# Operating Systems

**Lock (Chapter 28 ~ 29**)

Dr. Young-Woo Kwon

# Definitions

- **Race condition**: output of a concurrent program depends on the order of operations between threads
- **Mutual exclusion**: only one thread does a particular thing at a time
- **Critical section**: piece of code that only one thread can execute at once
- **Lock**: prevent someone from doing something

# Locks: The Basic Idea

- Ensure that any **critical section** executes as if it were a **single atomic instruction.**
  - An example: the canonical update of a shared variable

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

  - Add some code around the critical section

    ```
    1    lock_t mutex; // some globally-allocated lock 'mutex'
    2    …
    3    lock(&mutex);
    4    balance = balance + 1;
    5    unlock(&mutex);
    ```

# Locks: The Basic Idea

- Lock variable holds <u>the state of</u> the lock.
  - **available** (or **unlocked** or **free**)
    - No thread holds the lock.

  - **acquired** (or **locked** or **held**)
    - Exactly one thread holds the lock and presumably is in a critical section.

# The semantics of the lock()

- `lock()`
  - **Try to** acquire the lock.
  - If <u>no other thread holds</u> the lock, the thread will **acquire** the lock.
  - **Enter** the *critical section*.
    - This thread is said to be <u>**the owner of**</u> the lock.

  - Other threads are *prevented from* entering the critical section while the first thread that holds the lock is in there.

# Pthread Locks - mutex

- The name that the POSIX library uses for a <u>lock</u>.
  - Used to provide <span style="color:red">mutual exclusion</span> between threads.

```
1    pthread_mutex_t lock = PTHREAD_MUTEX_INITIALIZER;
2
3    Pthread_mutex_lock(&lock); // wrapper for pthread_mutex_lock()
4    balance = balance + 1;
5    Pthread_mutex_unlock(&lock);
```

  - We may be using *different locks* to protect *different variables* → Increase **concurrency** (a more **fine-grained** approach).

# Building A Lock

- <u>Efficient locks</u> provided mutual exclusion at **<span style="color:red">low cost</span>**.
- Building a lock need some help from the **hardware** and the **OS**.

# Evaluating locks – Basic criteria

- **Mutual exclusion**
  - Does the lock work, preventing multiple threads from entering *a critical section*?
- **Fairness**
  - Does each thread contending for the lock get a fair shot at acquiring it once it is free? (Starvation)
- **Performance**
  - The time overheads added by using the lock

# Why hardware support needed?

- **First attempt**: Using a *flag* denoting whether the lock is held or not.
  - The code below has problems.

```
1    typedef struct __lock_t { int flag; } lock_t;
2
3    void init(lock_t *mutex) {
4        // 0 → lock is available, 1 → held
5        mutex->flag = 0;
6    }
7
8    void lock(lock_t *mutex) {
9        while (mutex->flag == 1)  // TEST the flag
10               ;   // spin-wait (do nothing)
11        mutex->flag = 1;   // now SET it !
12    }
13
14   void unlock(lock_t *mutex) {
15       mutex->flag = 0;
16   }`
```

# Why hardware support needed? (Cont.)

– **Problem 1**: <span style="color:darkred">**No Mutual Exclusion**</span> (assume `flag=0` to begin)

| Thread1 | Thread2 |
|---|---|
| call `lock()` | |
| while (flag == 1) | |
| interrupt: switch to Thread 2 | |
| | call `lock()` |
| | while (flag == 1) |
| | flag = 1; |
| | interrupt: switch to Thread 1 |
| flag = 1;  // set flag to 1 (too!) | |

– **Problem 2**: <span style="color:darkred">**Spin-waiting**</span> wastes time waiting for another thread.

• So, we need an **atomic instruction** supported by <span style="color:darkred">**Hardware**</span>!

