barriers (finish) / deadlock

# last time

race conditions
    inconsistent results due to timing variation
    example: "lose" update due to reading value while update being
    computed

compilers, processors and memory access reordering
    order you write in C code [or even assembly] might not be order of
    accesses
    need special operations that gaurentee consistent order

locks for taking turns
    one thread can "hold" lock at a time
    lock operation waits for lock to be available (unlock'd)
    requires threads agree to get lock before using shared thing

barriers — advance threads in lock-step

# exercise

*pthread. barrier -init ()*

```
pthread_barrier_t barrier; int x = 0, y = 0;
void thread_one() {
    y = 10;
    pthread_barrier_wait(&barrier);          ← ①
    y = x + y;       y = 30
    pthread_barrier_wait(&barrier);          ← ②
    pthread_barrier_wait(&barrier);          ← ③
    printf("%d %d\n", x, y);
}                    50  30          ← after 3rd barrier
void thread_two() {
    x = 20;
    pthread_barrier_wait(&barrier);          ← ①
    pthread_barrier_wait(&barrier);          ← ②
    x = x + y;   X = 20 + 30  = 50
    pthread_barrier_wait(&barrier);          ← ③
}
```
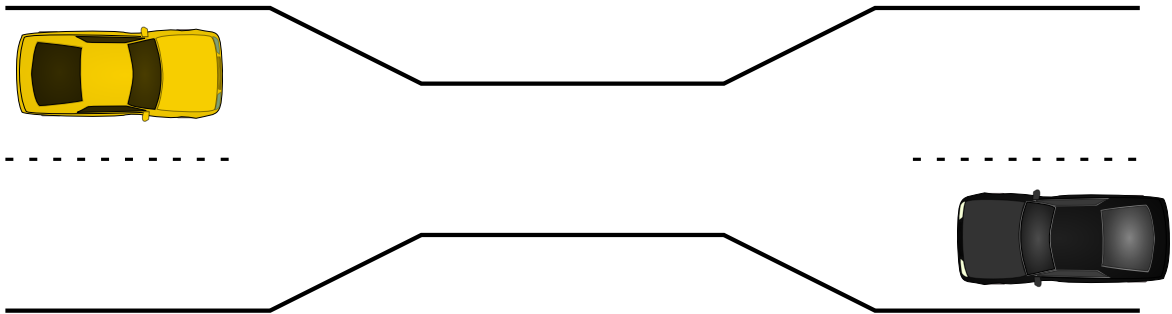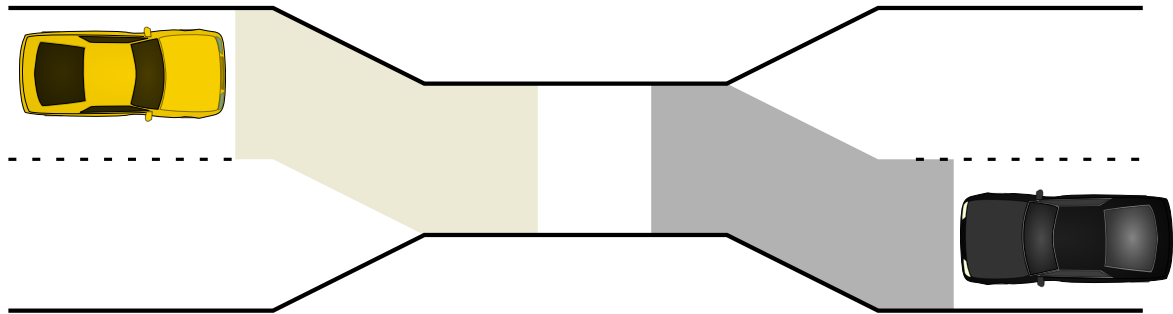
*main ()*
*barrier -init (2)*
*thread - one*
*thread - two*

