# last time

translation lookaside buffers
    special additional cache for last-level page table entries
    looked by virtual page number
    can practically be very small and therefore very fast


pthread API — pthread_create, pthread_join
    pthread join

# thread joining

pthread_join allows collecting thread return value

if you don't join joinable thread, then <span style="color:red">memory leak</span>!

# thread joining

pthread_join allows collecting thread return value

if you don't join joinable thread, then <span style="color:red">memory leak</span>!

avoiding memory leak?

always join...or

"detach" thread to make it not joinable

# pthread_detach

```
void *show_progress(void * ...) { ... }
void spawn_show_progress_thread() {
    pthread_t show_progress_thread;
    pthread_create(&show_progress_thread, NULL,
                   show_progress, NULL);

    /* instead of keeping pthread_t around to join thread later: */
    pthread_detach(show_progress_thread);
}

int main() {
    spawn_show_progress_thread();
    do_other_stuff();
    ...
}
```

detach = don't care about return value, etc.
system will deallocate when thread terminates

4

# starting threads detached

```
void *show_progress(void * ...) { ... }
void spawn_show_progress_thread() {
    pthread_t show_progress_thread;
    pthread_attr_t attrs;
    pthread_attr_init(&attrs);
    pthread_attr_setdetachstate(&attrs, PTHREAD_CREATE_DETACHED);
    pthread_create(&show_progress_thread, attrs,
                   show_progress, NULL);
    pthread_attr_destroy(&attrs);
}
```

# setting stack sizes

```c
void *show_progress(void * ...) { ... }
void spawn_show_progress_thread() {
    pthread_t show_progress_thread;
    pthread_attr_t attrs;
    pthread_attr_init(&attrs);
    pthread_attr_setstacksize(&attrs, 32 * 1024 /* bytes */);
    pthread_create(&show_progress_thread, attrs,
                   show_progress, NULL);
}
```
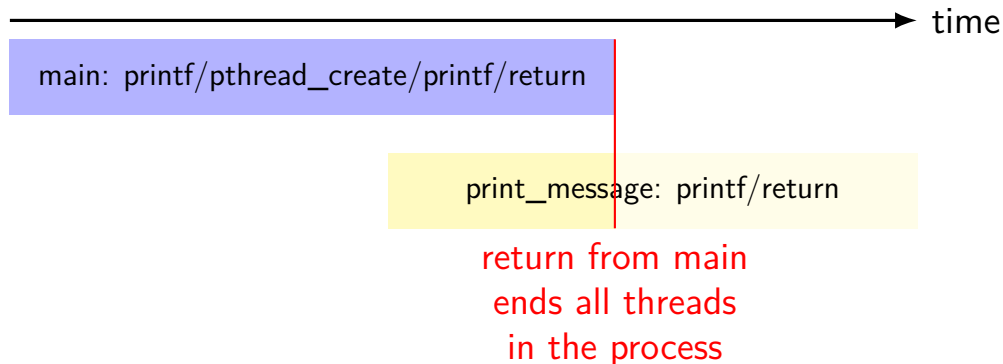
# a threading race

```
#include <pthread.h>
#include <stdio.h>
void *print_message(void *ignored_argument) {
    printf("In the thread\n"); return NULL;
}
int main() {
    printf("About to start thread\n");
    pthread_t the_thread;
    pthread_create(&the_thread, NULL, print_message, NULL);
    printf("Done starting thread\n");
    return 0;
}
```

My machine: outputs In the thread about 4% of the time.
What happened?

# a race

returning from main exits the entire process (all its threads)
> same as calling exit; not like other threads

race: main's return 0 or print_message's printf first?

# the correctness problem

two threads?

introduces *non-determinism*

which one runs first?

allows for "race condition" bugs

...to be avoided with synchronization constructs

# example application: ATM server

commands: withdraw, deposit

one correctness goal: don't lose money

# ATM server
(pseudocode)

```
ServerLoop() {
    while (true) {
        ReceiveRequest(&operation, &accountNumber, &amount);
        if (operation == DEPOSIT) {
            Deposit(accountNumber, amount);
        } else ...
    }
}
Deposit(accountNumber, amount) {
    account = GetAccount(accountNumber);
    account->balance += amount;
    SaveAccountUpdates(account);
}
```

# a threaded server?

```
Deposit(accountNumber, amount) {
    account = GetAccount(accountId);
    account−>balance += amount;
    SaveAccountUpdates(account);
}
```

maybe GetAccount/SaveAccountUpdates can be slow?
  read/write disk sometimes? contact another server sometimes?

maybe lots of requests to process?
  maybe real logic has more checks than Deposit()
  …

all reasons to handle multiple requests at once

$\rightarrow$ many threads all running the server loop

# multiple threads

```
main() {
    for (int i = 0; i < NumberOfThreads; ++i) {
        pthread_create(&server_loop_threads[i], NULL,
                       ServerLoop, NULL);
    }
    ...
}

ServerLoop() {
    while (true) {
        ReceiveRequest(&operation, &accountNumber, &amount);
        if (operation == DEPOSIT) {
            Deposit(accountNumber, amount);
        } else ...
    }
}
```

# the lost write

account−>balance += amount; (in two threads, same account)

........................................................................................................................

|                  Thread A                  |                  Thread B                  |

mov account−>balance, %rax
add amount, %rax

──────────────────────── context switch ────────────────────────

                                     mov account−>balance, %rax
                                     add amount, %rax

──────────────────────── context switch ────────────────────────

mov %rax, account−>balance

──────────────────────── context switch ────────────────────────

                                     mov %rax, account−>balance

# the lost write

account−>balance += amount; (in two threads, same account)

|             Thread A             |             Thread B             |
|----------------------------------|----------------------------------|

mov account−>balance, %rax
add amount, %rax

──────────────── context switch ────────────────
                                    mov account−>balance, %rax
                                    add amount, %rax
──────────────── context switch ────────────────
mov %rax, account−>balance
──────────────── context switch ────────────────
                                    mov %rax, account−>balance

lost write to balance

"winner" of the race

# the lost write

`account->balance += amount;` (in two threads, same account)

|  Thread A  |  Thread B  |
|---|---|

```
mov account->balance, %rax
add amount, %rax
```
──────────── context switch ────────────
```
                              mov account->balance, %rax
                              add amount, %rax
```
──────────── context switch ────────────
```
mov %rax, account->balance
```
──────────── context switch ────────────
```
                              mov %rax, account->balance
```

lost write to balance

"winner" of the race

lost track of thread A's money

# thinking about race conditions (1)

what are the possible values of $x$? (initially $x = y = 0$)

| Thread A | Thread B |
|----------|----------|
| $x \leftarrow 1$ | $y \leftarrow 2$ |

# thinking about race conditions (2)

possible values of $x$? (initially $x = y = 0$)

| Thread A | Thread B |
|---|---|
| $x \leftarrow y + 1$ | $y \leftarrow 2$ |
| | $y \leftarrow y \times 2$ |

# thinking about race conditions (2)

possible values of $x$? (initially $x = y = 0$)

| Thread A | Thread B |
|---|---|
| $x \leftarrow y + 1$ | $y \leftarrow 2$ |
| | $y \leftarrow y \times 2$ |

# thinking about race conditions (3)

what are the possible values of $x$?