– *test-and-set* instruction, also known as *atomic exchange*

# Test And Set (Atomic Exchange)

- An instruction to support the creation of simple locks

```
1    int TestAndSet(int *ptr, int new) {
2        int old = *ptr;    // fetch old value at ptr
3        *ptr = new;        // store 'new' into ptr
4        return old;        // return the old value
5    }
```

- **return**(testing) old value pointed to by the `ptr`.
- *Simultaneously* **update**(setting) said value to `new`.
- This sequence of operations is **performed atomically.**

# A Simple Spin Lock using test-and-set

```
1    typedef struct __lock_t {
2        int flag;
3    } lock_t;
4
5    void init(lock_t *lock) {
6        // 0 indicates that lock is available,
7        // 1 that it is held
8        lock->flag = 0;
9    }
10
11   void lock(lock_t *lock) {
12       while (TestAndSet(&lock->flag, 1) == 1)
13           ;            // spin-wait
14   }
15
16   void unlock(lock_t *lock) {
17       lock->flag = 0;
18   }
```

– **Note**: To work correctly on a single CPU, it requires <u>a preemptive scheduler</u>.
  • **Why?**

# Evaluating Spin Locks

- **Correctness**: <span style="color:red">yes</span>
  - The spin lock only allows a single thread to entry the critical section.

- **Fairness**: <span style="color:red">no</span>
  - Spin locks <u>don't provide any fairness</u> guarantees.
  - Indeed, a thread spinning may spin *forever*.

- **Performance**:
  - In the single CPU, performance overheads can be quire *painful*.
  - If **the number of threads roughly equals the number of CPUs**, spin locks work *reasonably well*.

# Compare-And-Swap (SPARC)

- Test whether the value at the address(`ptr`) is equal to `expected`.
  - *If so*, update the memory location pointed to by `ptr` with the `new` value.
  - *In either case*, return the actual value at that memory location.

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

**Compare-and-Swap hardware atomic instruction (C-style)**

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

**Spin lock with compare-and-swap**

# Compare-And-Exchange (x86)

- C-callable x86-version of compare-and-swap

```
1    char CompareAndSwap(int *ptr, int old, int new) {
2        unsigned char ret;
3
4        // Note that sete sets a 'byte' not the word
5        __asm__ __volatile__ (
6                " lock\n"
7                " cmpxchgl %2,%1\n"
8                " sete %0\n"
9                : "=q" (ret), "=m" (*ptr)
10               : "r" (new), "m" (*ptr), "a" (old)
11               : "memory");
12       return ret;
13   }
```

# Fetch-And-Add

- Atomically increment a value while returning the old value at a particular address.

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

**Fetch-And-Add Hardware atomic instruction (C-style)**

# Ticket Lock

- **Ticket lock** can be built with <u>fetch-and add</u>.
  - Ensure progress for all threads. → fairness

```c
1    typedef struct __lock_t {
2        int ticket;
3        int turn;
4    } lock_t;
5
6    void lock_init(lock_t *lock) {
7        lock->ticket = 0;
8        lock->turn = 0;
9    }
10
11   void lock(lock_t *lock) {
12       int myturn = FetchAndAdd(&lock->ticket);
13       while (lock->turn != myturn)
14               ; // spin
15   }
16   void unlock(lock_t *lock) {
17       FetchAndAdd(&lock->turn);
18   }
```

# So Much Spinning

- Hardware-based spin locks are simple and they work.

- In some cases, these solutions can be quite inefficient.
  - Any time a thread gets caught *spinning*, it **wastes an entire time slice** doing nothing but checking a value.

> **How To Avoid *Spinning*?**
> **We'll need OS Support too!**

# A Simple Approach: Just Yield

- When you are going to spin, **give up the CPU** to another thread.
  - OS system call moves the caller from the *running state* **to the** *ready state*.
  - The cost of a **context switch** can be substantial and the **starvation** problem still exists.

```
1    void init() {
2        flag = 0;
3    }
4
5    void lock() {
6        while (TestAndSet(&flag, 1) == 1)
7            yield(); // give up the CPU
8    }
9
10   void unlock() {
11       flag = 0;
12   }
```

**Lock with Test-and-set and Yield**

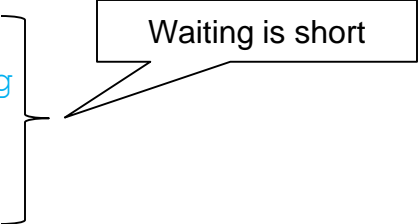# Using Queues: Sleeping Instead of Spinning

- **Queue** to keep track of which threads are <u>waiting</u> to enter the lock.
- `park()`
  - Put a calling thread to sleep
- `unpark(threadID)`
  - Wake a particular thread as designated by `threadID`.