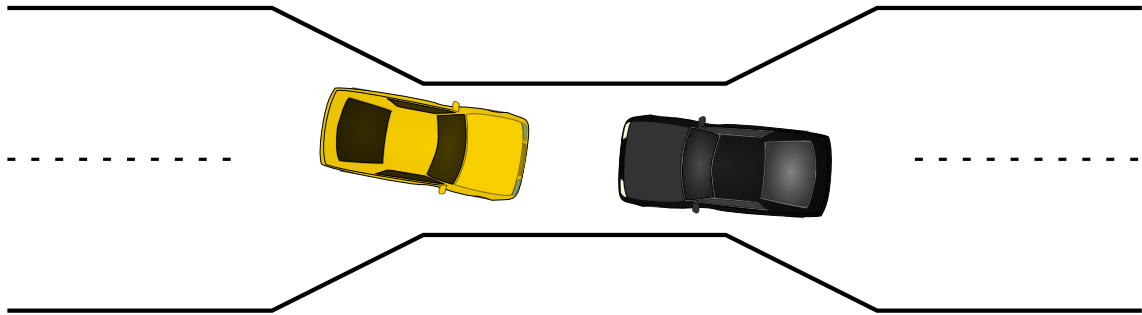output? (if both run at once, barrier set for 2 threads)
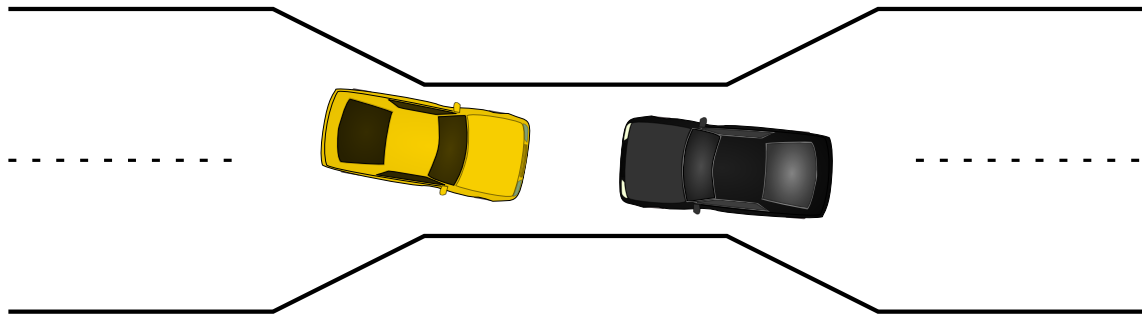
3

# 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); | |
| lock(&B->lock); | |
| | lock(&B->lock… |
| | (waiting for B lock) |
| (do move) | |
| unlock(&B->lock); | |
| | lock(&B->lock); |
| | lock(&A->lock… |
| unlock(&A->lock); | |
| | lock(&A->lock); |
| | (do move) |
| | unlock(&A->lock); |

# moving two files: unlucky timeline

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

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

# 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 |

# moving two files: unlucky timeline

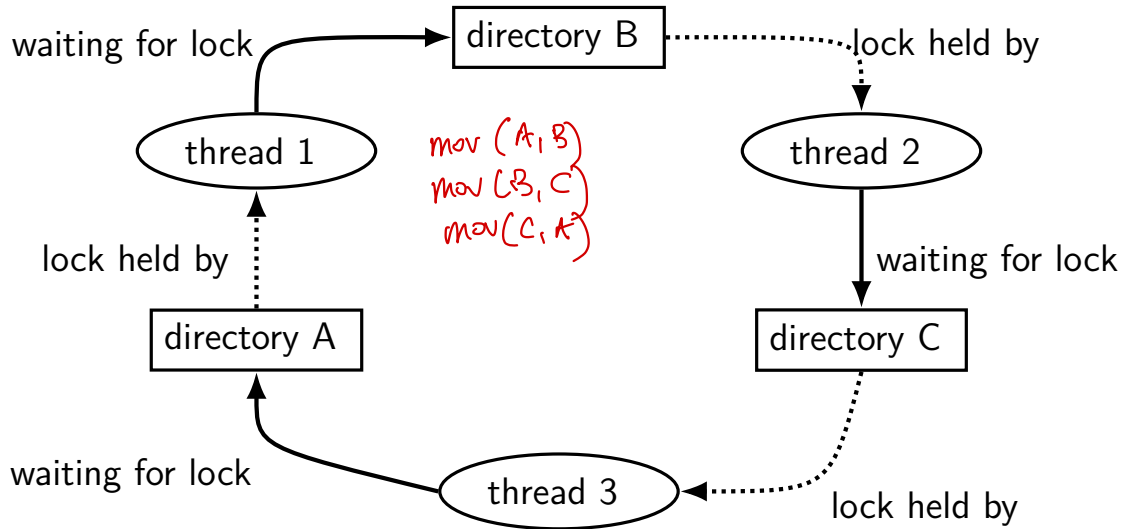| Thread 1 | Thread 2 |
|---|---|
| MoveFile(A, B, "foo") | MoveFile(B, A, "bar") |

