cond_broadcast(&finished_cv);
  pthread_mutex_unlock(&lock);
}
```

allow all waiters to proceed
(once we unlock the lock)

# WaitForFinish timeline 1

| WaitForFinish thread | Finish thread |
|---|---|
| mutex_lock(&lock) | |
| (thread has lock) | |
| | mutex_lock(&lock) |
| | (start waiting for lock) |
| while (!finished) ... | |
| cond_wait(&finished_cv, &lock); | |
| (start waiting for cv) | (done waiting for lock) |
| | finished = true |
| | cond_broadcast(&finished_cv) |
| (done waiting for cv) | |
| (start waiting for lock) | |
| | mutex_unlock(&lock) |
| (done waiting for lock) | |
| while (!finished) ... | |
| (finished now true, so return) | |
| mutex_unlock(&lock) | |

# WaitForFinish timeline 2

| WaitForFinish thread | Finish thread |
|---|---|
| | mutex_lock(&lock) |
| | finished = **true** |
| | cond_broadcast(&finished_cv) |
| | mutex_unlock(&lock) |
| mutex_lock(&lock) | |
| **while** (!finished) ... | |
| (finished now true, so return) | |
| mutex_unlock(&lock) | |

# why the loop

```
while (!finished) {
  pthread_cond_wait(&finished_cv, &lock);
}
```

we only broadcast if finished is true

so why check finished afterwards?

# why the loop

```
while (!finished) {
  pthread_cond_wait(&finished_cv, &lock);
}
```

we only `broadcast` if `finished` is true

so why check `finished` afterwards?

pthread_cond_wait manual page:
> "Spurious wakeups ... may occur."

spurious wakeup = `wait` returns even though nothing happened

# backup slides

# cache coherency states

extra information for each cache block
    overlaps with/replaces valid, dirty bits

stored in each cache

update states based on reads, writes and heard messages on bus

different caches may have different states for same block

# MSI state summary

**Modified**    value may be <span style="color:red">different than memory</span> *and* I am the only one who has it

**Shared**    value is the <span style="color:red">same as memory</span>

**Invalid**    I don't have the value; I will need to ask for it

# MSI scheme

| from state | hear read | hear write | read | write |
|---|---|---|---|---|
| Invalid | — | — | to Shared | to Modified |
| Shared | — | to Invalid | — | to Modified |
| Modified | to Shared | to Invalid | — | — |

blue: transition requires sending message on bus

# MSI scheme

| from state | hear read | hear write | read | write |
|---|---|---|---|---|
| Invalid | — | — | to Shared | to Modified |
| Shared | — | to Invalid | — | to Modified |
| Modified | to Shared | to Invalid | — | — |

blue: transition requires sending message on bus

example: write while Shared

    must send write — inform others with Shared state
    then change to Modified

# MSI scheme

| from state | hear read | hear write | read | write |
|---|---|---|---|---|
| Invalid | — | — | to Shared | to Modified |
| Shared | — | to Invalid | — | to Modified |
| Modified | to Shared | to Invalid | — | — |

blue: transition requires sending message on bus

example: write while Shared
    must send write — inform others with Shared state
    then change to Modified

example: hear write while Shared
    change to Invalid
    can send read later to get value from writer

example: write while Modified
    nothing to do — no other CPU can have a copy

# MSI example



| address | value | state |
|---------|-------|--------|
| 0xA300 | 100 | Shared |
| 0xC400 | 200 | Shared |
| 0xE500 | 300 | Shared |

| address | value | state |
|---------|-------|--------|
| 0x9300 | 172 | Shared |
| 0xA300 | 100 | Shared |
| 0xC500 | 200 | Shared |

CPU1    CPU2    MEM1

# MSI example

"CPU1 is writing 0xA3000"

maybe update memory?



CPU1

| address | value | state |
|---------|-------|-------|
| 0xA300 | ~~100~~101 | Modified |
| 0xC400 | 200 | Shared |
| 0xE500 | 300 | Shared |

CPU2

| address | value | state |
|---------|-------|-------|
| 0x9300 | 172 | Shared |
| ~~0xA300~~ | ~~100~~ | Invalid |
| 0xC500 | 200 | Shared |

MEM1

cache sees write: invalidate 0xA300

CPU1 writes 101 to 0xA300

# MSI example



nothing changed yet (writeback)

CPU1

| address | value | state |
|---------|-------|-------|
| 0xA300 | ~~101~~102 | Modified |
| 0xC400 | 200 | Shared |
| 0x... | | |

CPU2

| address | value | state |
|---------|-------|-------|
| 0x9300 | 172 | Shared |
| ~~0xA300~~ | ~~100~~ | Invalid |
| | | Shared |

MEM1

modified state — nothing communicated!
will "fix" later if there's a read

CPU1 writes 102 to 0xA300

# MSI example

"What is 0xA300?"



| CPU1 | | | CPU2 | | | MEM1 |

| address | value | state | address | value | state |
|---------|-------|----------|---------|-------|---------|
| 0xA300  | 102   | Modified | 0x9300  | 172   | Shared  |
| 0xC400  | 200   | Shared   | ~~0xA300~~ | ~~100~~ | Invalid |
| 0       |       |          |         |       | Shared  |

modified state — must update for CPU2!

CPU2 reads 0xA300

# MSI example

"Write 102 into 0xA300"



| address | value | state |   | address | value | state |
|---------|-------|-------|---|---------|-------|-------|
| 0xA300  | 102   | Shared |  | 0x9300  | 172   | Shared |
| 0xC400  | 200   | Shared |  | ~~0xA300~~ | ~~100~~ | Invalid |
| 0xE     |       |       |   |         |       | Shared |

written back to memory early
(could also become Invalid at CPU1)

CPU2 reads 0xA300

# MSI example



| address | value | state |
|---------|-------|-------|
| 0xA300  | 102   | Shared |
| 0xC400  | 200   | Shared |
| 0xE500  | 300   | Shared |

| address | value | state |
|---------|-------|-------|
| 0x9300  | 172   | Shared |
| ~~0xA300~~ | ~~100~~102 | Shared |
| 0xC500  | 200   | Shared |

# MSI: update memory

to write value (enter modified state), need to invalidate others

can avoid sending actual value (shorter message/faster)

"I am writing address $X$" versus "I am writing $Y$ to address $X$"

# MSI: on cache replacement/writeback

still happens — e.g. want to store something else

changes state to invalid

requires writeback if modified (= dirty bit)

# cache coherency exercise

modified/shared/invalid; all initially invalid; 32B blocks, 8B read/writes

    CPU 1: read 0x1000
    CPU 2: read 0x1000
    CPU 1: write 0x1000
    CPU 1: read 0x2000
    CPU 2: read 0x1000
    CPU 2: write 0x2008
    CPU 3: read 0x1008

Q1: final state of 0x1000 in caches?
    Modified/Shared/Invalid for CPU 1/2/3
    CPU 1:          CPU 2:          CPU 3:

Q2: final state of 0x2000 in caches?
    Modified/Shared/Invalid for CPU 1/2/3
    CPU 1:          CPU 2:          CPU 3:

# GCC: preventing reordering example (1)

