

Locks

Operating Systems

Based on: Three Easy Pieces by Arpaci-Dusseau

Moshe Sulamy

Tel-Aviv Academic College

- **Concurrency** issues
 - Execute a series of instructions **atomically**
 - With interrupts and concurrent processors
- Introducing: a **lock**
 - Critical section seemingly executes atomically

- **Lock variable**
 - Holds lock state
 - **Available** (or unlocked or **free**)
 - No thread holds the lock
 - **Acquired** (or **locked** or **held**)
 - Exactly one thread (**owner**) holds the lock
 - In a critical section

Basic Idea

- `lock()`
 - Try to acquire the lock
 - Will not return (or fail) if held by another thread
- `unlock()`
 - Lock is available again

Basic Idea

- Critical section:

```
balance = balance + 1;
```

- To use lock:

```
1 lock_t mutex; // lock variable
2 ...
3 lock(&mutex);
4 balance = balance + 1;
5 unlock(&mutex);
```

Pthread Locks

- POSIX library: **mutex (mutual exclusion)**
- Equivalent code:

```
1 pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;  
2 ...  
3 pthread_mutex_lock(&lock); // may fail!  
4 balance = balance + 1;  
5 pthread_mutex_unlock(&lock);
```

Pthread Locks

- POSIX library: **mutex** (**mutual exclusion**)
- Equivalent code:

```
1 pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;  
2 ...  
3 pthread_mutex_lock(&lock); // may fail!  
4 balance = balance + 1;  
5 pthread_mutex_unlock(&lock);
```

- Variable passed to lock and unlock
 - May use different locks for different sections
 - **Coarse-grained** locking: one big lock
 - **Fine-grained**: use various locks for different sections

Building A Lock

- Efficient locks provide mutual exclusion at low cost (overhead)
 - Support from hardware and the OS

How can we build an efficient lock?

Evaluating Locks

- **Mutual exclusion**

- At most one thread in the CS

- **Deadlock-freedom**

- Some thread eventually enters CS

- **Fairness (starvation-freedom)**

- Each thread eventually enters CS

- **Performance**

- Time overhead for using the lock
 - Single thread: overhead for grab & release
 - Multiple threads and CPUs

Controlling Interrupts

- Early solution: disable interrupts
 - For single-processor systems
 - No clock interrupt / context switch in critical section
- The negatives:

Controlling Interrupts

- Early solution: disable interrupts
 - For single-processor systems
 - No clock interrupt / context switch in critical section
- The negatives:
 - Trust arbitrary (greedy, malicious, or faulty) programs
 - Does not work on multiprocessors
 - Lost interrupts
- Used by OS

Just Using Loads/Stores

```
1 typedef struct __lock_t { int flag; } lock_t;
2
3 void init(lock_t* mutex) {
4     // 0: available, 1: locked
5     mutex->flag = 0;
6 }
7 void lock(lock_t* mutex) {
8     while (mutex->flag == 1)
9         ; // spin-wait
10    mutex->flag = 1;
11 }
12 void unlock(lock_t* mutex) {
13     mutex->flag = 0;
14 }
```