(initially $x = y = 0$)

| Thread A | Thread B |
|----------|----------|
| $x \leftarrow 1$ | $x \leftarrow 2$ |

# thinking about race conditions (2)

possible values of $x$? (initially $x = y = 0$)

| Thread A | Thread B |
|---|---|
| $x \leftarrow y + 1$ | $y \leftarrow 2$ |
| | $y \leftarrow y \times 2$ |

# atomic operation

*atomic operation* = operation that runs to completion or not at all

we will use these to let threads work together

most machines: loading/storing (aligned) words is atomic
    so can't get $3$ from $x \leftarrow 1$ and $x \leftarrow 2$ running in parallel
    aligned $\approx$ address of word is multiple of word size (typically done by
    compilers)

but some instructions are not atomic; examples:
    x86: integer add constant to memory location
    many CPUs: loading/storing values that cross cache blocks
        e.g. if cache blocks 0x40 bytes, load/store 4 byte from addr. 0x3E is not atomic

# lost adds (program)

```
.global update_loop
update_loop:
    addl $1, the_value // the_value (global variable) += 1
    dec %rdi            // argument 1 -= 1
    jg update_loop      // if argument 1 >= 0 repeat
    ret
```

```
int the_value;
extern void *update_loop(void *);
int main(void) {
    the_value = 0;
    pthread_t A, B;
    pthread_create(&A, NULL, update_loop, (void*) 1000000);
    pthread_create(&B, NULL, update_loop, (void*) 1000000);
    pthread_join(A, NULL); pthread_join(B, NULL);
    // expected result: 1000000 + 1000000 = 2000000
    printf("the_value = %d\n", the_value);
}
```

# lost adds (results)



the_value = ?

# but how?

probably not possible on single core
  exceptions can't occur in the middle of add instruction

…but 'add to memory' implemented with multiple steps
  still needs to load, add, store internally
  can be interleaved with what other cores do

# but how?

probably not possible on single core
> exceptions can't occur in the middle of add instruction

…but 'add to memory' implemented with multiple steps
> still needs to load, add, store internally
> can be interleaved with what other cores do

(and actually it's more complicated than that — we'll talk later)

# so, what is actually atomic

for now we'll assume: load/stores of 'words'
    (64-bit machine = 64-bits words)

in general: processor designer will tell you

their job to design caches, etc. to work as documented

# compilers move loads/stores (1)

```
void WaitForReady() {
    do {} while (!ready);
}
```

---

```
WaitForOther:
  movl ready, %eax   // eax <- other_ready
.L2:
  testl %eax, %eax
  je .L2                     // while (eax == 0) repeat
  ...
```

# compilers move loads/stores (1)

```
void WaitForReady() {
    do {} while (!ready);
}
```
---
```
WaitForOther:
  movl ready, %eax   // eax <- other_ready
.L2:
  testl %eax, %eax
  je .L2                      // while (eax == 0) repeat
  ...
```

# compilers move loads/stores (2)

```
void WaitForOther() {
    is_waiting = 1;
    do {} while (!other_ready);
    is_waiting = 0;
}
```

```
WaitForOther:
  // compiler optimization: don't set is_waiting to 1,
  // (why? it will be set to 0 anyway)
  movl other_ready, %eax  // eax <- other_ready
.L2:
  testl %eax, %eax
  je .L2                       // while (eax == 0) repeat
  ...
  movl $0, is_waiting  // is_waiting <- 0
```

# compilers move loads/stores (2)

```
void WaitForOther() {
    is_waiting = 1;
    do {} while (!other_ready);
    is_waiting = 0;
}
```

```
WaitForOther:
  // compiler optimization: don't set is_waiting to 1,
  // (why? it will be set to 0 anyway)
  movl other_ready, %eax  // eax <- other_ready
.L2:
  testl %eax, %eax
  je .L2                      // while (eax == 0) repeat
  ...
  movl $0, is_waiting  // is_waiting <- 0
```

# compilers move loads/stores (2)

```
void WaitForOther() {
    is_waiting = 1;
    do {} while (!other_ready);
    is_waiting = 0;
}
```

```
WaitForOther:
  // compiler optimization: don't set is_waiting to 1,
  // (why? it will be set to 0 anyway)
  movl other_ready, %eax   // eax <- other_ready
.L2:
  testl %eax, %eax
  je .L2                   // while (eax == 0) repeat
  ...
  movl $0, is_waiting  // is_waiting <- 0
```

# fixing compiler reordering?

isn't there a way to tell compiler not to do these optimizations?

yes, but that is <span style="color:red">still not enough</span>!

**processors** sometimes do this kind of reordering too (between cores)

# pthreads and reordering

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

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

pthread_create, pthread_join, other tools we'll talk about …
  basically: if pthreads is waiting for/starting something, no weird ordering

implementation part 1: prevent compiler reordering

implementation part 2: use special instructions
  example: x86 `mfence` instruction

# some definitions

**mutual exclusion**: ensuring only one thread does a particular thing at a time

   like checking for and, if needed, buying milk

# some definitions

**mutual exclusion**: ensuring only one thread does a particular thing at a time

   like checking for and, if needed, buying milk

**critical section**: code that exactly one thread can execute at a time

   result of critical section

# some definitions

**mutual exclusion**: ensuring only one thread does a particular thing at a time

> like checking for and, if needed, buying milk

**critical section**: code that exactly one thread can execute at a time

> result of critical section

**lock**: object only one thread can hold at a time

> interface for creating critical sections

# lock analogy

agreement: only change account balances while wearing this hat

normally hat kept on table

put on hat when editing balance

hopefully, only one person (= thread) can wear hat a time

need to wait for them to remove hat to put it on

# lock analogy

agreement: only change account balances while wearing this hat

normally hat kept on table

put on hat when editing balance

hopefully, only one person (= thread) can wear hat a time

need to wait for them to remove hat to put it on

"lock (or acquire) the lock" = get and put on hat

"unlock (or release) the lock" = put hat back on table

# the lock primitive

locks: an object with (at least) two operations:
    *acquire* or *lock* — wait until lock is free, then "grab" it
    *release* or *unlock* — let others use lock, wakeup waiters

typical usage: everyone acquires lock before using shared resource
    forget to acquire lock? weird things happen

```
Lock(account_lock);
balance += ...;
Unlock(account_lock);
```

# the lock primitive

locks: an object with (at least) two operations:
  *acquire* or *lock* — wait until lock is free, then "grab" it
  *release* or *unlock* — let others use lock, wakeup waiters

typical usage: everyone acquires lock before using shared resource
  forget to acquire lock? weird things happen

```
Lock(account_lock);
balance += ...;
Unlock(account_lock);
```

# waiting for lock?

when waiting — ideally:

not using processor (at least if waiting a while)