# Using Queues: Sleeping Instead of Spinning

```
typedef struct __lock_t {
  int flag;       // lock is acquired or not
  int guard;      // to protect the queue
  queue_t *q;
} lock_t;
```

# Using Queues: Sleeping Instead of Spinning

```
1    typedef struct __lock_t { int flag; int guard; queue_t *q; } lock_t;
2
3    void lock_init(lock_t *m) {
4        m->flag = 0;
5        m->guard = 0;
6        queue_init(m->q);
7    }
8
9    void lock(lock_t *m) {
10       while (TestAndSet(&m->guard, 1) == 1)
11           ; // acquire guard lock by spinning
12       if (m->flag == 0) {
13           m->flag = 1; // lock is acquired
14           m->guard = 0;
15       } else {
16           queue_add(m->q, gettid());
17           m->guard = 0;
18           park();
19       }
20   }
21   …
```

Waiting is short

**Lock With Queues, Test-and-set, Yield, And Wakeup**

# Using Queues: Sleeping Instead of Spinning

```
22   void unlock(lock_t *m) {
23       while (TestAndSet(&m->guard, 1) == 1)
24           ; // acquire guard lock by spinning
25       if (queue_empty(m->q))
26           m->flag = 0; // let go of lock; no one wants it
27       else
28           unpark(queue_remove(m->q)); // hold lock (for next thread!)
29       m->guard = 0;
30   }
```

**Lock With Queues, Test-and-set, Yield, And Wakeup (Cont.)**

# Wakeup/waiting race

```
1   typedef struct __lock_t { int flag; int guard; queue_t *q; } lock_t;
2
3   void lock_init(lock_t *m) {
4       m->flag = 0;
5       m->guard = 0;
6       queue_init(m->q);
7   }
8
9   void lock(lock_t *m) {
10      while (TestAndSet(&m->guard, 1) == 1)
11          ; // acquire guard lock by spinning
12      if (m->flag == 0) {
13          m->flag = 1; // lock is acquired
14          m->guard = 0;
15      } else {
16          queue_add(m->q, gettid());
17          m->guard = 0;
18          park();
19      }
20  }
21  …
```

*What if a lock holder releases the lock just before the thread B executes "park()"?*

- In case of releasing the lock (*thread A*) just before the call to `park()` (*thread B*) → Thread B would **sleep forever** (potentially).
- **Solaris** solves this problem by adding a third system call: `setpark()`.
  - By calling this routine, a thread can indicate it *is about to* `park`.
  - If it happens to be interrupted and another thread calls `unpark` before `park` is actually called, the subsequent `park` returns immediately instead of sleeping.

```
1           queue_add(m->q, gettid());
2           setpark(); // new code
3           m->guard = 0;
4           park();
```

