



## so far

building processes

virtual memory + virtual CPU (thread) + files (virtual storage)

connecting machines (networking, security)

caching — speedy memories

interfaces for controlling processes

revisiting threads

multi-threaded processes

passing values to/waiting for threads

# last time

reordering of loads/stores by compilers, processors

locks

- way to take turns

- don't actually restrict access to resourcej

barriers

- wait for every other thread

deadlock problem

# anonymous feedback

locks not covered in F2022 DSA1 due to end-of-semester stuff  
think we went into full detail last time

captions on recordings not there

some not-so-automatically imported automated (low-quality) captions  
was missing some I tried to import due to filename problem, probably  
fixes

## complimentary feedback

I really appreciate all of your patience in answering questions that students have, guidance on homework questions (that may seem easy to understand for you but may be a bit more challenging for students), and your responsiveness on piazza. I'm thankful for how you are willing to sit down with students in office hours to walk them through problems they are encountering step by step in an encouraging manner. It really makes me feel like I have all the support I need to be successful in this class, and it's something that not a lot of other teachers can offer, so I just wanted to express my gratitude amidst all the other negative complaints.

## Quiz: pthread\_join mistakes

```
pthread_join(pthread_t, void**)
```

I wrote code as if it was `pthread_join(pthread_t, void**)`

meant that:

needed to accept “most likely crash” as valid answer for Q2

in Q3 `pthread_joins` would not actually wait for threads to finish, so option F was correct!

wasn't my intention to test this kind of API detail

since it's something you'd notice by reading compiler warnings/errors

# quiz Q1

```
void foo() {
    int array[4] = {0, 0, 0, 0}; /* <--- on stack of foo() */
    ...
    int *q = &array[i]; /* <--- q points to stack of foo() */
    /* same as: int *q; q = &array[i]; */
    pthread_create(...., (void *) q); /* passing 'q' to thread */
    ...
}

void *thread_func(void *arg) {
    int *p;
    p = (int*) arg; /* <-- p = q, so p points to stack of foo() */
    *p = 1; /* <--- modifies what p points to */
    ...
}
```

## quiz Q4b

“assume that all code follows the convention that `dynamic_array` structs are only read or modified while holding the lock in the struct”

```
pthread_mutex_unlock(&to_add->lock);  
... /* <-- to_dad->data, to_add->size could change here */  
pthread_mutex_lock(&to_add->lock);  
/* but not while executing the below: */  
memcpy(dest->new_data + dest->size, to_add->data, to_add->size  
dest->size += to_add->size;
```

will either copy from new allocation or from original allocation  
*before new allocation is made*

should've written 'out-of-date allocation' instead of 'old allocation'  
in quiz option (or made it clearer that it needed to be a *problem*)

## quiz Q4c

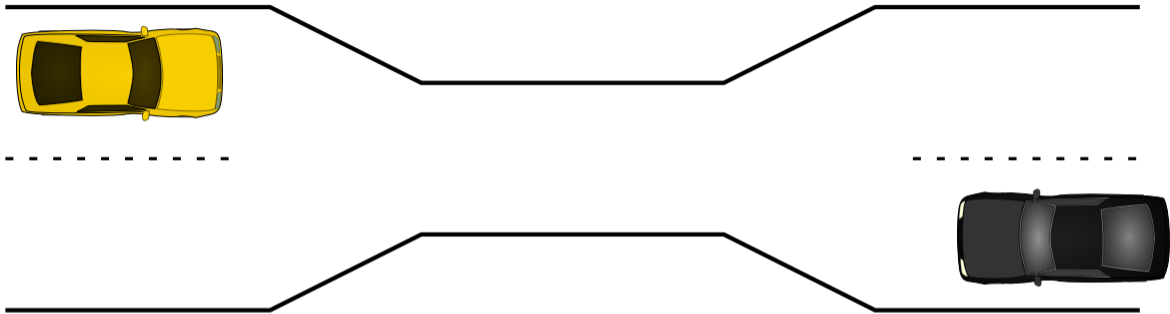
`append_to_array(dest, A) + append_to_array(dest, B)`

lock for `dest` is held during whole function

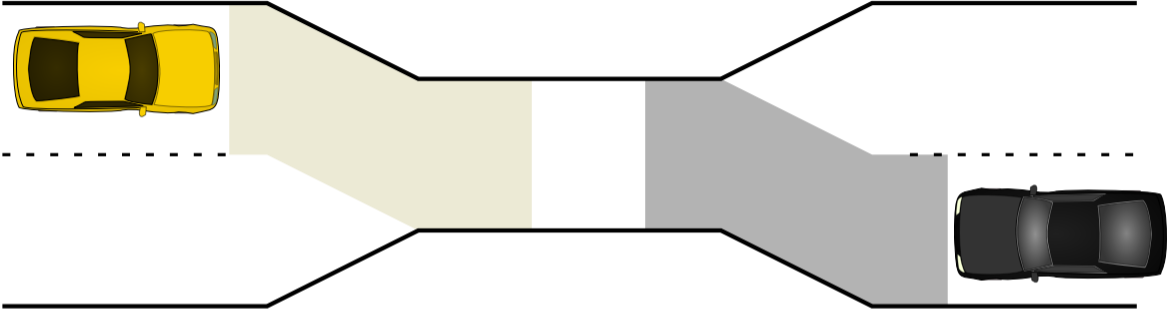
so the two calls take turns



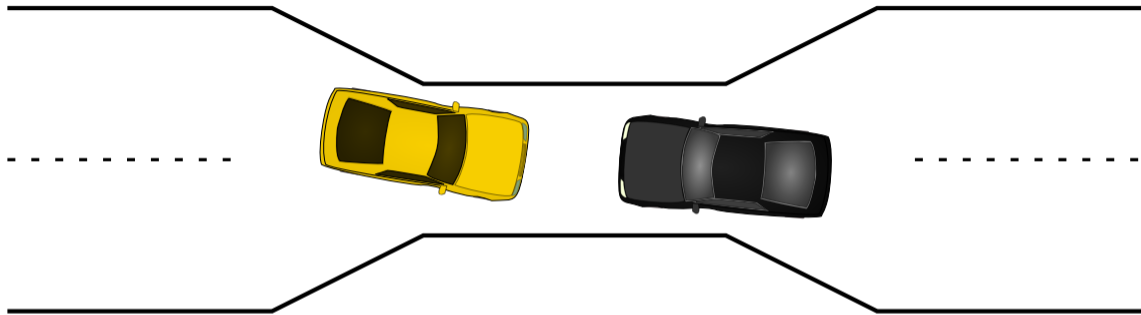
# the one-way bridge



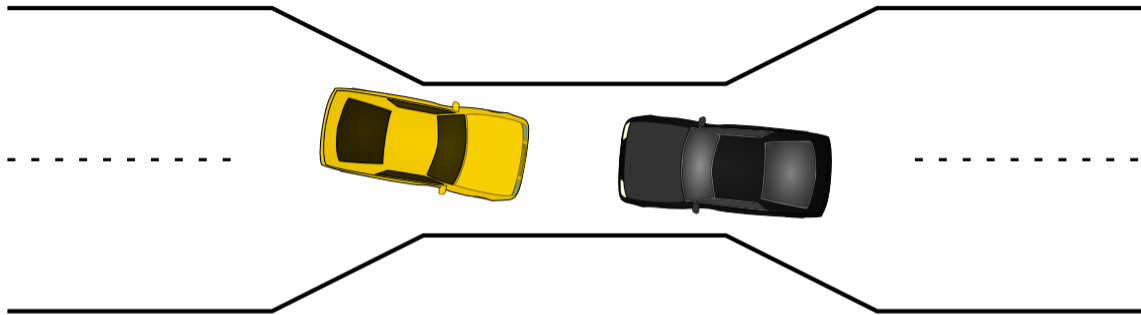
# the one-way bridge



# the one-way bridge



# the one-way bridge



## moving two files

```
struct Dir {
    mutex_t lock; HashMap entries;
};
void MoveFile(Dir *from_dir, Dir *to_dir, string filename) {
    mutex_lock(&from_dir->lock);
    mutex_lock(&to_dir->lock);

    Map_put(to_dir->entries, filename,
            Map_get(from_dir->entries, filename));
    Map_erase(from_dir->entries, filename);

    mutex_unlock(&to_dir->lock);
    mutex_unlock(&from_dir->lock);
}
```

Thread 1: MoveFile(A, B, "foo")

Thread 2: MoveFile(B, A, "bar")

# moving two files: lucky timeline (1)

## Thread 1

MoveFile(A, B, "foo")

---

lock(&A->lock);

lock(&B->lock);

(do move)

unlock(&B->lock);

unlock(&A->lock);

## Thread 2

MoveFile(B, A, "bar")

---

lock(&B->lock);

lock(&A->lock);

(do move)

unlock(&B->lock);

unlock(&A->lock);

## moving two files: lucky timeline (2)

### Thread 1

```
MoveFile(A, B, "foo")
```

---

```
lock(&A->lock);
```

```
lock(&B->lock);
```

```
(do move)
```

```
unlock(&B->lock);
```

```
unlock(&A->lock);
```

### Thread 2

```
MoveFile(B, A, "bar")
```

---

```
lock(&B->lock...
```

```
(waiting for B lock)
```

```
lock(&B->lock);
```

```
lock(&A->lock...
```

```
lock(&A->lock);
```

```
(do move)
```

```
unlock(&A->lock);
```

## moving two files: unlucky timeline

### Thread 1

```
MoveFile(A, B, "foo")
```

```
lock(&A->lock);
```

### Thread 2

```
MoveFile(B, A, "bar")
```

```
lock(&B->lock);
```



# moving two files: unlucky timeline

## Thread 1

```
MoveFile(A, B, "foo")
```

---

```
lock(&A->lock);
```

```
lock(&B->lock... stalled
```

```
(waiting for lock on B)
```

```
(waiting for lock on B)
```

## Thread 2

```
MoveFile(B, A, "bar")
```

---

```
lock(&B->lock);
```

```
lock(&A->lock... stalled
```

```
(waiting for lock on A)
```

# moving two files: unlucky timeline

## Thread 1

```
MoveFile(A, B, "foo")
```

---

```
lock(&A->lock);
```

```
lock(&B->lock... stalled
```

```
(waiting for lock on B)
```

```
(waiting for lock on B)
```

```
{do move} unreachable
```

```
unlock(&B->lock); unreachable
```

```
unlock(&A->lock); unreachable
```

## Thread 2

```
MoveFile(B, A, "bar")
```

---

```
lock(&B->lock);
```

```
lock(&A->lock... stalled
```

```
(waiting for lock on A)
```

```
{do move} unreachable
```

```
unlock(&A->lock); unreachable
```

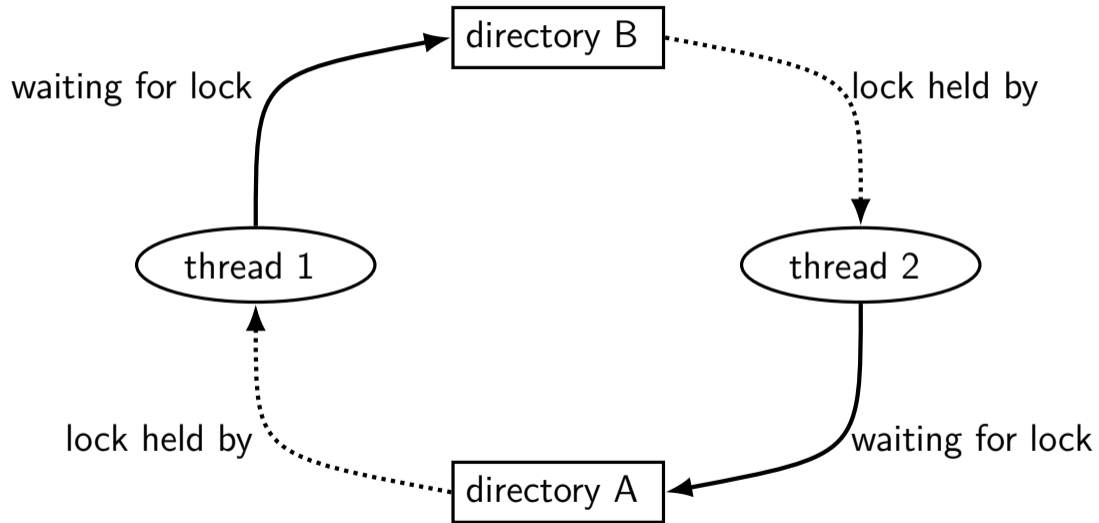
```
unlock(&B->lock); unreachable
```

# moving two files: unlucky timeline

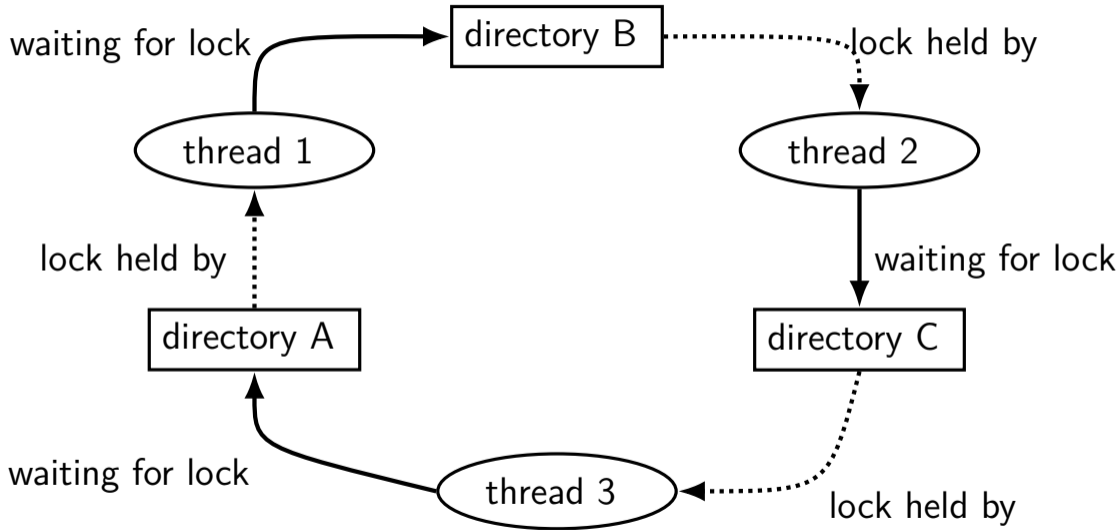
Thread 1	Thread 2
<code>MoveFile(A, B, "foo")</code>	<code>MoveFile(B, A, "bar")</code>
<code>lock(&amp;A-&gt;lock);</code>	
	<code>lock(&amp;B-&gt;lock);</code>
<code>lock(&amp;B-&gt;lock... stalled</code> (waiting for lock on B) (waiting for lock on B)	<code>lock(&amp;A-&gt;lock... stalled</code> (waiting for lock on A)
<del>(do move)</del> <code>unreachable</code>	<del>(do move)</del> <code>unreachable</code>
<del><code>unlock(&amp;B-&gt;lock);</code></del> <code>unreachable</code>	<del><code>unlock(&amp;A-&gt;lock);</code></del> <code>unreachable</code>
<del><code>unlock(&amp;A-&gt;lock);</code></del> <code>unreachable</code>	<del><code>unlock(&amp;B-&gt;lock);</code></del> <code>unreachable</code>

Thread 1 holds A lock, waiting for Thread 2 to release B lock

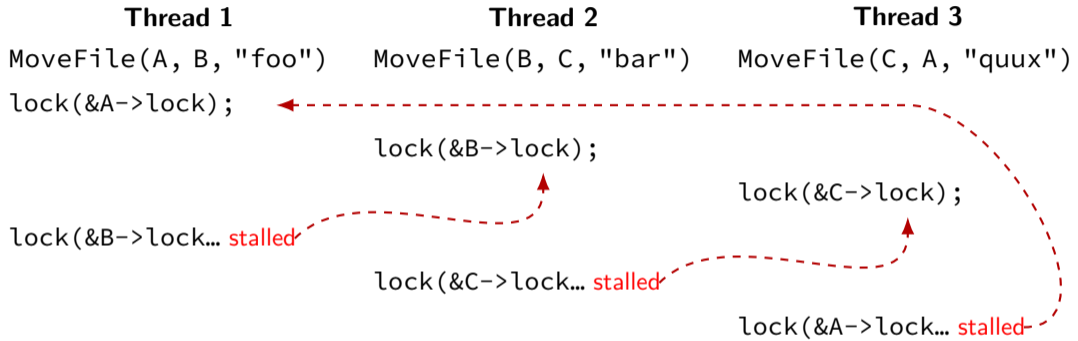
# moving two files: dependencies



# moving three files: dependencies



# moving three files: unlucky timeline



# deadlock with free space

## Thread 1

```
AllocateOrWaitFor(1 MB)
AllocateOrWaitFor(1 MB)
(do calculation)
Free(1 MB)
Free(1 MB)
```

## Thread 2

```
AllocateOrWaitFor(1 MB)
AllocateOrWaitFor(1 MB)
(do calculation)
Free(1 MB)
Free(1 MB)
```

2 MB of space — deadlock possible with unlucky order

# deadlock with free space (unlucky case)

## Thread 1

AllocateOrWaitFor(1 MB)

AllocateOrWaitFor(1 MB... stalled

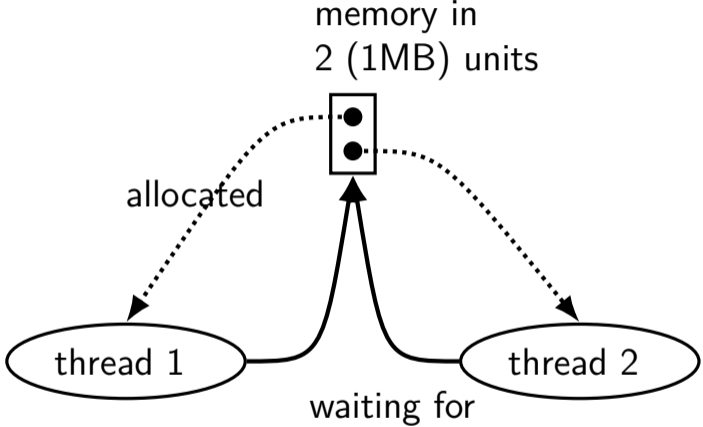
## Thread 2

AllocateOrWaitFor(1 MB)