```
void Alice() {
    int one = 1;
    __atomic_store(&note_from_alice, &one, __ATOMIC_SEQ_CST);
    do {
    } while (__atomic_load_n(&note_from_bob, __ATOMIC_SEQ_CST));
    if (no_milk) {++milk;}
}
```

---

```
Alice:
  movl $1, note_from_alice
  mfence
.L2:
  movl note_from_bob, %eax
  testl %eax, %eax
  jne .L2
  ...
```

# GCC: preventing reordering example (2)

```
void Alice() {
    note_from_alice = 1;
    do {
        __atomic_thread_fence(__ATOMIC_SEQ_CST);
    } while (note_from_bob);
    if (no_milk) {++milk;}
}
```

---

```
Alice:
  movl $1, note_from_alice  // note_from_alice ← 1
.L3:
  mfence  // make sure store is visible to other cores before
          // on x86: not needed on second+ iteration of loop
  cmpl $0, note_from_bob  // if (note_from_bob == 0) repeat fe
  jne .L3
  cmpl $0, no_milk
  ...
```

# C++: preventing reordering

to help implementing things like pthread_mutex_lock

C++ 2011 standard: *atomic* header, *std::atomic* class

prevent CPU reordering *and* prevent compiler reordering

also provide other tools for implementing locks (more later)

could also hand-write assembly code
    compiler can't know what assembly code is doing

# C++: preventing reordering example

```cpp
#include <atomic>
void Alice() {
    note_from_alice = 1;
    do {
        std::atomic_thread_fence(std::memory_order_seq_cst);
    } while (note_from_bob);
    if (no_milk) {++milk;}
}
```

---