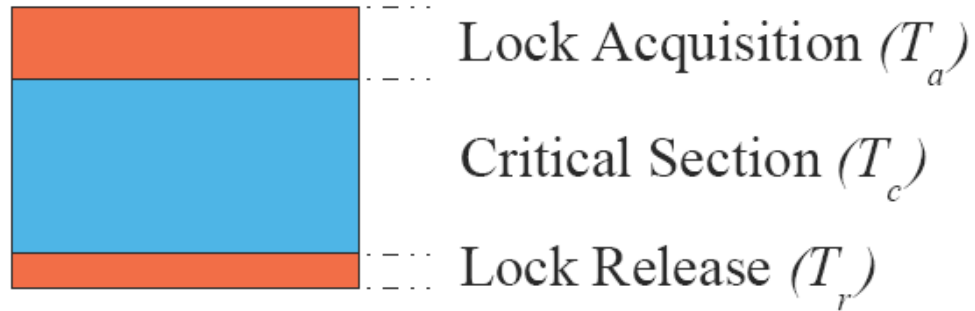**Code modification inside of `lock()`**

# Spinning vs Blocking

- Some lock implementations combine spinning and blocking locks
- Blocking has a cost
  - Shouldn't block if lock becomes available in less time than it takes to block
- Strategy: spin for time it would take to block
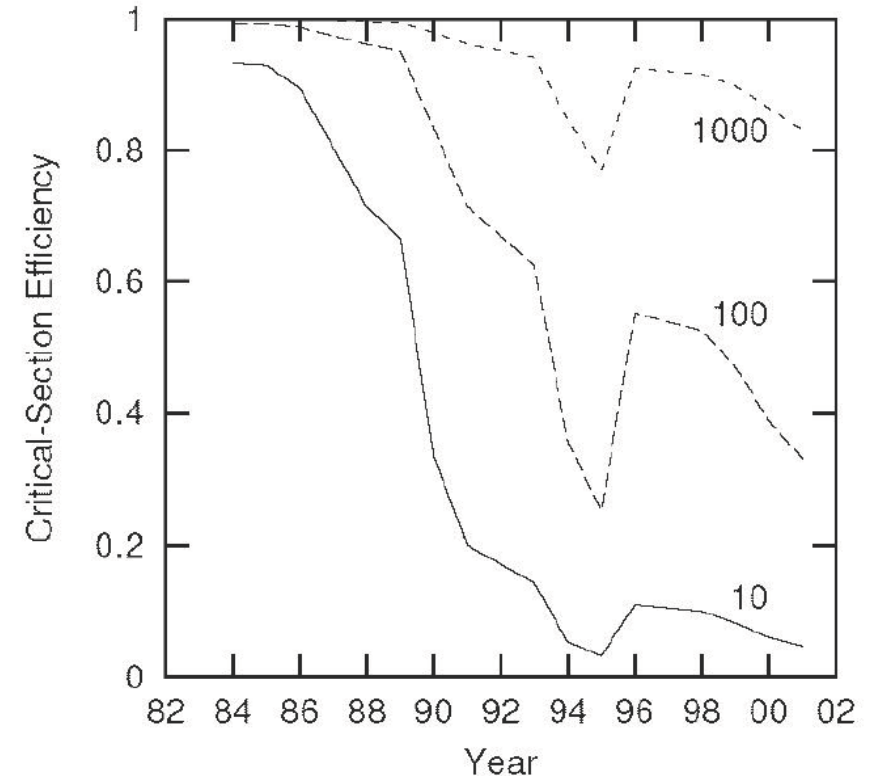  - Even in worst case, total cost for lock() is less than 2*block time

# Two-Phase Locks

- A two-phase lock realizes that spinning can be useful if the lock *is about to* be released.

  - **First phase**

    - The lock spins for a while, *hoping that* it can acquire the lock.

    - If the lock is not acquired during the first spin phase, a second phase is entered,

  - **Second phase**

    - The caller is put to sleep.

    - The caller is only woken up when the lock becomes free later.

# Critical Section Efficiency



Lock Acquisition $(T_a)$

Critical Section $(T_c)$

Lock Release $(T_r)$

$$Efficiency = \frac{T_c}{T_c + T_a + T_r}$$

- As processors get faster, CSE decreases because atomic instructions become relatively more expensive



*Source: McKenney, 2005*

# 29. Lock-based Concurrent Data Structures

# Lock-based Concurrent Data structure

- Adding locks to a data structure makes the structure **thread safe**.
  - How locks are added determine both the correctness and performance of the data structure.

# Example: Concurrent Counters without Locks

- Simple but not scalable

```c
1         typedef struct __counter_t {
2                 int value;
3         } counter_t;
4
5         void init(counter_t *c) {
6                 c->value = 0;
7         }
8
9         void increment(counter_t *c) {
10                c->value++;
11        }
12
13        void decrement(counter_t *c) {
14                c->value--;
15        }
16
17        int get(counter_t *c) {
18                return c->value;
19        }
```