AllocateOrWaitFor(1 MB... stalled



# free space: dependency graph



# deadlock with free space (lucky case)

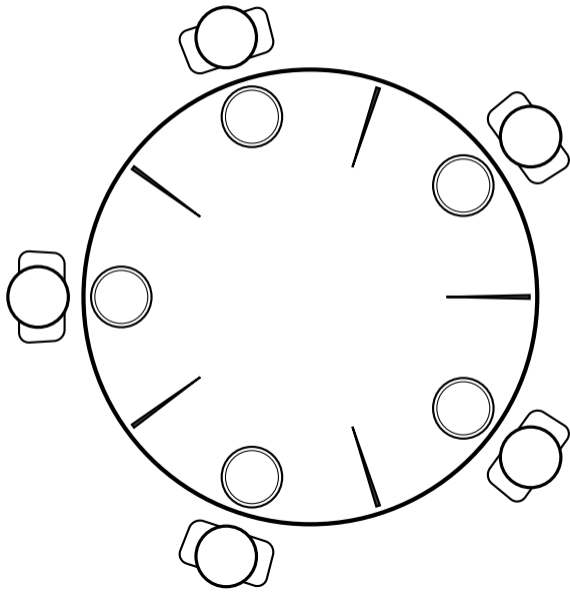
## Thread 1

```
AllocateOrWaitFor(1 MB)
AllocateOrWaitFor(1 MB)
(do calculation)
Free(1 MB);
Free(1 MB);
```

## Thread 2

```
AllocateOrWaitFor(1 MB)
AllocateOrWaitFor(1 MB)
(do calculation)
Free(1 MB);
Free(1 MB);
```

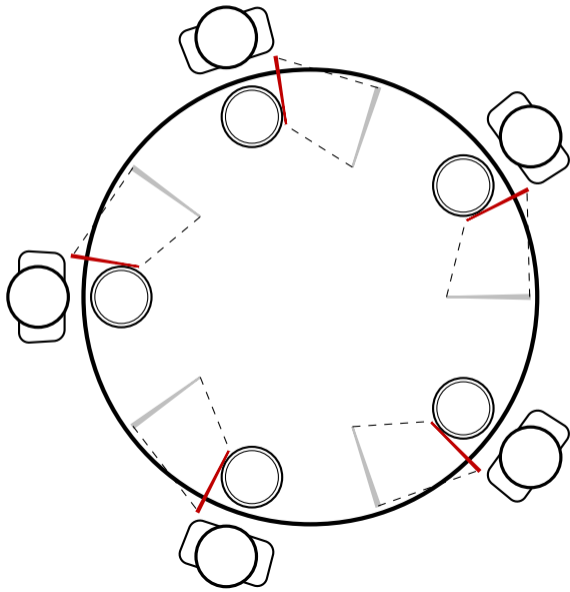
# dining philosophers



five philosophers either think or eat  
to eat:

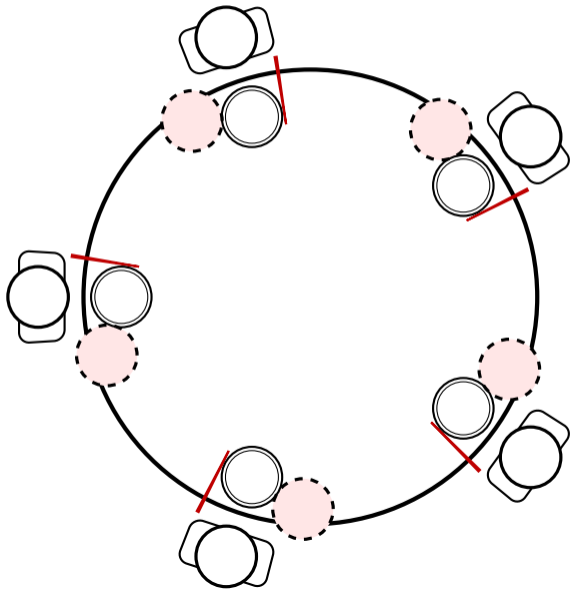
grab chopstick on left, then  
grab chopstick on right, then  
eat, then  
return chopsticks

# dining philosophers



everyone eats at the same time?  
grab left chopstick, then...

# dining philosophers



everyone eats at the same time?  
grab left chopstick, then  
try to grab right chopstick, ...  
we're at an impasse

# deadlock

deadlock — circular waiting for resources

resource = something needed by a thread to do work

locks

CPU time

disk space

memory

...

often non-deterministic in practice

most common example: **when acquiring multiple locks**

# deadlock

deadlock — circular waiting for **resources**

resource = something needed by a thread to do work

- locks

- CPU time

- disk space

- memory

- ...

often non-deterministic in practice

most common example: **when acquiring multiple locks**

# deadlock versus starvation

starvation: one+ unlucky (no progress), one+ lucky (yes progress)

example: low priority threads versus high-priority threads

deadlock: no one involved in deadlock makes progress



# deadlock versus starvation

starvation: one+ unlucky (no progress), one+ lucky (yes progress)

example: low priority threads versus high-priority threads

deadlock: no one involved in deadlock makes progress

starvation: once starvation happens, taking turns will resolve

low priority thread just needed a chance...

deadlock: once it happens, taking turns won't fix

# deadlock requirements

## mutual exclusion

one thread at a time can use a resource

## hold and wait

thread holding a resources waits to acquire *another* resource

## no preemption of resources

resources are only released voluntarily

thread trying to acquire resources can't 'steal'

## circular wait

there exists a set  $\{T_1, \dots, T_n\}$  of waiting threads such that

$T_1$  is waiting for a resource held by  $T_2$

$T_2$  is waiting for a resource held by  $T_3$

...

$T_n$  is waiting for a resource held by  $T_1$

# how is deadlock possible?

Given list: A, B, C, D, E

```
RemoveNode(LinkedListNode *node) {  
    pthread_mutex_lock(&node->lock);  
    pthread_mutex_lock(&node->prev->lock);  
    pthread_mutex_lock(&node->next->lock);  
    node->next->prev = node->prev; node->prev->next = node->next;  
    pthread_mutex_unlock(&node->next->lock); pthread_mutex_unlock(&node->p  
    pthread_mutex_unlock(&node->lock);  
}
```

Which of these (all run in parallel) can deadlock?

- A. RemoveNode(B) and RemoveNode(C)
- B. RemoveNode(B) and RemoveNode(D)
- C. RemoveNode(B) and RemoveNode(C) and RemoveNode(D)
- D. A and C
- E. B and C
- F. all of the above
- G. none of the above

## how is deadlock — solution

Remove B

lock B

lock A (prev)

wait to lock C (next)

Remove C

lock C

wait to lock B (prev)

---

With B and D — only overlap in in node C — no circular wait possible

# deadlock prevention techniques

## infinite resources

or at least enough that never run out

no *mutual exclusion*

## no shared resources

no *mutual exclusion*

## no waiting

“busy signal” — abort and (maybe) retry  
revoke/preempt resources

no *hold and wait/*  
*preemption*

acquire resources in **consistent order**

no *circular wait*

# deadlock prevention techniques

## infinite resources

or at least enough that never run out

no *mutual exclusion*

## no shared resources

no *mutual exclusion*

## no waiting

“busy signal” — abort and (maybe) retry  
revoke/preempt resources

no *hold and wait/*  
*preemption*

acquire resources in **consistent order**

no *circular wait*

# deadlock prevention techniques

## infinite resources

or at least enough that never run out

no *mutual exclusion*

## no shared resources

no *mutual exclusion*

## no waiting

“busy signal” — abort and (maybe) retry  
revoke/preempt resources

no *hold and wait/*  
*preemption*

acquire resources in **consistent order**

no *circular wait*

# deadlock prevention techniques

**infinite resources**

or at least enough that never run out

*no mutual exclusion*

memory allocation: malloc() fails rather than waiting (no deadlock)

locks: pthread\_mutex\_trylock fails rather than waiting

...

*exclusion*

**no waiting**

“busy signal” — abort and (maybe) retry

revoke/preempt resources

*no hold and wait/  
preemption*

acquire resources in **consistent order**

*no circular wait*



# stealing locks???

how do we make stealing locks possible

unclean: just kill the thread

    problem: inconsistent state?

clean: have code to undo partial operation

    some databases do this

won't go into detail in this class

# revokable locks?

```
try {  
    AcquireLock();  
    use shared data  
} catch (LockRevokedException le) {  
    undo operation hopefully?  
} finally {  
    ReleaseLock();  
}
```

# deadlock prevention techniques

**infinite resources**

or at least enough that never run out

*no mutual exclusion*

**no shared resources**

*no mutual exclusion*

**no waiting**

“**busy signal**” — **abort and (maybe) retry**  
revoke/preempt resources

*no hold and wait/  
preemption*

acquire resources in **consistent order**

*no circular wait*

# abort and retry limits?

abort-and-retry

pthread's mutexes:

`pthread_mutex_trylock`

`pthread_mutex_timedlock`

how many times will you retry?

## moving two files: abort-and-retry

```
struct Dir { mutex_t lock; HashMap entries; };  
void MoveFile(Dir *from_dir, Dir *to_dir, string filename) {  
    while (true) {  
        mutex_lock(&from_dir->lock);  
        if (mutex_trylock(&to_dir->lock) == LOCKED) break;  
        mutex_unlock(&from_dir->lock);  
    }  
}
```

```
Map_put(to_dir->entries, filename, Map_get(from_dir->entries, filename));  
from_dir->entries.erase(filename);
```

```
mutex_unlock(&to_dir->lock);  
mutex_unlock(&from_dir->lock);  
}
```

Thread 1: MoveFile(A, B, "foo"); Thread 2: MoveFile(B, A, "bar")

# moving two files: lots of bad luck?

## Thread 1

MoveFile(A, B, "foo")

---

lock(&A->lock) → LOCKED

trylock(&B->lock) → FAILED

unlock(&A->lock)

lock(&A->lock) → LOCKED

trylock(&B->lock) → FAILED

unlock(&A->lock)

## Thread 2

MoveFile(B, A, "bar")

---

lock(&B->lock) → LOCKED

trylock(&A->lock) → FAILED

unlock(&B->lock)

lock(&B->lock) → LOCKED

trylock(&A->lock) → FAILED

# livelock

livelock: keep aborting and retrying without end

like deadlock — no one's making progress  
potentially forever

unlike deadlock — threads are not waiting

# preventing livelock

make schedule random — e.g. random waiting after abort

make threads run one-at-a-time if lots of aborting

other ideas?



# deadlock prevention techniques

**infinite resources**

or at least enough that never run out

*no mutual exclusion*

**no shared resources**

*no mutual exclusion*

requires some way to undo partial changes to avoid errors  
common approach for databases

**no waiting**

...

“busy signal” — abort and (maybe) retry

revoke/preempt resources

*no hold and wait/  
preemption*

acquire resources in **consistent order**

*no circular wait*

# deadlock prevention techniques

## infinite resources

or at least enough that never run out

no *mutual exclusion*

## no shared resources

no *mutual exclusion*

## no waiting

“busy signal” — abort and (maybe) retry  
revoke/preempt resources

no *hold and wait/*  
*preemption*

acquire resources in **consistent order**

no *circular wait*

# acquiring locks in consistent order (1)

```
MoveFile(Dir* from_dir, Dir* to_dir, string filename) {  
    if (from_dir->path < to_dir->path) {  
        lock(&from_dir->lock);  
        lock(&to_dir->lock);  
    } else {  
        lock(&to_dir->lock);  
        lock(&from_dir->lock);  
    }  
    ...  
}
```

# acquiring locks in consistent order (1)

```
MoveFile(Dir* from_dir, Dir* to_dir, string filename) {  
    if (from_dir->path < to_dir->path) {  
        lock(&from_dir->lock);  
        lock(&to_dir->lock);  
    } else {  
        lock(&to_dir->lock);  
        lock(&from_dir->lock);  
    }  
    ...  
}
```

any ordering will do  
e.g. compare pointers

## acquiring locks in consistent order (2)

often by convention, e.g. Linux kernel comments:

```
/*  
 * ...  
 * Lock order:  
 *     contex.ldt_usr_sem  
 *     mmap_sem  
 *     context.lock  
 */
```

---

```
/*  
 * ...  
 * Lock order:  
 * 1. slab_mutex (Global Mutex)  
 * 2. node->list_lock  
 * 3. slab_lock(page) (Only on some arches and for debugging)  
 * ...  
 */
```

# deadlock prevention techniques

## infinite resources

or at least enough that never run out

*no mutual exclusion*

## no shared resources

*no mutual exclusion*

## no waiting

“busy signal” — abort and (maybe) retry  
revoke/preempt resources

*no hold and wait/  
preemption*

acquire resources in **consistent order**

*no circular wait*

# deadlock detection

why? debugging or fix deadlock by aborting operations

idea: search for cyclic dependencies

# detecting deadlocks on locks

let's say I want to detect deadlocks that only involve mutexes  
goal: help programmers debug deadlocks

...by modifying my threading library:

```
struct Thread {  
    ... /* stuff for implementing thread */  
    /* what extra fields go here? */  
};  
  
struct Mutex {  
    ... /* stuff for implementing mutex */  
    /* what extra fields go here? */  
};
```



# deadlock detection

why? debugging or fix deadlock by aborting operations

idea: search for cyclic dependencies

need:

- list of all contended resources

- what thread is waiting for what?

- what thread 'owns' what?

## aside: divisible resources

deadlock is possible with divisible resources like memory,...

example: suppose 6MB of RAM for threads total:

- thread 1 has 2MB allocated, waiting for 2MB

- thread 2 has 2MB allocated, waiting for 2MB

- thread 3 has 1MB allocated, waiting for keypress

cycle: thread 1 waiting on memory owned by thread 2?

not a deadlock — thread 3 can still finish

- and after it does, thread 1 or 2 can finish

## aside: divisible resources

deadlock is possible with divisible resources like memory,...

example: suppose 6MB of RAM for threads total:

- thread 1 has 2MB allocated, waiting for 2MB

- thread 2 has 2MB allocated, waiting for 2MB

- thread 3 has 1MB allocated, waiting for keypress

cycle: thread 1 waiting on memory owned by thread 2?

not a deadlock — thread 3 can still finish

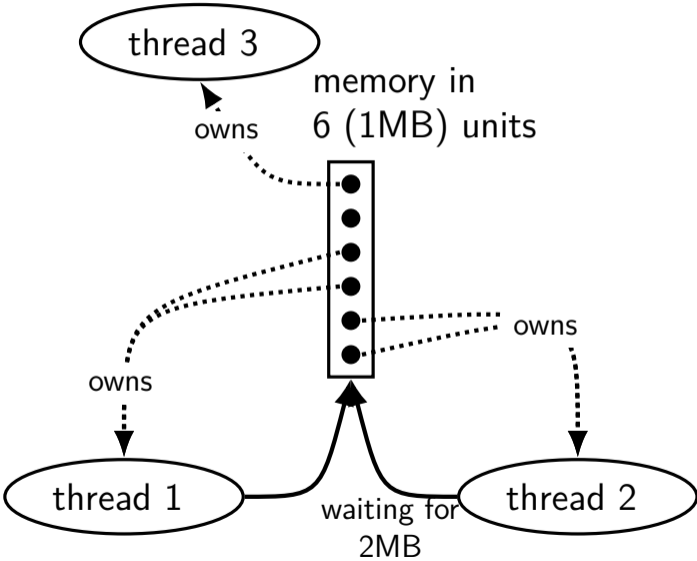
- and after it does, thread 1 or 2 can finish

...but would be deadlock

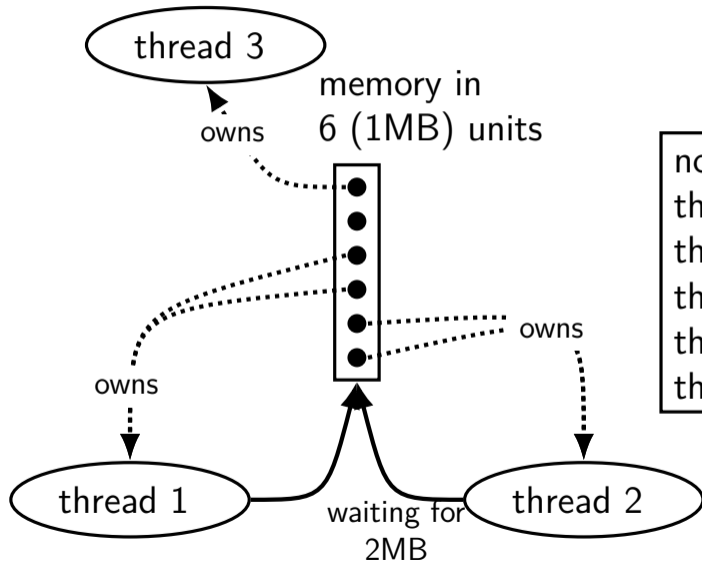
- ...if thread 3 waiting lock held by thread 1

- ...with 5MB of RAM

# divisible resources: not deadlock

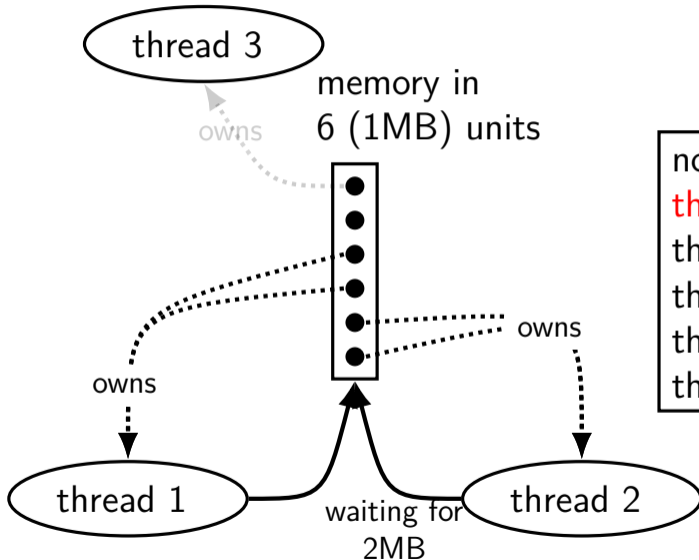


# divisible resources: not deadlock



not deadlock:  
thread 3 finishes  
then thread 1 can get memory  
then thread 1 finishes  
then thread 2 can get resources  
then thread 2 can finish

