# 29

# Lock-based Concurrent Data Structures

Before moving beyond locks, we'll first describe how to use locks in some common data structures. Adding locks to a data structure to make it usable by threads makes the structure **thread safe**. Of course, exactly how such locks are added determines both the correctness and performance of the data structure. And thus, our challenge:

CRUX: HOW TO ADD LOCKS TO DATA STRUCTURES
When given a particular data structure, how should we add locks to it, in order to make it work correctly? Further, how do we add locks such that the data structure yields high performance, enabling many threads to access the structure at once, i.e., **concurrently**?

Of course, we will be hard pressed to cover all data structures or all methods for adding concurrency, as this is a topic that has been studied for years, with (literally) thousands of research papers published about it. Thus, we hope to provide a sufficient introduction to the type of thinking required, and refer you to some good sources of material for further inquiry on your own. We found Moir and Shavit's survey to be a great source of information [MS04].

## 29.1 Concurrent Counters

One of the simplest data structures is a counter. It is a structure that is commonly used and has a simple interface. We define a simple non-concurrent counter in Figure 29.1.

### Simple But Not Scalable

As you can see, the non-synchronized counter is a trivial data structure, requiring a tiny amount of code to implement. We now have our next challenge: how can we make this code **thread safe**? Figure 29.2 shows how we do so.

1

```
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  }
```

Figure 29.1: **A Counter Without 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  }
```

Figure 29.2: **A Counter With Locks**

This concurrent counter is simple and works correctly. In fact, it fol-
lows a design pattern common to the simplest and most basic concurrent
data structures: it simply adds a single lock, which is acquired when call-
ing a routine that manipulates the data structure, and is released when
returning from the call. In this manner, it is similar to a data structure
built with **monitors** [BH73], where locks are acquired and released auto-
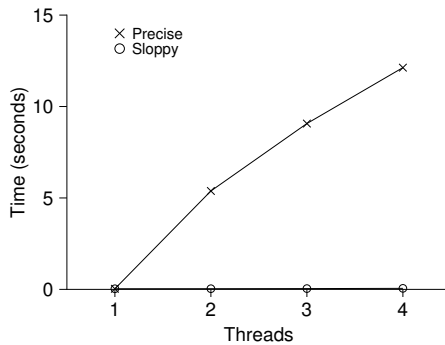matically as you call and return from object methods.

Figure 29.3: **Performance of Traditional vs. Sloppy Counters**

At this point, you have a working concurrent data structure. The problem you might have is performance. If your data structure is too slow, you'll have to do more than just add a single lock; such optimizations, if needed, are thus the topic of the rest of the chapter. Note that if the data structure is *not* too slow, you are done! No need to do something fancy if something simple will work.

To understand the performance costs of the simple approach, we run a benchmark in which each thread updates a single shared counter a fixed number of times; we then vary the number of threads. Figure 29.3 shows the total time taken, with one to four threads active; each thread updates the counter one million times. This experiment was run upon an iMac with four Intel 2.7 GHz i5 CPUs; with more CPUs active, we hope to get more total work done per unit time.

From the top line in the figure (labeled *precise*), you can see that the performance of the synchronized counter scales poorly. Whereas a single thread can complete the million counter updates in a tiny amount of time (roughly 0.03 seconds), having two threads each update the counter one million times concurrently leads to a massive slowdown (taking over 5 seconds!). It only gets worse with more threads.

Ideally, you'd like to see the threads complete just as quickly on multiple processors as the single thread does on one. Achieving this end is called **perfect scaling**; even though more work is done, it is done in parallel, and hence the time taken to complete the task is not increased.

### Scalable Counting

Amazingly, researchers have studied how to build more scalable counters for years [MS04]. Even more amazing is the fact that scalable counters matter, as recent work in operating system performance analysis has shown [B+10]; without scalable counting, some workloads running on

| 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 | $5 \rightarrow 0$ | 1 | 3 | 4 | 5 (from $L_1$) |
| 7 | 0 | 2 | 4 | $5 \rightarrow 0$ | 10 (from $L_4$) |

Figure 29.4: **Tracing the Sloppy Counters**

Linux suffer from serious scalability problems on multicore machines.

Though many techniques have been developed to attack this problem, we'll now describe one particular approach. The idea, introduced in recent research [B+10], is known as a **sloppy counter**.

The sloppy counter works by representing a single logical counter via numerous *local* physical counters, one per CPU core, as well as a single *global* counter. Specifically, on a machine with four CPUs, there are four local counters and one global one. In addition to these counters, there are also locks: one for each local counter, and one for the global counter.

The basic idea of sloppy counting is as follows. When a thread running on a given core wishes to increment the counter, it increments its local counter; access to this local counter is synchronized via the corresponding local lock. Because each CPU has its own local counter, threads across CPUs can update local counters without contention, and thus counter updates are scalable.

However, to keep the global counter up to date (in case a thread wishes to read its value), the local values are periodically transferred to the global counter, by acquiring the global lock and incrementing it by the local counter's value; the local counter is then reset to zero.

How often this local-to-global transfer occurs is determined by a threshold, which we call $S$ here (for sloppiness). The smaller $S$ is, the more the counter behaves like the non-scalable counter above; the bigger $S$ is, the more scalable the counter, but the further off the global value might be from the actual count. One could simply acquire all the local locks and the global lock (in a specified order, to avoid deadlock) to get an exact value, but that is not scalable.

To make this clear, let's look at an example (Figure 29.4). In this example, the threshold $S$ is set to $5$, and there are threads on each of four CPUs updating their local counters $L_1$ ... $L_4$. The global counter value ($G$) is also shown in the trace, with time increasing downward. At each time step, a local counter may be incremented; if the local value reaches the threshold $S$, the local value is transferred to the global counter and the local counter is reset.

The lower line in Figure 29.3 (labeled *sloppy*, on page 3) shows the performance of sloppy counters with a threshold $S$ of $1024$. Performance is excellent; the time taken to update the counter four million times on four processors is hardly higher than the time taken to update it one million times on one processor.