Just Using Loads/Stores

- No mutual exclusion



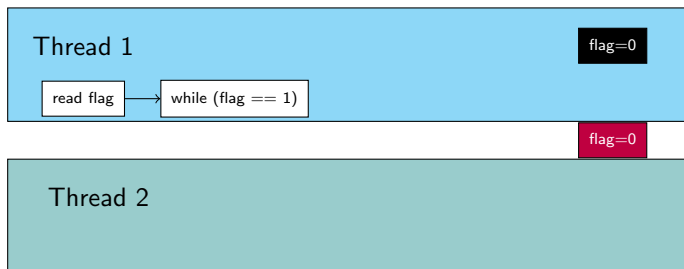
Just Using Loads/Stores

- No mutual exclusion



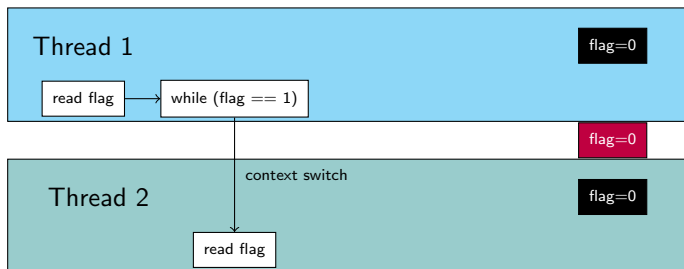
Just Using Loads/Stores

- No mutual exclusion



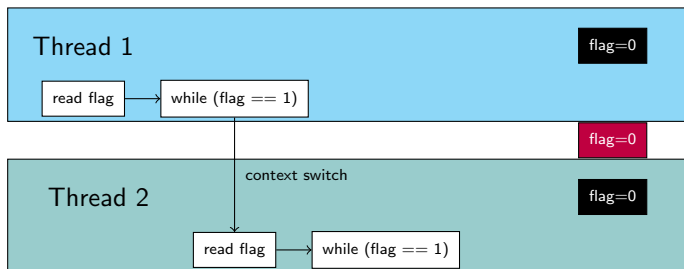
Just Using Loads/Stores

- No mutual exclusion



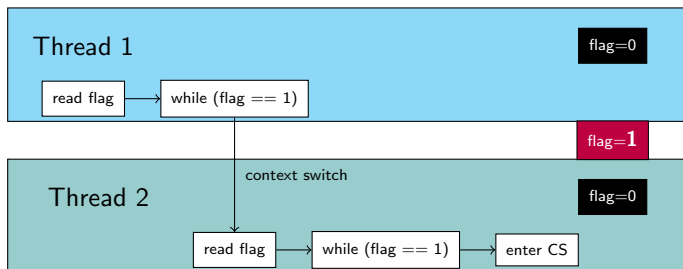
Just Using Loads/Stores

- No mutual exclusion



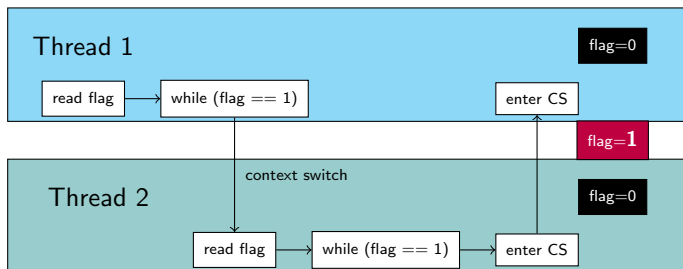
Just Using Loads/Stores

- No mutual exclusion



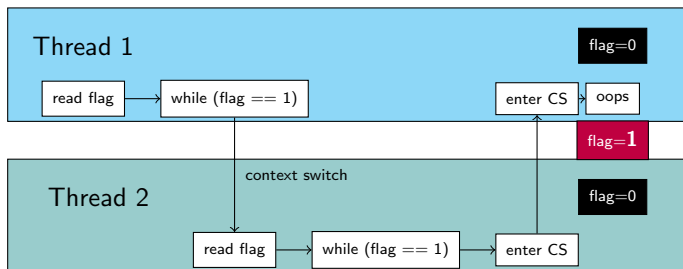
Just Using Loads/Stores

- No mutual exclusion



Just Using Loads/Stores

- No mutual exclusion



Test-And-Set

- Hardware support: a new instruction **test-and-set**
 - Update value and return previous, **atomically**
- Defined as:

```
1 int TestAndSet(int* old_ptr, int new) {  
2     int old = *old_ptr;  
3     *old_ptr = new;  
4     return old;  
5 }
```

New Spin Lock

```
1 typedef struct __lock_t { int flag; } lock_t;
2
3 void init(lock_t* mutex) {
4     // 0: available, 1: locked
5     mutex->flag = 0;
6 }
7 void lock(lock_t* mutex) {
8     while (TestAndSet(&mutex->flag, 1))
9         ; // spin-wait
10    mutex->flag = 1;
11 }
12 void unlock(lock_t* mutex) {
13     mutex->flag = 0;
14 }
```

Evaluating Spin Locks

- Correctness (**mutual exclusion**)?

Evaluating Spin Locks

- Correctness (**mutual exclusion**)? yes
- **Deadlock-freedom?**

Evaluating Spin Locks

- Correctness (**mutual exclusion**)? yes
- **Deadlock-freedom**? yes
- **Fairness**?