# divisible resources: not deadlock



not deadlock:

**thread 3 finishes**

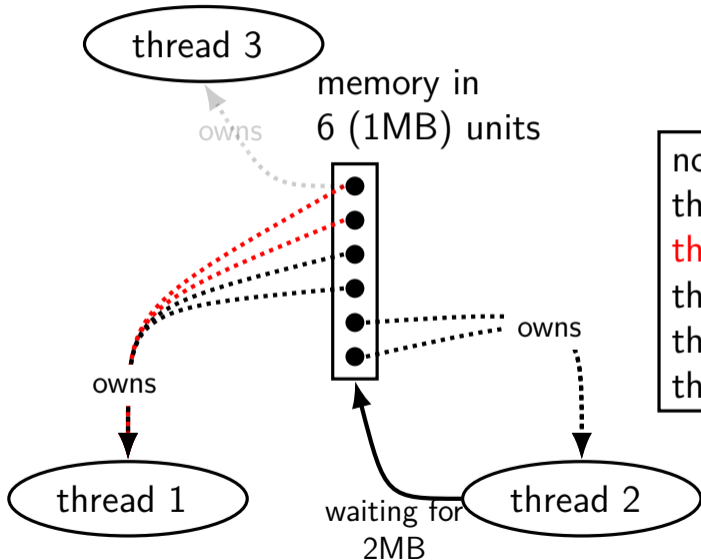
then thread 1 can get memory

then thread 1 finishes

then thread 2 can get resources

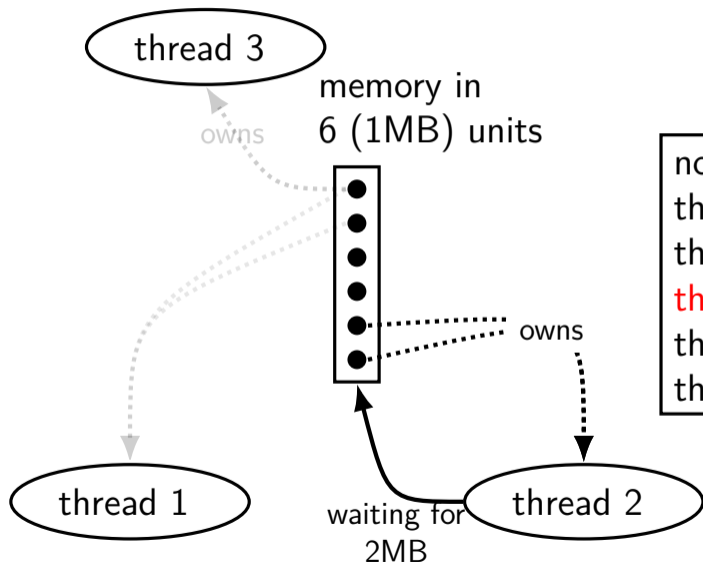
then thread 2 can finish

# divisible resources: not deadlock



not deadlock:  
thread 3 finishes  
then thread 1 can get memory  
then thread 1 finishes  
then thread 2 can get resources  
then thread 2 can finish

# divisible resources: not deadlock



not deadlock:

thread 3 finishes

then thread 1 can get memory

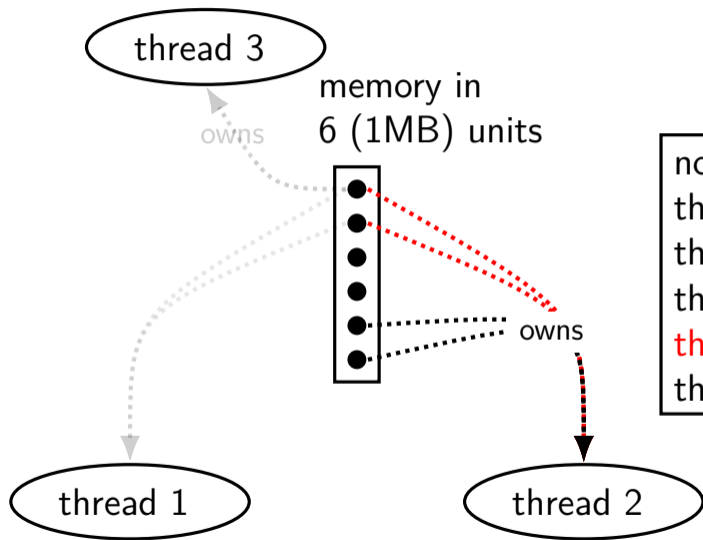
**then thread 1 finishes**

then thread 2 can get resources

then thread 2 can finish

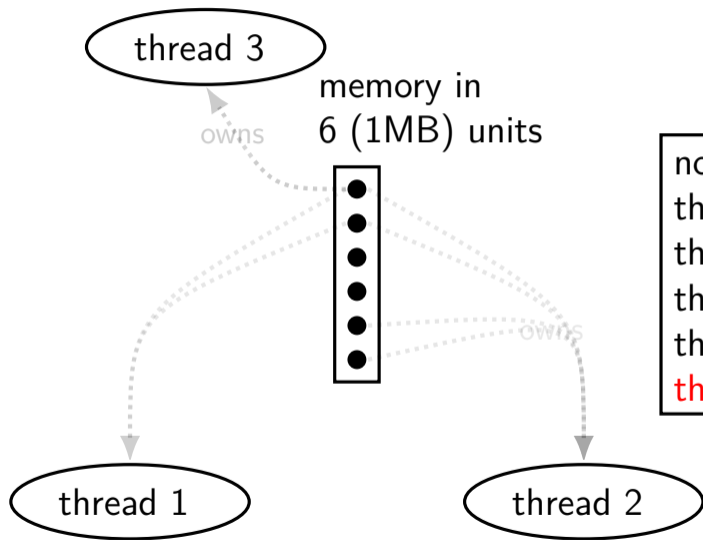


# divisible resources: not deadlock



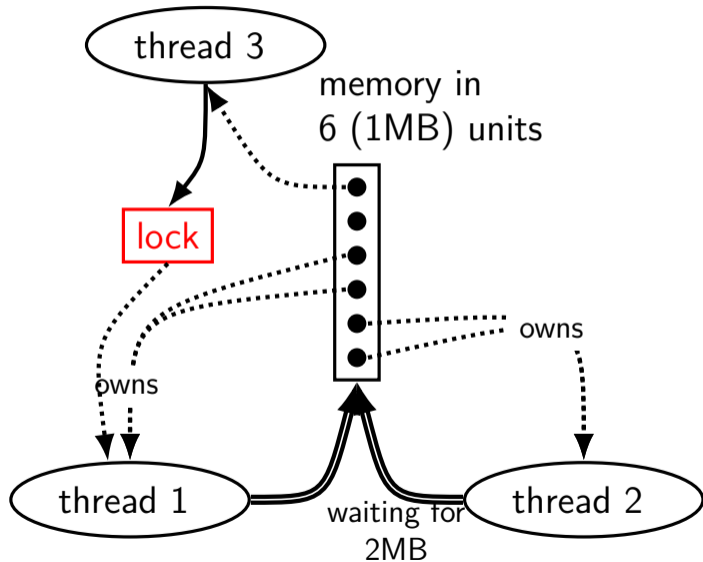
not deadlock:  
thread 3 finishes  
then thread 1 can get memory  
then thread 1 finishes  
then thread 2 can get resources  
then thread 2 can finish

# divisible resources: not deadlock

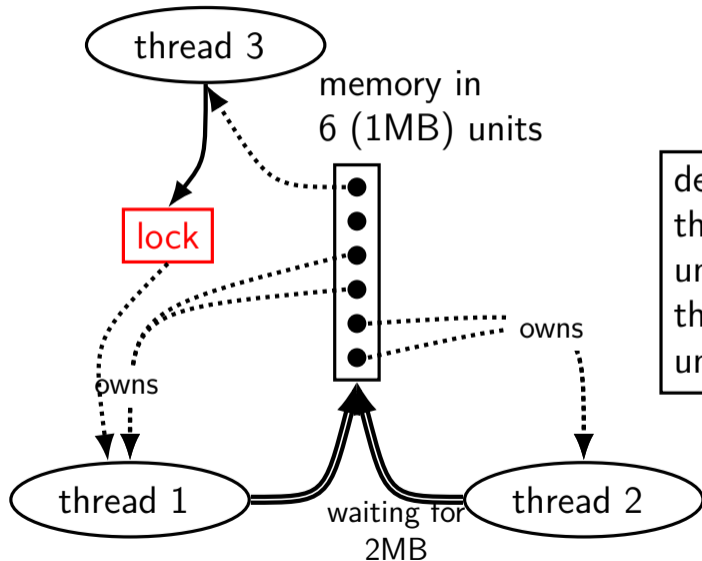


not deadlock:  
thread 3 finishes  
then thread 1 can get memory  
then thread 1 finishes  
then thread 2 can get resources  
**then thread 2 can finish**

# divisible resources: is deadlock



# divisible resources: is deadlock



deadlock:

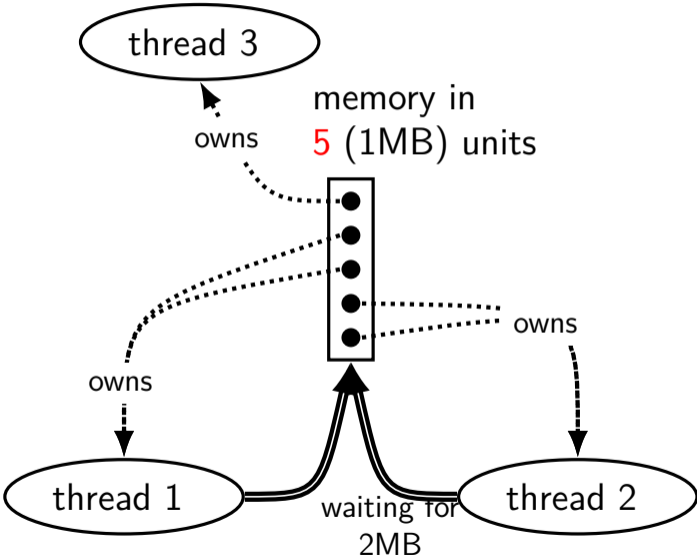
thread 3 can't finish

until thread 1 releases lock, but

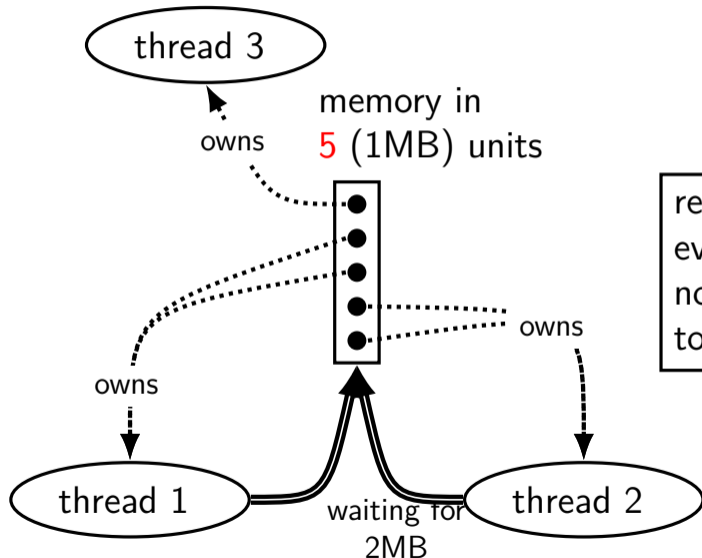
thread 1 can't finish

until thread 3 releases memory

# divisible resources: is deadlock

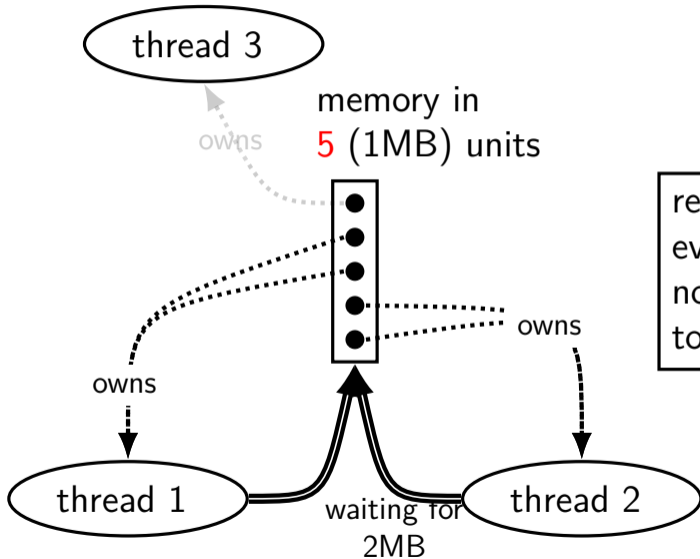


# divisible resources: is deadlock



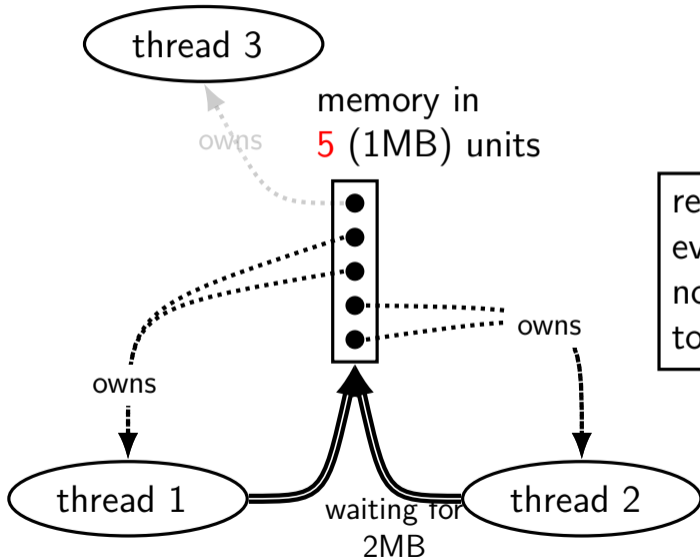
reducing memory: deadlock:  
even after thread 3 finishes  
no way for thread 1+2  
to get what they want

# divisible resources: is deadlock



reducing memory: deadlock:  
even after thread 3 finishes  
no way for thread 1+2  
to get what they want

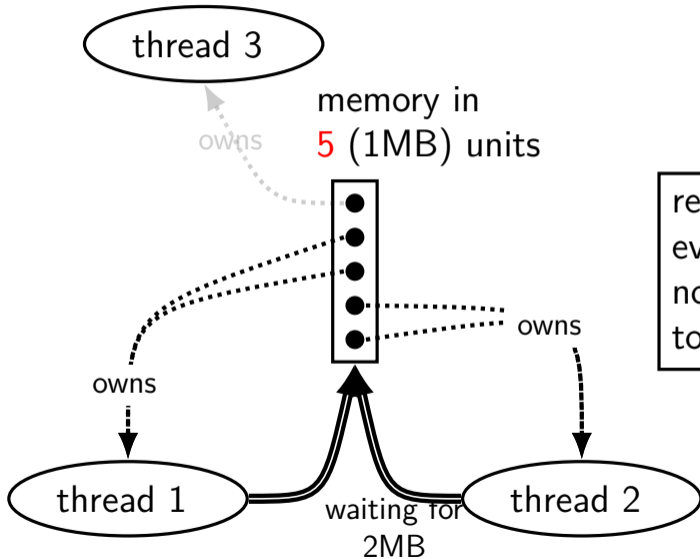
# divisible resources: is deadlock



reducing memory: deadlock:  
even after thread 3 finishes  
no way for thread 1+2  
to get what they want

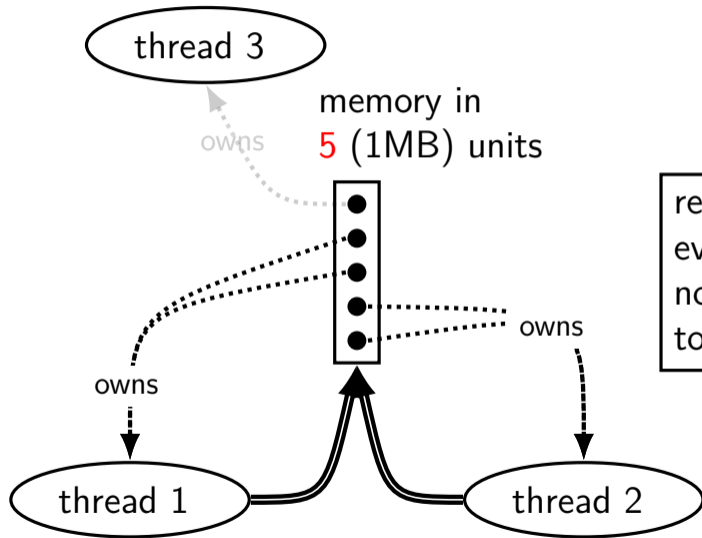


# divisible resources: is deadlock



reducing memory: deadlock:  
even after thread 3 finishes  
no way for thread 1+2  
to get what they want

# divisible resources: is deadlock



reducing memory: deadlock:  
even after thread 3 finishes  
no way for thread 1+2  
to get what they want

# deadlock detection with divisible resources

for each resource: track which threads have those resources

for each thread: resources they are waiting for

repeatedly:

- find a thread where all the resources it needs are available

- remove that thread and mark the resources it has as free — it can complete now!

either: all threads eliminated *or* found deadlock

# aside: deadlock detection in reality

requires:

- instrumenting contended resources

- “undo” to get out of deadlock

common example: for locks in a database

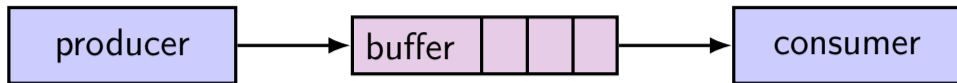
- database typically has customized locking code