```
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->threshold = threshold;
13      c->global = 0;
14      pthread_mutex_init(&c->glock, NULL);
15      int i;
16      for (i = 0; i < NUMCPUS; i++) {
17          c->local[i] = 0;
18          pthread_mutex_init(&c->llock[i], NULL);
19      }
20  }
21
22  // update: usually, just grab local lock and update local amount
23  //         once local count has risen by 'threshold', grab global
24  //         lock and transfer local values to it
25  void update(counter_t *c, int threadID, int amt) {
26      int cpu = threadID % NUMCPUS;
27      pthread_mutex_lock(&c->llock[cpu]);
28      c->local[cpu] += amt;                  // assumes amt > 0
29      if (c->local[cpu] >= c->threshold) {  // transfer to global
30          pthread_mutex_lock(&c->glock);
31          c->global += c->local[cpu];
32          pthread_mutex_unlock(&c->glock);
33          c->local[cpu] = 0;
34      }
35      pthread_mutex_unlock(&c->llock[cpu]);
36  }
37
38  // get: just return global amount (which may not be perfect)
39  int get(counter_t *c) {
40      pthread_mutex_lock(&c->glock);
41      int val = c->global;
42      pthread_mutex_unlock(&c->glock);
43      return val; // only approximate!
44  }
```

Figure 29.5: **Sloppy Counter Implementation**

Figure 29.6 shows the importance of the threshold value $S$, with four threads each incrementing the counter 1 million times on four CPUs. If $S$ is low, performance is poor (but the global count is always quite accurate); if $S$ is high, performance is excellent, but the global count lags (by at most the number of CPUs multiplied by $S$). This accuracy/performance trade-off is what sloppy counters enables.

A rough version of such a sloppy counter is found in Figure 29.5. Read it, or better yet, run it yourself in some experiments to better understand how it works.
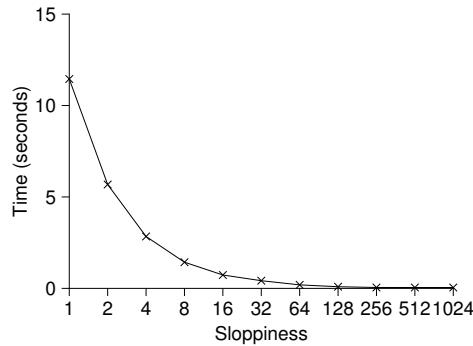
Figure 29.6: **Scaling Sloppy Counters**

## 29.2 Concurrent Linked Lists

We next examine a more complicated structure, the linked list. Let's start with a basic approach once again. For simplicity, we'll omit some of the obvious routines that such a list would have and just focus on concurrent insert; we'll leave it to the reader to think about lookup, delete, and so forth. Figure 29.7 shows the code for this rudimentary data structure.

As you can see in the code, the code simply acquires a lock in the insert routine upon entry, and releases it upon exit. One small tricky issue arises if `malloc()` happens to fail (a rare case); in this case, 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; a recent study of Linux kernel patches found that a huge fraction of bugs (nearly 40%) are found on such rarely-taken code paths (indeed, this observation sparked some of our own research, in which we removed all memory-failing paths from a Linux file system, resulting in a more robust system [S+11]).

Thus, a challenge: can we rewrite the insert and lookup routines to remain correct under concurrent insert but avoid the case where the failure path also requires us to add the call to unlock?

The answer, in this case, is yes. Specifically, we can rearrange the code a bit so that the lock and release only surround the actual critical section in the insert code, and that a common exit path is used in the lookup code. The former works because part of the lookup actually need not be locked; assuming that `malloc()` itself is thread-safe, each thread can call into it without worry of race conditions or other concurrency bugs. Only when updating the shared list does a lock need to be held. See Figure 29.8 for the details of these modifications.