Evaluating Spin Locks

- Correctness (**mutual exclusion**)? yes
- **Deadlock-freedom**? yes
- **Fairness**? no
- Performance?
 - Single CPU:

Evaluating Spin Locks

- Correctness (**mutual exclusion**)? yes
- **Deadlock-freedom**? yes
- **Fairness**? no
- Performance?
 - Single CPU:painful
 - Owner thread is preempted, all $N - 1$ others spin-wait needlessly
 - Multiple CPUs:

Evaluating Spin Locks

- Correctness (**mutual exclusion**)? yes
- **Deadlock-freedom**? yes
- **Fairness**? no
- Performance?
 - Single CPU: painful
 - Owner thread is preempted, all $N - 1$ others spin-wait needlessly
 - Multiple CPUs: reasonably well

Compare-And-Swap

- Another hardware primitive: **compare-and-swap**
- Compare to `expected`, update only if equal, return previous
- Defined as:

```
1 int CompareAndSwap(int* ptr, int expected, int new) {  
2     int original = *ptr;  
3     if (original == expected)  
4         *ptr = new;  
5     return original;  
6 }
```

Compare-And-Swap

- Spin-lock with CAS:

```
1 void lock(lock_t* lock) {  
2     while (CompareAndSwap(&mutex->flag, 0, 1) == 1)  
3         ; // spin  
4 }
```

- Fairness? performance?

Compare-And-Swap

- Spin-lock with CAS:

```
1 void lock(lock_t* lock) {  
2     while (CompareAndSwap(&mutex->flag, 0, 1) == 1)  
3         ; // spin  
4 }
```

- Fairness? performance?
 - Pretty much the same

Fetch-And-Add

- Final hardware primitive: **fetch-and-add**
- Atomically increment a value and return old value
- Defined as:

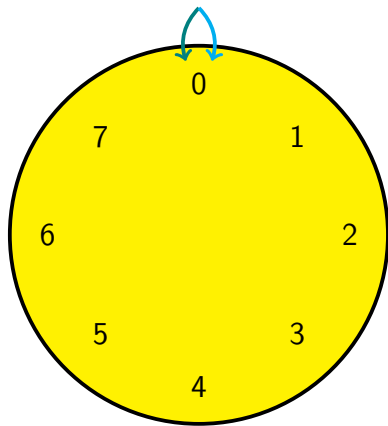
```
1 int FetchAndAdd(int* ptr) {  
2     int old = *ptr;  
3     *ptr = old + 1;  
4     return old;  
5 }
```


Fetch-And-Add

- We can now build a fair **ticket lock**:

```
1  typedef struct __lock_t {
2      int ticket;
3      int turn;
4  } lock_t;
5
6  void init(lock_t* lock) {
7      lock->ticket = 0;
8      lock->turn = 0;
9  }
10 void lock(lock_t* lock) {
11     int myturn = FetchAndAdd(&lock->ticket);
12     while (lock->turn != myturn)
13         ; // spin
14 }
15 void unlock(lock_t* lock) {
16     lock->turn = lock->turn + 1;
17 }
```

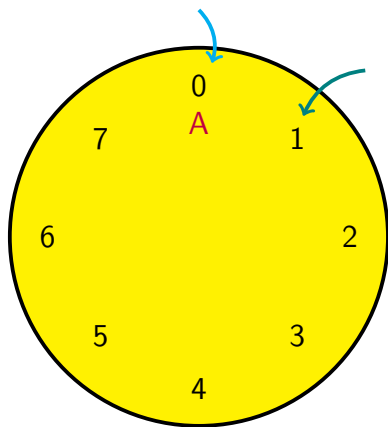
Ticket Lock



Ticket
Turn

Ticket Lock

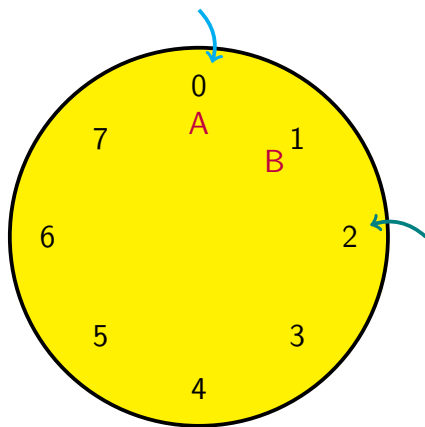
- A: lock(), gets ticket 0 & runs



Ticket
Turn

Ticket Lock

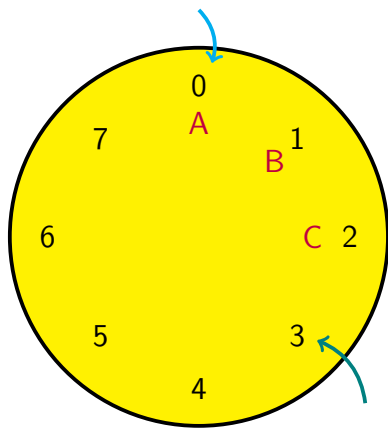
- A: lock(), gets ticket 0 & runs
- B: lock(), gets ticket 1, spins