- “undo” exists as side-effect of code for handling power/disk failures

related idea: *avoid* deadlock with detection on “what if” scenario

- see Banker’s algorithm

## example: producer/consumer

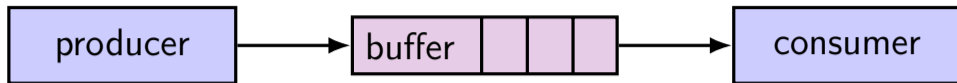


shared buffer (queue) of fixed size

one or more producers inserts into queue

one or more consumers removes from queue

## example: producer/consumer



shared buffer (queue) of fixed size

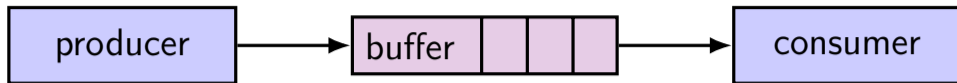
one or more producers inserts into queue

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 → compiler → assembler → 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)

pthread: build your own: provides you locks + condition variables



# monitor idea

a monitor

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

# monitor idea

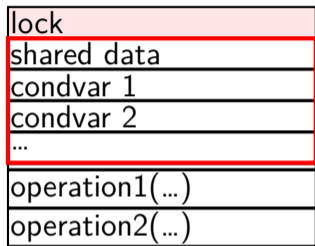
a monitor

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

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

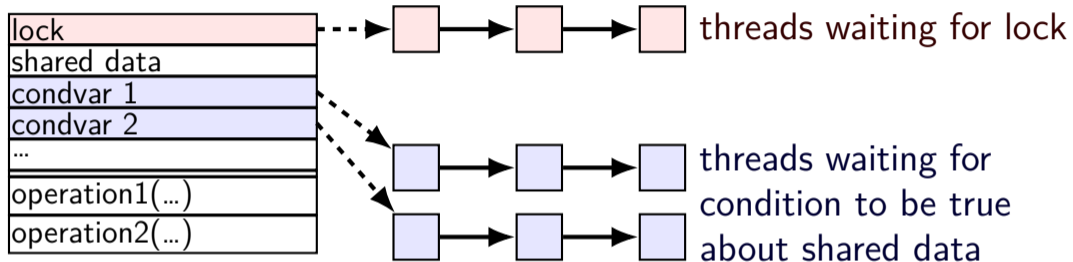
# monitor idea

a monitor



# monitor idea

a monitor



# condvar operations

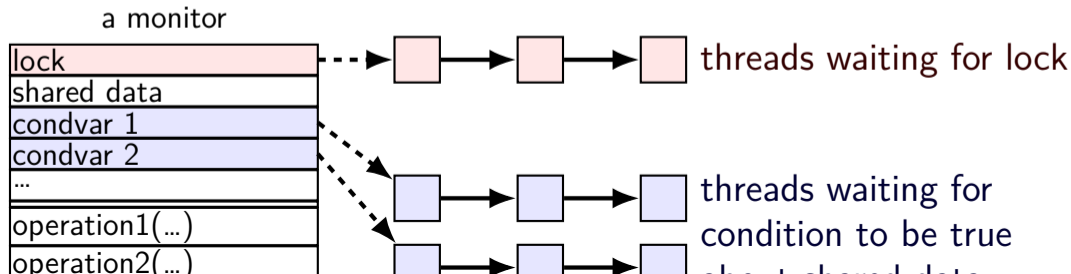
condvar operations:

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

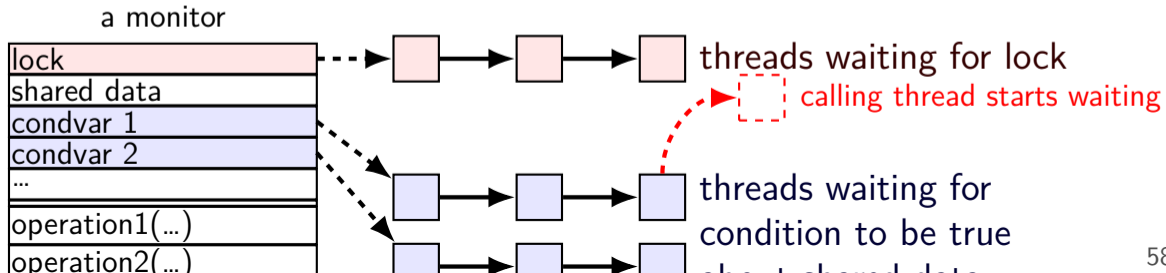
condvar operations:

**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



# condvar operations

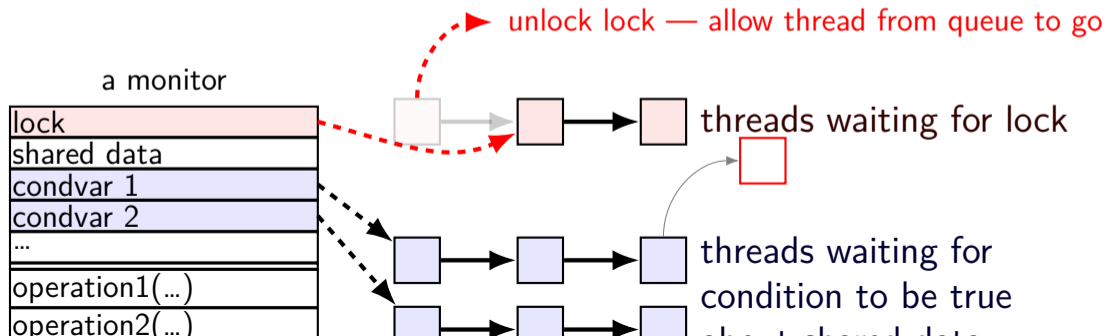
condvar operations:

`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

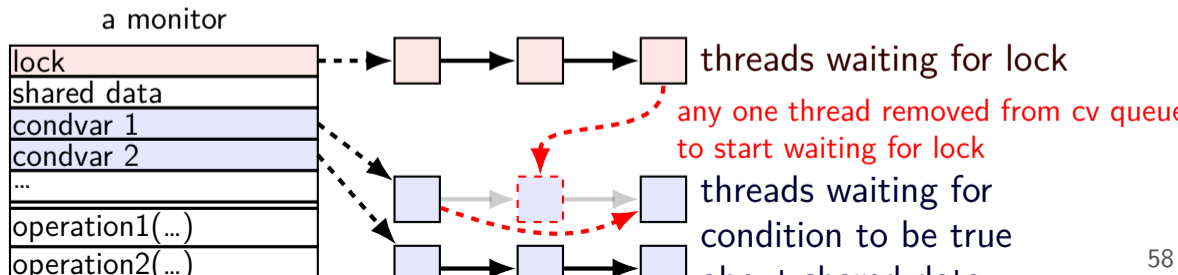
condvar operations:

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



# 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);
}
```

acquire lock before  
reading or writing finished

```
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);
}
```

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

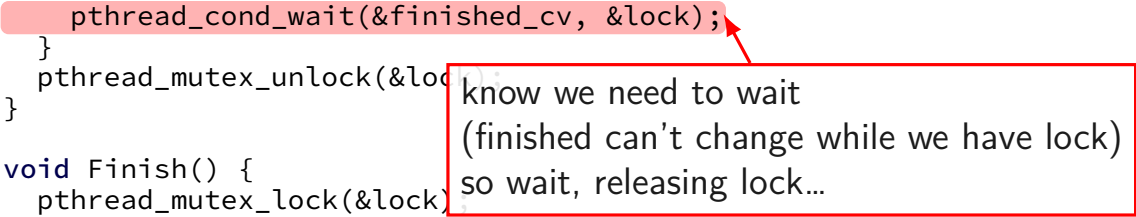
```
void Finish() {
    pthread_mutex_lock(&lock);
    finished = true;
    pthread_cond_broadcast(&finished_cv);
    pthread_mutex_unlock(&lock);
}
```

# 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);  
}
```



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

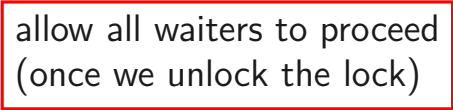
# pthread cv usage

```
// 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
<code>mutex_lock(&amp;lock)</code> (thread has lock)	
	<code>mutex_lock(&amp;lock)</code> (start waiting for lock)
<code>while (!finished) ...</code> <code>cond_wait(&amp;finished_cv, &amp;lock);</code> (start waiting for cv)	(done waiting for lock)
	<code>finished = true</code> <code>cond_broadcast(&amp;finished_cv)</code>
(done waiting for cv) (start waiting for lock)	
(done waiting for lock) <code>while (!finished) ...</code> (finished now true, so return)	<code>mutex_unlock(&amp;lock)</code>

## WaitForFinish timeline 2

WaitForFinish thread	Finish thread
	<code>mutex_lock(&amp;lock)</code> <code>finished = true</code> <code>cond_broadcast(&amp;finished_cv)</code> <code>mutex_unlock(&amp;lock)</code>
<code>mutex_lock(&amp;lock)</code> <code>while (!finished) ...</code> (finished now true, so return) <code>mutex_unlock(&amp;lock)</code>	



## 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

# unbounded buffer producer/consumer

```
pthread_mutex_t lock;  
pthread_cond_t data_ready;  
UnboundedQueue buffer;
```

```
Produce(item) {  
    pthread_mutex_lock(&lock);  
    buffer.enqueue(item);  
    pthread_cond_signal(&data_ready);  
    pthread_mutex_unlock(&lock);  
}  
Consume() {  
    pthread_mutex_lock(&lock);  
    while (buffer.empty()) {  
        pthread_cond_wait(&data_ready, &lock);  
    }  
    item = buffer.dequeue();  
    pthread_mutex_unlock(&lock);  
    return item;  
}
```

# unbounded buffer producer/consumer

```
pthread_mutex_t lock;  
pthread_cond_t data_ready;  
UnboundedQueue buffer;
```

```
Produce(item) {  
    pthread_mutex_lock(&lock);  
    buffer.enqueue(item);  
    pthread_cond_signal(&data_ready);  
    pthread_mutex_unlock(&lock);  
}  
Consume() {  
    pthread_mutex_lock(&lock);  
    while (buffer.empty()) {  
        pthread_cond_wait(&data_ready, &lock);  
    }  
    item = buffer.dequeue();  
    pthread_mutex_unlock(&lock);  
    return item;  
}
```

rule: never touch buffer  
without acquiring lock

otherwise: what if two threads  
simultaneously en/dequeue?  
(both use same array/linked list entry?)  
(both reallocate array?)

# unbounded buffer producer/consumer

```
pthread_mutex_t lock;  
pthread_cond_t data_ready;  
UnboundedQueue buffer;
```

```
Produce(item) {  
    pthread_mutex_lock(&lock);  
    buffer.enqueue(item);  
    pthread_cond_signal(&data_ready);  
    pthread_mutex_unlock(&lock);  
}
```

```
Consume() {  
    pthread_mutex_lock(&lock);  
    while (buffer.empty()) {  
        pthread_cond_wait(&data_ready, &lock);  
    }  
    item = buffer.dequeue();  
    pthread_mutex_unlock(&lock);  
    return item;  
}
```

check if empty  
if so, dequeue

okay because have lock  
other threads **cannot** dequeue here

}

# unbounded buffer producer/consumer

```
pthread_mutex_t lock;  
pthread_cond_t data_ready;  
UnboundedQueue buffer;
```

```
Produce(item) {  
    pthread_mutex_lock(&lock);  
    buffer.enqueue(item);  
    pthread_cond_signal(&data_ready);  
    pthread_mutex_unlock(&lock);  
}
```

wake one Consume thread  
*if any are waiting*



```
Consume() {  
    pthread_mutex_lock(&lock);  
    while (buffer.empty()) {  
        pthread_cond_wait(&data_ready, &lock);  
    }  
    item = buffer.dequeue();  
    pthread_mutex_unlock(&lock);  
    return item;  
}
```

# unbounded buffer producer/consumer

```
pthread_mutex_t lock;  
pthread_cond_t data_ready;  
UnboundedQueue buffer;
```

```
Produce(item) {  
    pthread_mutex_lock(&lock);  
    buffer.enqueue(item);  
    pthread_cond_signal(&data_ready);  
    pthread_mutex_unlock(&lock);  
}
```

```
Consume() {  
    pthread_mutex_lock(&lock);  
    while (buffer.empty()) {  
        pthread_cond_wait(&data_ready, &lock);  
    }  
    item = buffer.dequeue();  
    pthread_mutex_unlock(&lock);  
    return item;  
}
```

Thread 1

Produce()
...lock
...enqueue
...signal
...unlock

Thread 2

Consume()
...lock
...empty? no
...dequeue
...unlock
return

0 iterations: Produce() called before Consume()  
1 iteration: Produce() signalled, probably  
2+ iterations: spurious wakeup or ...?

# unbounded buffer producer/consumer

```
pthread_mutex_t lock;  
pthread_cond_t data_ready;  
UnboundedQueue buffer;
```

```
Produce(item) {  
    pthread_mutex_lock(&lock);  
    buffer.enqueue(item);  
    pthread_cond_signal(&data_ready);  
    pthread_mutex_unlock(&lock);  
}
```

```
Consume() {  
    pthread_mutex_lock(&lock);  
    while (buffer.empty()) {  
        pthread_cond_wait(&data_ready, &lock);  
    }  
    item = buffer.dequeue();  
    pthread_mutex_unlock(&lock);  
    return item;  
}
```

Thread 1	Thread 2
	Consume()
	...lock
	...empty? yes
	...unlock/start wait
Produce()	waiting for data_ready
...lock	
...enqueue	
...signal	stop wait
...unlock	lock
	...empty? no
	...dequeue
	...unlock
	return

0 iterations: Produce() called before Consume()  
1 iteration: Produce() signalled, probably  
2+ iterations: spurious wakeup or ...?

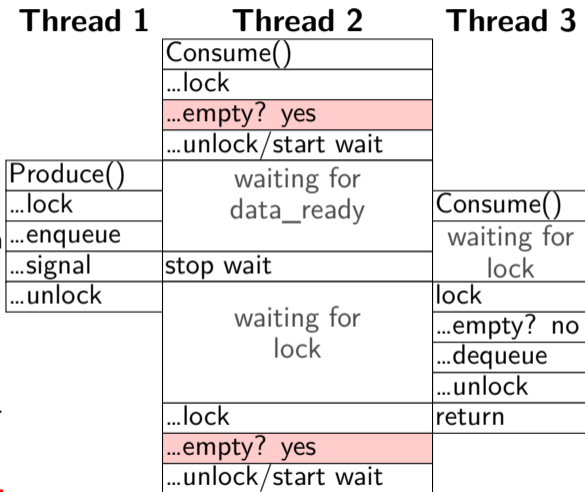


# unbounded buffer producer/consumer

```
pthread_mutex_t lock;  
pthread_cond_t data_ready;  
UnboundedQueue buffer;
```

```
Produce(item) {  
    pthread_mutex_lock(&lock);  
    buffer.enqueue(item);  
    pthread_cond_signal(&data_ready);  
    pthread_mutex_unlock(&lock);  
}
```

```
Consume() {  
    pthread_mutex_lock(&lock);  
    while (buffer.empty()) {  
        pthread_cond_wait(&data_ready,  
                          &lock);  
    }  
    item = buffer.dequeue();  
    pthread_mutex_unlock(&lock);  
    return item;  
}
```



0 iterations: Produce() called before Consume()  
1 iteration: Produce() signalled, probably  
**2+ iterations: spurious wakeup or ...?**

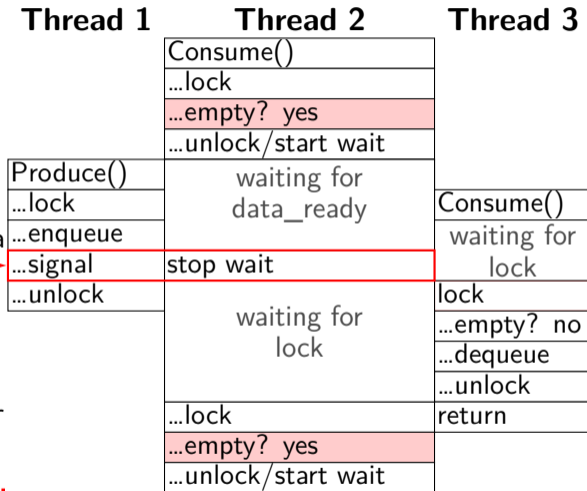
# unbounded buffer producer/consumer

```
pthread_mutex_t lock;
pthread_cond_t data_ready;
UnboundedQueue buffer;
```

in pthreads: signalled thread not guaranteed to hold lock next

alternate design: signalled thread gets lock next called "Hoare scheduling" not done by pthreads, Java, ...