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

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

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

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

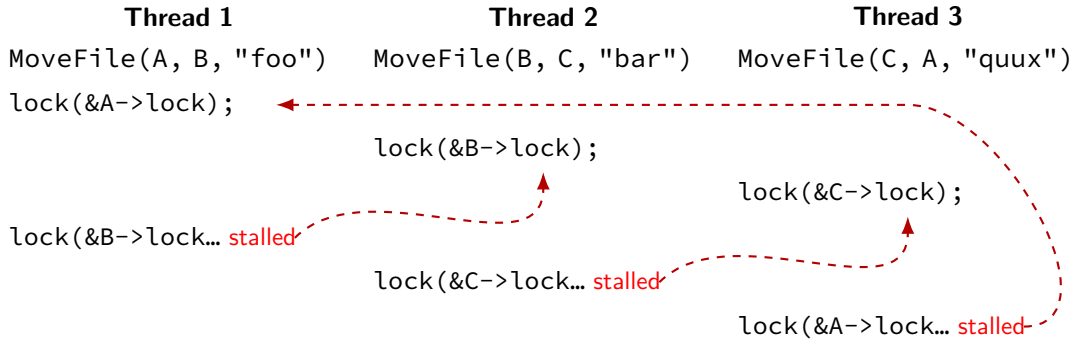Thread 1 holds A lock, waiting for Thread 2 to release B lock

