# monitors (con't) / deadlock

# last time

correction re: binary semaphores with counting

monitor pattern: condition variable for each reason to wait
    loop checking reason + wait on CV
    broadcast/signal when need to wait might have changed

monitor sloppiness: spurious wakeups, signal/broadcast more than needed...

reader/writer locks

pattern: count waiters for waiter priority

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
```

lock to protect shared state

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
unsigned int readers, writers;
```

state: number of active readers, writers

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
unsigned int readers, writers;
/* condition, signal when writers becomes 0 */
cond_t ok_to_read_cv;
/* condition, signal when readers + writers becomes 0 */
cond_t ok_to_write_cv;
```

conditions to wait for (no readers or writers, no writers)

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
unsigned int readers, writers;
/* condition, signal when writers becomes 0 */
cond_t ok_to_read_cv;
/* condition, signal when readers + writers becomes 0 */
cond_t ok_to_write_cv;
```

```
ReadLock() {
  mutex_lock(&lock);
  while (writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);
  }
  ++readers;
  mutex_unlock(&lock);
}
ReadUnlock() {
  mutex_lock(&lock);
  --readers;
  if (readers == 0) {
    cond_signal(&ok_to_write_cv);
  }
  mutex_unlock(&lock);
}
```

```
WriteLock() {
  mutex_lock(&lock);
  while (readers + writers != 0) {
    cond_wait(&ok_to_write_cv);
  }
  ++writers;
  mutex_unlock(&lock);
}
WriteUnlock() {
  mutex_lock(&lock);
  --writers;
  cond_signal(&ok_to_write_cv);
  cond_broadcast(&ok_to_read_cv);
  mutex_unlock(&lock);
}
```

broadcast — wakeup all readers when no writers

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
unsigned int readers, writers;
/* condition, signal when writers becomes 0 */
cond_t ok_to_read_cv;
/* condition, signal when readers + writers becomes 0 */
cond_t ok_to_write_cv;
ReadLock() {                        WriteLock() {
  mutex_lock(&lock);                  mutex_lock(&lock);
  while (writers != 0) {              while (readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);   cond_wait(&ok_to_write_cv);
  }                                   }
  ++readers;                          ++writers;
  mutex_unlock(&lock);                mutex_unlock(&lock);
}                                   }
ReadUnlock() {                      WriteUnlock() {
  mutex_lock(&lock);                  mutex_lock(&lock);
  --readers;                          --writers;
  if (readers == 0) {                 cond_signal(&ok_to_write_cv);
    cond_signal(&ok_to_write_cv);     cond_broadcast(&ok_to_read_cv);
  }                                   mutex_unlock(&lock);
  mutex_unlock(&lock);              }
}
```

wakeup a single writer when no readers or writers

# rwlocks with monitors (attempt 1)

```
mutex_t lock;
unsigned int readers, writers;
/* condition, signal when writers becomes 0 */
cond_t ok_to_read_cv;
/* condition, signal when readers + writers becomes 0 */
cond_t ok_to_write_cv;
ReadLock() {                              WriteLock() {
  mutex_lock(&lock);                        mutex_lock(&lock);
  while (writers != 0) {                     while (readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);          cond_wait(&ok_to_write_cv);
  }                                          }
  ++readers;                                 ++writers;
  mutex_unlock(&lock);                       mutex_unlock(&lock);
}                                         }
ReadUnlock() {                            WriteUnlock() {
  mutex_lock(&lock);                        mutex_lock(&lock);
  --readers;                                --writers;
  if (readers == 0) {                       cond_signal(&ok_to_write_cv);
    cond_signal(&ok_to_write_cv);           cond_broadcast(&ok_to_read_cv);
  }                                         mutex_unlock(&lock);
  mutex_unlock(&lock);                    }
}
```

problem: wakeup readers first or writer first?

this solution: wake them all up and they fight! inefficient!

# reader/writer-priority

policy question: writers first or readers first?
   writers-first: no readers go when writer waiting
   readers-first: no writers go when reader waiting

previous implementation: whatever randomly happens
   writers signalled first, maybe gets lock first?
   …but non-determinstic in pthreads

can make explicit decision

# writer-priority (1)

```
mutex_t lock; cond_t ok_to_read_cv; cond_t ok_to_write_cv;
int readers = 0, writers = 0;
int waiting_writers = 0;
ReadLock() {                          WriteLock() {
  mutex_lock(&lock);                    mutex_lock(&lock);
  while (writers != 0                   ++waiting_writers;
        && waiting_writers != 0) {      while (readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);     cond_wait(&ok_to_write_cv, &lock);
  }                                     }
  ++readers;                            --waiting_writers;
  mutex_unlock(&lock);                  ++writers;
}                                       mutex_unlock(&lock);
                                      }
ReadUnlock() {
  mutex_lock(&lock);                  WriteUnlock() {
  --readers;                            mutex_lock(&lock);
  if (readers == 0) {                   --writers;
    cond_signal(&ok_to_write_cv);       if (waiting_writers != 0) {
  }                                       cond_signal(&ok_to_write_cv);
  mutex_unlock(&lock);                  } else {
}                                         cond_broadcast(&ok_to_read_cv);
                                        }
                                        mutex_unlock(&lock);
                                      }
```

# writer-priority (1)