OS can context switch to other programs

# pthread mutex

```
#include <pthread.h>

pthread_mutex_t account_lock;
pthread_mutex_init(&account_lock, NULL);
    // or: pthread_mutex_t account_lock =
    //                  PTHREAD_MUTEX_INITIALIZER;
...
pthread_mutex_lock(&account_lock);
balance += ...;
pthread_mutex_unlock(&account_lock);
```

# exercise

```
pthread_mutex_t lock1 = PTHREAD_MUTEX_INITIALIZER;
pthread_mutex_t lock2 = PTHREAD_MUTEX_INITIALIZER;
string one = "init one", two = "init two";
void ThreadA() {
    pthread_mutex_lock(&lock1);
    one = "one in ThreadA";  // (A1)
    pthread_mutex_unlock(&lock1);
    pthread_mutex_lock(&lock2);
    two = "two in ThreadA";  // (A2)
    pthread_mutex_unlock(&lock2);
}
void ThreadB() {
    pthread_mutex_lock(&lock1);
    one = "one in ThreadB";  // (B1)
    pthread_mutex_lock(&lock2);
    two = "two in ThreadB";  // (B2)
    pthread_mutex_unlock(&lock2);
    pthread_mutex_unlock(&lock1);
}
```

possible values of one/two after A+B run?

# exercise (alternate 1)

```
pthread_mutex_t lock1 = PTHREAD_MUTEX_INITIALIZER;
pthread_mutex_t lock2 = PTHREAD_MUTEX_INITIALIZER;
string one = "init one", two = "init two";
void ThreadA() {
    pthread_mutex_lock(&lock2);
    two = "two in ThreadA";  // (A2)
    pthread_mutex_unlock(&lock2);
    pthread_mutex_lock(&lock1);
    one = "one in ThreadA";  // (A1)
    pthread_mutex_unlock(&lock1);
}
void ThreadB() {
    pthread_mutex_lock(&lock1);
    one = "one in ThreadB";  // (B1)
    pthread_mutex_lock(&lock2);
    two = "two in ThreadB";  // (B2)
    pthread_mutex_unlock(&lock2);
    pthread_mutex_unlock(&lock1);
}
```

possible values of one/two after A+B run?

# exercise (alternate 2)

```
pthread_mutex_t lock1 = PTHREAD_MUTEX_INITIALIZER;
pthread_mutex_t lock2 = PTHREAD_MUTEX_INITIALIZER;
string one = "init one", two = "init two";
void ThreadA() {
    pthread_mutex_lock(&lock2);
    two = "two in ThreadA";  // (A2)
    pthread_mutex_unlock(&lock2);
    pthread_mutex_lock(&lock1);
    one = "one in ThreadA";  // (A1)
    pthread_mutex_unlock(&lock1);
}
void ThreadB() {
    pthread_mutex_lock(&lock1);
    one = "one in ThreadB";  // (B1)
    pthread_mutex_unlock(&lock1);
    pthread_mutex_lock(&lock2);
    two = "two in ThreadB";  // (B2)
    pthread_mutex_unlock(&lock2);
}
```

possible values of one/two after A+B run?

# POSIX mutex restrictions

pthread_mutex rule:  unlock from same thread you lock in

does this actually matter?

depends on how pthread_mutex is implemented

# preview: general sync

lots of coordinating threads beyond locks/barriers

will talk about two general tools later:
    monitors/condition variables
    semaphores

big added feature: wait for arbitrary thing to happen

# a bad idea

one bad idea to wait for an event:

```
pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER; bool ready = false;
void WaitForReady() {
    pthread_mutex_lock(&lock);
    do {
        pthread_mutex_unlock(&lock);
        /* only time MarkReady() can run */
        pthread_mutex_lock(&lock);
    } while (!ready);
    pthread_mutex_unlock(&lock);
}
void MarkReady() {
    pthread_mutex_lock(&lock);
    ready = true;
    pthread_mutex_unlock(&lock);
}
```

wastes processor time; MarkReady can stall waiting for unlock
window

# beyond locks

in practice: want more than locks for synchronization

for waiting for arbtirary events (without CPU-hogging-loop):
    monitors
    semaphores

for common synchornization patterns:
    barriers
    reader-writer locks

higher-level interface:
    transactions

# barriers

compute minimum of 100M element array with 2 processors

algorithm:

compute minimum of 50M of the elements on each CPU
>   one thread for each CPU

wait for all computations to finish

take minimum of all the minimums

# barriers

compute minimum of 100M element array with 2 processors

algorithm:

compute minimum of 50M of the elements on each CPU
>   one thread for each CPU

wait for all computations to finish

take minimum of all the minimums

# barriers API

barrier.Initialize(NumberOfThreads)

barrier.Wait() — return after all threads have waited

idea: multiple threads perform computations in parallel

threads wait for all other threads to call Wait()

# barrier: waiting for finish

```
barrier.Initialize(2);
```

Thread 0

```
partial_mins[0] =
    /* min of first
       50M elems */;

barrier.Wait();


total_min = min(
    partial_mins[0],
    partial_mins[1]
);
```

Thread 1

```
partial_mins[1] =
    /* min of last
       50M elems */
barrier.Wait();
```

# barriers: reuse

```
results[0][0] = getInitial(0);        results[0][1] = getInitial(1);
barrier.Wait();                       barrier.Wait();

results[1][0] =                       results[1][1] =
    computeFrom(                          computeFrom(
        results[0][0],                        results[0][0],
        results[0][1]                         results[0][1]
    );                                    );
barrier.Wait();                       barrier.Wait();

results[2][0] =                       results[2][1] =
    computeFrom(                          computeFrom(
        results[1][0],                        results[1][0],
        results[1][1]                         results[1][1]
    );                                    );
```

# barriers: reuse

Thread 0

```
results[0][0] = getInitial(0);
barrier.Wait();

results[1][0] =
    computeFrom(
        results[0][0],
        results[0][1]
    );
barrier.Wait();

results[2][0] =
    computeFrom(
        results[1][0],
        results[1][1]
    );
```

Thread 1

```
results[0][1] = getInitial(1);
barrier.Wait();

results[1][1] =
    computeFrom(
        results[0][0],
        results[0][1]
    );
barrier.Wait();

results[2][1] =
    computeFrom(
        results[1][0],
        results[1][1]
    );
```

# barriers: reuse

|                 Thread 0                  |                 Thread 1                  |

```
          Thread 0                              Thread 1
results[0][0] = getInitial(0);     results[0][1] = getInitial(1);
barrier.Wait();                    barrier.Wait();

results[1][0] =                    results[1][1] =
    computeFrom(                       computeFrom(
        results[0][0],                     results[0][0],
        results[0][1]                      results[0][1]
    );                                 );
barrier.Wait();                    barrier.Wait();

results[2][0] =                    results[2][1] =
    computeFrom(                        computeFrom(
        results[1][0],                      results[1][0],
        results[1][1]                       results[1][1]
    );                                  );
```

# pthread barriers

```
pthread_barrier_t barrier;
pthread_barrier_init(
    &barrier,
    NULL /* attributes */,
    numberOfThreads
);
...
...
pthread_barrier_wait(&barrier);
```