```
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  }
```

Figure 29.7: **Concurrent Linked List**

As for the lookup routine, it is a simple code transformation to jump out of the main search loop to a single return path. Doing so again reduces the number of lock acquire/release points in the code, and thus decreases the chances of accidentally introducing bugs (such as forgetting to unlock before returning) into the code.

### Scaling Linked Lists

Though we again have a basic concurrent linked list, once again we are in a situation where it does not scale particularly well. One technique that researchers have explored to enable more concurrency within a list is

```
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   }
```

Figure 29.8: **Concurrent Linked List: Rewritten**

something called **hand-over-hand locking** (a.k.a. **lock coupling**) [MS04].
    The idea is pretty simple. Instead of having a single lock for the entire
list, you instead add a lock per node of the list. When traversing the
list, the code first grabs the next node's lock and then releases the current
node's lock (which inspires the name hand-over-hand).
    Conceptually, a hand-over-hand linked list makes some sense; it en-
ables a high degree of concurrency in list operations. However, in prac-
tice, it is hard to make such a structure faster than the simple single lock
approach, as the overheads of acquiring and releasing locks for each node
of a list traversal is prohibitive. Even with very large lists, and a large
number of threads, the concurrency enabled by allowing multiple on-
going traversals is unlikely to be faster than simply grabbing a single
lock, performing an operation, and releasing it. Perhaps some kind of hy-
brid (where you grab a new lock every so many nodes) would be worth
investigating.

TIP: MORE CONCURRENCY ISN'T NECESSARILY FASTER
If the scheme you design adds a lot of overhead (for example, by acquiring and releasing locks frequently, instead of once), the fact that it is more concurrent may not be important. Simple schemes tend to work well, especially if they use costly routines rarely. Adding more locks and complexity can be your downfall. All of that said, there is one way to really know: build both alternatives (simple but less concurrent, and complex but more concurrent) and measure how they do. In the end, you can't cheat on performance; your idea is either faster, or it isn't.

TIP: BE WARY OF LOCKS AND CONTROL FLOW
A general design tip, which is useful in concurrent code as well as elsewhere, is to be wary of control flow changes that lead to function returns, exits, or other similar error conditions that halt the execution of a function. Because many functions will begin by acquiring a lock, allocating some memory, or doing other similar stateful operations, when errors arise, the code has to undo all of the state before returning, which is error-prone. Thus, it is best to structure code to minimize this pattern.

## 29.3 Concurrent Queues

As you know by now, there is always a standard method to make a concurrent data structure: add a big lock. For a queue, we'll skip that approach, assuming you can figure it out.

Instead, we'll take a look at a slightly more concurrent queue designed by Michael and Scott [MS98]. The data structures and code used for this queue are found in Figure 29.9 on the following page.

If you study this code carefully, you'll notice that there are two locks, one for the head of the queue, and one for the tail. The goal of these two locks is to enable concurrency of enqueue and dequeue operations. In the common case, the enqueue routine will only access the tail lock, and dequeue only the head lock.

One trick used by Michael and Scott is to add a dummy node (allocated in the queue initialization code); this dummy enables the separation of head and tail operations. Study the code, or better yet, type it in, run it, and measure it, to understand how it works deeply.

Queues are commonly used in multi-threaded applications. However, the type of queue used here (with just locks) often does not completely meet the needs of such programs. A more fully developed bounded queue, that enables a thread to wait if the queue is either empty or overly full, is the subject of our intense study in the next chapter on condition variables. Watch for it!