```
mutex_t lock; cond_t ok_to_read_cv; cond_t ok_to_write_cv;
int readers = 0, writers = 0;
int waiting_writers = 0;
ReadLock() {                          WriteLock() {
  mutex_lock(&lock);                    mutex_lock(&lock);
  while (writers != 0                   ++waiting_writers;
         && waiting_writers != 0) {     while (readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);     cond_wait(&ok_to_write_cv, &lock);
  }                                     }
  ++readers;                            --waiting_writers;
  mutex_unlock(&lock);                  ++writers;
}                                       mutex_unlock(&lock);
                                      }
ReadUnlock() {
  mutex_lock(&lock);                  WriteUnlock() {
  --readers;                            mutex_lock(&lock);
  if (readers == 0) {                   --writers;
    cond_signal(&ok_to_write_cv);       if (waiting_writers != 0) {
  }                                       cond_signal(&ok_to_write_cv);
  mutex_unlock(&lock);                  } else {
}                                         cond_broadcast(&ok_to_read_cv);
                                        }
                                        mutex_unlock(&lock);
                                      }
```

5

# writer-priority (1)

```
mutex_t lock; cond_t ok_to_read_cv; cond_t ok_to_write_cv;
int readers = 0, writers = 0;
int waiting_writers = 0;
ReadLock() {                        WriteLock() {
  mutex_lock(&lock);                  mutex_lock(&lock);
  while (writers != 0                 ++waiting_writers;
         && waiting_writers != 0) {   while (readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);   cond_wait(&ok_to_write_cv, &lock);
  }                                   }
  ++readers;                          --waiting_writers;
  mutex_unlock(&lock);                ++writers;
}                                     mutex_unlock(&lock);
                                    }

ReadUnlock() {
  mutex_lock(&lock);                WriteUnlock() {
  --readers;                          mutex_lock(&lock);
  if (readers == 0) {                 --writers;
    cond_signal(&ok_to_write_cv);     if (waiting_writers != 0) {
  }                                     cond_signal(&ok_to_write_cv);
  mutex_unlock(&lock);                } else {
}                                       cond_broadcast(&ok_to_read_cv);
                                      }
                                      mutex_unlock(&lock);
                                    }
```

5

# reader-priority (1)

```
...
int waiting_readers = 0;
ReadLock() {                          WriteLock() {
  mutex_lock(&lock);                    mutex_lock(&lock);
  ++waiting_readers;                    while (waiting_readers +
  while (writers != 0) {                        readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);     cond_wait(&ok_to_write_cv);
  }                                     }
  --waiting_readers;                    ++writers;
  ++readers;                            mutex_unlock(&lock);
  mutex_unlock(&lock);                }
}                                     WriteUnlock() {
                                        mutex_lock(&lock);
ReadUnlock() {                          --writers;
  ...                                   if (waiting_readers == 0) {
  if (waiting_readers == 0) {             cond_signal(&ok_to_write_cv);
    cond_signal(&ok_to_write_cv);       } else {
  }                                       cond_broadcast(&ok_to_read_cv);
}                                       }
                                        mutex_unlock(&lock);
                                      }
```

# reader-priority (1)

```
...
int waiting_readers = 0;
ReadLock() {                          WriteLock() {
  mutex_lock(&lock);                    mutex_lock(&lock);
  ++waiting_readers;                    while (waiting_readers +
  while (writers != 0) {                       readers + writers != 0) {
    cond_wait(&ok_to_read_cv, &lock);     cond_wait(&ok_to_write_cv);
  }                                     }
  --waiting_readers;                    ++writers;
  ++readers;                            mutex_unlock(&lock);
  mutex_unlock(&lock);                }
}                                     WriteUnlock() {
                                        mutex_lock(&lock);
ReadUnlock() {                          --writers;
  ...                                   if (waiting_readers == 0) {
  if (waiting_readers == 0) {             cond_signal(&ok_to_write_cv);
    cond_signal(&ok_to_write_cv);       } else {
  }                                       cond_broadcast(&ok_to_read_cv);
}                                       }
                                        mutex_unlock(&lock);
                                      }
```

# choosing orderings?

can use monitors to implement lots of lock policies

want $X$ to go first/last — add extra variables
    (number of waiters, even lists of items, etc.)

need way to write condition "you can go now"
    e.g. writer-priority: readers can go if no writer waiting

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|----------|----------|----------|----------|---|---|----|
|          |          |          |          | 0 | 0 | 0  |

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|---|---|---|---|---|---|---|
| | | | | 0 | 0 | 0 |

# simulation of reader/write lock

writer-priority version

$W$ = writers, $R$ = readers, $WW$ = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|----------|----------|----------|----------|---|---|----|
| | | | | 0 | 0 | 0 |
| ReadLock | | | | 0 | 1 | 0 |

```
mutex_lock(&lock);
while (writers != 0 && waiting_writers != 0) {
  cond_wait(&ok_to_read_cv, &lock);
}
++readers;
mutex_unlock(&lock);
```

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|---|---|---|---|---|---|---|
| | | | | 0 | 0 | 0 |
| ReadLock | | | | 0 | 1 | 0 |
| (reading) | ReadLock | | | 0 | 2 | 0 |

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|---|---|---|---|---|---|---|
| | | | | 0 | 0 | 0 |
| ReadLock | | | | 0 | 1 | 0 |
| (reading) | ReadLock | | | 0 | 2 | 0 |
| (reading) | (reading) | WriteLock wait | | 0 | 2 | 1 |

```
mutex_lock(&lock);
++waiting_writers;
while (readers + writers != 0) {
  cond_wait(&ok_to_write_cv, &lock);
}
```

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|---|---|---|---|---|---|---|
| | | | | 0 | 0 | 0 |
| ReadLock | | | | 0 | 1 | 0 |
| (reading) | ReadLock | | | 0 | 2 | 0 |
| (reading) | (reading) | WriteLock wait | | 0 | 2 | 1 |
| (reading) | (reading) | WriteLock wait | ReadLock wait | 0 | 2 | 1 |

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|----------|----------|----------|----------|---|---|----|
| | | | | 0 | 0 | 0 |
| ReadLock | | | | 0 | 1 | 0 |
| (reading) | ReadLock | | | 0 | 2 | 0 |
| (reading) | (reading) | WriteLock wait | | 0 | 2 | 1 |
| (reading) | (read | wait | ReadLock wait | 0 | 2 | 1 |
| ReadUnlock | ← | wait | ReadLock wait | 0 | 1 | 1 |