```
Alice:
  movl $1, note_from_alice  // note_from_alice ← 1
.L2:
  mfence  // make sure store visible on/from other cores
  cmpl $0, note_from_bob  // if (note_from_bob == 0) repeat fence
  jne .L2
  cmpl $0, no_milk
  ...
```

# C++ atomics: no reordering

```
std::atomic<int> note_from_alice, note_from_bob;
void Alice() {
    note_from_alice.store(1);
    do {
    } while (note_from_bob.load());
    if (no_milk) {++milk;}
}
```

---

```
Alice:
  movl $1, note_from_alice
  mfence
.L2:
  movl note_from_bob, %eax
  testl %eax, %eax
  jne .L2
  ...
```

# GCC: built-in atomic functions

used to implement std::atomic, etc.

predate std::atomic

builtin functions starting with `__sync` and `__atomic`

these are what xv6 uses

# aside: some x86 reordering rules

each core sees its own loads/stores in order
> (if a core stores something, it can always load it back)

stores *from other cores* appear in a consistent order
> (but a core might observe its own stores too early)

*causality*:
*if* a core reads X=a and (after reading X=a) writes Y=b,
*then* a core that reads Y=b cannot later read X=older value than a

# how do you do anything with this?

difficult to reason about what modern CPU's reordering rules do

typically: don't depend on details, instead:

special instructions with stronger (and simpler) ordering rules
   often same instructions that help with implementing locks in other ways

special instructions that restrict ordering of instructions around them ("fences")
   loads/stores can't cross the fence

# xv6 spinlock: debugging stuff

```c
void acquire(struct spinlock *lk) {
  ...
  if(holding(lk))
    panic("acquire")
  ...
  // Record info about lock acquisition for debugging.
  lk->cpu = mycpu();
  getcallerpcs(&lk, lk->pcs);
}
void release(struct spinlock *lk) {
  if(!holding(lk))
    panic("release");

  lk->pcs[0] = 0;
  lk->cpu = 0;
  ...
}
```

# xv6 spinlock: debugging stuff

```
void acquire(struct spinlock *lk) {
  ...
  if(holding(lk))
    panic("acquire")
  ...
  // Record info about lock acquisition for debugging.
  lk->cpu = mycpu();
  getcallerpcs(&lk, lk->pcs);
}
void release(struct spinlock *lk) {
  if(!holding(lk))
    panic("release");

  lk->pcs[0] = 0;
  lk->cpu = 0;
  ...
}
```

# xv6 spinlock: debugging stuff

```
void acquire(struct spinlock *lk) {
  ...
  if(holding(lk))
    panic("acquire")
  ...
  // Record info about lock acquisition for debugging.
  lk->cpu = mycpu();
  getcallerpcs(&lk, lk->pcs);
}
void release(struct spinlock *lk) {
  if(!holding(lk))
    panic("release");

  lk->pcs[0] = 0;
  lk->cpu = 0;
  ...
}
```

# xv6 spinlock: debugging stuff

```
void acquire(struct spinlock *lk) {
  ...
  if(holding(lk))
    panic("acquire")
  ...
  // Record info about lock acquisition for debugging.
  lk->cpu = mycpu();
  getcallerpcs(&lk, lk->pcs);
}
void release(struct spinlock *lk) {
  if(!holding(lk))
    panic("release");

  lk->pcs[0] = 0;
  lk->cpu = 0;
  ...
}
```

# fetch-and-add with CAS (1)

```
compare-and-swap(address, old_value, new_value) {
    if (memory[address] == old_value) {
        memory[address] = new_value;
        return true;
    } else {
        return false;
    }
}
```

---

```
long my_fetch_and_add(long *pointer, long amount) { ... }
```

implementation sketch:

    fetch value from pointer old
    compute in temporary value result of addition new
    try to change value at pointer from old to new
    [compare-and-swap]
    if not successful, repeat

# fetch-and-add with CAS (2)

```
long my_fetch_and_add(long *p, long amount) {
    long old_value;
    do {
        old_value = *p;
    } while (!compare_and_swap(p, old_value, old_value + amount);
    return old_value;
}
```

## exercise: append to singly-linked list

ListNode is a singly-linked list

assume: threads *only* append to list (no deletions, reordering)

use `compare-and-swap(pointer, old, new)`:
    atomically change `*pointer` from `old` to `new`
    return true if successful
    return false (and change nothing) if `*pointer` is not `old`