```
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
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  }
```

Figure 29.9: **Michael and Scott Concurrent Queue**

## 29.4 Concurrent Hash Table

We end our discussion with a simple and widely applicable concurrent data structure, the hash table. We'll focus on a simple hash table that does not resize; a little more work is required to handle resizing, which we leave as an exercise for the reader (sorry!).

This concurrent hash table is straightforward, is built using the concurrent lists we developed earlier, and works incredibly well. The reason

```
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  }
```

Figure 29.10: **A Concurrent Hash Table**

for its good performance is that instead of having a single lock for the en-
tire structure, it uses a lock per hash bucket (each of which is represented
by a list). Doing so enables many concurrent operations to take place.

Figure 29.11 shows the performance of the hash table under concur-
rent updates (from 10,000 to 50,000 concurrent updates from each of four
threads, on the same iMac with four CPUs). Also shown, for the sake
of comparison, is the performance of a linked list (with a single lock).
As you can see from the graph, this simple concurrent hash table scales
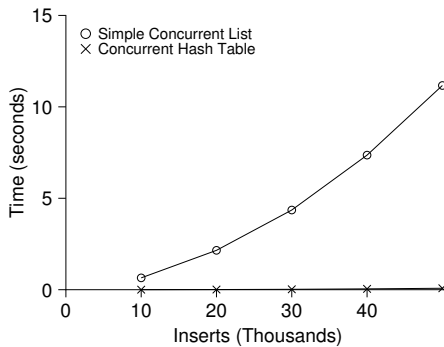magnificently; the linked list, in contrast, does not.



Figure 29.11: **Scaling Hash Tables**

> TIP: AVOID PREMATURE OPTIMIZATION (KNUTH'S LAW)
> When building a concurrent data structure, start with the most basic approach, which is to add a single big lock to provide synchronized access. By doing so, you are likely to build a *correct* lock; if you then find that it suffers from performance problems, you can refine it, thus only making it fast if need be. As **Knuth** famously stated, "Premature optimization is the root of all evil."
>
> Many operating systems utilized a single lock when first transitioning to multiprocessors, including Sun OS and Linux. In the latter, this lock even had a name, the **big kernel lock** (**BKL**). For many years, this simple approach was a good one, but when multi-CPU systems became the norm, only allowing a single active thread in the kernel at a time became a performance bottleneck. Thus, it was finally time to add the optimization of improved concurrency to these systems. Within Linux, the more straightforward approach was taken: replace one lock with many. Within Sun, a more radical decision was made: build a brand new operating system, known as Solaris, that incorporates concurrency more fundamentally from day one. Read the Linux and Solaris kernel books for more information about these fascinating systems [BC05, MM00].

## 29.5  Summary

We have introduced a sampling of concurrent data structures, from counters, to lists and queues, and finally to the ubiquitous and heavily-used hash table. We have learned a few important lessons along the way: to be careful with acquisition and release of locks around control flow changes; that enabling more concurrency does not necessarily increase performance; that performance problems should only be remedied once they exist. This last point, of avoiding **premature optimization**, is central to any performance-minded developer; there is no value in making something faster if doing so will not improve the overall performance of the application.