```
pthread_cond_wait(&data_r
}
item = buffer.dequeue();
pthread_mutex_unlock(&lock);
return item;
}
```



0 iterations: Produce() called before Consume()  
 1 iteration: Produce() signalled, probably  
 2+ iterations: spurious wakeup or ...?

# Hoare versus Mesa monitors

Hoare-style monitors

signal 'hands off' lock to awoken thread

Mesa-style monitors

any eligible thread gets lock next  
(maybe some other idea of priority?)

every current threading library I know of does Mesa-style

# backup slides



# backup slides

# pipe() deadlock

**BROKEN** example:

```
int child_to_parent_pipe[2], parent_to_child_pipe[2];
pipe(child_to_parent_pipe); pipe(parent_to_child_pipe);
if (fork() == 0) {
    /* child */
    write(child_to_parent_pipe[1], buffer, HUGE_SIZE);
    read(parent_to_child_pipe[0], buffer, HUGE_SIZE);
    exit(0);
} else {
    /* parent */
    write(parent_to_child_pipe[1], buffer, HUGE_SIZE);
    read(child_to_parent_pipe[0], buffer, HUGE_SIZE);
}
```

This will **hang forever** (if `HUGE_SIZE` is big enough).

# deadlock waiting

child writing to pipe waiting for free buffer space

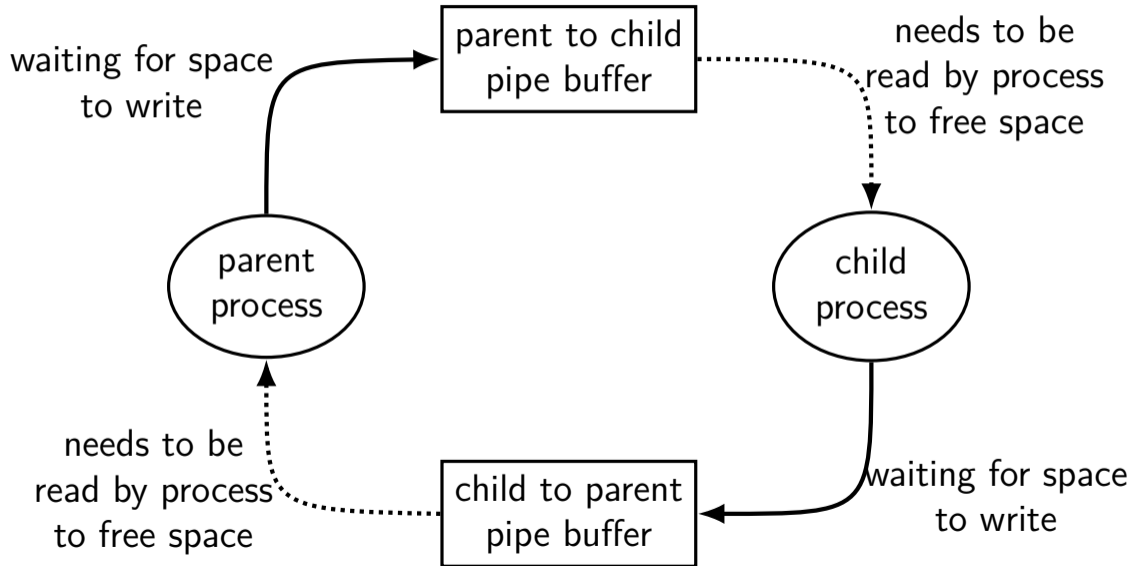
...which will not be available until parent reads

parent writing to pipe waiting for free buffer space

...which will not be available until child reads



# circular dependency



# allocating all at once?

for resources like disk space, memory

figure out maximum allocation **when starting thread**

“only” need conservative estimate

only start thread if those resources are available

okay solution for embedded systems?

# deadlock with free space

## Thread 1

AllocateOrWaitFor(1 MB)

AllocateOrWaitFor(1 MB)

(do calculation)

Free(1 MB)

Free(1 MB)

## Thread 2

AllocateOrWaitFor(1 MB)

AllocateOrWaitFor(1 MB)

(do calculation)

Free(1 MB)

Free(1 MB)

2 MB of space — deadlock possible with unlucky order

# deadlock with free space (unlucky case)

## Thread 1

AllocateOrWaitFor(1 MB)

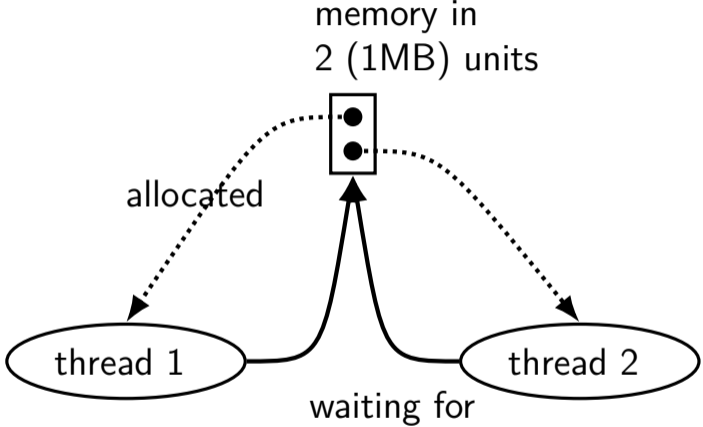
AllocateOrWaitFor(1 MB... stalled

## Thread 2

AllocateOrWaitFor(1 MB)

AllocateOrWaitFor(1 MB... stalled

# free space: dependency graph



# deadlock with free space (lucky case)

## Thread 1

```
AllocateOrWaitFor(1 MB)
AllocateOrWaitFor(1 MB)
(do calculation)
Free(1 MB);
Free(1 MB);
```

## Thread 2

```
AllocateOrWaitFor(1 MB)
AllocateOrWaitFor(1 MB)
(do calculation)
Free(1 MB);
Free(1 MB);
```

# AllocateOrFail

## Thread 1

AllocateOrFail(1 MB)

AllocateOrFail(1 MB) **fails!**

Free(1 MB) **(cleanup after failure)**

## Thread 2

AllocateOrFail(1 MB)

AllocateOrFail(1 MB) **fails!**

Free(1 MB) **(cleanup after failure)**

okay, now what?

give up?

both try again? — maybe this will keep happening? (called **livelock**)

try one-at-a-time? — gaurenteed to work, but tricky to implement

# AllocateOrSteal

## Thread 1

AllocateOrSteal(1 MB)

AllocateOrSteal(1 MB)

(do work)

## Thread 2

AllocateOrSteal(1 MB)

Thread killed to free 1MB

problem: can one actually implement this?

problem: can one kill thread and keep system in consistent state?



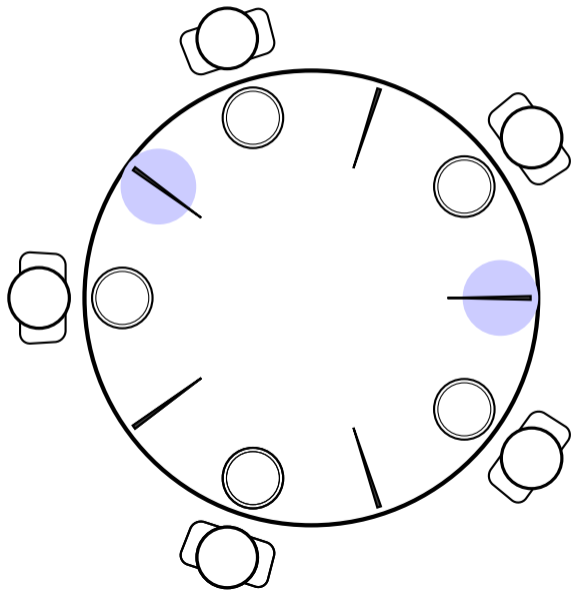
# fail/steal with locks

pthread provides `pthread_mutex_trylock` — “lock or fail”

some databases implement *revocable locks*

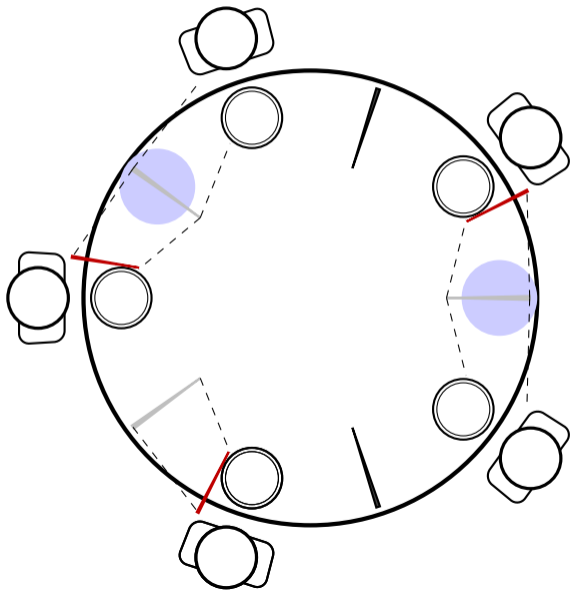
do equivalent of throwing exception in thread to ‘steal’ lock  
need to carefully arrange for operation to be cleaned up

# dining philosophers — ordering



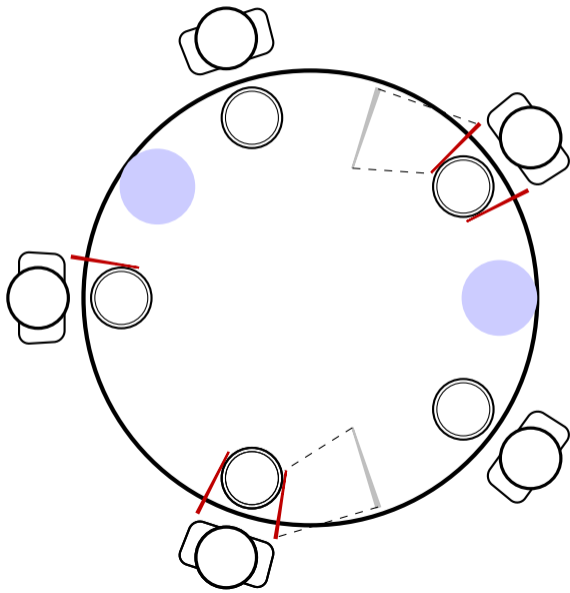
mark some chopsticks places  
rule: grab from marked place first  
only grab other chopstick after that

# dining philosophers — ordering



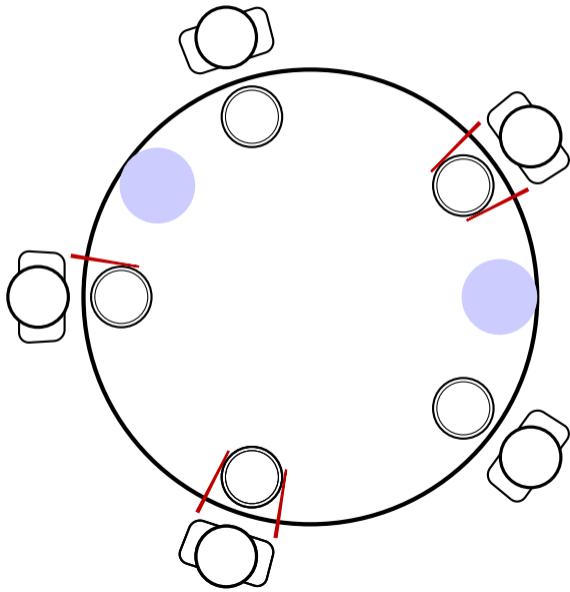
mark some chopsticks places  
rule: grab from marked place first  
only grab other chopstick after that

# dining philosophers — ordering



mark some chopsticks places  
rule: grab from marked place first  
only grab other chopstick after that

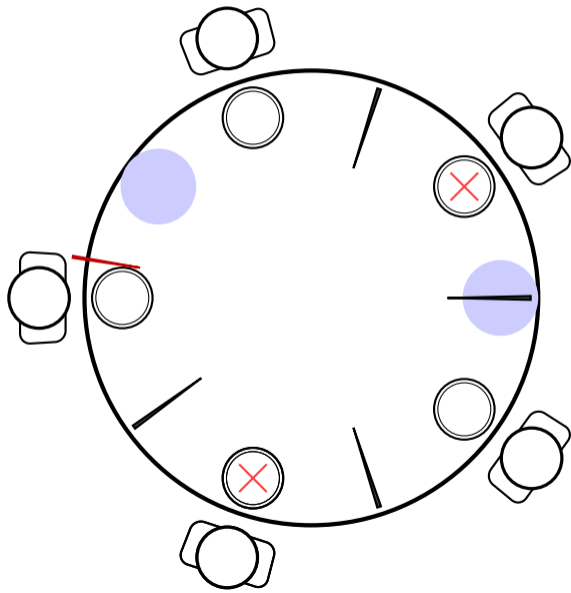
# dining philosophers — ordering



mark some chopsticks places  
rule: grab from marked place first  
only grab other chopstick after that

avoids circular dependency,  
means everyone else  
**eventually** gets a turn

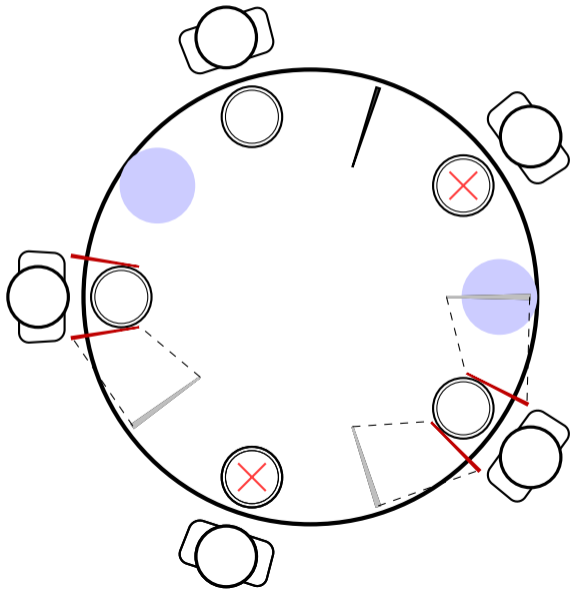
# dining philosophers — ordering



mark some chopsticks places  
rule: grab from marked place first  
only grab other chopstick after that

avoids circular dependency,  
means everyone else  
**eventually** gets a turn

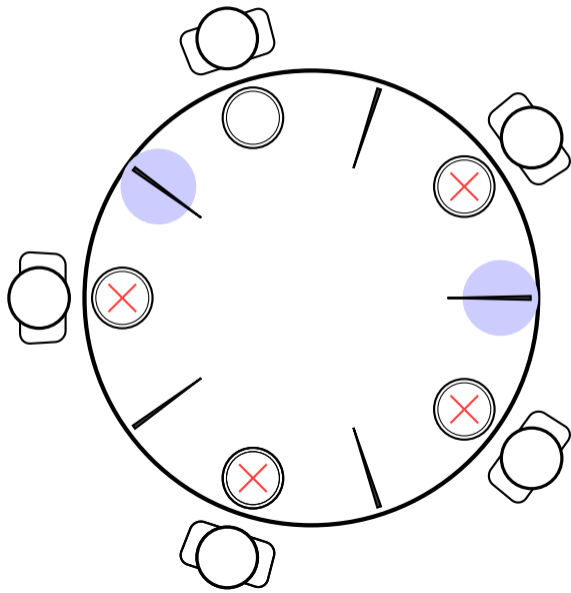
# dining philosophers — ordering



mark some chopsticks places  
rule: grab from marked place first  
only grab other chopstick after that

avoids circular dependency,  
means everyone else  
**eventually** gets a turn

## dining philosophers — ordering

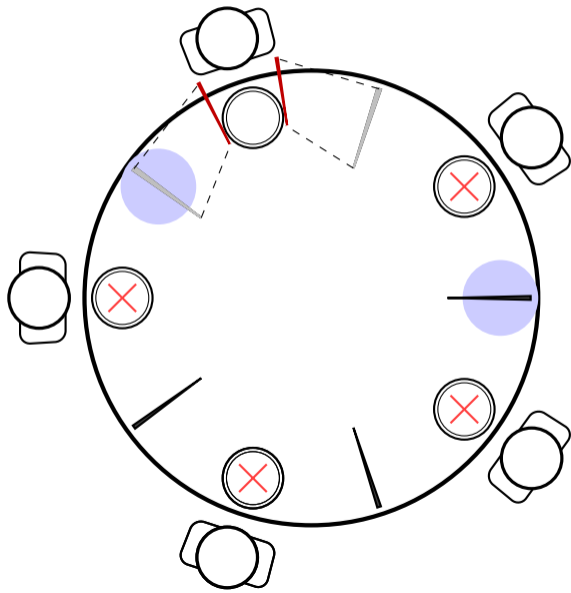


mark some chopsticks places  
rule: grab from marked place first  
only grab other chopstick after that

avoids circular dependency,  
means everyone else  
**eventually** gets a turn



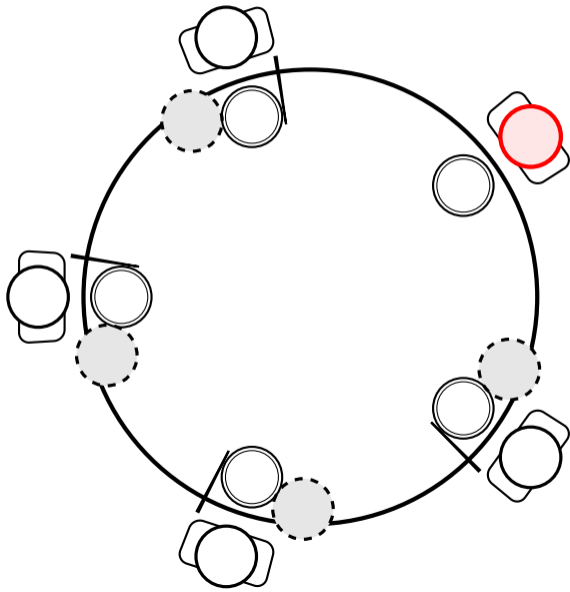
# dining philosophers — ordering



mark some chopsticks places  
rule: grab from marked place first  
only grab other chopstick after that

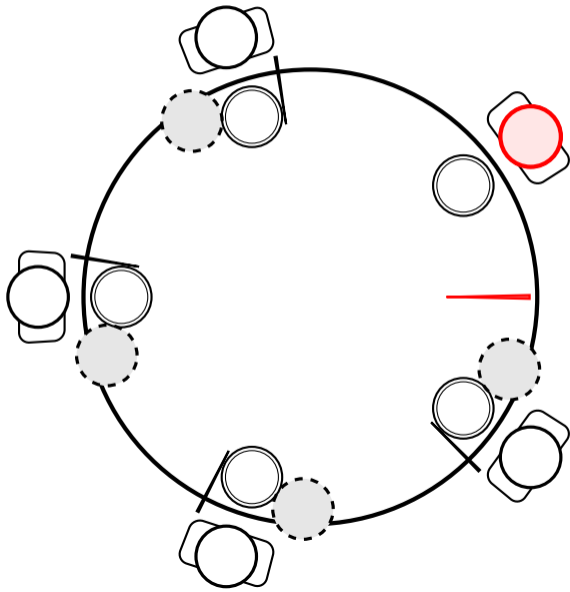
avoids circular dependency,  
means everyone else  
**eventually** gets a turn

# dining philosophers — aborting



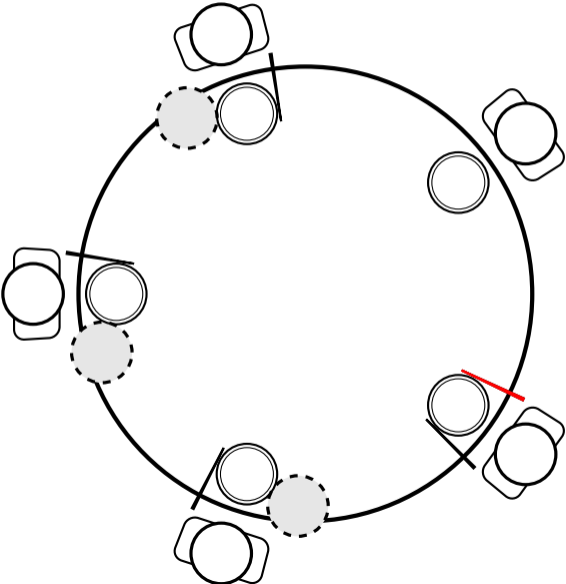
dining philosopher  
what if someone's impatient  
just gives up instead of waiting

# dining philosophers — aborting



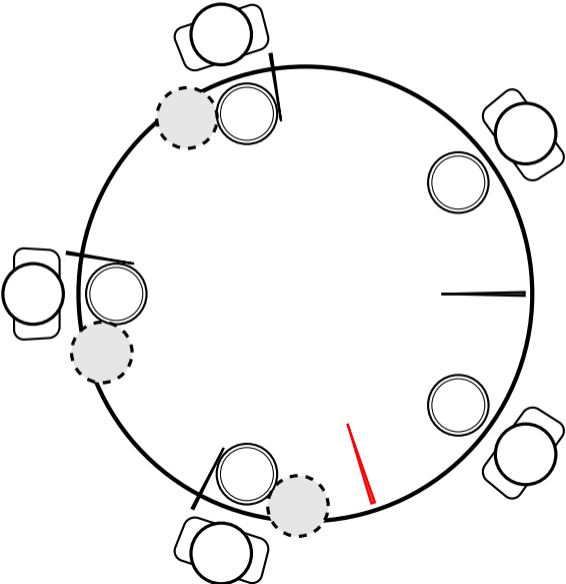
dining philosopher  
what if someone's impatient  
just gives up instead of waiting

# dining philosophers — aborting



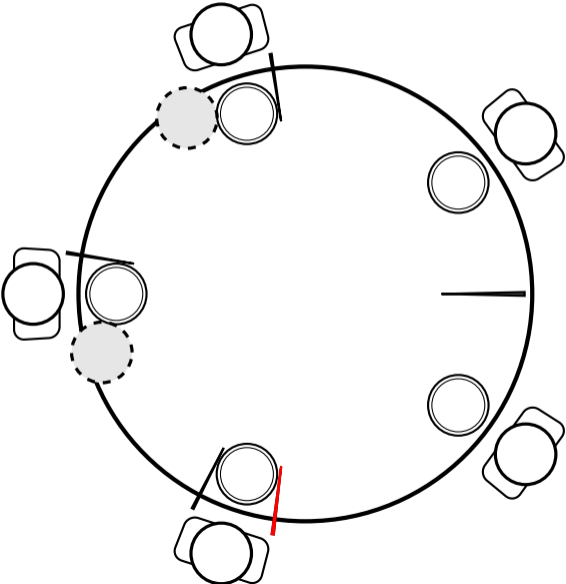
now everyone else can eat

# dining philosophers — aborting



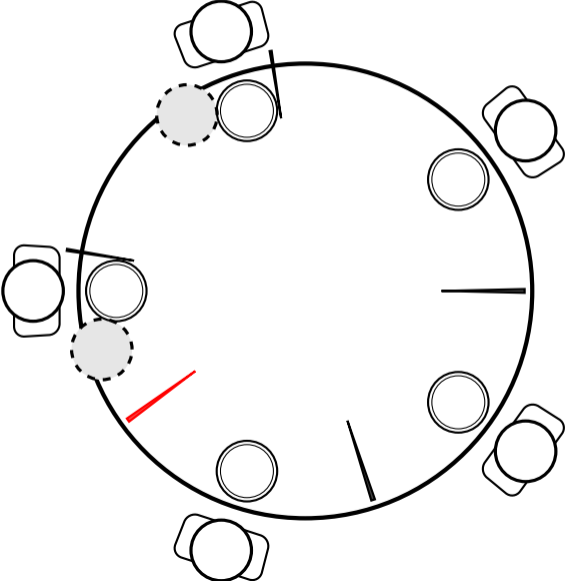
now everyone else can eat

# dining philosophers — aborting



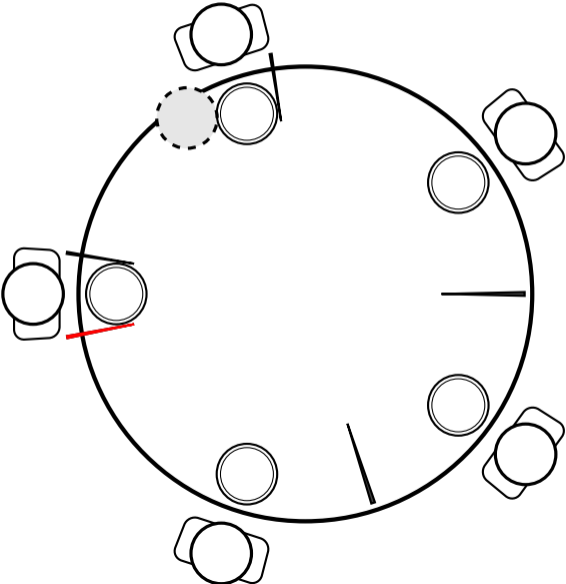
now everyone else can eat

# dining philosophers — aborting



now everyone else can eat

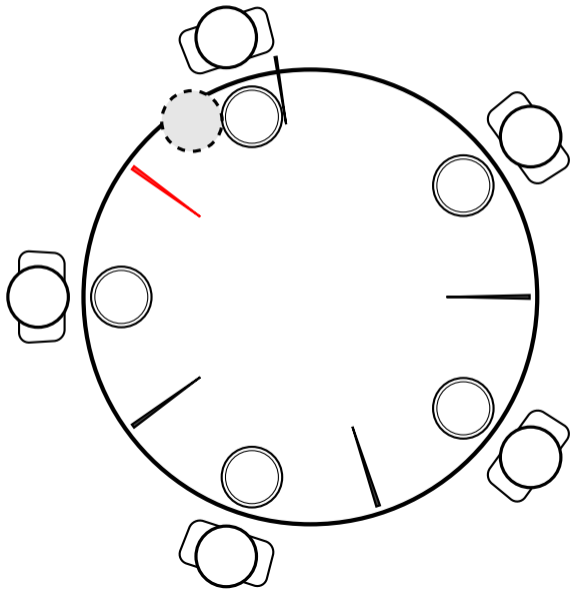
# dining philosophers — aborting



now everyone else can eat

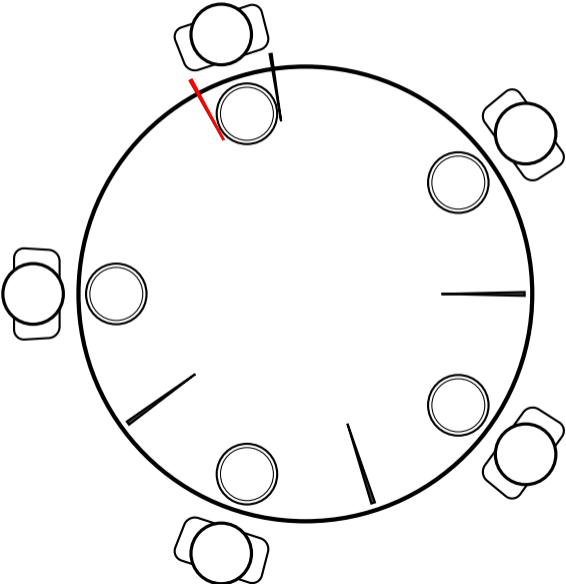


# dining philosophers — aborting



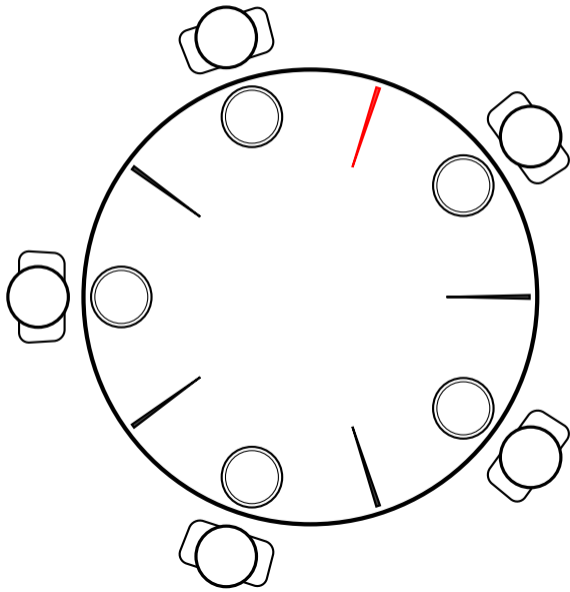
now everyone else can eat

# dining philosophers — aborting



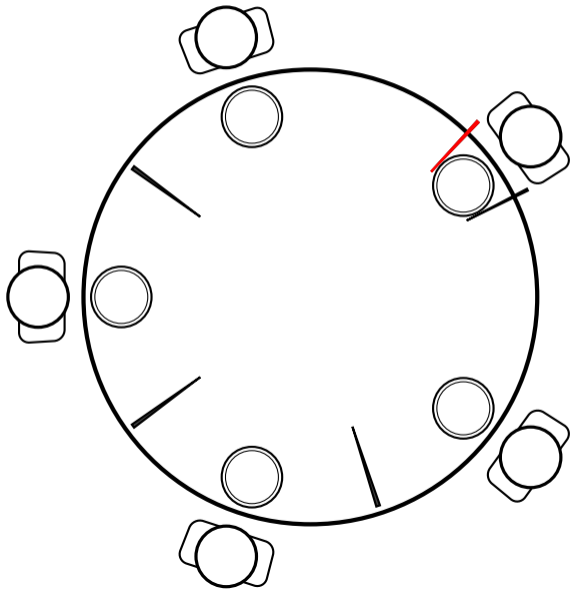
now everyone else can eat

# dining philosophers — aborting



now everyone else can eat

# dining philosophers — aborting



and person who gave up  
might succeed later

## using deadlock detection for prevention

suppose you know the *maximum resources* a process could request

make decision **when starting process** (“*admission control*”)

# using deadlock detection for prevention

suppose you know the *maximum resources* a process could request

make decision **when starting process** (“*admission control*”)

ask “what if every process was waiting for maximum resources”  
including the one we’re starting

would it cause deadlock? then **don’t let it start**

called Banker’s algorithm



# backup slides



## 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);
```