```
void append_to_list(ListNode *head, ListNode *new_last_node) {
    ...
}
```

# spinlock problems

lock abstraction is not powerful enough
>  lock/unlock operations don't handle "wait for event"
>  common thing we want to do with threads
>  solution: other synchronization abstractions

spinlocks waste CPU time more than needed
>  want to run another thread instead of infinite loop
>  solution: lock implementation integrated with scheduler

spinlocks can send a lot of messages on the shared bus
>  more efficient atomic operations to implement locks

# ping-ponging



| address | value | state | | address | value | state | | address | value | state |
|---------|-------|-------|--|---------|-------|-------|--|---------|-------|-------|
| lock | locked | Modified | | lock | --- | Invalid | | lock | --- | Invalid |

CPU1    CPU2    CPU3    MEM1

# ping-ponging

"I want to modify `lock`?"



| address | value | state | address | value | state | address | value | state |
|---------|-------|-------|---------|-------|-------|---------|-------|-------|
| lock | --- | Invalid | lock | locked | Modified | lock | --- | Invalid |

CPU2 read-modify-writes lock
(to see it is still locked)

# ping-ponging

"I want to modify `lock`"



| address | value | state |
|---------|-------|-------|
| lock | --- | Invalid |

| address | value | state |
|---------|-------|-------|
| lock | --- | Invalid |

| address | value | state |
|---------|-------|-------|
| lock | locked | Modified |

CPU3 read-modify-writes lock
(to see it is still locked)

# ping-ponging

"I want to modify `lock`?"



| address | value | state |
|---------|-------|-------|
| lock | --- | Invalid |

CPU1

| address | value | state |
|---------|-------|-------|
| lock | locked | Modified |

CPU2

| address | value | state |
|---------|-------|-------|
| lock | --- | Invalid |

CPU3

MEM1

CPU2 read-modify-writes lock
(to see it is still locked)

# ping-ponging

"I want to modify `lock`"

| CPU1 | | | | CPU2 | | | | CPU3 | | | | MEM1 |

| address | value | state |
|---------|-------|-------|
| lock | --- | Invalid |

| address | value | state |
|---------|-------|-------|
| lock | --- | Invalid |

| address | value | state |
|---------|-------|-------|
| lock | locked | Modified |

CPU3 read-modify-writes lock
(to see it is still locked)

# ping-ponging

"I want to modify `lock`"



| address | value | state |
|---------|-------|-------|
| lock | unlocked | Modified |

| address | value | state |
|---------|-------|-------|
| lock | --- | Invalid |

| address | value | state |
|---------|-------|-------|
| lock | | Invalid |

CPU1 sets lock to unlocked

# ping-ponging

"I want to modify `lock`"



| address | value | state |
|---------|-------|-------|
| lock | --- | Invalid |

| address | value | state |
|---------|-------|-------|
| lock | locked | Modified |

| address | value | state |
|---------|-------|-------|
| lock | | Invalid |

some CPU (this example: CPU2) acquires lock

# ping-ponging

test-and-set problem: cache block "ping-pongs" between caches
    each waiting processor reserves block to modify
    could maybe wait until it determines modification needed — but not
    typical implementation

each transfer of block sends messages on bus

…so bus can't be used for real work
    like what the processor with the lock is doing

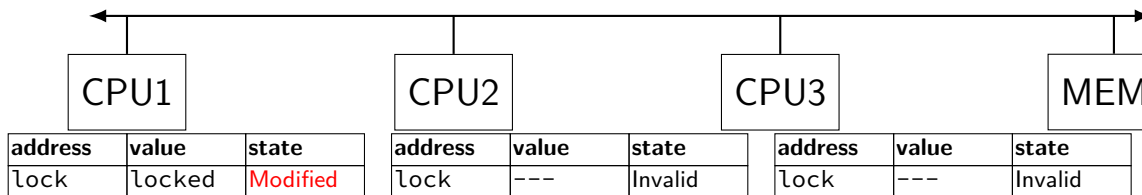# test-and-test-and-set (pseudo-C)