Of course, we have just scratched the surface of high performance structures. See Moir and Shavit's excellent survey for more information, as well as links to other sources [MS04]. In particular, you might be interested in other structures (such as B-trees); for this knowledge, a database class is your best bet. You also might be interested in techniques that don't use traditional locks at all; such **non-blocking data structures** are something we'll get a taste of in the chapter on common concurrency bugs, but frankly this topic is an entire area of knowledge requiring more study than is possible in this humble book. Find out more on your own if you are interested (as always!).

# References

[B+10] "An Analysis of Linux Scalability to Many Cores"
Silas Boyd-Wickizer, Austin T. Clements, Yandong Mao, Aleksey Pesterev, M. Frans Kaashoek,
Robert Morris, Nickolai Zeldovich
OSDI '10, Vancouver, Canada, October 2010
*A great study of how Linux performs on multicore machines, as well as some simple solutions.*

[BH73] "Operating System Principles"
Per Brinch Hansen, Prentice-Hall, 1973
Available: http://portal.acm.org/citation.cfm?id=540365
*One of the first books on operating systems; certainly ahead of its time. Introduced monitors as a
concurrency primitive.*

[BC05] "Understanding the Linux Kernel (Third Edition)"
Daniel P. Bovet and Marco Cesati
O'Reilly Media, November 2005
*The classic book on the Linux kernel. You should read it.*

[L+13] "A Study of Linux File System Evolution"
Lanyue Lu, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau, Shan Lu
FAST '13, San Jose, CA, February 2013
*Our paper that studies every patch to Linux file systems over nearly a decade. Lots of fun findings in
there; read it to see! The work was painful to do though; the poor graduate student, Lanyue Lu, had to
look through every single patch by hand in order to understand what they did.*

[MS98] "Nonblocking Algorithms and Preemption-safe Locking on Multiprogrammed Shared-
memory Multiprocessors"
M. Michael and M. Scott
Journal of Parallel and Distributed Computing, Vol. 51, No. 1, 1998
*Professor Scott and his students have been at the forefront of concurrent algorithms and data structures
for many years; check out his web page, numerous papers, or books to find out more.*

[MS04] "Concurrent Data Structures"
Mark Moir and Nir Shavit
In Handbook of Data Structures and Applications
(Editors D. Metha and S.Sahni)
Chapman and Hall/CRC Press, 2004
Available: www.cs.tau.ac.il/˜shanir/concurrent-data-structures.pdf
*A short but relatively comprehensive reference on concurrent data structures. Though it is missing
some of the latest works in the area (due to its age), it remains an incredibly useful reference.*

[MM00] "Solaris Internals: Core Kernel Architecture"
Jim Mauro and Richard McDougall
Prentice Hall, October 2000
*The Solaris book. You should also read this, if you want to learn in great detail about something other
than Linux.*

[S+11] "Making the Common Case the Only Case with Anticipatory Memory Allocation"
Swaminathan Sundararaman, Yupu Zhang, Sriram Subramanian,
Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau
FAST '11, San Jose, CA, February 2011
*Our work on removing possibly-failing calls to malloc from kernel code paths. The idea is to allocate all
potentially needed memory before doing any of the work, thus avoiding failure deep down in the storage
stack.*

# Condition Variables

Thus far we have developed the notion of a lock and seen how one can be properly built with the right combination of hardware and OS support. Unfortunately, locks are not the only primitives that are needed to build concurrent programs.

In particular, there are many cases where a thread wishes to check whether a **condition** is true before continuing its execution. For example, a parent thread might wish to check whether a child thread has completed before continuing (this is often called a `join()`); how should such a wait be implemented? Let's look at Figure 30.1.