```
mutex_lock(&lock);
--readers;
if (readers == 0)
    ...
```

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|----------|----------|----------|----------|---|---|-----|
|          |          |          |          | 0 | 0 | 0 |
| ReadLock |          |          |          | 0 | 1 | 0 |
| (reading) | ReadLock |          |          | 0 | 2 | 0 |
| (reading) | (reading) | WriteLock wait |     | 0 | 2 | 1 |
| (reading) | (reading) | WriteLock wait | ReadLock wait | 0 | 2 | 1 |
| ReadUnlock | (reading) | Write |          | 1 | 1 |  |
|          | ReadUnlock | W... |          | 0 | 1 |  |

```
mutex_lock(&lock);
--readers;
if (readers == 0)
  cond_signal(&ok_to_write_cv)
mutex_unlock(&lock);
```

8

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|----------|----------|----------|----------|---|---|----|
| | | | | 0 | 0 | 0 |
| ReadLock | | | | 0 | 1 | 0 |
| (reading) | Read | | | 0 | 2 | 0 |
| (reading) | (rea | | | 0 | 2 | 1 |
| (reading) | (rea | | it | 0 | 2 | 1 |
| ReadUnlock | (reading) | WriteLock wait | ReadLock wait | 0 | 1 | 1 |
| | ReadUnlock | WriteLock wait | ReadLock wait | 0 | 0 | 1 |
| | | WriteLock | ReadLock wait | 1 | 0 | 0 |

```
while (readers + writers != 0) {
    cond_wait(&ok_to_write_cv, &lock);
}
--waiting_writers; ++writers;
mutex_unlock(&lock);
```

8

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|----------|----------|----------|----------|---|---|----|
|          |          |          |          | 0 | 0 | 0  |
| ReadLock |          |          |          | 0 | 1 | 0  |
| (reading) | ReadLock |          |          | 0 | 2 | 0  |
| (reading) | (reading) | WriteLock wait |     | 0 | 2 | 1  |
| (reading) | (reading) | WriteLock wait | ReadLock wait | 0 | 2 | 1 |
| ReadUnlock | (reading) | WriteLock wait | ReadLock wait | 0 | 1 | 1 |
|          | ReadUnlock | WriteLock wait | ReadLock wait | 0 | 0 | 1 |
|          |          | WriteLock | ReadLock wait | 1 | 0 | 0 |
|          |          | (read+writing) | ReadLock wait | 1 | 0 | 0 |

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|---|---|---|---|---|---|---|
| | | | | 0 | 0 | 0 |
| ReadLock | | | | 0 | 1 | 0 |
| (reading) | ReadLock | | | 0 | 2 | 0 |
| (reading) | (readi... | | | 0 | 2 | 1 |
| (reading) | (readi... | | wait | 0 | 2 | 1 |
| ReadUnlock | (readi... | | wait | 0 | 1 | 1 |
| | ReadUn... | | wait | 0 | 0 | 1 |
| | | WriteLock | ReadLock wait | 1 | 0 | 0 |
| | | (read+writing) | ReadLock wait | 1 | 0 | 0 |
| | | WriteUnlock | ReadLock wait | 0 | 0 | 0 |

```
mutex_lock(&lock);
if (waiting_writers != 0) {
  cond_signal(&ok_to_write_cv);
} else {
  cond_broadcast(&ok_to_read_cv);
}
```

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|----------|----------|----------|----------|---|---|-----|
|          |          |          |          | 0 | 0 | 0 |
| ReadLock |          |          |          | 0 | 1 | 0 |
| (reading) | ReadLock |         |          | 0 | 2 | 0 |
| (reading) | (reading) | `while (writers != 0 && waiting_writers != 0) {` | | | | |
| (reading) | (reading) | `    cond_wait(&ok_to_read_cv, &lock);` | | | | |
| (reading) | (reading) | `}` | | | | |
| ReadUnlock | (reading) | `++readers;` | | | | |
|          | ReadUnlock | `mutex_unlock(&lock);` | | | | |
|          |          | WriteLock | ReadLock wait | 1 | 0 | 0 |
|          |          | (read+writing) | ReadLock wait | 1 | 0 | 0 |
|          |          | WriteUnlock | ReadLock wait | 0 | 0 | 0 |
|          |          |          | ReadLock | 0 | 1 | 0 |

8

# simulation of reader/write lock

writer-priority version

W = writers, R = readers, WW = waiting_writers

| reader 1 | reader 2 | writer 1 | reader 3 | W | R | WW |
|---|---|---|---|---|---|---|
| | | | | 0 | 0 | 0 |
| ReadLock | | | | 0 | 1 | 0 |
| (reading) | ReadLock | | | 0 | 2 | 0 |
| (reading) | (reading) | WriteLock wait | | 0 | 2 | 1 |
| (reading) | (reading) | WriteLock wait | ReadLock wait | 0 | 2 | 1 |
| ReadUnlock | (reading) | WriteLock wait | ReadLock wait | 0 | 1 | 1 |
| | ReadUnlock | WriteLock wait | ReadLock wait | 0 | 0 | 1 |
| | | WriteLock | ReadLock wait | 1 | 0 | 0 |
| | | (read+writing) | ReadLock wait | 1 | 0 | 0 |
| | | WriteUnlock | ReadLock wait | 0 | 0 | 0 |
| | | | ReadLock | 0 | 1 | 0 |

# rwlock exercise

suppose there are multiple waiting writers

which one gets waken up first?
  whichever gets signal'd or gets lock first

could instead keep in order they started waiting

exercise: what extra information should we track?
  hint: we might need an array

```
mutex_t lock; cond_t ok_to_read_cv, ok_to_write_cv;
int readers, writers, waiting_writers;
```

# rwlock exercise solution?

list of waiting writes?