# Example: Concurrent Counters with Locks

- Add a **single lock**.
  - The lock is acquired when calling a routine that manipulates the data structure.
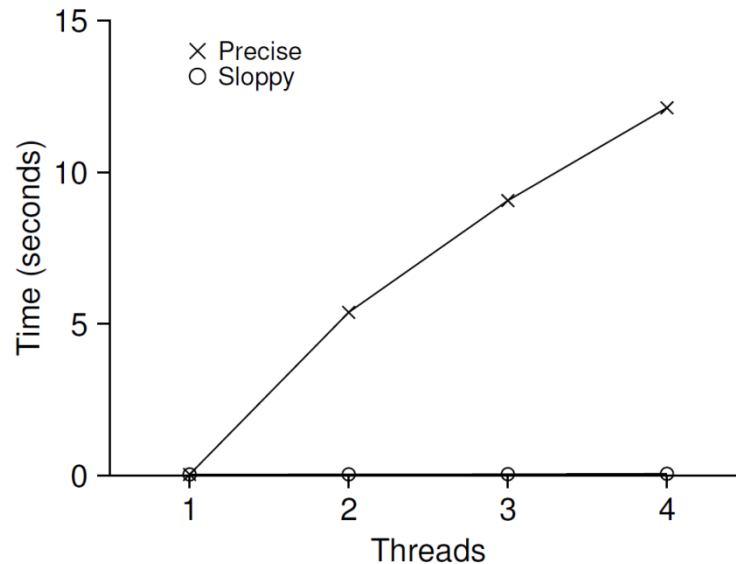
```
1       typedef struct __counter_t {
2               int value;
3               pthread_lock_t lock;
4       } counter_t;
5
6       void init(counter_t *c) {
7               c->value = 0;
8               Pthread_mutex_init(&c->lock, NULL);
9       }
10
11      void increment(counter_t *c) {
12              Pthread_mutex_lock(&c->lock);
13              c->value++;
14              Pthread_mutex_unlock(&c->lock);
15      }
16
```

# Example: Concurrent Counters with Locks

```
(Cont.)
17      void decrement(counter_t *c) {
18              Pthread_mutex_lock(&c->lock);
19              c->value--;
20              Pthread_mutex_unlock(&c->lock);
21      }
22
23      int get(counter_t *c) {
24              Pthread_mutex_lock(&c->lock);
25              int rc = c->value;
26              Pthread_mutex_unlock(&c->lock);
27              return rc;
28      }
```

# The performance costs of the simple approach

- Each thread updates a single shared counter.
  - Each thread updates the counter one million times.



**Performance of Traditional vs. Sloppy Counters**

**Synchronized counter scales poorly.**

# Perfect Scaling

- Even though more work is done, it is **done in parallel**.
- The time taken to complete the task is *not increased*.

# Sloppy counter

- The sloppy counter works by representing …
  - A single **logical counter** via numerous local physical counters, on per CPU core
  - A single **global counter**
  - There are **lock**s:
    - One for each local counter and one for the global counter
- Example: on a machine with four CPUs
  - Four local counters
  - One global counter

# The basic idea of sloppy counting

- When a thread running on a core wishes to increment the counter.
  - It increment its local counter.
  - Each CPU has its own local counter:
    - Threads across CPUs can update local counters *without contention*.
    - Thus counter updates are scalable.
  - The local values are periodically transferred to the global counter.
    - Acquire the global lock
    - Increment it by the local counter's value
    - The local counter is then reset to zero.