# exercise

```
pthread_barrier_t barrier; int x = 0, y = 0;
void thread_one() {
    y = 10;
    pthread_barrier_wait(&barrier);
    y = x + y;
    pthread_barrier_wait(&barrier);
    pthread_barrier_wait(&barrier);
    printf("%d %d\n", x, y);
}
void thread_two() {
    x = 20;
    pthread_barrier_wait(&barrier);
    pthread_barrier_wait(&barrier);
    x = x + y;
    pthread_barrier_wait(&barrier);
}
```

output? (if both run at once, barrier set for 2 threads)

# life homework (pseudocode)

```
for (int time = 0; time < MAX_ITERATIONS; ++time) {
    for (int y = 0; y < size; ++y) {
        for (int x = 0; x < size; ++x) {
            to_grid(x, y) = computeValue(from_grid, x, y);
        }
    }
    swap(from_grid, to_grid);
}
```

# life homework

compute grid of values for time $t$ from grid for time $t - 1$
  compute new value at $i, j$ based on surrounding values

parallel version: produce parts of grid in different threads

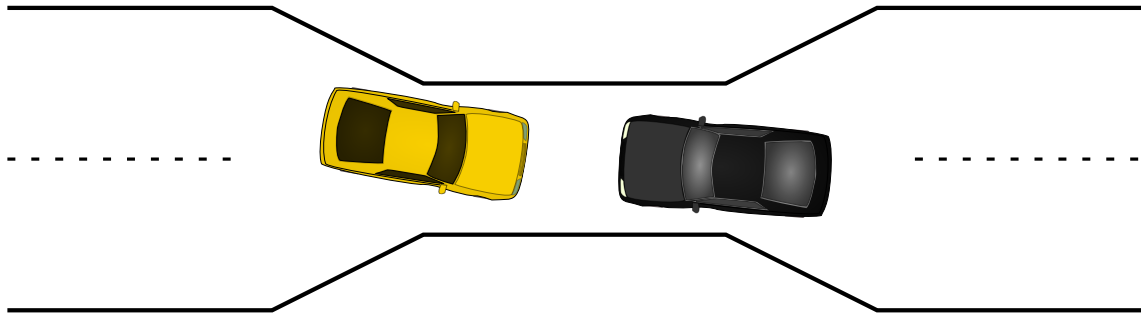use barriers to finish time $t$ before going to time $t + 1$

# 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<br>MoveFile(A, B, "foo") | Thread 2<br>MoveFile(B, A, "bar") |
|---|---|
| `lock(&A->lock);` | |
| `lock(&B->lock);` | |
| (do move) | |
| `unlock(&B->lock);` | |
| `unlock(&A->lock);` | |
| | `lock(&B->lock);` |
| | `lock(&A->lock);` |
| | (do move) |
| | `unlock(&B->lock);` |
| | `unlock(&A->lock);` |

# moving two files: lucky timeline (2)

| **Thread 1**<br>MoveFile(A, B, "foo") | **Thread 2**<br>MoveFile(B, A, "bar") |
|---|---|
| lock(&A->lock);<br>lock(&B->lock); | |
| | lock(&B->lock…<br>(waiting for B lock) |
| (do move)<br>unlock(&B->lock); | |
| | lock(&B->lock);<br>lock(&A->lock… |
| unlock(&A->lock); | |
| | lock(&A->lock);<br>(do move)<br>unlock(&A->lock);<br>unlock(&B->lock); |

# moving two files: unlucky timeline

| Thread 1 | Thread 2 |
|---|---|
| MoveFile(A, B, "foo") | MoveFile(B, A, "bar") |

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

# moving two files: unlucky timeline

| Thread 1 | Thread 2 |
|---|---|
| MoveFile(A, B, "foo") | MoveFile(B, A, "bar") |

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

                            lock(&B->lock);