```
struct WaitingWriter {
    cond_t cv;
    bool ready;
};
Queue<WaitingWriter*> waiting_writers;

WriteLock(...) {
  ...
  if (need to wait) {
    WaitingWriter self;
    self.ready = false;
    ...
    while(!self.ready) {
        pthread_cond_wait(&self.cv, &lock);
    }
  }
  ...
}
```

# rwlock exercise solution?

dedicated writing thread with queue
     (DoWrite~Produce; WritingThread~Consume)

```
ThreadSafeQueue<WritingTask*> waiting_writes;
WritingThread() {
    while (true) {
        WritingTask* task = waiting_writer.Dequeue();
        WriteLock();
        DoWriteTask(task);
        task.done = true;
        cond_broadcast(&task.cv);
    }
}
DoWrite(task) {
    // instead of WriteLock(); DoWriteTask(...); WriteUnlock()
    WritingTask task = ...;
    waiting_writes.Enqueue(&task);
    while (!task.done) { cond_wait(&task.cv); }
}
```

# the one-way bridge

# the one-way bridge

# the one-way bridge

# the one-way bridge

# dining philosophers



five philosophers either think or eat
to eat, grab chopsticks on either side

# 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

# 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[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

# moving two files

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

  to_dir->entries[filename] = 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: 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); |
| | 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<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) |

# 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…` <span style="color:red">stalled</span>

<span style="color:gray">(waiting for lock on B)</span>

<span style="color:gray">(waiting for lock on B)</span>               `lock(&A->lock…` <span style="color:red">stalled</span>

                                                    <span style="color:gray">(waiting for lock on A)</span>

~~(do move)~~ <span style="color:red">unreachable</span>                         ~~(do move)~~ <span style="color:red">unreachable</span>

~~unlock(&B->lock);~~ <span style="color:red">unreachable</span>           ~~unlock(&A->lock);~~ <span style="color:red">unreachable</span>

~~unlock(&A->lock);~~ <span style="color:red">unreachable</span>           ~~unlock(&B->lock);~~ <span style="color:red">unreachable</span>

# 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

# moving three files: dependencies

# moving three files: unlucky timeline



**Thread 1**
MoveFile(A, B, "foo")
lock(&A->lock);

lock(&B->lock… stalled

**Thread 2**
MoveFile(B, C, "bar")
lock(&B->lock);

lock(&C->lock… stalled

**Thread 3**
MoveFile(C, A, "quux")

lock(&C->lock);

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**

```
AllocateOrWaitFor(1 MB)

AllocateOrWaitFor(1 MB… stalled
```

**Thread 2**

```
AllocateOrWaitFor(1 MB)

AllocateOrWaitFor(1 MB… stalled
```

# 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);`                          |

# 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: <span style="color:red">when acquiring multiple locks</span>

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

# deadlock prevention techniques

**infinite resources**
    or at least enough that never run out

no *mutual exclusion*

**no shared resources**

no *mutual exclusion*

**no waiting** (e.g. abort and retry)
    "busy signal"

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*

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

**no waiting** (e.g. abort and retry)          no *hold and wait*/
    "busy signal"          *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** (e.g. abort and retry)      no *hold and wait*/
    "busy signal"      *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*

**<span style="color:red">no waiting</span>** (e.g. abort and retry)
    "busy signal"

no *hold and wait*/
*preemption*

acquire resources in **consistent order**

no *circular wait*

request **all resources at once**

no *hold and wait*

# AllocateOrFail

| **Thread 1** | **Thread 2** |
|---|---|