```
1   void *child(void *arg) {
2       printf("child\n");
3       // XXX how to indicate we are done?
4       return NULL;
5   }
6
7   int main(int argc, char *argv[]) {
8       printf("parent: begin\n");
9       pthread_t c;
10      Pthread_create(&c, NULL, child, NULL); // create child
11      // XXX how to wait for child?
12      printf("parent: end\n");
13      return 0;
14  }
```

Figure 30.1: **A Parent Waiting For Its Child**

What we would like to see here is the following output:

```
parent: begin
child
parent: end
```

We could try using a shared variable, as you see in Figure 30.2. This solution will generally work, but it is hugely inefficient as the parent spins and wastes CPU time. What we would like here instead is some way to put the parent to sleep until the condition we are waiting for (e.g., the child is done executing) comes true.

1

```
1    volatile int done = 0;
2
3    void *child(void *arg) {
4        printf("child\n");
5        done = 1;
6        return NULL;
7    }
8
9    int main(int argc, char *argv[]) {
10       printf("parent: begin\n");
11       pthread_t c;
12       Pthread_create(&c, NULL, child, NULL); // create child
13       while (done == 0)
14           ; // spin
15       printf("parent: end\n");
16       return 0;
17   }
```

Figure 30.2: **Parent Waiting For Child: Spin-based Approach**

THE CRUX: HOW TO WAIT FOR A CONDITION
In multi-threaded programs, it is often useful for a thread to wait for some condition to become true before proceeding. The simple approach, of just spinning until the condition becomes true, is grossly inefficient and wastes CPU cycles, and in some cases, can be incorrect. Thus, how should a thread wait for a condition?

## 30.1 Definition and Routines

To wait for a condition to become true, a thread can make use of what is known as a **condition variable**. A **condition variable** is an explicit queue that threads can put themselves on when some state of execution (i.e., some **condition**) is not as desired (by **waiting** on the condition); some other thread, when it changes said state, can then wake one (or more) of those waiting threads and thus allow them to continue (by **signaling** on the condition). The idea goes back to Dijkstra's use of "private semaphores" [D68]; a similar idea was later named a "condition variable" by Hoare in his work on monitors [H74].

To declare such a condition variable, one simply writes something like this: pthread_cond_t c;, which declares c as a condition variable (note: proper initialization is also required). A condition variable has two operations associated with it: wait() and signal(). The wait() call is executed when a thread wishes to put itself to sleep; the signal() call is executed when a thread has changed something in the program and thus wants to wake a sleeping thread waiting on this condition. Specifically, the POSIX calls look like this:

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