lock(&B->lock… stalled
(waiting for lock on B)     lock(&A->lock… stalled
(waiting for lock on B)     (waiting for lock on A)
```

# moving two files: unlucky timeline

| Thread 1 MoveFile(A, B, "foo") | Thread 2 MoveFile(B, A, "bar") |
|---|---|
| lock(&A->lock); | |
| | lock(&B->lock); |
| lock(&B->lock… stalled | |
| (waiting for lock on B) | lock(&A->lock… stalled |
| (waiting for lock on B) | (waiting for lock on A) |
| ~~(do move)~~ unreachable | ~~(do move)~~ unreachable |
| ~~unlock(&B->lock);~~ unreachable | ~~unlock(&A->lock);~~ unreachable |
| ~~unlock(&A->lock);~~ unreachable | ~~unlock(&B->lock);~~ unreachable |

# moving two files: unlucky timeline

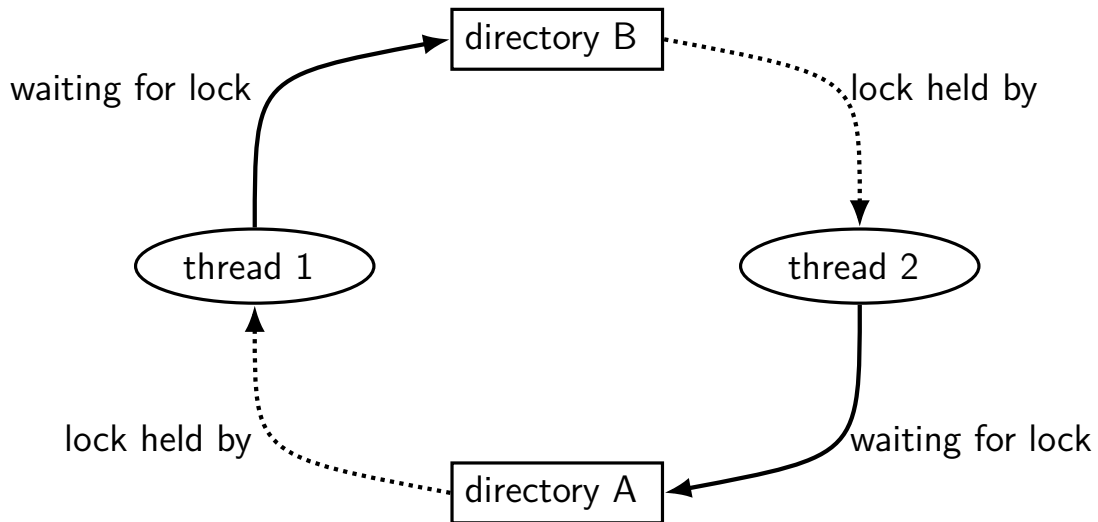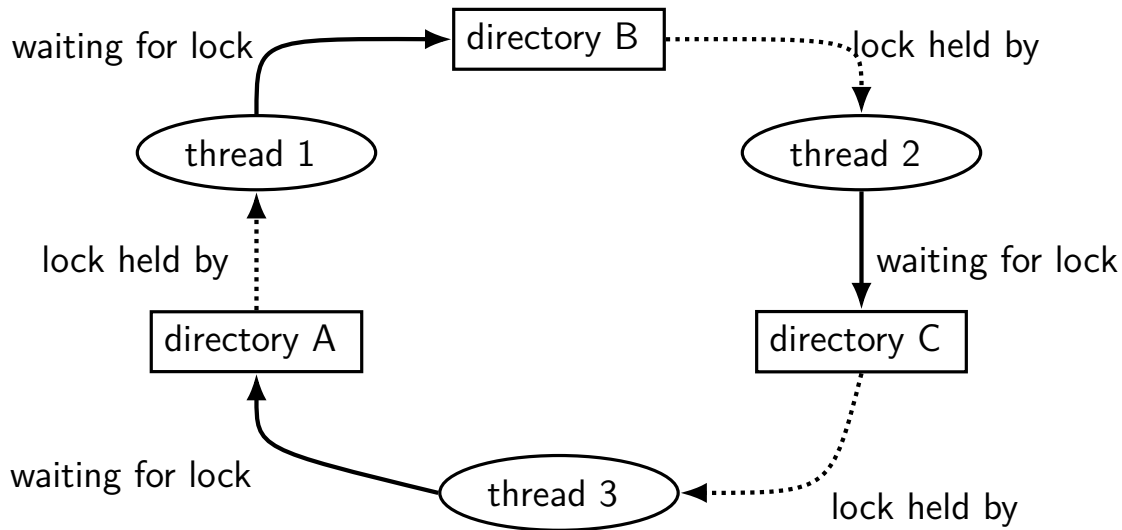| **Thread 1**<br>MoveFile(A, B, "foo") | **Thread 2**<br>MoveFile(B, A, "bar") |
|---|---|
| lock(&A->lock); | |
| | lock(&B->lock); |
| lock(&B->lock… stalled | |
| (waiting for lock on B) | lock(&A->lock… stalled |
| (waiting for lock on B) | (waiting for lock on A) |
| | |
| (do move) unreachable | (do move) unreachable |
| unlock(&B->lock); unreachable | unlock(&A->lock); unreachable |
| unlock(&A->lock); unreachable | unlock(&B->lock); unreachable |

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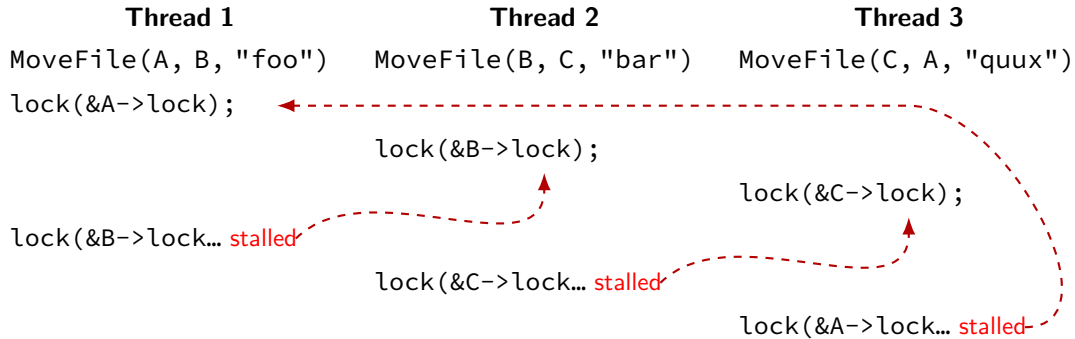Thread 2 holds B lock, waiting for Thread 1 to release A lock

# moving two files: dependencies



directory B

waiting for lock

lock held by

thread 1

thread 2

lock held by

waiting for lock

directory A

# moving three files: dependencies

# moving three files: unlucky timeline

| Thread 1 | Thread 2 | Thread 3 |
|---|---|---|
| MoveFile(A, B, "foo") | MoveFile(B, C, "bar") | MoveFile(C, A, "quux") |

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

                    lock(&B->lock);

                                        lock(&C->lock);

lock(&B->lock… stalled

                    lock(&C->lock… stalled

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

# deadlock with free space

| Thread 1 | Thread 2 |
|---|---|
| AllocateOrWaitFor(1 MB) | AllocateOrWaitFor(1 MB) |
| AllocateOrWaitFor(1 MB) | AllocateOrWaitFor(1 MB) |
| (do calculation) | (do calculation) |
| Free(1 MB) | Free(1 MB) |
| Free(1 MB) | Free(1 MB) |

2 MB of space — deadlock possible with unlucky order

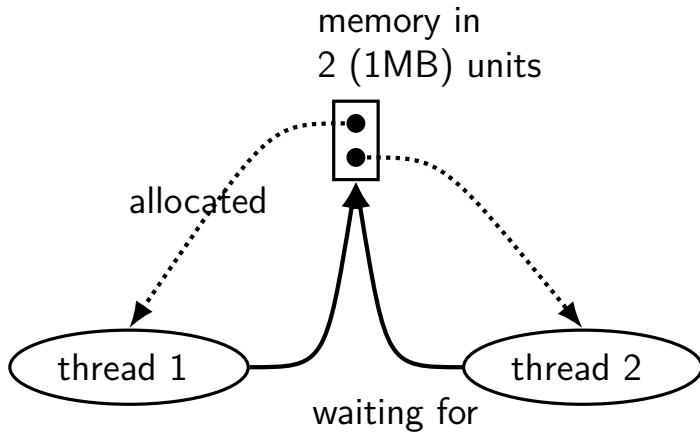# deadlock with free space (unlucky case)

| Thread 1 | Thread 2 |
|---|---|
| `AllocateOrWaitFor(1 MB)` | |
| | `AllocateOrWaitFor(1 MB)` |
| `AllocateOrWaitFor(1 MB…` stalled | |
| | `AllocateOrWaitFor(1 MB…` stalled |

# free space: dependency graph



memory in
2 (1MB) units

allocated

thread 1

thread 2

waiting for

# deadlock with free space (lucky case)

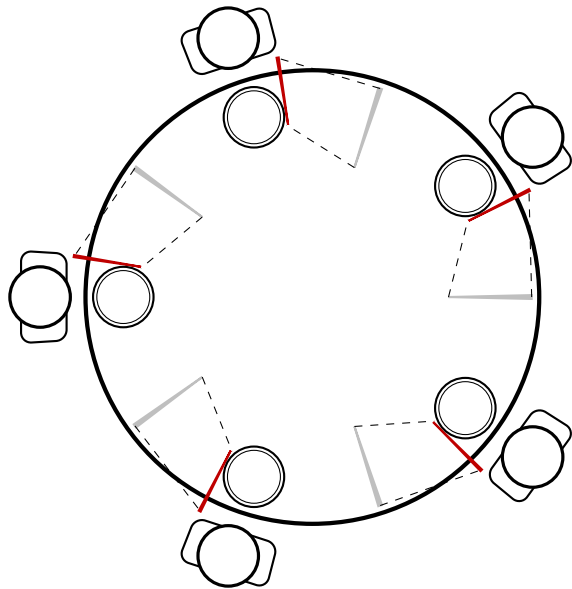| Thread 1 | Thread 2 |
|---|---|
| AllocateOrWaitFor(1 MB) | |
| AllocateOrWaitFor(1 MB) | |
| (do calculation) | |
| Free(1 MB); | |
| Free(1 MB); | |
| | AllocateOrWaitFor(1 MB) |
| | AllocateOrWaitFor(1 MB) |
| | (do calculation) |
| | Free(1 MB); |
| | Free(1 MB); |

# dining philosophers
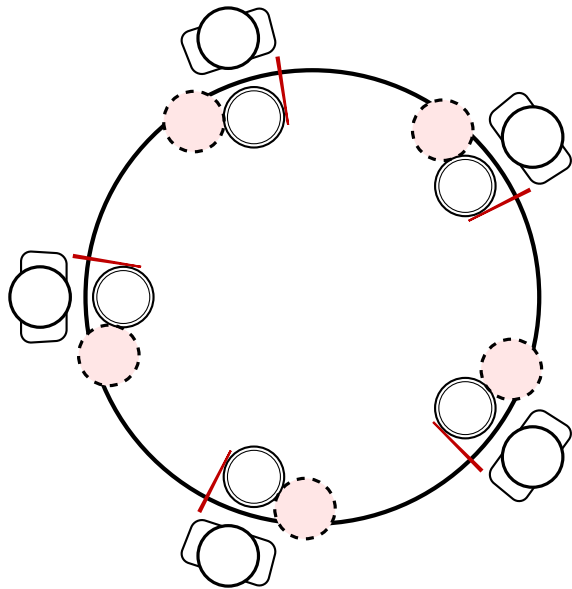


five philosophers either think or eat
to eat:
grab chopstick on left, then
grba chopstick on right, then
then eat, then
return chopsticks

# 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 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 $\{T_1, \ldots, T_n\}$ 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

# 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    no *hold and wait/*
    revoke/preempt resources         *preemption*

acquire resources in **consistent order**      no *circular wait*

request **all resources at once**      no *hold and 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      no *hold and wait/*
    revoke/preempt resources        *preemption*

acquire resources in **consistent order**        no *circular wait*

request **all resources at once**        no *hold and 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*

request **all resources at once**

no *hold and 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
> problem: retry how many times? no bound on number of tries needed
> …

*exclusion*

**no waiting**
    "busy signal" — abort and (maybe) retry      no *hold and wait*/
    revoke/preempt resources                      *preemption*

acquire resources in **consistent order**      no *circular wait*

request **all resources at once**      no *hold and 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        no *hold and wait*/
    revoke/preempt resources                *preemption*


acquire resources in **consistent order**        no *circular wait*