```
AllocateOrFail(1 MB)
```

                                   `AllocateOrFail(1 MB)`

`AllocateOrFail(1 MB)` <span style="color:red">fails!</span>

                                   `AllocateOrFail(1 MB)` <span style="color:red">fails!</span>

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

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


okay, now what?

    give up?

    both try again? — maybe this will keep happening? (called <span style="color:red">livelock</span>)

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

# AllocateOrSteal

| **Thread 1** | **Thread 2** |
|---|---|
| `AllocateOrSteal(1 MB)` | |
| | `AllocateOrSteal(1 MB)` |
| `AllocateOrSteal(1 MB)` | Thread killed to free 1MB |
| (do work) | |

problem: can one actually implement this?

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

# fail/steal with locks

pthreads 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

# livelock

abort-and-retry

how many times will you retry?

# moving two files: abort-and-retry

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

  to_dir->entries[filename] = 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<br>MoveFile(A, B, "foo") | Thread 2<br>MoveFile(B, A, "bar") |
| --- | --- |
| trylock(&A->lock) → LOCKED | |
| | trylock(&B->lock) → LOCKED |
| trylock(&B->lock) → FAILED | |
| | trylock(&A->lock) → FAILED |
| unlock(&A->lock) | |
| | unlock(&B->lock) |
| trylock(&A->lock) → LOCKED | |
| | trylock(&B->lock) → LOCKED |
| trylock(&B->lock) → FAILED | |
| | trylock(&A->lock) → FAILED |
| unlock(&A->lock) | |
| | unlock(&B->lock) |

# livelock

like deadlock — no one's making progress
> potentially forever

unlike deadlock — threads are trying

...but keep aborting and retrying

# preventing livelock

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

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

other ideas?

# stealing locks???

how do we make stealing locks possible

# revokable locks

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

# Linux out-of-memory killer

Linux by default *overcommits* memory
    tell processes they have more memory than is available
    (some recommend disabling this feature)

problem: what if wrong?
    could wait for program to finish, free memory…
    but could be waiting forever because of deadlock

solution: kill a process
    (and try to choose one that's not important)

# database transactions

databases operations organized into *transactions*
>     happens all at once or not at all

until transaction is *committed*, not finalized

code to undo transaction in case it's not okay

database deadlock solution: invoke undo transaction code

...then rerun transaction later

# deadlock prevention techniques

**infinite resources**
    or at least enough that never run out

no *mutual exclusion*

**no shared resources**

no *mutual exclusion*

**no waiting** (e.g. abort and retry)
    "busy signal"

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** (e.g. abort and retry)
    "busy signal"

no *hold and wait*/
*preemption*

acquire resources in **consistent order**

no *circular wait*

request **all resources at once**

no *hold and wait*

# allocating all at once?

for resources like disk space, memory

figure out maximum allocation <span style="color:red">when starting thread</span>
   "only" need conservative estimate

only start thread if those resources are available

okay solution for embedded systems?

# deadlock detection

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

idea: search for cyclic dependencies

need:

list of all contended resources
what thread is waiting for what?
what thread 'owns' what?

# aside: deadlock detection in reality

instrument all contended resources?
    add tracking of who locked what
    modify every lock implementation — no simple spinlocks?
    some tricky cases: e.g. what about counting semaphores?

doing something useful on deadlock?
    want way to "undo" partially done operations

...but done for some applications

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

# resource allocation graphs

nodes: resources or threads

edge thread→resource: thread waiting for resource

edge resource→thread: resource is "owned" by thread
    holds lock on
    will be deallocated by
    …

# resource allocate graphs

# searching for cycles

cycle $\rightarrow$ deadlock happened!

finding cycles: recall 2150 topological sort (maybe???)

# divided resources

what about resources like memory?

allocating 1MB of memory:
    thread 'owns' the 1MB, but…
    another thread can use can use any other 1MB

want to track all of memory together

"partial ownership"
    locked half the memory

# dividable/interchangeable resources

# deadlock detection

cycle-finding not enough

new idea: try to simulate progress
    anything not waiting releases resources (as it finishes)
    anything waiting on only free resources no one else wants takes resources

see if everything gets resources eventually

# deadlock detection (with variable resources)

(pseudocode)

```
class Resources { map<ResourceType, int> amounts; ... };
Resources free_resources;
map<Thread, Resources> requested;
map<Thread, Resources> owned;
```

# deadlock detection (with variable resources)

(pseudocode)

```
class Resources { map<ResourceType, int> amounts; ... };
Resources free_resources;
map<Thread, Resources> requested;
map<Thread, Resources> owned;

...
do { done = true;
  for (Thread t : all threads with owned or requested resources) {
    // if everything requested is free, finish
    if (requested[t] <= free_resources ) {
      requested[t] = no_resources;
      free_resources += owned[t];
      owned[t] = no_resources;
      done = false;
    }
  }
} while (!done);
if (owned.size() > 0) { DeadlockDetected() }
```

# deadlock detection (with variable resources)

(pseudocode)

```
class Resources { map<ResourceType, int> amounts; ... };
Resources free_resources;
map<Thread, Resources> requested;
map<Thread, Resources> owned;
...
do { done =
  for (Thread
    // if everything requested is free, finish
    if (requested[t] <= free_resources ) {
      requested[t] = no_resources;
      free_resources += owned[t];
      owned[t] = no_resources;
      done = false;
    }
  }
} while (!done);
if (owned.size() > 0) { DeadlockDetected() }
```

$\leq$ — free resources include *everything* being requested
(enough memory, disk, each lock requested, etc.)
note: not requesting anything right now? — always true

# deadlock detection (with variable resources)

(pseudocode)
```
class Resources { map<ResourceType, int> amounts; ... };
Resources free_resources;
map<Thread, Resources> requested;
map<Thread, Resources> owned;

...
do { done = true;
  for (Thread t : all       │assume requested resources taken    │resources) {
    // if everything re      │then everything taken released      │
    if (requested[t] <= free_resources ) {
      requested[t] = no_resources;
      free_resources += owned[t];
      owned[t] = no_resources;
      done = false;
    }
  }
} while (!done);
if (owned.size() > 0) { DeadlockDetected() }
```

# deadlock detection (with variable resources)

(pseudocode)

```
class Resources { map<ResourceType, int> amounts; ... };
Resources free_resources;
map<Thread, Resources> requested;
map<Thread, Resources> owned;

...
do { done = true;
  for (Thread t : all threads with owned or requested resources) {
    // if everything requested is free, finish
    if (requested[t] <= free_resources ) {
      requested[t] = no_resources;
```

keep going until nothing changes