# The basic idea of sloppy counting (Cont.)

- How often the local-to-global transfer occurs is determined by a threshold, $S$ (sloppiness).
  - The smaller $S$:
    - The more the counter behaves like the *non-scalable counter*.
  - The bigger $S$:
    - The more scalable the counter.
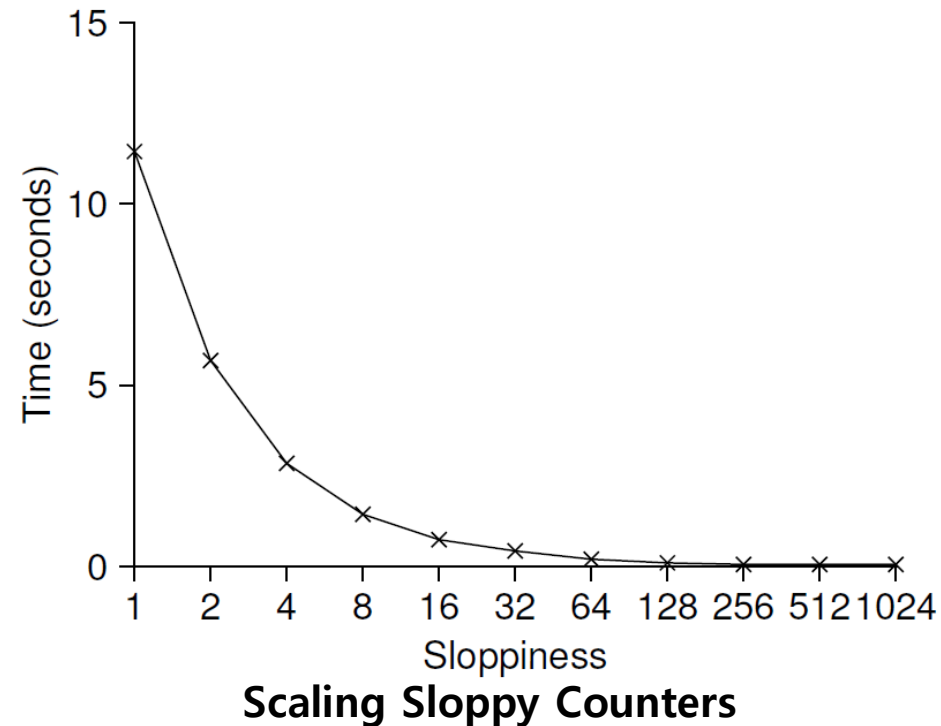    - The further off the global value might be from the *actual count*.

# Sloppy counter example

- Tracing the Sloppy Counters
  - The threshold S is set to 5.
  - There are threads on each of four CPUs
  - Each thread updates their local counters $L_1 \ldots L_4$.

| Time | L1 | L2 | L3 | L4 | G |
|------|------|------|------|------|------|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 1 | 1 | 0 |
| 2 | 1 | 0 | 2 | 1 | 0 |
| 3 | 2 | 0 | 3 | 1 | 0 |
| 4 | 3 | 0 | 3 | 2 | 0 |
| 5 | 4 | 1 | 3 | 3 | 0 |
| 6 | 5 → 0 | 1 | 3 | 4 | 5 (from ) |
| 7 | 0 | 2 | 4 | 5 → 0 | 10 (from ) |

# Importance of the threshold value $S$

- Each of four threads increments a counter 1 million times on four CPUs.
  - Low S → Performance is **poor**, The global count is always quite **accurate**.
  - High S → Performance is **excellent**, The global count **lags**.



**Scaling Sloppy Counters**

# Sloppy Counter Implementation

```c
1       typedef struct __counter_t {
2           int global;             // global count
3           pthread_mutex_t glock;  // global lock
4           int local[NUMCPUS];     // local count (per cpu)
5           pthread_mutex_t llock[NUMCPUS]; // ... and locks
6           int threshold;          // update frequency
7       } counter_t;
8
9       // init: record threshold, init locks, init values
10      //       of all local counts and global count
11      void init(counter_t *c, int threshold) {
12          c->thres hold = threshold;
13
14          c->global = 0;
15          pthread_mutex_init(&c->glock, NULL);
16
17          int i;
18          for (i = 0; i < NUMCPUS; i++) {
19              c->local[i] = 0;
20              pthread_mutex_init(&c->llock[i], NULL);
21          }
22      }
23
```

# Sloppy Counter Implementation (Cont.)

```
(Cont.)
24      // update: usually, just grab local lock and update local amount
25      //        once local count has risen by 'threshold', grab global
26      //        lock and transfer local values to it
27      void update(counter_t *c, int threadID, int amt) {
28          pthread_mutex_lock(&c->llock[threadID]);
29          c->local[threadID] += amt;        // assumes amt > 0
30          if (c->local[threadID] >= c->threshold) { // transfer to global
31              pthread_mutex_lock(&c->glock);
32              c->global += c->local[threadID];
33              pthread_mutex_unlock(&c->glock);
34              c->local[threadID] = 0;
35          }
36          pthread_mutex_unlock(&c->llock[threadID]);
37      }
38
39      // get: just return global amount (which may not be perfect)
40      int get(counter_t *c) {
41          pthread_mutex_lock(&c->glock);
42          int val = c->global;
43          pthread_mutex_unlock(&c->glock);
44          return val;    // only approximate!
45      }
```