# moving two files: dependencies

# moving three files: dependencies

# moving three files: unlucky timeline



| Thread 1 | Thread 2 | Thread 3 |
|----------|----------|----------|

```
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) | *1* | AllocateOrWaitFor(1 MB) | *2* |
| AllocateOrWaitFor(1 MB) | *3rd* | AllocateOrWaitFor(1 MB) | *4th* |
| → (do calculation) | *not avail* | (do calculation) | *not aval* |
| Free(1 MB) | | Free(1 MB) | |
| Free(1 MB) | | Free(1 MB) | |

2 MB of space — deadlock possible with unlucky order

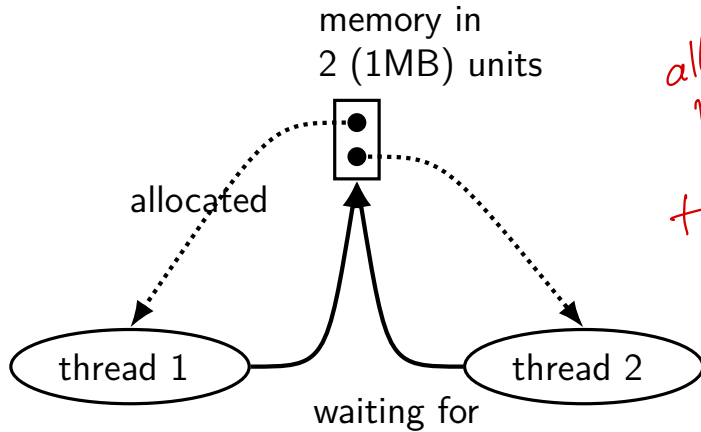# deadlock with free space (unlucky case)

| 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

waiting for

thread 2

allocation
held by... First MB

waiting

thr1

waiting
for

thr2

holds
alloc

Second MB

# 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); |

# lab next week

applying solutions to deadlock to classic *dining philosphers* problem

# dining philosophers



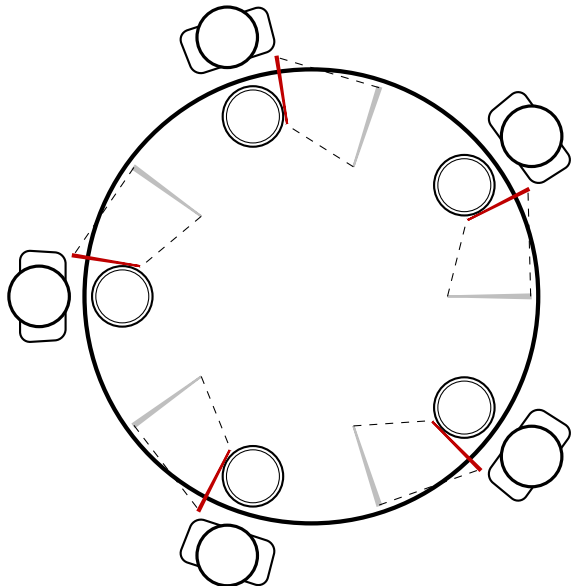five philosophers either think or eat
to eat:
chopstix[tid]
grab chopstick on left, then
grba chopstick on right, then
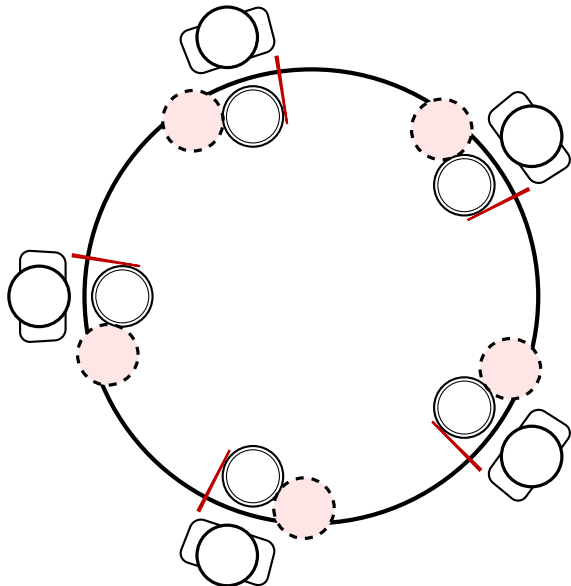then eat, then    chopstix[tid+1]
return chopsticks

# 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

# how is deadlock — solution

| Remove B | Remove C |
|----------|----------|
| lock B | lock C |
| lock A (prev) | wait to lock B (prev) |
| wait to lock C (next) | |

With B and D — only overlap in in node C — no circular wait possible
(thread can't be waiting while holding something other thread wants)

Rem B ⊂ deadlock ⟶ Rem C ⊂ deadlock ⟶ Rem D

Rem B

lock B
lock A

lock C

Rem D

lock D
lock C
lock E

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

**<span style="color:red">no shared resources</span>**            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*

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

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

**no waiti**
> "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                    no *hold and wait*/
  revoke/preempt resources                                   *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);
  }
  ...
}
```

Thr1 : mov (A, B) → lock A
lock B

Thr2: mov (B, A) → lock A
lock B

Thr1
lock A
lock B
mov A→B

Thr 2
stall on lock A
acq A
lock B
mov B→A

# 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

or lexicographic sort
numerical sort

# 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** *avoid*
    or at least enough that never run out                    no *mutual exclusion*

**no shared resources** *avoid*                              no *mutual exclusion*

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

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

# backup slides

# 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                    Thread 1

```
partial_mins[0] =
    /* min of first
       50M elems */;                partial_mins[1] =
                                        /* min of last
barrier.Wait();                           50M elems */
                                   barrier.Wait();


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

# 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]
    );
```

51

# pthread barriers

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