```
      owned[t] = no_resources;
      done = false;
    }
  }
} while (!done);
if (owned.size() > 0) { DeadlockDetected() }
```

# finding no deadlock (take free)



waiting on
two units

resource A — 3 units

thread 1

thread 2

owned by

resource B — 1 unit

waiting on
one unit

# finding no deadlock (take free)

# finding no deadlock (take free)



owns 2 units
(was waiting on)

resource A — 3 units

thread 1

algorithm: if all resources
requested by thread free
simulate: it gets ownership of them

owned by

resource B — 1 unit

waiting on
one unit

# finding no deadlock (take free)



owns 2 units
(was waiting on)

resource A — 3 units

thread 1

thread 2

algorithm: if thread not
requesting more resources
simulate: it finishes

owned by

resource B — 1 unit

waiting on
one unit

# finding no deadlock (take free)



owns 2 units
(was waiting on)

resource A — 3 units

thread 1

thread 2

algorithm: if thread not
requesting more resources
simulate: it finishes

owned by

resource B — 1 unit

waiting on
one unit

# finding no deadlock (take free)



owns 2 units
(was waiting on)

resource A — 3 units

algorithm: if all resources
requested by thread free
simulate: it gets ownership of them

thread 2

owned by

resource B — 1 unit

owns
(was waiting)

# finding no deadlock (take free)



owns 2 units
(was waiting on)

resource A — 3 units

thread 1

thread 2

algorithm: if thread not
requesting more resources
simulate: it finishes

owned by

resource B — 1 unit

owns
(was waiting)

# finding no deadlock (take free)



owns 2 units
(was waiting on)

resource A — 3 units

thread 1

thread 2

simulate all threads finishing? no deadlock

owned by

resource B — 1 unit

owns
(was waiting)

# using deadlock detection for prevention

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

make decision <span style="color:red">when starting process</span> ("*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 Baker's algorithm

# recovering from deadlock?

what if it's too late?

kill a thread involved in the deadlock?
   hopefully won't mess things up???

tell owner to release a resource
   need code written to do this???

same concept as locks you can steal

# additional threading topics (if time)

queuing spinlocks: ticket spinlocks?

Linux kernel support for user locks: futexes?

fast synchronization for read-mostly data: read-copy-update?

# threads are hard

get synchronization wrong? weird things happen

...and only sometimes

are there better ways to handle the same problems?
     concurrency — multiple things at once
     parallelism — same thing, use more cores/etc.

# beyond threads: event based programming

writing server that servers multiple clients?
> e.g. multiple web browsers at a time

maybe don't really need multiple processors/cores
> one network, not that fast

idea: one thread handles multiple connections

# beyond threads: event based programming

writing server that servers multiple clients?
   e.g. multiple web browsers at a time

maybe don't really need multiple processors/cores
   one network, not that fast

idea: one thread handles multiple connections

issue: read from/write to multiple streams at once?

# event loops

```
while (true) {
    event = WaitForNextEvent();
    switch (event.type) {
    case NEW_CONNECTION:
        handleNewConnection(event); break;
    case CAN_READ_DATA_WITHOUT_WAITING:
        connection = LookupConnection(event.fd);
        handleRead(connection);
        break;
    case CAN_WRITE_DATA_WITHOUT_WAITING:
        connection = LookupConnection(event.fd);
        handleWrite(connection);
        break;
        ...
    }
}
```

# some single-threaded processing code

```
void ProcessRequest(int fd) {
    while (true) {
        char command[1024] = {};
        size_t comamnd_length = 0;
        do {
            ssize_t read_result =
                read(fd, command + com
                    sizeof(command) -
            if (read_result <= 0) handle_error();
            command_length += read_result;
        } while (command[command_length - 1] != '\n');
        if (IsExitCommand(command)) { return; }
        char response[1024];
        computeResponse(response, commmand);
        size_t total_written = 0;
        while (total_written < sizeof(response)) {
            ...
        }
    }
}
```

```
class Connection {
    int fd;
    char command[1024];
    size_t command_length;
    char response[1024];
    size_t total_written;
    ...
};
```

# some single-threaded processing code

```
void ProcessRequest(int fd) {
    while (true) {
        char command[1024] = {};
        size_t comamnd_length = 0;
        do {
            ssize_t read_result =
                read(fd, command + com
                    sizeof(command) -
            if (read_result <= 0) handle_error();
            command_length += read_result;
        } while (command[command_length - 1] != '\n');
        if (IsExitCommand(command)) { return; }
        char response[1024];
        computeResponse(response, commmand);
        size_t total_written = 0;
        while (total_written < sizeof(response)) {
            ...
        }
    }
}
```

```
class Connection {
    int fd;
    char command[1024];
    size_t command_length;
    char response[1024];
    size_t total_written;
    ...
};
```

# as event code

```
handleRead(Connection *c) {
    ssize_t read_result =
        read(fd, c->command + command_length,
            sizeof(command) - c->command_length);
    if (read_result <= 0) handle_error();
    c->command_length += read_result;

    if (c->command[c->command_length - 1] == '\n') {
        computeResponse(c->response, c->command);
        if (IsExitCommand(command)) {
            FinishConnection(c);
        }
        StopWaitingToRead(c->fd);
        StartWaitingToWrite(c->fd);
    }
}
```

# as event code

```
handleRead(Connection *c) {
    ssize_t read_result =
        read(fd, c->command + command_length,
            sizeof(command) - c->command_length);
    if (read_result <= 0) handle_error();
    c->command_length += read_result;

    if (c->command[c->command_length - 1] == '\n') {
        computeResponse(c->response, c->command);
        if (IsExitCommand(command)) {
            FinishConnection(c);
        }
        StopWaitingToRead(c->fd);
        StartWaitingToWrite(c->fd);
    }
}
```

# POSIX support for event loops

select and poll functions
    take list(s) of file descriptors to read and to write
    wait for them to be read/writeable without waiting
    (or for new connections associated with them, etc.)

many OS-specific extensions/improvements/alternatives:
    examples: Linux epoll, Windows IO completion ports
    better ways of managing list of file descriptors
    do read/write when ready instead of just returning when reading/writing
    is okay

# message passing

instead of having variables, locks between threads...

send messages between threads/processes

what you need anyways between machines
    big 'supercomputers' = really many machines together

arguably an easier model to program
    can't have locking issues

## message passing API

core functions: Send(toId, data)/Recv(fromId, data)

simplest version: functions wait for other processes/threads
extensions: send/recv at same time, multiple messages at once, don't
wait, etc.

