

Complex

CO3

Session - 3

COURSE NAME: OPERATING SYSTEMS
COURSE CODE: 23CS2104R/A

LOCKS, LOCKED DATA STRUCTURES, CONDITION VARIABLES, MUTEX.

Simple

AIM OF THE SESSION

To familiarize students with the basic concept of locks.

INSTRUCTIONAL OBJECTIVES

This Session is designed to:

1. Demonstrate what is meant by Lock.
2. Describe the types of Locks.
3. Describe the Locked Data Structures.
4. Describe the Concurrent Counters with Locks.

LEARNING OUTCOMES

At the end of this session, you should be able to:

1. Defines what is a lock.
2. Describe Test and Set and Compare and Swap.
3. Summarize the Role of Mutex Locks.

SYNCHRONIZATION HARDWARE

- Many systems provide hardware support for critical section code
- Uniprocessors – could disable interrupts
 - Currently running code would execute without preemption
 - Generally too inefficient on multiprocessor systems
 - Operating systems using this not broadly scalable
- Modern machines provide special atomic hardware instructions
 - Atomic = non-interruptible
 - Either test memory word and set value
 - Or swap contents of two memory words

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

```
lock_t mutex; // some globally-allocated lock 'mutex'  
...  
lock(&mutex);  
balance = balance + 1; unlock(&mutex);
```

CONTD..

- Lock variable holds the state of 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 no other thread holds the lock, the thread will **acquire** the lock.
 - **Enter** the critical section.
 - This thread is said to be the owner of the lock.
 - Other threads are prevented from entering the critical section while the first thread that holds the lock is in there.

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

CONTROLLING INTERRUPTS

- **Disable Interrupts** for critical sections
- One of the earliest solutions used to provide mutual exclusion
- Invented for single-processor systems.

```
1  void lock() {  
2      DisableInterrupts();  
3  }  
4  void unlock() {  
5      EnableInterrupts();  
6  }
```


LOCKS

- Test and Set
- Compare and Swap
- Load-Linked and Store-Conditional
- Fetch-And-Add
- Ticket Lock
- Futex Lock

TEST AND SET (ATOMIC EXCHANGE)

- An instruction to support the creation of simple locks
- **return**(testing) old value pointed to by the ptr.
- Simultaneously **update**(setting) said value to new.
- This sequence of operations is performed atomically.

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

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, 1 that it is held
7      lock->flag = 0;
8  }
9
10 void lock(lock_t *lock) {
11     while (TestAndSet(&lock->flag, 1) == 1)
12         ; // spin-wait (do nothing)
13 }
14
15 void unlock(lock_t *lock) {
16     lock->flag = 0;
17 }
```

COMPARE-AND-SWAP

- 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.

```
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

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

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

- Spin lock with compare-and-swap

```
1  int LoadLinked(int *ptr) {  
2      return *ptr;  
3  }  
4  
5  int StoreConditional(int *ptr, int value) {  
6      if (no one has updated *ptr since the LoadLinked to this address) {  
7          *ptr = value;  
8          return 1; // success!  
9      } else {  
10         return 0; // failed to update  
11     }  
12 }
```

The store-conditional only succeeds if no intermittent store to the address has taken place.

- **success:** return 1 and update the value at ptr to value.
- **fail:** the value at ptr is not updated and 0 is returned.

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 fetch-and add.
- Ensure progress for all threads. → fairness

```
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
17 void unlock(lock_t *lock) {
18     lock->turn = lock->turn + 1;
19 }
```


USING QUEUES: SLEEPING INSTEAD OF SPINNING

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

WAKEUP/WAITING RACE

- 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.

FUTEX

- Linux provides a futex (is similar to Solaris's park and unpark). More functionality goes into the kernel.
- `futex_wait(address, expected)`
 - Put the calling thread to sleep
 - If the value at address is not equal to expected, the call returns immediately.
- `futex_wake(address)`
 - \nexists Wake one thread that is waiting on the queue.

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.