```
acquire(int *the_lock) {
    do {
        while (ATOMIC_READ(the_lock) == 0) { /* try again */ }
    } while (ATOMIC_TEST_AND_SET(the_lock) == ALREADY_SET);
}
```
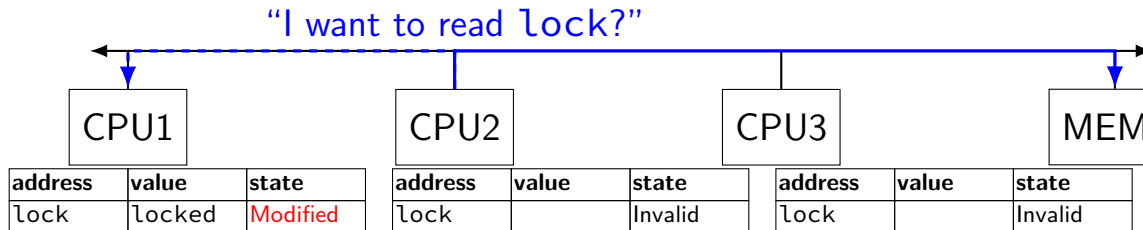
# test-and-test-and-set (assembly)

```
acquire:
    cmp $0, the_lock        // test the lock non-atomically
            // unlike lock xchg --- keeps lock in Shared state!
    jne acquire             // try again (still locked)
    // lock possibly free
    // but another processor might lock
    // before we get a chance to
    // ... so try wtih atomic swap:
    movl $1, %eax           // %eax ← 1
    lock xchg %eax, the_lock  // swap %eax and the_lock
            // sets the_lock to 1
            // sets %eax to prior value of the_lock
    test %eax, %eax         // if the_lock wasn't 0 (someone else
    jne acquire             //   try again
    ret
```

# less ping-ponging



| address | value | state |
|---------|--------|----------|
| lock | locked | Modified |

| address | value | state |
|---------|-------|---------|
| lock | --- | Invalid |

| address | value | state |
|---------|-------|---------|
| lock | --- | Invalid |

CPU1  CPU2  CPU3  MEM

# less ping-ponging

"I want to read `lock`?"



| address | value | state |
| --- | --- | --- |
| lock | locked | Modified |

| address | value | state |
| --- | --- | --- |
| lock | | Invalid |

| address | value | state |
| --- | --- | --- |
| lock | | Invalid |

CPU2 reads lock
(to see it is still locked)

# less ping-ponging

"set lock to locked"



| address | value | state |
|---------|-------|-------|
| lock | locked | Shared |

| address | value | state |
|---------|-------|-------|
| lock | locked | Shared |

| address | value | state |
|---------|-------|-------|
| lock | | Invalid |

CPU1 writes back lock value,
then CPU2 reads it

# less ping-ponging



"I want to read `lock`"

| address | value | state |
|---------|-------|-------|
| lock | locked | Shared |

CPU1

| address | value | state |
|---------|-------|-------|
| lock | locked | Shared |

CPU2

| address | value | state |
|---------|-------|-------|
| lock | locked | Shared |

CPU3

MEM

CPU3 reads lock
(to see it is still locked)

# less ping-ponging



| address | value | state |
|---------|--------|--------|
| lock | locked | Shared |

| address | value | state |
|---------|--------|--------|
| lock | locked | Shared |

| address | value | state |
|---------|--------|--------|
| lock | locked | Shared |

CPU2, CPU3 continue to read lock from cache
no messages on the bus

# less ping-ponging

"I want to modify `lock`"



| address | value | state |
|---------|-------|-------|
| lock | unlocked | Modified |

| address | value | state |
|---------|-------|-------|
| lock | --- | Invalid |

| address | value | state |
|---------|-------|-------|
| lock | --- | Invalid |

CPU1 · CPU2 · CPU3 · MEM

CPU1 sets lock to unlocked

# less ping-ponging

"I want to modify `lock`"



| address | value | state |
|---------|-------|-------|
| lock | | Modified |

| address | value | state |
|---------|-------|-------|
| lock | | Invalid |

| address | value | state |
|---------|-------|-------|
| lock | | Invalid |

CPU1    CPU2    CPU3    MEM

some CPU (this example: CPU2) acquires lock
(CPU1 writes back value, then CPU2 reads + modifies it)

# couldn't the read-modify-write instruction…

notice that the value of the lock isn't changing…

and keep it in the shared state

maybe — but extra step in "common" case
(swapping different values)

# more room for improvement?

can still have a lot of attempts to modify locks after unlocked

there other spinlock designs that avoid this
ticket locks
MCS locks
…

# MSI extensions

real cache coherency protocols sometimes more complex:

separate tracking modifications from whether other caches have copy

send values directly between caches (maybe skip write to memory)

send messages only to cores which might care (no shared bus)