```
if (thread_id == 0) {
    for (int i = 1; i < MAX_THREAD; ++i) {
        Send(i, getWorkForThread(i));
    }
    for (int i = 1; i < MAX_THREAD; ++i) {
        WorkResult result;
        Recv(i, &result);
        handleResultForThread(i, result);
    }
} else {
    WorkInfo work;
    Recv(0, &work);
    Send(0, ComputeResultFor(work));
```

# message passing game of life



divide grid
like you would for normal threads

each process stores cells
in that part of grid

(no shared memory!)

# message passing game of life



process 3 only needs values
of cells around its area
(values of cells adjacent to
the ones it computes)

# message passing game of life



small slivers of
other process's cells needed

solution: process 2, 4
send messages with cells every iterat

# message passing game of life



some of process 3's cells
also needed by process 2/4

so process 3 also sends messages

# message passing game of life



one possible pseudocode:
all even processes send messages
(while odd receives), then
all odd processes send messages
(while even receives)

# message passing game of life



one possible pseudocode:
all even processes send messages
(while odd receives), then
all odd processes send messages
(while even receives)

# backup slides

# fairer spinlocks

so far — everything on spinlocks
>    mutexes, condition variables — built with spinlocks

spinlocks are pretty 'unfair'
>    where fair = get lock if waiting longest

last CPU that held spinlock more likely to get it again
>    already has the lock in its cache…

but there are many other ways to spinlocks…

## ticket spinlocks

```
unsigned int serving_number;
unsigned int next_number;

Lock() {
    // "take a number"
    unsigned int my_number = atomic_read_and_increme
    // wait until "now serving" that number
    while (atomic_read(&serving_number) != my_number
        /* do nothing */
    }
    // MISSING: code to prevent reordering reads/wri
}

Unlock() {
```

# ticket spinlocks and cache contention

still have contention to write next_number

...but no retrying writes!
> should limit 'ping-ponging'?

threads loop performing a read repeatedly while waiting
> value will be broadcasted to all processors
> 'free' if using a bus
> not-so-free if another way of connecting CPUs

# beyond ticket spinlocks

Linux kernel used to use ticket spinlocks

now uses variant of MCS spinlocks — locks have linked-list queue!
    careful use of atomic operations to modify queue

still try

goal: even less contention
    unlocking value doesn't require broadcasting to all CPUs
    each processor waits on its own cache block

# Linux futexes

futex — **f**ast **u**serspace mu**tex**

goal: implement waiting like 'proper' mutexes, but…

don't enter kernel mode most of the time

challenge: can't acquire lock to call scheduler from user mode

# futex operations

```
futex(&lock_value, FUTEX_WAIT, expected_value, ...);
```

check if lock_value is expected_value
    if not — return immediately
    otherwise, sleep until it futex(…, FUTEX_WAKE is called

```
futex(&lock_value, FUTEX_WAKE, num_processes);
```

wakeup up to num_processes which called FUTEX_WAIT

# mutexes with futexes