LOCK-BASED CONCURRENT DATA STRUCTURES

- 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

```
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

```
1  typedef struct __counter_t {
2      int          value;
3      pthread_mutex_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
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 }
```

SLOPPY COUNTER

- The sloppy counter works by representing ...
 - A single **logical counter** via numerous local physical counters, one per CPU core
 - A single **global counter**
 - There are **locks**:
 - 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 increments its local counter.
 - Each CPU has its 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.

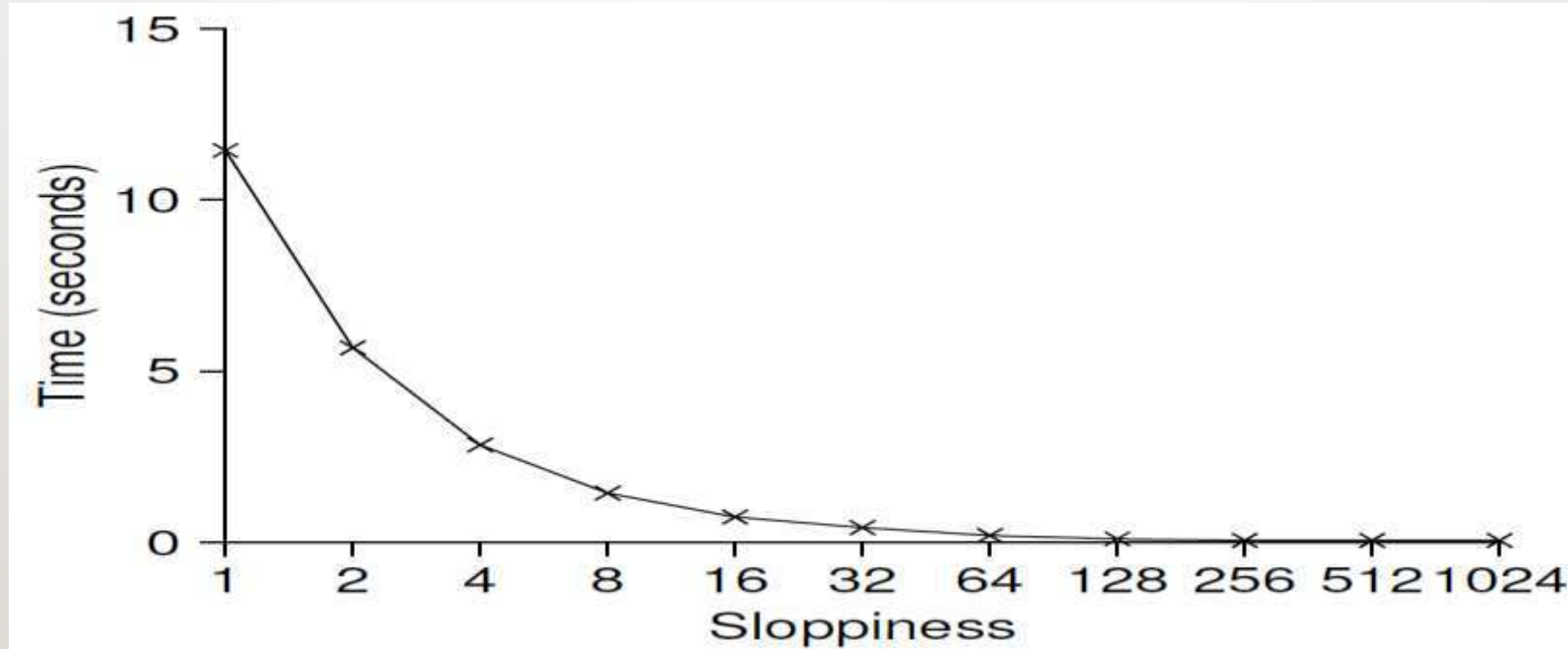
SLOPPY COUNTER EXAMPLE

- **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 $L_1 \dots L_4$.

Time	L_1	L_2	L_3	L_4	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	1	3	4	5 (from L_1)
7	0	2	4	5	10 (from L_4)

IMPORTANCE OF THE THRESHOLD VALUE S

- Each 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**.



CONCURRENT LINKED LIST:

```
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 }
17
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 }
32
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 release the lock before failing the insert.
 - This kind of exceptional control flow has been shown to be quite error prone.
 - **Solution:** The lock and release *only surround* the actual critical section in the insert code

CONCURRENT LINKED LIST: REFACTORED INSERT

```
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
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 grab 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 can be *prohibitive*.
 - Perhaps a hybrid (where you grab a new lock every so many nodes)

approach?

MICHAEL AND SCOTT CONCURRENT QUEUES

- Queues uses enqueue / dequeue operations only
- 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.)

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

(Cont.)

CONCURRENT QUEUES (CONT.)

(Cont.)

```
21      void Queue_Enqueue(queue_t *q, int value) {
22          node_t *tmp = malloc(sizeof(node_t));
23          assert(tmp != NULL && "Malloc failed");
24
25          tmp->value = value;
26          tmp->next = NULL;
27
28          pthread_mutex_lock(&q->tailLock);
29          q->tail->next = tmp;
30          q->tail = tmp;
31          pthread_mutex_unlock(&q->tailLock);
32      }
```

(Cont.)

CONCURRENT QUEUES (CONT.)

(Cont.)

```
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); return -1;
39             // 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, simplest hash, slow searches
 - Built using the concurrent lists
 - It uses a lock per hash bucket each of which is represented by a list.

CONCURRENT HASH TABLE

```
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      }
```

CONDITION VARIABLES

- There are many cases where a thread wishes to check whether a
 - **condition** is true before continuing its execution.
- Example: A parent thread might wish to check whether a child thread has *completed*.
 - This is often called a `join()`.

HOW TO WAIT FOR A CONDITION

- Condition variable
 - **Waiting** on the condition
 - An explicit queue that threads can put themselves on when some state of execution is not as desired.
 - **Signaling** on the condition
 - Some other thread, *when it changes said state*, can wake one of those waiting threads and allow them to continue.

DEFINITION AND ROUTINES

- Declare condition variable

```
pthread_cond_t c;
```

- Proper initialization is required.

Operation (the POSIX calls)

```
pthread_cond_wait(pthread_cond_t *c,  
pthread_mutex_t *m);  
pthread_cond_signal(pthread_cond_t *c);
```


DEFINITION AND ROUTINES

- The wait() call takes a mutex as a parameter.
- ⚡ The wait() call **release the lock** and put the calling thread to sleep.
- ⚡ When the thread wakes up, it must **re-acquire the lock**.

PARENT WAITING FOR THE CHILD

```
1      int done = 0;
2      pthread_mutex_t m = PTHREAD_MUTEX_INITIALIZER;
3      pthread_cond_t c = PTHREAD_COND_INITIALIZER;
4
5      void thr_exit() { Pthread_mutex_lock(&m);
6                      done = 1;
7                      Pthread_cond_signal(&c);
8                      Pthread_mutex_unlock(&m);
9
10     }
11
12     void *child(void *arg) {
13         printf("child\n");
14         thr_exit();
15         return NULL;
16     }
17
18     void thr_join() { Pthread_mutex_lock(&m);
19                     while (done == 0)
20                         Pthread_cond_wait(&c,
21                                         Pthread_mutex_unlock(&m);
22
23     }
24
```

PARENT WAITING FOR CHILD: USE A CONDITION VARIABLE

```
(cont.  
)    in main(int argc, char *argv[]) {  
25    t  
26        printf("parent: begin\n");  
27        pthread_t p;  
28        Pthread_create(&p, NULL,  
29        child, NULL);  
30        thr_join();  
31        printf("parent: end\n");  
32        return 0;  
33    }
```