# life homework even/odd

naive way has an operation that needs locking:

```
for (int time = 0; time < MAX_ITERATIONS; ++time) {  
    ... compute to_grid ...  
    swap(from_grid, to_grid);  
}
```

but this alternative needs less locking:

```
Grid grids[2];  
for (int time = 0; time < MAX_ITERATIONS; ++time) {  
    from_grid = &grids[time % 2];  
    to_grid = &grids[(time % 2) + 1];  
    ... compute to_grid ...  
}
```

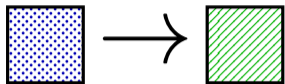
# life homework even/odd

naive way has an operation that needs locking:

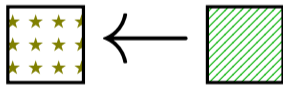
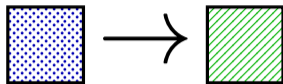
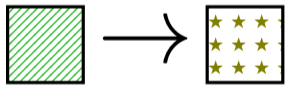
```
for (int time = 0; time < MAX_ITERATIONS; ++time) {  
    ... compute to_grid ...  
    swap(from_grid, to_grid);  
}
```

but this alternative needs less locking:

```
Grid grids[2];  
for (int time = 0; time < MAX_ITERATIONS; ++time) {  
    from_grid = &grids[time % 2];  
    to_grid = &grids[(time % 2) + 1];  
    ... compute to_grid ...  
}
```



swap



## 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 the_lock
    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 t

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)
                        // if lock was already locked, %eax is 1, so it will be swapped with 0
                        // of t

test %eax, %eax
jne acquire
ret
```

if lock was already locked retry  
“spin” until lock is released elsewhere

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)
```

```
test %eax, %eax
jne acquire
ret
```

release lock by setting it to 0 (not taken)  
allows looping acquire to finish

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)
```

```
test %eax, %eax
jne acquire
ret
```

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

release:

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

## 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 is
```

## exercise: spin wait

finished: .quad 0

ThreadFinish:

```
-----A-----  
ret
```

ThreadWaitForFinish:

```
-----B-----  
lock xchg %eax, finished  
cmp $0, %eax  
__C_ ThreadWaitForFinish  
ret
```

*/\* or without using a writing instr*

```
mov %eax, finished  
mfence  
cmp $0, %eax  
je ThreadWaitForFinish  
ret
```

A. mfence; mov \$1, finished

B. mov \$1, finished; mfence

C. mov \$0, %eax E. je

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 as waiting
        /* xv6: myproc()->state = SLEEPING; */
        UnlockSpinlock(&m->guard_spinlock);
        run scheduler (context switch)
    } 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
        mark thread as no longer waiting
        /* 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 as waiting
        /* xv6: myproc()->state = SLEEPING; */
        UnlockSpinlock(&m->guard_spinlock);
        run scheduler (context switch)
    } 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
        mark thread as no longer waiting
        /* 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;
};
```

subtly: if UnlockMutex runs here on another core  
need to make sure scheduler on the other core doesn't switch to thread  
while it is still running (would 'clone' thread/mess up registers)

```
LockMutex(Mutex *m) {
    LockSpinlock(&m->guard_spinlock);
    if (m->lock_taken) {
        put current thread on m->wait_queue
        mark current thread as waiting
        /* xv6: myproc()->state = SLEEPING; */
        UnlockSpinlock(&m->guard_spinlock);
        run scheduler (context switch)
    } 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
        mark thread as no longer waiting
        /* 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;
};
```

```
LockMutex(Mutex *m) {
    LockSpinlock(&m->guard_spinlock);
    if (m->lock_taken) {
        put current thread on m->wait_queue
        mark current thread as waiting
        /* xv6: myproc()->state = SLEEPING; */
        UnlockSpinlock(&m->guard_spinlock);
        run scheduler (context switch)
    } 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
        mark thread as no longer waiting
        /* 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)	
start LockMutex		
acquire spinlock		
discover lock taken		
enqueue thread A		
thread A set not runnable		
release spinlock	start UnlockMutex	
	thread A set runnable	
	finish UnlockMutex	
	run scheduler	
	scheduler switches to A	
	...with old version of registers	
thread A runs scheduler		...
...finally saving registers		...

Linux soln.: track 'thread running' separately from 'thread

# mutex and scheduler subtly

core 0 (thread A)	core 1 (thread B)	
start LockMutex		
acquire spinlock		
discover lock taken		
enqueue thread A		
thread A set not runnable		
release spinlock	start UnlockMutex	
	thread A set runnable	
	finish UnlockMutex	
	run scheduler	
	scheduler switches to A	
	...with old version of registers	
thread A runs scheduler		...
...finally saving registers		...