```
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() {
6        Pthread_mutex_lock(&m);
7        done = 1;
8        Pthread_cond_signal(&c);
9        Pthread_mutex_unlock(&m);
10   }
11
12   void *child(void *arg) {
13       printf("child\n");
14       thr_exit();
15       return NULL;
16   }
17
18   void thr_join() {
19       Pthread_mutex_lock(&m);
20       while (done == 0)
21           Pthread_cond_wait(&c, &m);
22       Pthread_mutex_unlock(&m);
23   }
24
25   int main(int argc, char *argv[]) {
26       printf("parent: begin\n");
27       pthread_t p;
28       Pthread_create(&p, NULL, child, NULL);
29       thr_join();
30       printf("parent: end\n");
31       return 0;
32   }
```

Figure 30.3: **Parent Waiting For Child: Use A Condition Variable**

We will often refer to these as wait() and signal() for simplicity. One thing you might notice about the wait() call is that it also takes a mutex as a parameter; it assumes that this mutex is locked when wait() is called. The responsibility of wait() is to release the lock and put the calling thread to sleep (atomically); when the thread wakes up (after some other thread has signaled it), it must re-acquire the lock before returning to the caller. This complexity stems from the desire to prevent certain race conditions from occurring when a thread is trying to put itself to sleep. Let's take a look at the solution to the join problem (Figure 30.3) to understand this better.

There are two cases to consider. In the first, the parent creates the child thread but continues running itself (assume we have only a single processor) and thus immediately calls into thr_join() to wait for the child thread to complete. In this case, it will acquire the lock, check if the child is done (it is not), and put itself to sleep by calling wait() (hence releasing the lock). The child will eventually run, print the message "child", and call thr_exit() to wake the parent thread; this code just grabs the lock, sets the state variable done, and signals the parent thus waking it. Finally, the parent will run (returning from wait() with the lock held), unlock the lock, and print the final message "parent: end".

In the second case, the child runs immediately upon creation, sets
done to 1, calls signal to wake a sleeping thread (but there is none, so
it just returns), and is done. The parent then runs, calls thr_join(), sees
that done is 1, and thus does not wait and returns.

One last note: you might observe the parent uses a while loop instead
of just an if statement when deciding whether to wait on the condition.
While this does not seem strictly necessary per the logic of the program,
it is always a good idea, as we will see below.

To make sure you understand the importance of each piece of the
thr_exit() and thr_join() code, let's try a few alternate implemen-
tations. First, you might be wondering if we need the state variable done.
What if the code looked like the example below? Would this work?

```
1   void thr_exit() {
2       Pthread_mutex_lock(&m);
3       Pthread_cond_signal(&c);
4       Pthread_mutex_unlock(&m);
5   }
6
7   void thr_join() {
8       Pthread_mutex_lock(&m);
9       Pthread_cond_wait(&c, &m);
10      Pthread_mutex_unlock(&m);
11  }
```

Unfortunately this approach is broken. Imagine the case where the
child runs immediately and calls thr_exit() immediately; in this case,
the child will signal, but there is no thread asleep on the condition. When
the parent runs, it will simply call wait and be stuck; no thread will ever
wake it. From this example, you should appreciate the importance of
the state variable done; it records the value the threads are interested in
knowing. The sleeping, waking, and locking all are built around it.

Here is another poor implementation. In this example, we imagine
that one does not need to hold a lock in order to signal and wait. What
problem could occur here? Think about it!

```
1   void thr_exit() {
2       done = 1;
3       Pthread_cond_signal(&c);
4   }
5
6   void thr_join() {
7       if (done == 0)
8           Pthread_cond_wait(&c);
9   }
```

The issue here is a subtle race condition. Specifically, if the parent calls
thr_join() and then checks the value of done, it will see that it is 0 and
thus try to go to sleep. But just before it calls wait to go to sleep, the parent
is interrupted, and the child runs. The child changes the state variable
done to 1 and signals, but no thread is waiting and thus no thread is
woken. When the parent runs again, it sleeps forever, which is sad.

TIP: ALWAYS HOLD THE LOCK WHILE SIGNALING
Although it is strictly not necessary in all cases, it is likely simplest and
best to hold the lock while signaling when using condition variables. The
example above shows a case where you *must* hold the lock for correct-
ness; however, there are some other cases where it is likely OK not to, but
probably is something you should avoid. Thus, for simplicity, **hold the
lock when calling signal**.

The converse of this tip, i.e., hold the lock when calling wait, is not just
a tip, but rather mandated by the semantics of wait, because wait always
(a) assumes the lock is held when you call it, (b) releases said lock when
putting the caller to sleep, and (c) re-acquires the lock just before return-
ing. Thus, the generalization of this tip is correct: **hold the lock when
calling signal or wait**, and you will always be in good shape.

Hopefully, from this simple join example, you can see some of the ba-
sic requirements of using condition variables properly. To make sure you
understand, we now go through a more complicated example: the **pro-
ducer/consumer** or **bounded-buffer** problem.

## 30.2 The Producer/Consumer (Bounded Buffer) Problem

The next synchronization problem we will confront in this chapter is
known as the **producer/consumer** problem, or sometimes as the **bounded
buffer** problem, which was first posed by Dijkstra [D72]. Indeed, it was
this very producer/consumer problem that led Dijkstra and his co-workers
to invent the generalized semaphore (which can be used as either a lock
or a condition variable) [D01]; we will learn more about semaphores later.

Imagine one or more producer threads and one or more consumer
threads. Producers generate data items and place them in a buffer; con-
sumers grab said items from the buffer and consume them in some way.

This arrangement occurs in many real systems. For example, in a
multi-threaded web server, a producer puts HTTP requests into a work
queue (i.e., the bounded buffer); consumer threads take requests out of
this queue and process them.

A bounded buffer is also used when you pipe the output of one pro-
gram into another, e.g., `grep foo file.txt | wc -l`. This example
runs two processes concurrently; `grep` writes lines from `file.txt` with
the string `foo` in them to what it thinks is standard output; the UNIX
shell redirects the output to what is called a UNIX pipe (created by the
**pipe** system call). The other end of this pipe is connected to the stan-
dard input of the process `wc`, which simply counts the number of lines in
the input stream and prints out the result. Thus, the `grep` process is the
producer; the `wc` process is the consumer; between them is an in-kernel
bounded buffer; you, in this example, are just the happy user.