PARENT WAITING FOR CHILD: USE A CONDITION VARIABLE

- **Parent:**
 - Create the child thread and continue running itself.
 - Call into `thr_join()` to wait for the child thread to complete.
 - $\cancel{\phi}$ Acquire the lock
 - $\cancel{\phi}$ Check if the child is done
 - $\cancel{\phi}$ Put itself to sleep by calling `wait()`
 - $\cancel{\phi}$ Release the lock

PARENT WAITING FOR CHILD: USE A CONDITION VARIABLE

- **Child:**
 - Print the message “child”
 - Call `thr_exit()` to wake the parent thread
 - \nexists Grab the lock
 - \nexists Set the state variable done
 - \nexists Signal the parent thus waking it.

THE PRODUCER / CONSUMER (BOUND BUFFER) PROBLEM

- **Producer**
 - Produce data items
 - Wish to place data items in a buffer
- **Consumer**
 - Grab data items out of the buffer consume them in some way
- Example: Multi-threaded web server
 - A producer puts HTTP requests in to a work queue
 - Consumer threads take requests out of this queue and process them

BOUNDED BUFFER

- A bounded buffer is used when you pipe the output of one program into another.
- Example: `grep foo file.txt | wc -l`
 - ¢ The `grep` process is the producer.
 - ¢ The `wc` process is the consumer.
 - ¢ Between them is an in-kernel bounded buffer.
- Bounded buffer is a Shared resource à **Synchronized access** is required.

THE SINGLE BUFFER PRODUCER/CONSUMER SOLUTION

- Use two condition variables and while
- **Producer** threads wait on the condition empty, and signals fill.
- **Consumer** threads wait on fill and signal empty.

```
1      cond_t empty, fill;
2      mutex_t mutex;
3
4      void *producer(void *arg) {  int
5          for      i, loops = (int) arg;
6              (i = 0; i < loops; i++) {
7                  Pthread_mutex_lock(&mutex);  while
8                      (count == 1)
9                      Pthread_cond_wait(&empty, &mutex);  put(i);
10                     Pthread_cond_signal(&fill);
11                     Pthread_mutex_unlock(&mutex);
12
13                 }
14     }
15
```


THE SINGLE BUFFER PRODUCER/CONSUMER SOLUTION

```
17 void *consumer(void *arg) {
18     int i;
19     for (i = 0; i < loops; i++) {
20         Pthread_mutex_lock(&mutex);           // c1
21         while (count == 0)                     // c2
22             Pthread_cond_wait(&cond, &mutex); // c3
23         int tmp = get();                       // c4
24         Pthread_cond_signal(&cond);           // c5
25         Pthread_mutex_unlock(&mutex);         // c6
26         printf("%d\n", tmp);
27     }
28 }
```

THE FINAL PRODUCER/CONSUMER SOLUTION(MULTIPLE BUFFERS)

- More **concurrency** and **efficiency** à Add more buffer slots.
- Allow concurrent production or consuming to take place.
- Reduce context switches.

```
1      int buffer[MAX];  int
2      fill = 0;  int use =
3      0;  int count = 0;
4
5
6      void put(int value) {  buffer[fill] =
7          value;  fill = (fill + 1) % MAX;
8          count++;
9      }
10
11     int
12         get() {
13         int tmp = buffer[use];  use =
14         (use + 1) % MAX;  count--;
15         return tmp;
16
17     }
```

The Final Put and Get Routines

THE FINAL PRODUCER/CONSUMER SOLUTION (CONT.)

```
1      cond_t empty, fill;
2      mutex_t mutex;
3
4      void *producer(void *arg) {
5          int i, loops = (int) arg;
6          for (i = 0; i < loops; i++) {
7              Pthread_mutex_lock(&mutex);           // p1
8              while (count == MAX)                 // p2
9                  Pthread_cond_wait(&empty, &mutex); // p3
10             put(i);                               // p4
11             Pthread_cond_signal(&fill);           // p5
12             Pthread_mutex_unlock(&mutex);         // p6
13         }
14     }
15
16     void *consumer(void *arg) {
17         int i, loops = (int) arg;
18         for (i = 0; i < loops; i++) {
19             Pthread_mutex_lock(&mutex);           // c1
20             while (count == 0)                    // c2
21                 Pthread_cond_wait(&fill,          // c3
22                                 &mutex);
23             int tmp = get();                       // c4
```

```
(Cont.)  
23         Pthread_cond_signal(&empty);           // c5  
24         Pthread_mutex_unlock(&mutex);          // c6  
25         printf("%d\n", tmp);  
26     }  
27 }
```

The Final Working Solution (Cont.)