request **all resources at once**        no *hold and wait*

# 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      no *hold and wait /*
    revoke/preempt resources      *preemption*

acquire resources in **consistent order**      no *circular wait*

request **all resources at once**      no *hold and 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*

request **all resources at once**

no *hold and 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      no *hold and wait*/
    revoke/preempt resources                          *preemption*

acquire resources in **consistent order**             no *circular wait*

request **all resources at once**                         no *hold and wait*

# monitors/condition variables

locks for mutual exclusion

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

related data structures


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

# monitor idea

a monitor

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

# monitor idea

a monitor



| lock |
| shared data |
| condvar 1 |
| condvar 2 |
| ... |

| operation1(...) |
| operation2(...) |

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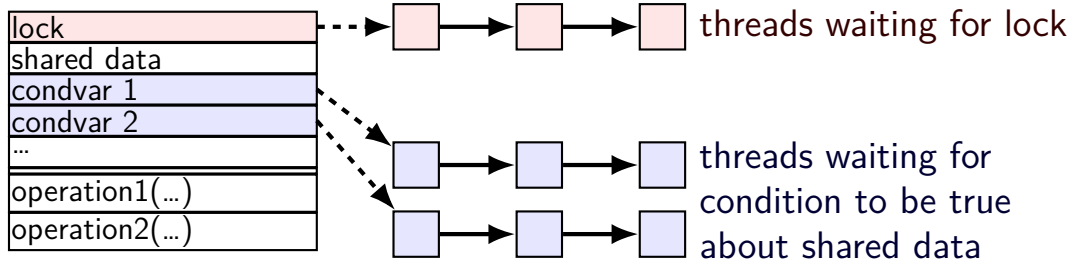before accessing
any part of monitor's stuff

# monitor idea

a monitor

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



threads waiting for lock

# monitor idea

a monitor



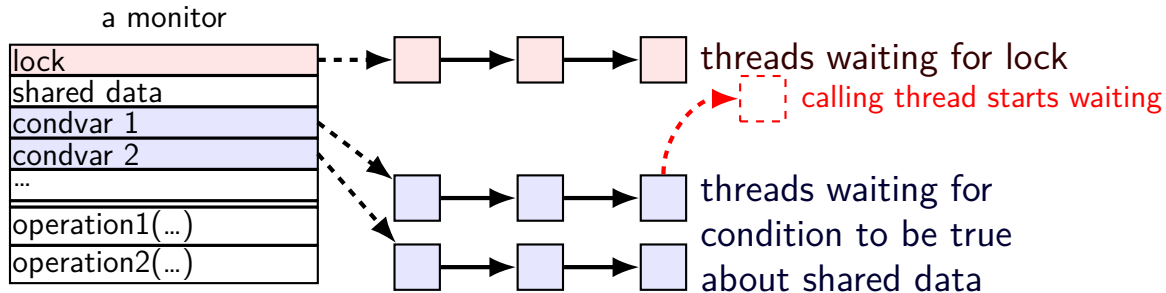| | |
|---|---|
| lock | threads waiting for lock |
| shared data | |
| condvar 1 | |
| condvar 2 | threads waiting for |
| ... | condition to be true |
| operation1(...) | about shared data |
| operation2(...) | |

# condvar operations

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

a monitor

# condvar operations
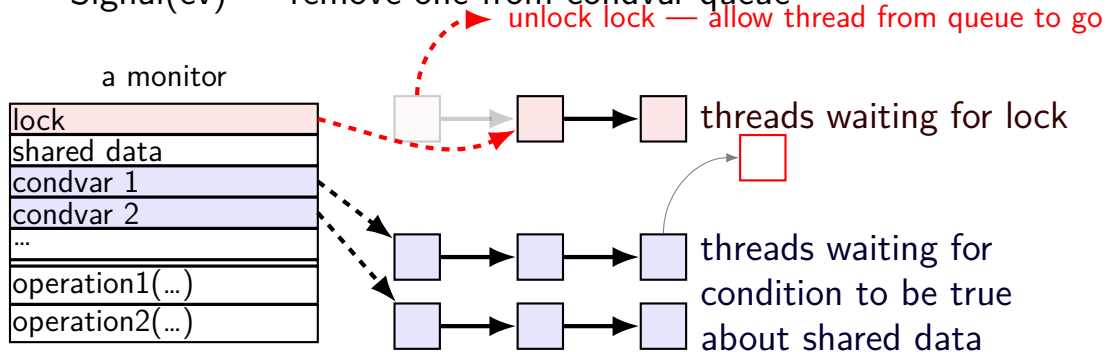
condvar operations:

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

…and reacquire lock before returning

Broadcast(cv) — remove all from condvar queue

Signal(cv) — remove one from condvar queue



a monitor

threads waiting for lock

calling thread starts waiting

threads waiting for condition to be true about shared data

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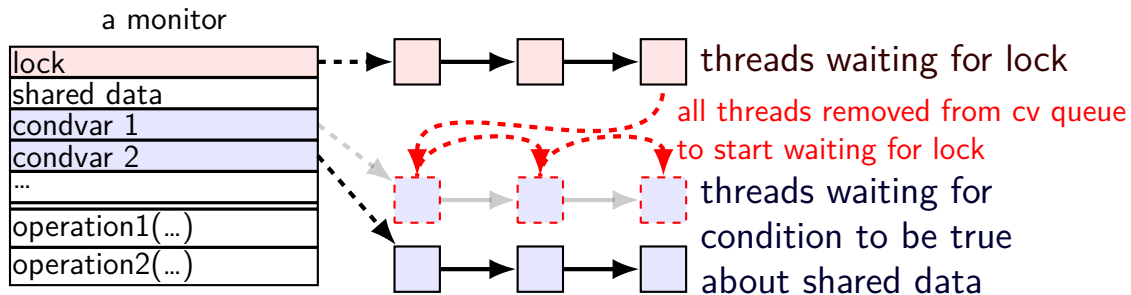operation1(…)
operation2(…)

# condvar operations

condvar operations:

Wait(cv, lock) — unlock lock, add current thread to cv queue
...and reacquire lock before returning

Broadcast(cv) — remove all from condvar queue

Signal(cv) — remove one from condvar queue
→ unlock lock — allow thread from queue to go

a monitor

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

threads waiting for lock

threads waiting for
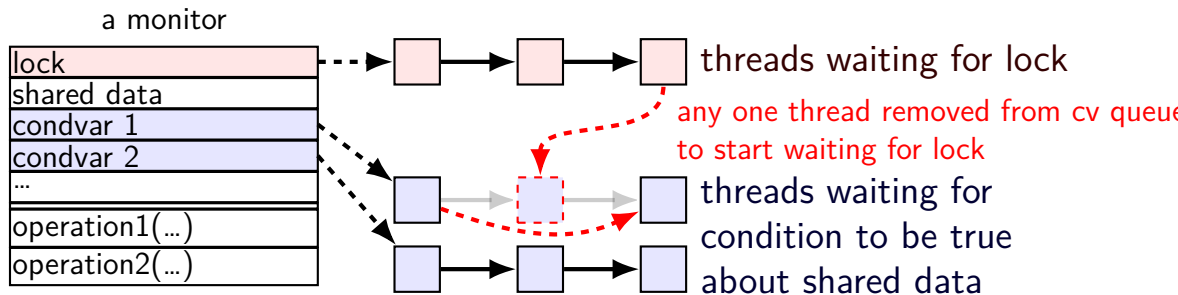condition to be true
about shared data

# condvar operations

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



a monitor

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

threads waiting for lock

all threads removed from cv queue
to start waiting for lock

threads waiting for
condition to be true
about shared data

# condvar operations

condvar operations:

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

...and reacquire lock before returning

Broadcast(cv) — remove all from condvar queue

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



a monitor

threads waiting for lock

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

threads waiting for
condition to be true
about shared data

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

# pthread cv usage

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

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

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

# pthread cv usage

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

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

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

acquire lock before
reading or writing `finished`

# pthread cv usage

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

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

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

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

# pthread cv usage

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

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

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

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

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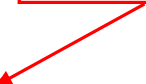
```
// MISSING: init calls, etc.
pthread_mutex_t lock;
bool finished;      // data, only accessed with after acquiring lock
pthread_cond_t finished_cv;  // to wait for 'finished' to be true