Ticket
Turn

Ticket Lock

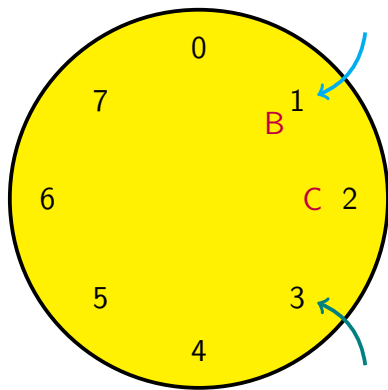
- A: lock(), gets ticket 0 & runs
- B: lock(), gets ticket 1, spins
- C: lock(), gets ticket 2, spins



Ticket
Turn

Ticket Lock

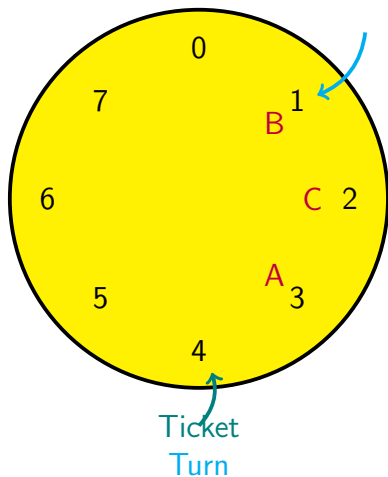
- A: lock(), gets ticket 0 & runs
- B: lock(), gets ticket 1, spins
- C: lock(), gets ticket 2, spins
- A: unlock(), turn++, B runs



Ticket
Turn

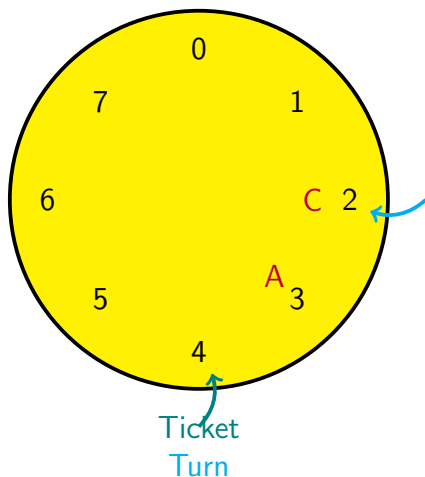
Ticket Lock

- A: lock(), gets ticket 0 & runs
- B: lock(), gets ticket 1, spins
- C: lock(), gets ticket 2, spins
- A: unlock(), turn++, B runs
- A: lock(), gets ticket 3, spins



Ticket Lock

- A: lock(), gets ticket 0 & runs
- B: lock(), gets ticket 1, spins
- C: lock(), gets ticket 2, spins
- A: unlock(), turn++, B runs
- A: lock(), gets ticket 3, spins
- B: unlock(), turn++, C runs
- ...



Decker's Algorithm

- What about a lock without hardware support?

```
1  int flag[2];      // wants to grab lock?
2  int turn;         // whose turn?
3
4  void init() {
5      flag[0] = flag[1] = 0;
6      turn = 0;
7  }
8  void lock(int self) {
9      flag[self] = 1;
10     turn = 1 - self;    // let other run
11     while ((flag[1-self] == 1) && (turn == 1-self))
12         ; // spin-wait
13 }
14 void unlock(int self) {
15     flag[self] = 0;
16 }
```

- Various issues → concurrency course

Too Much Spinning

- Locks so far used spinning
 - Quite inefficient
- Consider N threads
 - Thread 1 grabs lock, $N - 1$ threads waiting for lock
 - Timer interrupt \rightarrow context switch
 - $N - 1$ threads execute, waste $N - 1$ time slices
- Solution? hardware again!