- p2: **A producer** only sleeps if all buffers are currently filled.
- c2: **A consumer** only sleeps if all buffers are currently empty.
- **Note:** The Producer Consumer Problem **with Multiple Buffer Slots** is practically demonstrated in the next Chapter using semaphores and mutex.

MUTEX LOCK

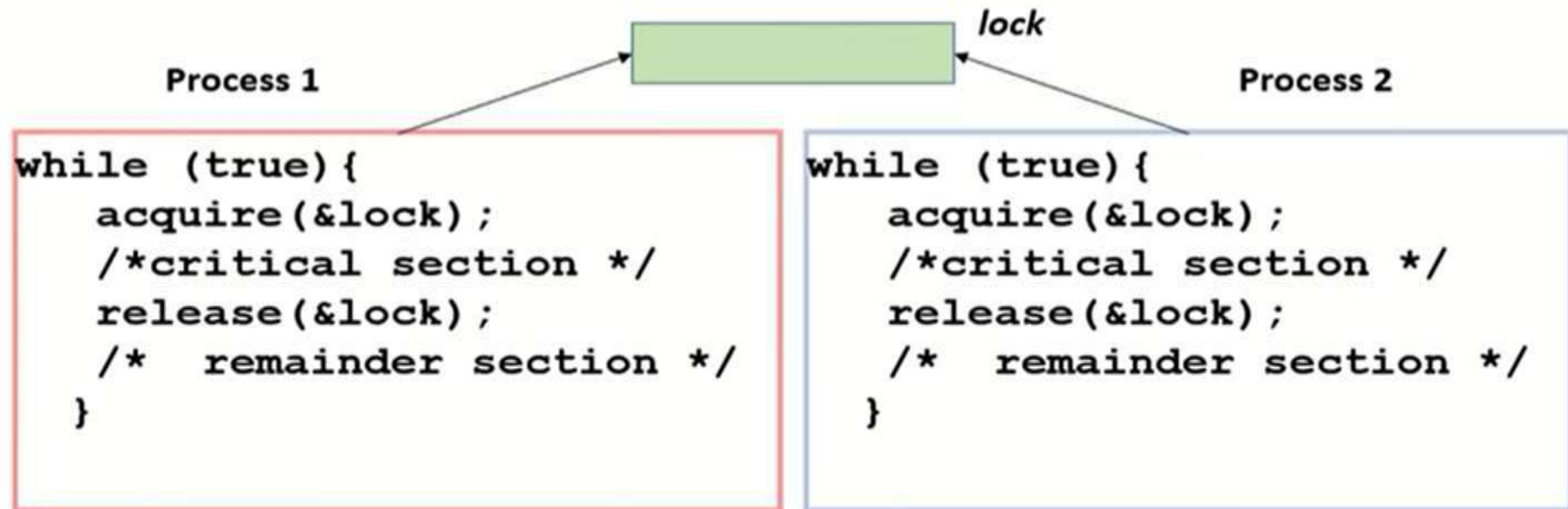
- A lock variable provides the simplest synchronization mechanism for processes.
 1. Its a **software mechanism** implemented in user mode, i.e. no support required from the Operating System.
 2. Its a busy waiting solution (keeps the CPU busy even when its technically waiting).
 3. It can be used for more than two processes.
 4. When Lock = 0 implies critical section is vacant (initial value) and Lock = 1 implies critical section occupied.