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

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

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

# WaitForFinish timeline 1

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

# WaitForFinish timeline 2

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

# why the loop

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

we only broadcast if finished is true

so why check finished afterwards?

# why the loop

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

we only `broadcast` if `finished` is true

so why check `finished` afterwards?

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

spurious wakeup = `wait` returns even though nothing happened

# 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
simulatenously 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 can**not** 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);
}
Consume() {
    pthread_mutex_lock(&lock);
    while (buffer.empty()) {
        pthread_cond_wait(&data_ready, &lock);
    }
    item = buffer.dequeue();
    pthread_mutex_unlock(&lock);
    return item;
}
```

wake one Consume thread
*if any are waiting*

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

| Thread 1 | Thread 2 | Thread 3 |
|---|---|---|
| | Consume() | |
| | ...lock | |
| | ...empty? yes | |
| | ...unlock/start wait | |
| Produce() | waiting for data_ready | |
| ...lock | | Consume() |
| ...enqueue | | waiting for lock |
| ...signal | stop wait | |
| ...unlock | waiting for lock | lock |
| | | ...empty? no |
| | | ...dequeue |
| | | ...unlock |
| | ...lock | return |
| | ...empty? yes | |
| | ...unlock/start wait | |

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

85

# unbounded buffer producer/consumer

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

in pthreads: signalled thread not
    gaurenteed 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;
}
```

| Thread 1 | Thread 2 | Thread 3 |
|---|---|---|
| | Consume() | |
| | ...lock | |
| | ...empty? yes | |
| | ...unlock/start wait | |
| Produce() | waiting for data_ready | |
| ...lock | | Consume() |
| ...enqueue | | waiting for lock |
| ...signal | stop wait | |
| ...unlock | | lock |
| | waiting for lock | ...empty? no |
| | | ...dequeue |
| | | ...unlock |
| | ...lock | return |
| | ...empty? yes | |
| | ...unlock/start wait | |

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

85

# 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

# bounded buffer producer/consumer

```
pthread_mutex_t lock;
pthread_cond_t data_ready; pthread_cond_t space_ready;
BoundedQueue buffer;
Produce(item) {
    pthread_mutex_lock(&lock);
    while (buffer.full()) { pthread_cond_wait(&space_ready, &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_cond_signal(&space_ready);
    pthread_mutex_unlock(&lock);
    return item;
}
```
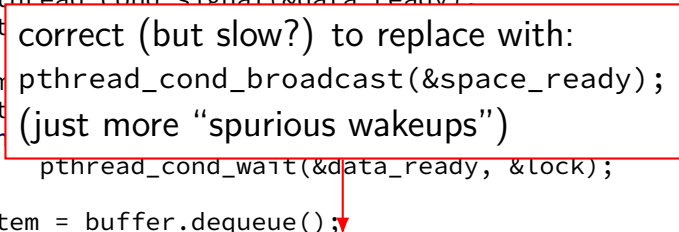
# bounded buffer producer/consumer

```
pthread_mutex_t lock;
pthread_cond_t data_ready; pthread_cond_t space_ready;
BoundedQueue buffer;
Produce(item) {
    pthread_mutex_lock(&lock);
    while (buffer.full()) { pthread_cond_wait(&space_ready, &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_cond_signal(&space_ready);
    pthread_mutex_unlock(&lock);
    return item;
}
```