Linux soln.: track 'thread running' separately from 'thread

# 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

# 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() {  
    disable interrupts  
}
```

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

# naive interrupt enable/disable (1)

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

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

## naive interrupt enable/disable (2)

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

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

## naive interrupt enable/disable (2)

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

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

## naive interrupt enable/disable (2)

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

# C++ containers and locking

can you use a vector from multiple threads?

...question: how is it implemented?

# C++ containers and locking

can you use a vector from multiple threads?

...question: how is it implemented?

- dynamically allocated array
- reallocated on size changes

# C++ containers and locking

can you use a vector from multiple threads?

...question: how is it implemented?

- dynamically allocated array
- reallocated on size changes

can access from multiple threads ...as long as not  
append/erase/etc.?

assuming it's implemented like we expect...

- but can we really depend on that?

- e.g. could shrink internal array after a while with no expansion save memory?



# C++ standard rules for containers

multiple threads can **read anything at the same time**

can only read element **if no other thread is modifying it**

can safely **add/remove elements if no other threads** are accessing container

(sometimes can safely add/remove in extra cases)

exception: vectors of bools — can't safely read and write at same time

might be implemented by putting multiple bools in one int

# 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 f
    jne .L3
    cmpl $0, no_milk
```

# exercise: fetch-and-add with compare-and-swap

exercise: implement fetch-and-add with compare-and-swap

```
compare_and_swap(address, old_value, new_value) {  
    if (memory[address] == old_value) {  
        memory[address] = new_value;  
        return true;    // x86: set ZF flag  
    } else {  
        return false;  // x86: clear ZF flag  
    }  
}
```

# solution

```
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;  
    }  
}
```

## xv6 spinlock: acquire

```
void
acquire(struct spinlock *lk)
{
    pushcli(); // disable interrupts to avoid deadlock.
    ...
    // The xchg is atomic.
    while(xchg(&lk->locked, 1) != 0)
        ;

    // Tell the C compiler and the processor to not move loads or stores
    // past this point, to ensure that the critical section's memory
    // references happen after the lock is acquired.
    __sync_synchronize();
    ...
}
```

# xv6 spinlock: acquire

```
void
acquire(struct spinlock *lk)
{
    pushcli(); // disable interrupts to avoid deadlock.
    ...
    // The xchg is atomic.
    while(xchg(&lk->locked, 1) != 0)
        ;

    // don't let us be interrupted after while have the lock
    // problem: interruption might try to do something with the lock
    ...but that can never succeed until we release the lock
    ...
    ...but we won't release the lock until interruption finishes
}
```

*or store  
memory*

# xv6 spinlock: acquire

```
void
acquire(struct spinlock *lk)
{
    pushcli(); // disable interrupts to avoid deadlock.
    ...
    // The xchg is atomic.
    while(xchg(&lk->locked, 1) != 0)
        ;

    // Tell the C compiler and the processor to not move loads or stores
    // past this point, to ensure that the critical section's memory
    // references happen after the lock is acquired.
    __sync_synchronize();
    ...
}
```

xchg wraps the lock xchg instruction  
same loop as before



## xv6 spinlock: acquire

```
void
acquire(struct spinlock *lk)
{
    pushcli(); // disable interrupts to avoid deadlock.
    ...
    // The xchg is atomic.
    while(xchg(&lk->locked, 1) != 0)
        ;

    // Tell the C compiler and the processor to not move loads or stores
    // past this point, to ensure that the critical section's memory
    // avoid load store reordering (including by compiler)
    -- on x86, xchg alone is enough to avoid processor's reordering
    .. (but compiler may need more hints)
}
```

# xv6 spinlock: release

```
void
```

```
release(struct spinlock *lk)
```

```
...
```

```
// Tell the C compiler and the processor to not move loads or stores
```

```
// past this point, to ensure that all the stores in the critical
```

```
// section are visible to other cores before the lock is released.
```

```
// Both the C compiler and the hardware may re-order loads and
```

```
// stores; __sync_synchronize() tells them both not to.
```

```
__sync_synchronize();
```

```
// Release the lock, equivalent to lk->locked = 0.
```

```
// This code can't use a C assignment, since it might
```

```
// not be atomic. A real OS would use C atomics here.
```

```
asm volatile("movl $0, %0" : "+m" (lk->locked) : );
```

```
popcli();
```

```
}
```

# xv6 spinlock: release

```
void
```

```
release(struct spinlock *lk)
```

```
...
```

```
// Tell the C compiler and the processor to not move loads or stores
```

```
// past this point, to ensure that all the stores in the critical
```

```
// section are visible to other cores before the lock is released
```

```
// Both the C compiler and the hardware may re-order loads and
```

```
// stores; __sync_synchronize() tells them both not to.
```

```
__sync_synchronize();
```

```
// Release the lock, equivalent to lk->locked = 0.
```

```
// This code can't use a C assignment, since it might
```

```
// not
```

```
asm volatile
```

turns into instruction to tell processor not to reorder  
plus tells compiler not to reorder

```
popcli(),
```

```
}
```

# xv6 spinlock: release

```
void  
release(struct spinlock *lk)
```

```
...
```

```
// Tell the C compiler and the processor to not move loads or stores  
// past this point, to ensure that all the stores in the critical  
// section are visible to other cores before the lock is released.  
// Both the C compiler and the hardware may re-order loads and  
// stores; __sync_synchronize() tells them both not to.  
__sync_synchronize();
```

```
// Release the lock, equivalent to lk->locked = 0.  
// This code can't use a C assignment, since it might  
// not be atomic. A real OS would use C atomics here.
```

```
asm volatile("movl $0, %0" : "+m" (lk->locked) : );
```

turns into mov of constant 0 into lk->locked

```
popcli(),
```

```
}
```

# xv6 spinlock: release

```
void  
release(struct spinlock *lk)
```

```
...
```

```
// Tell the C compiler and the processor to not move loads or stores  
// past this point, to ensure that all the stores in the critical  
// section are visible to other cores before the lock is released.  
// Both the C compiler and the hardware may re-order loads and  
// stores; __sync_synchronize() tells them both not to.
```

```
__sync_synchronize();
```

```
// Release the lock, equivalent to lk->locked = 0.  
// This code can't use a C assignment, since it might  
// not be atomic. A real OS would use C atomics here.
```

```
asm volatile("movl $0, %0" : "=r" (lk->locked) : : "memory");
```

reenable interrupts (taking nested locks into account)

```
popcli(),
```

```
}
```

# 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) {  
    ...  
}
```



# append to singly-linked list

```
/* assumption: other threads may be appending to list,  
 *           but nodes are not being removed, reordered, etc.  
 */  
void append_to_list(ListNode *head, ListNode *new_last_node) {  
    memory_ordering_fence();  
    ListNode *current_last_node;  
    do {  
        current_last_node = head;  
        while (current_last_node->next) {  
            current_last_node = current_last_node->next;  
        }  
    } while (  
        !compare_and_swap(&current_last_node->next,  
                          NULL, new_last_node)  
    );  
}
```

# some common atomic operations (1)

```
// x86: emulate with exchange  
test_and_set(address) {  
    old_value = memory[address];  
    memory[address] = 1;  
    return old_value != 0; // e.g. set ZF flag  
}
```

```
// x86: xchg REGISTER, (ADDRESS)  
exchange(register, address) {  
    temp = memory[address];  
    memory[address] = register;  
    register = temp;  
}
```

## some common atomic operations (2)

```
// x86: mov OLD_VALUE, %eax; lock cmpxchg NEW_VALUE, (ADDRESS)  
compare-and-swap(address, old_value, new_value) {  
    if (memory[address] == old_value) {  
        memory[address] = new_value;  
        return true;    // x86: set ZF flag  
    } else {  
        return false;  // x86: clear ZF flag  
    }  
}
```

```
// x86: lock xaddl REGISTER, (ADDRESS)  
fetch-and-add(address, register) {  
    old_value = memory[address];  
    memory[address] += register;  
    register = old_value;  
}
```

# common atomic operation pattern

try to do operation, ...

detect if it failed

if so, repeat

atomic operation does “try and see if it failed” part

# 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 **different than memory** *and* I am the only one who has it

**Shared** value is the **same as memory**

**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

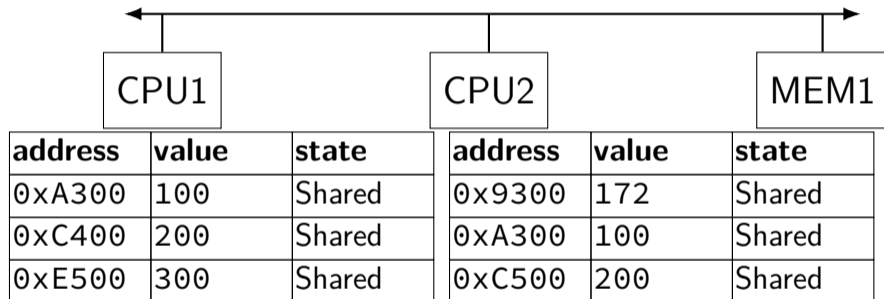
example: hear write while Shared

change to Invalid  
can send read later to get value from writer

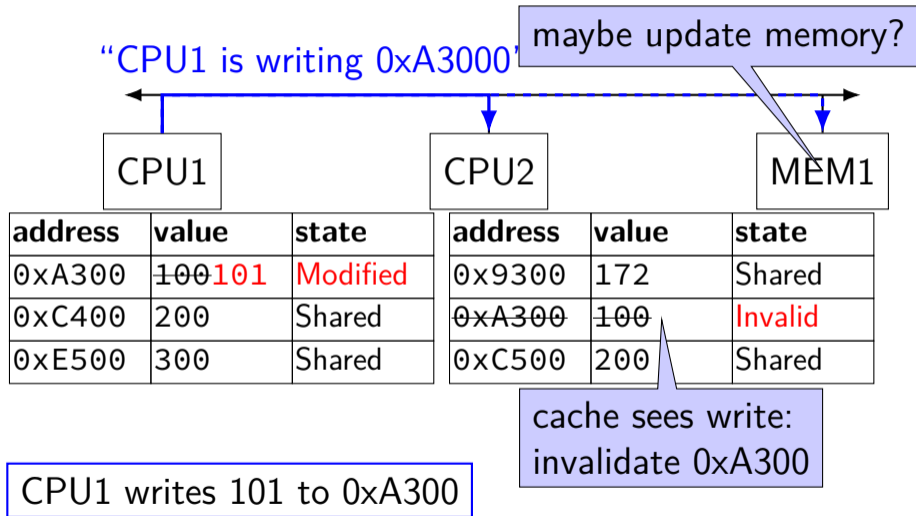
example: write while Modified

nothing to do — no other CPU can have a copy

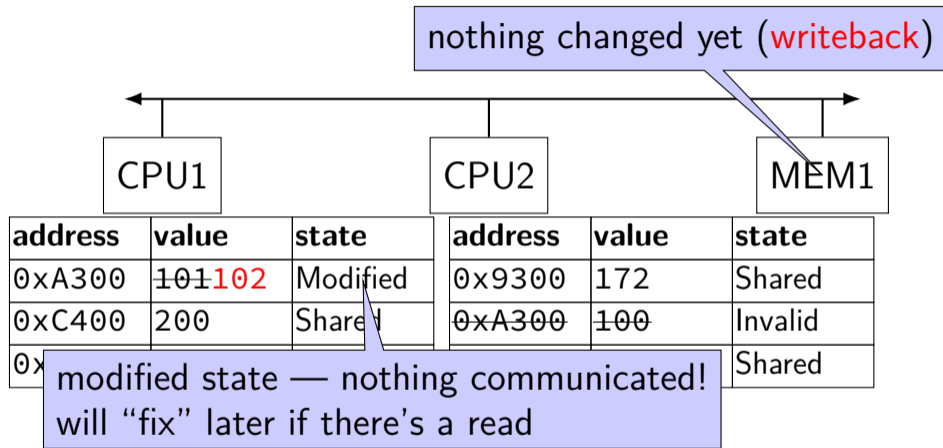
# MSI example



# MSI example

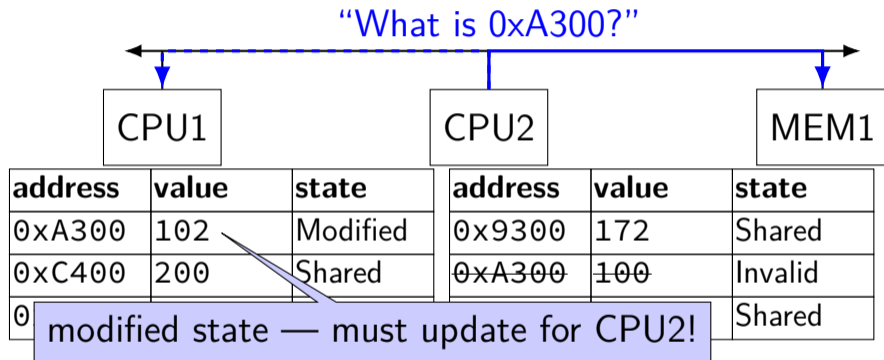


# MSI example



CPU1 writes 102 to 0xA300

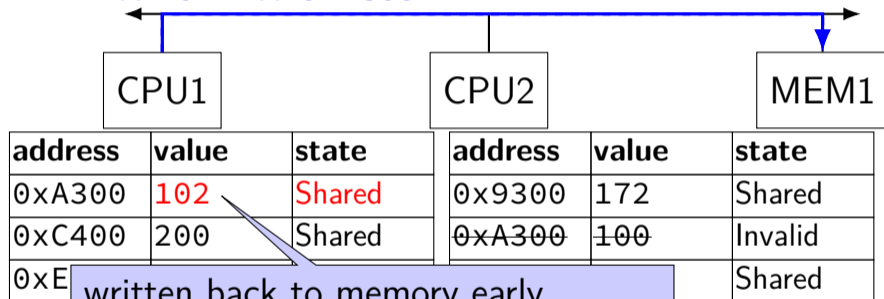
# MSI example



CPU2 reads 0xA300

# MSI example

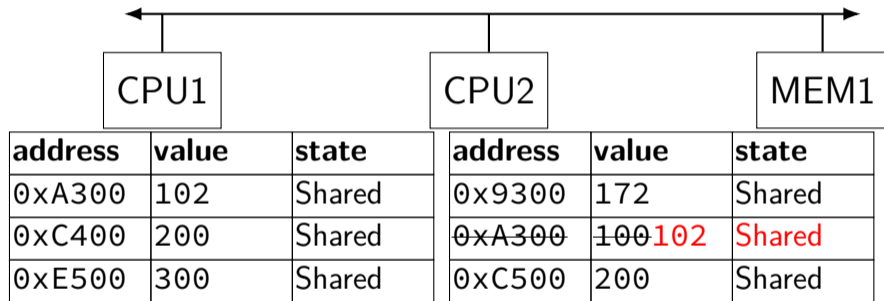
“Write 102 into 0xA300”



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

CPU2 reads 0xA300

# MSI example



## 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

# cache coherency exercise solution

action	0x1000-0x101f			0x2000-0x201f		
	CPU 1	CPU 2	CPU 3	CPU 1	CPU 2	CPU
	I	I	I	I	I	I
CPU 1: read 0x1000	S	I	I	I	I	I
CPU 2: read 0x1000	S	S	I	I	I	I
CPU 1: write 0x1000	M	I	I	I	I	I
CPU 1: read 0x2000	M	I	I	S	I	I
CPU 2: read 0x1000	S	S	I	S	I	I
CPU 2: write 0x2008	S	S	I	I	M	I
CPU 3: read 0x1008	S	S	S	I	M	I

# why load/store reordering?

fast processor designs can execute instructions out of order

goal: do something instead of waiting for slow memory accesses, etc.

more on this later in the semester

# 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

```
#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

# 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

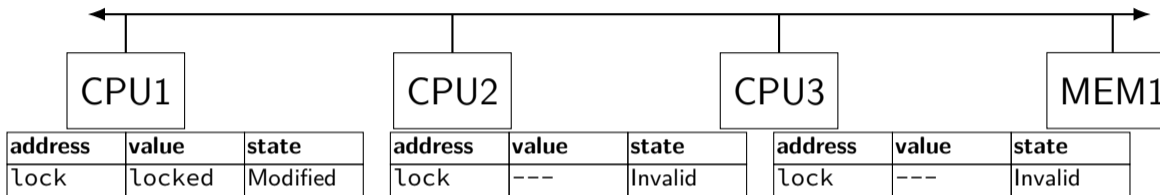
want to run another thread instead of infinite loop

solution: lock implementation integrated with scheduler

spinlocks can send a lot of messages on the shared bus

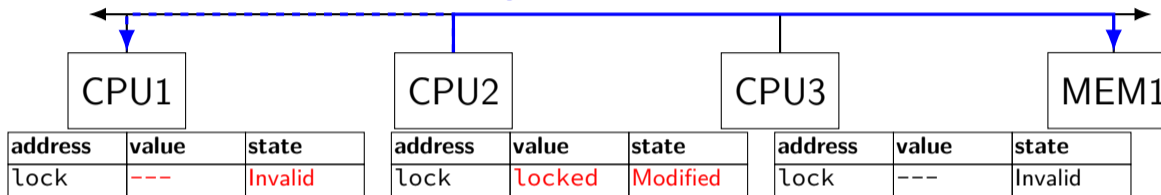
more efficient atomic operations to implement locks

# ping-ponging



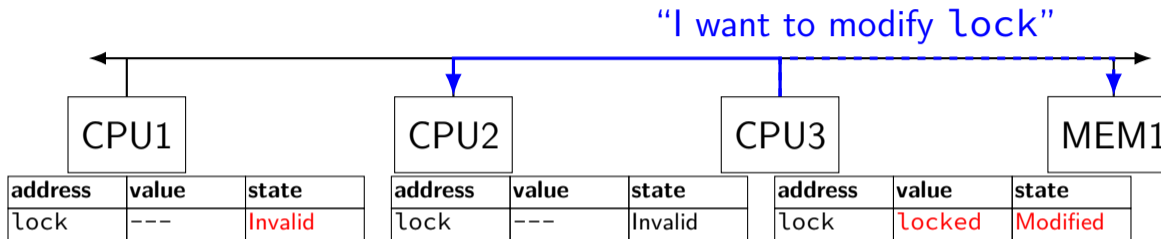
# ping-ponging

“I want to modify lock?”



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

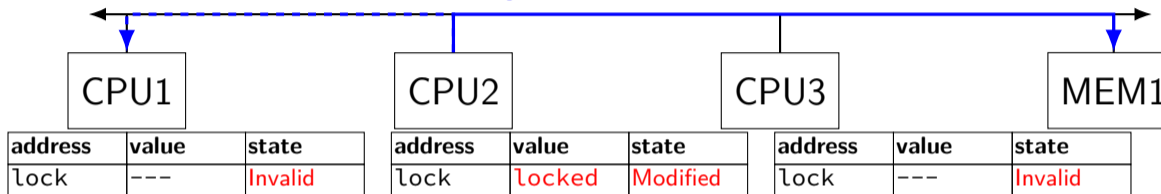
# ping-ponging



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

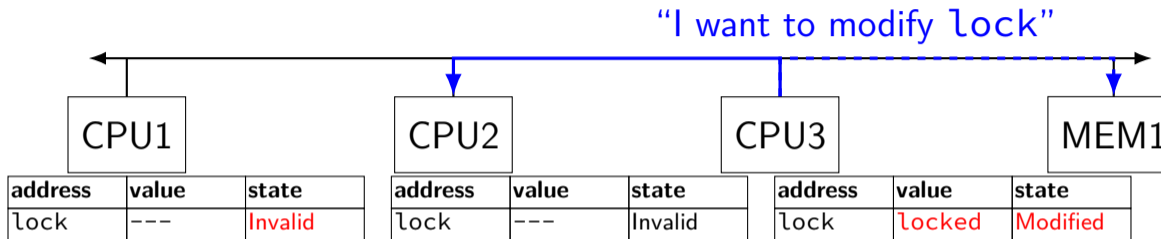
# ping-ponging

“I want to modify lock?”