```
int lock_value; // UNLOCKED or LOCKED_NO_WAITERS or LOCKED_WAITERS
Lock() {
retry:
    if (CompareAndSwap(&lock_value, UNLOCKED, LOCKED_NO_WAITERS) ==
        /* acquired lock */
        return;
    } else if (CompareAndSwap(&lock_value, LOCKED_NO_WAITERS, LOCKED
        futex(&lock_value, FUTEX_WAIT, LOCKED_WAITERS, ...);
    }
    goto retry;
}
Unlock() {
    if (CompareAndSwap(&lock_value, LOCKED_NO_WAITERS, UNLOCKED) ==
        return;
    } else {
        lock_value = UNLOCKED;
        futex(&lock_value, FUTEX_WAKE, 1, ...);
    }
}
```

# implementing futex_wait

hashtable: address $\rightarrow$ queue of waiting threads

use hashtable to look-up queue

lock queue

check value hasn't changed
    if so abort, releasing lock

add thread to queue

set thread as WAITING (not runnable)

unlock queue

call scheduler

# read-copy-update (high-level overview)

idea: read-mostly data structure

when reading:
    read normally <span style="color:red">via shared pointer</span>

when writing:
    make a copy
    atomically update the <span style="color:red">shared pointer</span>
    delete the old version <span style="color:red">eventually</span>

tricky part: when is it safe to delete old version
    implementation: <span style="color:red">scheduler integration</span>

# RCU operations

read lock — record: "I am reading now"

read unlock — record: "I am done reading now"

publish — atomically update pointer

after publish: wait until
    all threads currently running have context switched
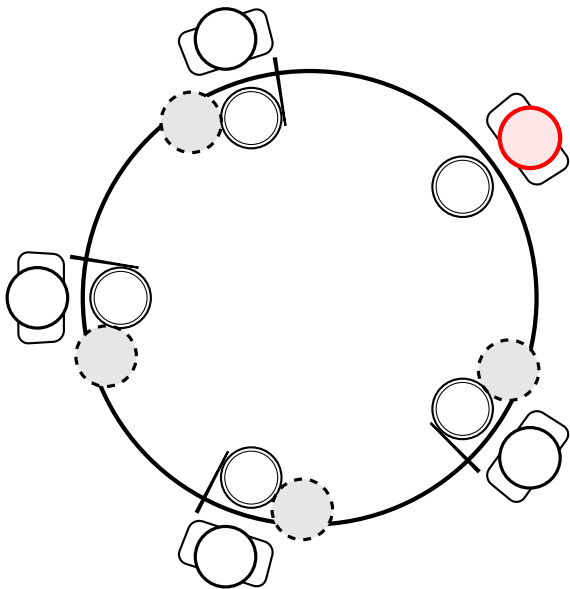    …and none of them set the "I am reading now" bit

# dining philosophers — ordering



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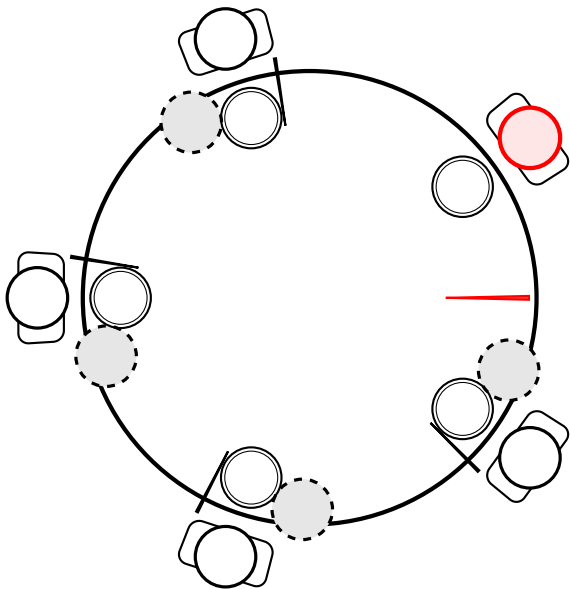rule: grab from marked place first
only grab other chopstick after that

# 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

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



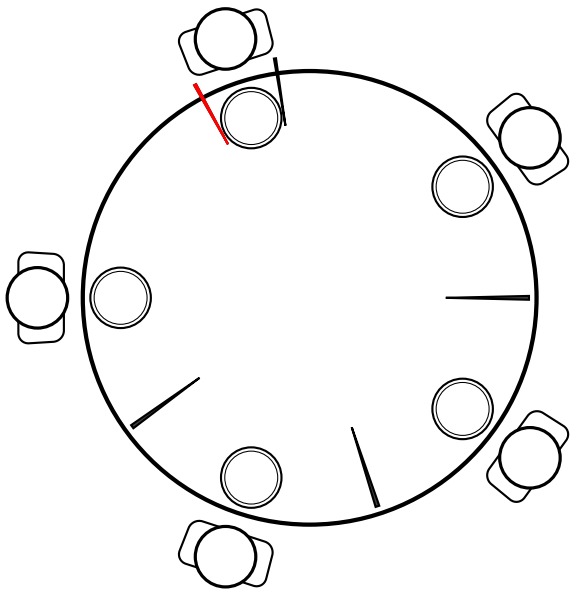now everyone else can eat

# dining philosophers — aborting



now everyone else can eat

# dining philosophers — aborting



now everyone else can eat

# dining philosophers — aborting



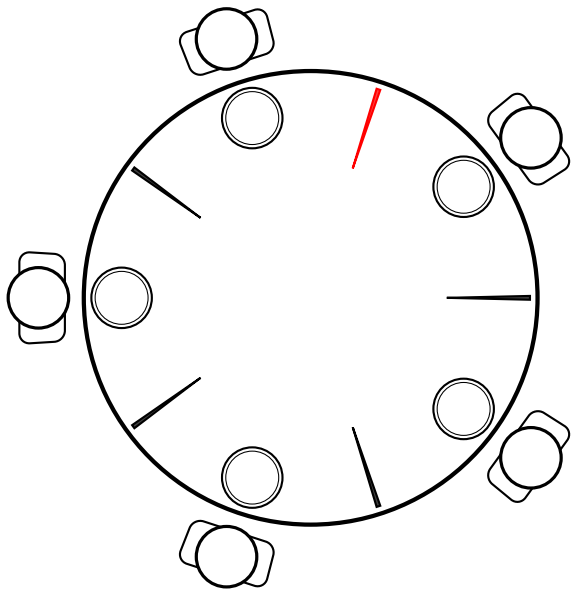now everyone else can eat

# dining philosophers — aborting



now everyone else can eat

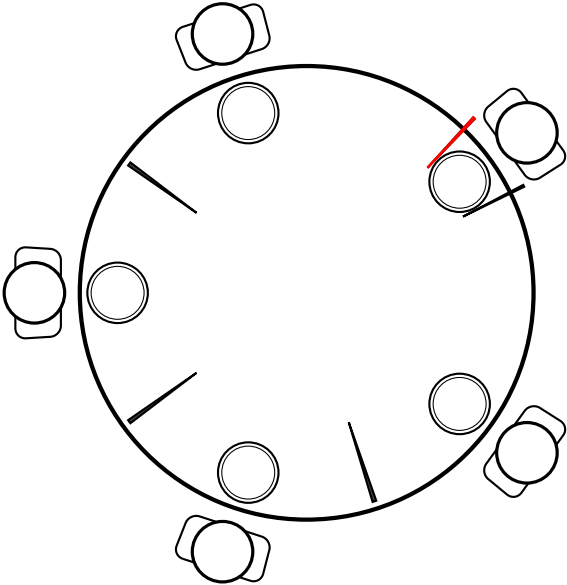# dining philosophers — aborting



now everyone else can eat

# dining philosophers — aborting



now everyone else can eat

# dining philosophers — aborting



and person who gave up
might succeed later