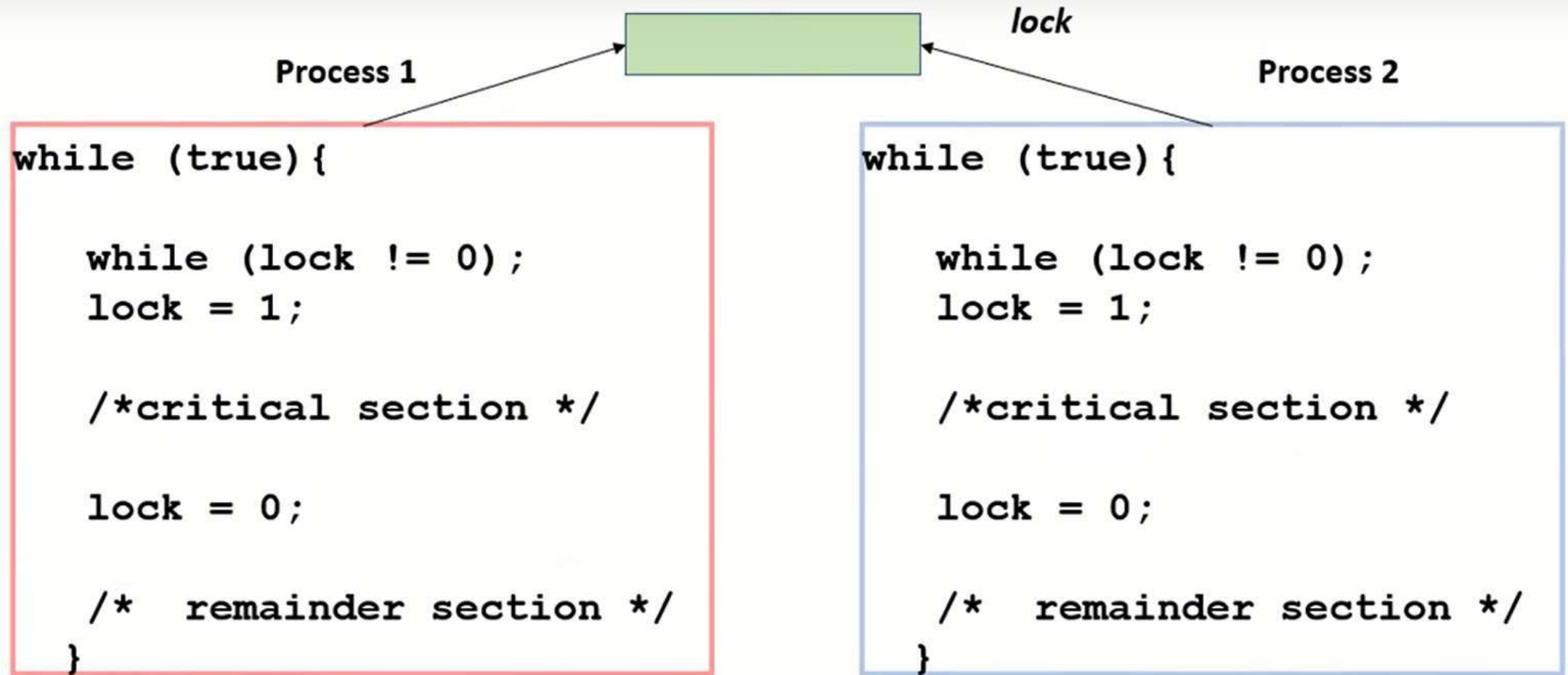
- Basic operations on a **lock**:
- **acquire**: mark the lock as owned by the current thread; if some other thread already owns the lock then first wait until the lock is free.
- **release**: mark the lock as free (it must currently be owned by the calling thread).

```
do {  
    acquire lock  
    critical section  
    release lock  
    remainder section  
} while
```

Mutex Locks

- **mutex** lock (**mut**ual **ex**clusion)





THANK YOU



Team – Operating System