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

# ping-ponging

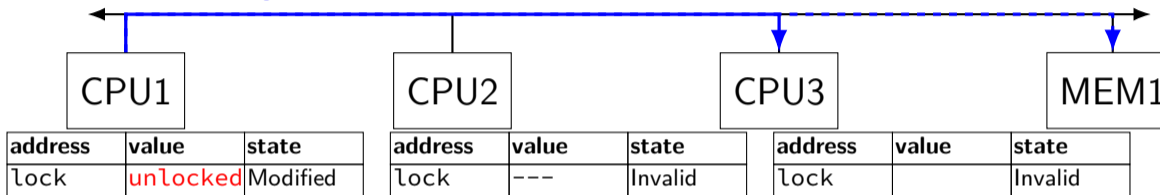


CPU3 read-modify-writes lock  
(to see it is still locked)



# ping-ponging

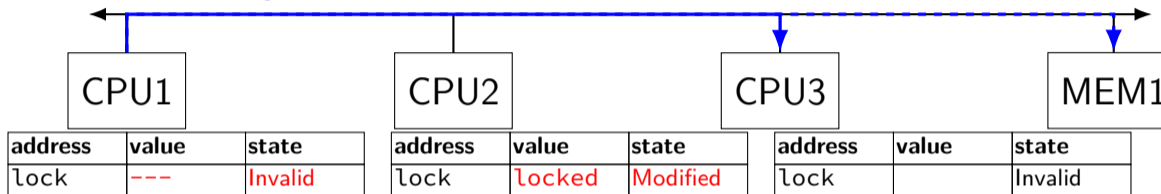
“I want to modify lock”



CPU1 sets lock to unlocked

# ping-ponging

“I want to modify lock”



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)

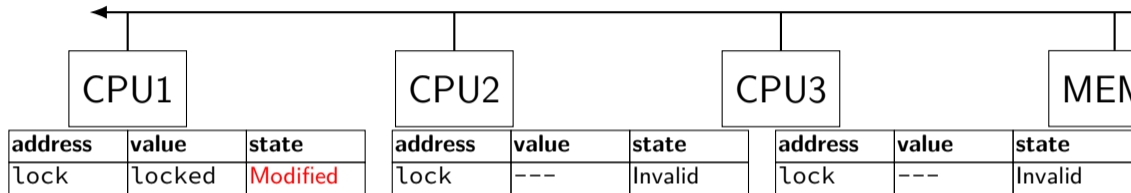
```
acquire(int *the_lock) {  
    do {  
        while (ATOMIC_READ(the_lock) == 0) { /* try again */ }  
    } while (ATOMIC_TEST_AND_SET(the_lock) == ALREADY_SET);  
}
```

# test-and-test-and-set (assembly)

acquire:

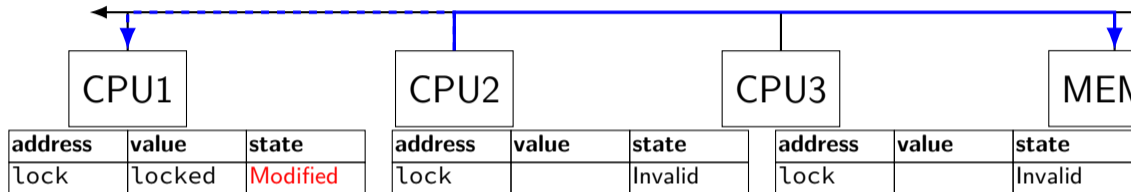
```
    cmp $0, the_lock           // test the lock non-atomically
                                // unlike lock xchg --- keeps lock in Shared state!
    jne acquire                // try again (still locked)
    // lock possibly free
    // but another processor might lock
    // before we get a chance to
    // ... so try with 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



# less ping-ponging

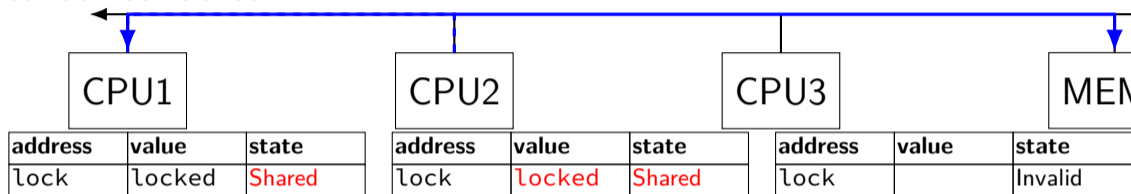
“I want to read lock?”



CPU2 reads lock  
(to see it is still locked)

# less ping-ponging

“set lock to locked”

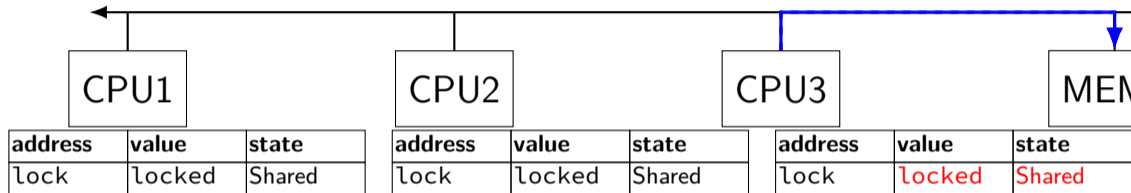


CPU1 writes back lock value,  
then CPU2 reads it



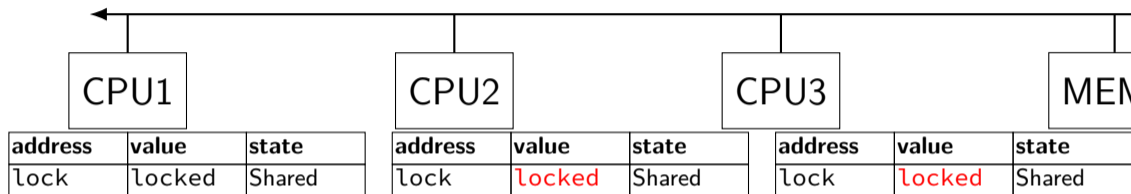
# less ping-ponging

“I want to read lock”



CPU3 reads lock  
(to see it is still locked)

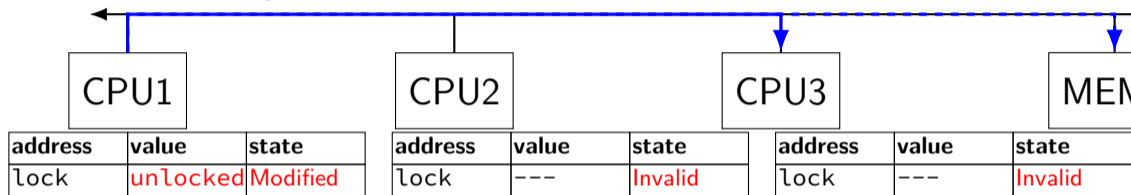
# less ping-ponging



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

# less ping-ponging

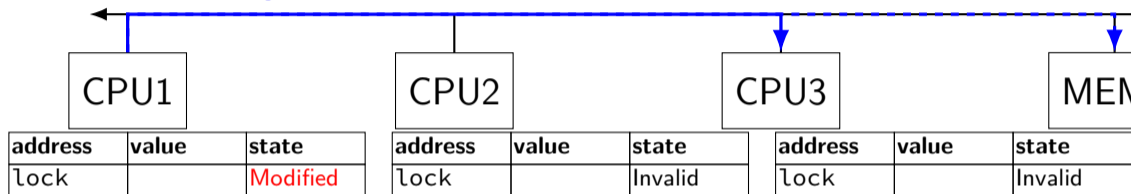
“I want to modify lock”



CPU1 sets lock to unlocked

# less ping-ponging

“I want to modify lock”



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)

## too much milk

roommates Alice and Bob want to keep fridge stocked with milk:

time	Alice	Bob
3:00	look in fridge. no milk	
3:05	leave for store	
3:10	arrive at store	look in fridge. no milk
3:15	buy milk	leave for store
3:20	return home, put milk in fridge	arrive at store
3:25		buy milk
3:30		return home, put milk in fridge

how can Alice and Bob coordinate better?



# too much milk “solution” 1 (algorithm)

leave a note: “I am buying milk”

place before buying, remove after buying

don't try buying if there's a note

≈ setting/checking a variable (e.g. “note = 1”)

with atomic load/store of variable

```
if (no milk) {  
    if (no note) {  
        leave note;  
        buy milk;  
        remove note;  
    }  
}
```

# too much milk “solution” 1 (algorithm)

leave a note: “I am buying milk”

place before buying, remove after buying

don't try buying if there's a note

≈ setting/checking a variable (e.g. “note = 1”)

with atomic load/store of variable

```
if (no milk) {  
    if (no note) {  
        leave note;  
        buy milk;  
        remove note;  
    }  
}
```

exercise: why doesn't this work?

# too much milk “solution” 1 (timeline)

**Alice**

```
if (no milk) {  
  if (no note) {  
  
    leave note;  
    buy milk;  
    remove note;  
  }  
}
```

**Bob**

```
if (no milk) {  
  if (no note) {  
  
    leave note;  
    buy milk;  
    remove note;  
  }  
}
```

## too much milk “solution” 2 (algorithm)

intuition: leave note when buying or checking if need to buy

```
leave note;
if (no milk) {
    if (no note) {
        buy milk;
    }
}
remove note;
```

## too much milk: “solution” 2 (timeline)

### Alice

```
leave note;  
if (no milk) {  
    if (no note) {  
        buy milk;  
    }  
}  
remove note;
```

## too much milk: “solution” 2 (timeline)

**Alice**

```
leave note;
```

```
if (no milk) {
```

```
    if (no note) { ← but there's always a note
```

```
        buy milk;
```

```
    }
```

```
}
```

```
remove note;
```

## too much milk: “solution” 2 (timeline)

**Alice**

```
leave note;
```

```
if (no milk) {
```

```
    if (no note) { ← but there's always a note
```

```
        buy milk;
```

```
    }
```

...will never buy milk (twice or once)

```
}
```

```
remove note;
```

## “solution” 3: algorithm

intuition: label notes so Alice knows which is hers (and vice-versa)

computer equivalent: separate noteFromAlice and noteFromBob variables

### Alice

```
leave note from Alice;  
if (no milk) {  
    if (no note from Bob) {  
        buy milk  
    }  
}  
remove note from Alice;
```

### Bob

```
leave note from Bob;  
if (no milk) {  
    if (no note from Alice)  
        buy milk  
    }  
}  
remove note from Bob;
```



# too much milk: “solution” 3 (timeline)

**Alice**

```
leave note from Alice
```

```
if (no milk) {
```

```
    if (no note from Bob) {
```

```
        buy milk
```

```
    }
```

```
}
```

```
remove note from Alice
```

**Bob**

```
leave note from Bob
```

```
if (no milk) {
```

```
    if (no note from Alice) {
```

```
        buy milk
```

```
    }
```

```
}
```

```
remove note from Bob
```

# too much milk: is it possible

is there a solutions with writing/reading notes?

≈ loading/storing from shared memory

yes, but it's not very elegant

## too much milk: solution 4 (algorithm)

### Alice

```
leave note from Alice
while (note from Bob) {
    do nothing
}
if (no milk) {
    buy milk
}
remove note from Alice
```

### Bob

```
leave note from Bob
if (no note from Alice) {
    if (no milk) {
        buy milk
    }
}
remove note from Bob
```

## too much milk: solution 4 (algorithm)

### Alice

```
leave note from Alice
while (note from Bob) {
    do nothing
}
if (no milk) {
    buy milk
}
remove note from Alice
```

### Bob

```
leave note from Bob
if (no note from Alice) {
    if (no milk) {
        buy milk
    }
}
remove note from Bob
```

exercise (hard): prove (in)correctness

## too much milk: solution 4 (algorithm)

### Alice

```
leave note from Alice
while (note from Bob) {
    do nothing
}
if (no milk) {
    buy milk
}
remove note from Alice
```

### Bob

```
leave note from Bob
if (no note from Alice) {
    if (no milk) {
        buy milk
    }
}
remove note from Bob
```

exercise (hard): prove (in)correctness

## too much milk: solution 4 (algorithm)

### Alice

```
leave note from Alice
while (note from Bob) {
    do nothing
}
if (no milk) {
    buy milk
}
remove note from Alice
```

exercise (hard): prove (in)correctness

exercise (hard): extend to three people

### Bob

```
leave note from Bob
if (no note from Alice) {
    if (no milk) {
        buy milk
    }
}
remove note from Bob
```

# Peterson's algorithm

general version of solution

see, e.g., Wikipedia

we'll use special hardware support instead

# 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



# 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

aside: this instruction did not exist in the original x86  
so x86 uses something older that's equivalent

# 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)

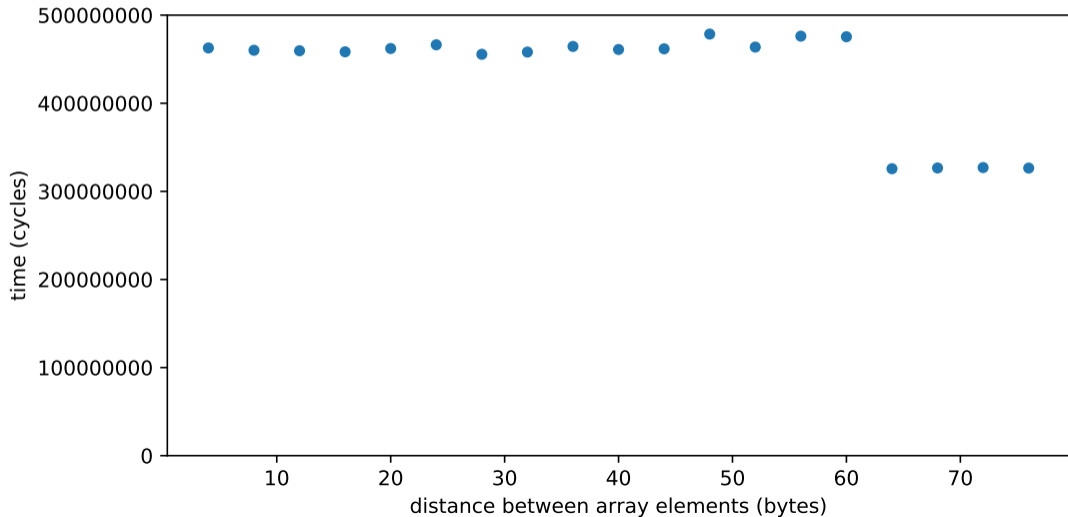
```
void *sum_up(void *raw_dest) {
    int *dest = (int *) raw_dest;
    for (int i = 0; i < 64 * 1024 * 1024; ++i) {
        *dest += data[i];
    }
}
```

```
__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

# exercise (1)

```
int values[1024];
int results[2];
void *sum_front(void *ignored_argument) {
    results[0] = 0;
    for (int i = 0; i < 512; ++i)
        results[0] += values[i];
    return NULL;
}
void *sum_back(void *ignored_argument) {
    results[1] = 0;
    for (int i = 512; i < 1024; ++i)
        results[1] += values[i];
    return NULL;
}
int sum_all() {
    pthread_t sum_front_thread, sum_back_thread;
    pthread_create(&sum_front_thread, NULL, sum_front, NULL);
    pthread_create(&sum_back_thread, NULL, sum_back, NULL);
    pthread_join(sum_front_thread, NULL);
    pthread_join(sum_back_thread, NULL);
    return results[0] + results[1];
}
```

## exercise (2)

```
struct ThreadInfo { int *values; int start; int end; int result };
void *sum_thread(void *argument) {
    ThreadInfo *my_info = (ThreadInfo *) argument;
    int sum = 0;
    for (int i = my_info->start; i < my_info->end; ++i) {
        my_info->result += my_info->values[i];
    }
    return NULL;
}

int sum_all(int *values) {
    ThreadInfo info[2]; pthread_t thread[2];
    for (int i = 0; i < 2; ++i) {
        info[i].values = values; info[i].start = i*512; info[i].end = (i+1)*512;
        pthread_create(&threads[i], NULL, sum_thread, (void *) &info[i]);
    }
    for (int i = 0; i < 2; ++i)
        pthread_join(threads[i], NULL);
    return info[0].result + info[1].result;
}
```

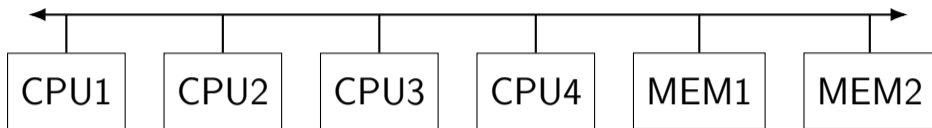
# connecting CPUs and memory

multiple processors, common memory

how do processors communicate with memory?



# shared bus



one possible design

we'll revisit later when we talk about I/O

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

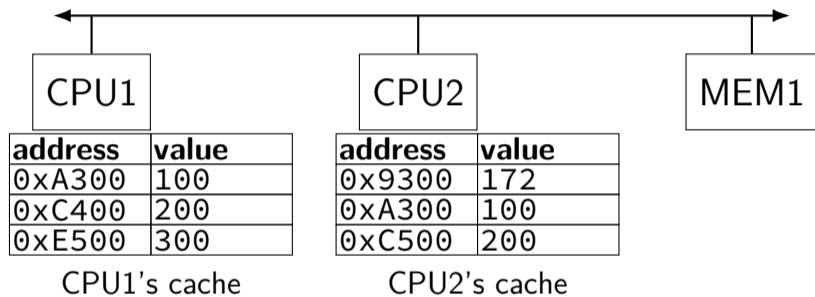
remember caches?

memory is **pretty slow**

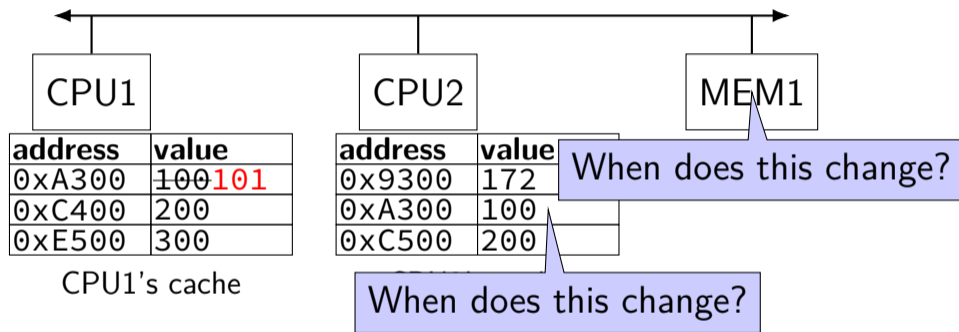
each CPU wants to keep local copies of memory

what happens when multiple CPUs cache same memory?

# the cache coherency problem



# the cache coherency problem



CPU1 writes 101 to 0xA300?