Too Much Spinning

- When you are going to spin, give up CPU
 - `yield()`: system call to change caller state
 - From **running** to **ready**
 - Essentially **deschedules** itself

```
1 void lock(lock_t* mutex) {  
2     while (TestAndSet(&mutex->flag, 1))  
3         yield(); // give up the CPU  
4     mutex->flag = 1;  
5 }
```

Too Much Spinning

- When you are going to spin, give up CPU
 - `yield()`: system call to change caller state
 - From **running** to **ready**
 - Essentially **deschedules** itself

```
1 void lock(lock_t* mutex) {  
2     while (TestAndSet(&mutex->flag, 1))  
3         yield(); // give up the CPU  
4     mutex->flag = 1;  
5 }
```

- Still costly
 - $N - 1$ system calls and context switches
 - Does not handle **starvation**

Using Queues

- Use **queue** to keep track of threads waiting for lock
- In **Solaris**: `park()`
 - Put calling thread to sleep
 - Until another thread calls `unpark(threadID)`

Using Queues

- Use **queue** to keep track of threads waiting for lock
- In **Solaris**: `park()`
 - Put calling thread to sleep
 - Until another thread calls `unpark(threadID)`

```
1 typedef struct __lock_t {
2     int flag;
3     int guard;
4     queue_t *q;
5 } lock_t;
6
7 void init(lock_t *m) {
8     m->flag = 0;
9     m->guard = 0;
10    queue_init(m->q);
11 }
```

Using Queues

```
1 void lock(lock_t *m) {
2     while (TestAndSet(&m->guard, 1) == 1)
3         ; // acquire guard lock by spinning
4     if (m->flag == 0) {
5         m->flag = 1;
6         m->guard = 0;
7     } else {
8         queue_add(m->q, getpid());
9         m->guard = 0;
10        park();
11    }
12 }
13 void unlock(lock_t *m) {
14     while (TestAndSet(&m->guard, 1) == 1)
15         ; // acquire guard lock by spinning
16     if (queue_empty(m->q))
17         m->flag = 0;
18     else
19         unpark(queue_remove(m->q));
20     m->guard = 0;
21 }
```

Using Queues

- Race condition:
 - 1 Thread adds itself to queue, releases lock
 - 2 Context switch just before call to `park()`
 - 3 Owner thread is done, removes new item from queue
 - 4 Calls `unpark()` for new thread: not yet parked!

Using Queues

- Also in Solaris: `setpark()`
 - Thread about to call `park()`
 - If interrupted and `unpark()` is called for it:
 - Subsequent `park()` returns immediately

```
1 void lock(lock_t *m) {  
2     ...  
3     } else {  
4         setpark(); // <- new code  
5         queue_add(m->q, gettid());  
6         m->guard = 0;  
7         park();  
8     }  
9 }
```

Using Queues: Different OS

- Support details vary between OS
- Linux: **futex**
 - Similar to Solaris
 - `futex_wait(address, expected)`
 - Puts calling thread to sleep if address is equal to expected
 - `futex_wake(address)`
 - Wakes one thread waiting on address

Using Queues: Different OS

- Snippet from POSIX thread library:

```
1 void mutex_lock(int *mutex) {
2     int v;
3     // Bit 31 was clear, we got the mutex (fastpath)
4     if (atomic_bit_test_set(mutex, 31) == 0)
5         return;
6     atomic_increment(mutex);
7     while (1) {
8         if (atomic_bit_test_set(mutex, 31) == 0) {
9             atomic_decrement(mutex);
10            return;
11        }
12        v = *mutex;
13        if (v >= 0)
14            continue;
15        futex_wait(mutex, v);
16    }
17 }
18 void mutex_unlock(int *mutex) {
19     if (atomic_add_zero(mutex, 0x80000000))
20         return; // zero iff no other interested threads
21
22     // there are other threads waiting
23     futex_wake(mutex);
24 }
```

Two-Phase Locks

- Hybrid approach: **two-phase lock**
 - Spinning can be useful
 - Particularly if lock is about to be released
- **First phase:** lock spins for a while
- **Second phase:** caller put to sleep, wakes up when lock becomes free

Summary

- **Lock**

- Execute a series of actions **atomically**
 - Evaluated by: **Mutual exclusion, Deadlock-freedom, fairness, performance**
 - POSIX library: **mutex, futex**
- Disabling interrupts: problematic, used by OS
- Hardware support: **test&set, compare&swap, fetch&add**
- Spin-locks: TAS lock & CAS lock
 - Avoid spinning with `yield()`
- Fairness: **ticket lock** or queue lock
- **Condition variables**