# Concurrent Linked Lists

```c
1       // basic node structure
2       typedef struct __node_t {
3               int key;
4               struct __node_t *next;
5       } node_t;
6
7       // basic list structure (one used per list)
8       typedef struct __list_t {
9               node_t *head;
10              pthread_mutex_t lock;
11      } list_t;
12
13      void List_Init(list_t *L) {
14              L->head = NULL;
15              pthread_mutex_init(&L->lock, NULL);
16      }
```

# Concurrent Linked Lists (Cont.)

```c
18      int List_Insert(list_t *L, int key) {
19              pthread_mutex_lock(&L->lock);
20              node_t *new = malloc(sizeof(node_t));
21              if (new == NULL) {
22                      perror("malloc");
23                      pthread_mutex_unlock(&L->lock);
24                      return -1; // fail
25              }
26              new->key = key;
27              new->next = L->head;
28              L->head = new;
29              pthread_mutex_unlock(&L->lock);
30              return 0; // success
31      }

33      int List_Lookup(list_t *L, int key) {
34              pthread_mutex_lock(&L->lock);
35              node_t *curr = L->head;
36              while (curr) {
37                      if (curr->key == key) {
38                              pthread_mutex_unlock(&L->lock);
39                              return 0; // success
40                      }
41                      curr = curr->next;
42              }
43              pthread_mutex_unlock(&L->lock);
44              return -1; // failure
45      }
```

# Concurrent Linked Lists (Cont.)

- The code **acquires** a lock in the insert routine upon entry.
- The code **releases** the lock upon exit.
  - If `malloc()` happens to *fail*, the code must also <u>release the lock</u> before failing the insert.
  - This kind of exceptional control flow has been shown to be <span style="color:red">quite error prone</span>.
  - **Solution**: The lock and release *only surround* the actual critical section in the insert code

# Concurrent Linked List: Rewritten

```
1        void List_Init(list_t *L) {
2                L->head = NULL;
3                pthread_mutex_init(&L->lock, NULL);
4        }
5
6        void List_Insert(list_t *L, int key) {
7                // synchronization not needed
8                node_t *new = malloc(sizeof(node_t));
9                if (new == NULL) {
10                       perror("malloc");
11                       return;
12               }
13               new->key = key;
14
15               // just lock critical section
16               pthread_mutex_lock(&L->lock);
17               new->next = L->head;
18               L->head = new;
19               pthread_mutex_unlock(&L->lock);
20       }
21
```

# Concurrent Linked List: Rewritten (Cont.)

```
(Cont.)
22        int List_Lookup(list_t *L, int key) {
23                int rv = -1;
24                pthread_mutex_lock(&L->lock);
25                node_t *curr = L->head;
26                while (curr) {
27                        if (curr->key == key) {
28                                rv = 0;
29                                break;
30                        }
31                        curr = curr->next;
32                }
33                pthread_mutex_unlock(&L->lock);
34                return rv; // now both success and failure
35        }
```

# Scaling Linked List

- Hand-over-hand locking (lock coupling)
  - Add **a lock per node** of the list instead of having a single lock for the entire list.
  - When traversing the list,
    - First grabs the next node's lock.
    - And then releases the current node's lock.
  - **Enable a high degree of concurrency in list operations.**
    - However, in practice, the overheads of acquiring and releasing locks for each node of a list traversal is *prohibitive*.

# Michael and Scott Concurrent Queues

- There are two locks.
  - One for the **head** of the queue.
  - One for the **tail**.
  - The goal of these two locks is to **enable concurrency of *enqueue* and *dequeue* operations.**
- Add a dummy node
  - Allocated in the queue initialization code
  - **Enable the separation of head and tail operations**

# Concurrent Queues (Cont.)

```c
1       typedef struct __node_t {
2               int value;
3               struct __node_t *next;
4       } node_t;
5
6       typedef struct __queue_t {
7               node_t *head;
8               node_t *tail;
9               pthread_mutex_t headLock;
10              pthread_mutex_t tailLock;
11      } queue_t;
12
13      void Queue_Init(queue_t *q) {
14              node_t *tmp = malloc(sizeof(node_t));
15              tmp->next = NULL;
16              q->head = q->tail = tmp;
17              pthread_mutex_init(&q->headLock, NULL);
18              pthread_mutex_init(&q->tailLock, NULL);
19      }
20
```