# bounded buffer producer/consumer

```
pthread_mutex_t lock;
pthread_cond_t data_ready; pthread_cond_t space_ready;
BoundedQueue buffer;
Produce(item) {
    pthread_mutex_lock(&lock);
    while (buffer.full()) { pthread_cond_wait(&space_ready, &lock); }
    buffer.enqueue(item);
    pthread_cond_signal(&data_ready);
    pt
}                correct (but slow?) to replace with:
Consum
    pt           pthread_cond_broadcast(&space_ready);
    wh
                 (just more "spurious wakeups")
        pthread_cond_wait(&data_ready, &lock);
    }
    item = buffer.dequeue();
    pthread_cond_signal(&space_ready);
    pthread_mutex_unlock(&lock);
    return item;
}
```

# bounded buffer producer/consumer

```
pthread_mutex_t lock;
pthread_cond_t data_ready; pthread_cond_t space_ready;
BoundedQueue buffer;
Produce(item) {
    pthread_mutex_lock(&lock);
    while (buffer.full()) { pthread_cond_wait(&space_ready, &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_cond_signal(&space_ready);
    pthread_mutex_unlock(&lock);
    return item;
}
```

correct but slow to replace
`data_ready` and `space_ready`
with 'combined' condvar `ready`
and use broadcast
(just more "spurious wakeups")

# monitor pattern

```
pthread_mutex_lock(&lock);
while (!condition A) {
    pthread_cond_wait(&condvar_for_A, &lock);
}
... /* manipulate shared data, changing other conditions */
if (set condition A) {
    pthread_cond_broadcast(&condvar_for_A);
    /* or signal, if only one thread cares */
}
if (set condition B) {
    pthread_cond_broadcast(&condvar_for_B);
    /* or signal, if only one thread cares */
}
...
pthread_mutex_unlock(&lock)
```

# monitors rules of thumb

never touch shared data without holding the lock

keep lock held for entire operation:
    verifying condition (e.g. buffer not full) *up to and including*
    manipulating data (e.g. adding to buffer)

create condvar for every kind of scenario waited for

always write loop calling cond_wait to wait for condition X

broadcast/signal condition variable every time you change X

# monitors rules of thumb

never touch shared data without holding the lock

keep lock held for entire operation:
> verifying condition (e.g. buffer not full) *up to and including*
> manipulating data (e.g. adding to buffer)

create condvar for every kind of scenario waited for

always write loop calling cond_wait to wait for condition X

broadcast/signal condition variable every time you change X

correct but slow to…
> broadcast when just signal would work
> broadcast or signal when nothing changed
> use one condvar for multiple conditions

# mutex/cond var init/destroy

```
pthread_mutex_t mutex;
pthread_cond_t cv;
pthread_mutex_init(&mutex, NULL);
pthread_cond_init(&cv, NULL);
// --OR--
pthread_mutex_t mutex = PTHREAD_MUTEX_INITIALIZER;
pthread_cond_t cv = PTHREAD_COND_INITIALIZER;

// and when done:
...
pthread_cond_destroy(&cv);
pthread_mutex_destroy(&mutex);
```

## wait for both finished

```
// MISSING: init calls, etc.
pthread_mutex_t lock;
bool finished[2];
pthread_cond_t both_finished_cv;

void WaitForBothFinished() {
  pthread_mutex_lock(&lock);
  while (_____) {
    pthread_cond_wait(&both_finished_cv, &lock);
  }
  pthread_mutex_unlock(&lock);
}

void Finish(int index) {
  pthread_mutex_lock(&lock);
  finished[index] = true;

  _____
  pthread_mutex_unlock(&lock);
}
```

# wait for both finished

```
// MISSING: init calls, etc.
pthread_mutex_t lock;
bool finished[2];
pthread_cond_t both_finished_cv;

void WaitForBothFinished() {
  pthread_mutex_lock(&lock);
  while (_____) {
    pthread_cond_wait(&both_finished_cv, &lock);
  }
  pthread_mutex_unlock(&lock);
}

void Finish(int index) {
  pthread_mutex_lock(&lock);
  finished[index] = true;
  _____
  pthread_mutex_unlock(&lock);
}
```

A. finished[0] && finished[1]
B. finished[0] || finished[1]
C. !finished[0] || !finished[1]
D. finished[0] != finished[1]
E. something else

# wait for both finished

```
// MISSING: init calls, etc.
pthread_mutex_t lock;
bool finished[2];
pthread_cond_t both_fini

void WaitForBothFinished
  pthread_mutex_lock(&lo
  while (_____
    pthread_cond_wait(&both_finished_cv, &lock);
  }
  pthread_mutex_unlock(&lock);
}

void Finish(int index) {
  pthread_mutex_lock(&lock);
  finished[index] = true;
  _____
  pthread_mutex_unlock(&lock);
}
```

A. pthread_cond_signal(&both_finished_cv)
B. pthread_cond_broadcast(&both_finished_cv)
C. if (finished[1−index])
        pthread_cond_singal(&both_finished_cv);
D. if (finished[1−index])
        pthread_cond_broadcast(&both_finished_cv);
E. something else

# monitor exercise: barrier

suppose we want to implement a one-use barrier; fill in blanks:

```
struct BarrierInfo {
    pthread_mutex_t lock;
    int total_threads; // initially total # of threads
    int number_reached; // initially 0
    _____
};
void BarrierWait(BarrierInfo *b) {
    pthread_mutex_lock(&b->lock);
    ++b->number_reached;
    if (b->number_reached == b->total_threads) {
        _____
    } else {
        _____
        _____
    }
    pthread_mutex_unlock(&b->lock);
}
```

## transactions

transaction: set of operations that occurs atomically

idea: something higher-level handles locking, etc.:

```
BeginTransaction();
int FromOldBalance = GetBalance(FromAccount);
int ToOldBalance = GetBalance(ToAccount);
SetBalance(FromAccount, FromOldBalance - 100);
SetBalance(ToAccount, FromOldBalance + 100);
EndTransaction();
```

idea: library/database/etc. makes "transaction" happens all at once

# consistency / durability

"happens all at once" = could mean:

locking to make sure no other operations interfere (consistency)

making sure on crash, no partial transaction seen (durability)

(some systems provide both, some provide only one)

we'll just talk about implementing consistency

# implementing consistency: simple

simplest idea: only one run transaction at a time

# implementing consistency: locking

everytime something read/written: acquire associated lock

on end transaction: release lock

if deadlock: undo everything, go back to BeginTransaction(), retry
    how to undo?
    one idea: keep list of writes instead of writing
    apply writes only at EndTransaction()

# implementing consistency: locking

everytime something read/written: acquire associated lock

on end transaction: release lock

if deadlock: undo everything, go back to BeginTransaction(), retry
    how to undo?
    one idea: keep list of writes instead of writing
    apply writes only at EndTransaction()

# implementing consistency: optimistic

on read: copy version # for value read

on write: record value to be written, but don't write yet

on end transaction:
  acquire locks on everything
  make sure values read haven't been changed since read

if they have changed, just retry transaction