# Concurrent Queues (Cont.)

```c
21          void Queue_Enqueue(queue_t *q, int value) {
22                  node_t *tmp = malloc(sizeof(node_t));
23                  assert(tmp != NULL);
24                  tmp->value = value;
25                  tmp->next = NULL;
26
27                  pthread_mutex_lock(&q->tailLock);
28                  q->tail->next = tmp;
29                  q->tail = tmp;
30                  pthread_mutex_unlock(&q->tailLock);
31          }
32
33          int Queue_Dequeue(queue_t *q, int *value) {
34                  pthread_mutex_lock(&q->headLock);
35                  node_t *tmp = q->head;
36                  node_t *newHead = tmp->next;
37                  if (newHead == NULL) {
38                          pthread_mutex_unlock(&q->headLock);
39                          return -1; // queue was empty
40                  }
41                  *value = newHead->value;
42                  q->head = newHead;
43                  pthread_mutex_unlock(&q->headLock);
44                  free(tmp);
45                  return 0;
46          }
```
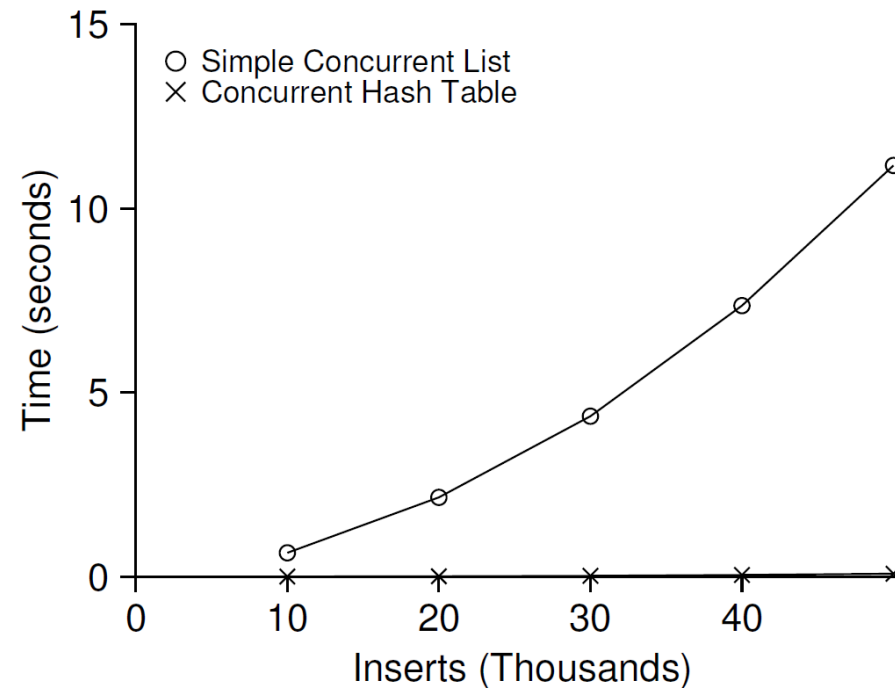
# Concurrent Hash Table

- Focus on a simple hash table
  - The hash table does not resize.
  - Built using the concurrent lists
  - It uses a **lock per hash bucket** each of which is represented by *a list*.

# Concurrent Hash Table

```c
1       #define BUCKETS (101)
2
3       typedef struct __hash_t {
4               list_t lists[BUCKETS];
5       } hash_t;
6
7       void Hash_Init(hash_t *H) {
8               int i;
9               for (i = 0; i < BUCKETS; i++) {
10                      List_Init(&H->lists[i]);
11              }
12      }
13
14      int Hash_Insert(hash_t *H, int key) {
15              int bucket = key % BUCKETS;
16              return List_Insert(&H->lists[bucket], key);
17      }
18
19      int Hash_Lookup(hash_t *H, int key) {
20              int bucket = key % BUCKETS;
21              return List_Lookup(&H->lists[bucket], key);
22      }
```

# Performance of Concurrent Hash Table

- From 10,000 to 50,000 concurrent updates from each of four threads.



**The simple concurrent hash table scales magnificently.**