```
1    int buffer;
2    int count = 0; // initially, empty
3
4    void put(int value) {
5        assert(count == 0);
6        count = 1;
7        buffer = value;
8    }
9
10   int get() {
11       assert(count == 1);
12       count = 0;
13       return buffer;
14   }
```

Figure 30.4: **The Put And Get Routines (Version 1)**

```
1    void *producer(void *arg) {
2        int i;
3        int loops = (int) arg;
4        for (i = 0; i < loops; i++) {
5            put(i);
6        }
7    }
8
9    void *consumer(void *arg) {
10       int i;
11       while (1) {
12           int tmp = get();
13           printf("%d\n", tmp);
14       }
15   }
```

Figure 30.5: **Producer/Consumer Threads (Version 1)**

Because the bounded buffer is a shared resource, we must of course require synchronized access to it, lest[1] a race condition arise. To begin to understand this problem better, let us examine some actual code.

The first thing we need is a shared buffer, into which a producer puts data, and out of which a consumer takes data. Let's just use a single integer for simplicity (you can certainly imagine placing a pointer to a data structure into this slot instead), and the two inner routines to put a value into the shared buffer, and to get a value out of the buffer. See Figure 30.4 for details.

Pretty simple, no? The put() routine assumes the buffer is empty (and checks this with an assertion), and then simply puts a value into the shared buffer and marks it full by setting count to 1. The get() routine does the opposite, setting the buffer to empty (i.e., setting count to 0) and returning the value. Don't worry that this shared buffer has just a single entry; later, we'll generalize it to a queue that can hold multiple entries, which will be even more fun than it sounds.

Now we need to write some routines that know when it is OK to access the buffer to either put data into it or get data out of it. The conditions for

---

[1]This is where we drop some serious Old English on you, and the subjunctive form.

```
1    cond_t  cond;
2    mutex_t mutex;
3
4    void *producer(void *arg) {
5        int i;
6        for (i = 0; i < loops; i++) {
7            Pthread_mutex_lock(&mutex);          // p1
8            if (count == 1)                      // p2
9                Pthread_cond_wait(&cond, &mutex); // p3
10           put(i);                              // p4
11           Pthread_cond_signal(&cond);          // p5
12           Pthread_mutex_unlock(&mutex);        // p6
13       }
14   }
15
16   void *consumer(void *arg) {
17       int i;
18       for (i = 0; i < loops; i++) {
19           Pthread_mutex_lock(&mutex);          // c1
20           if (count == 0)                      // c2
21                Pthread_cond_wait(&cond, &mutex); // c3
22           int tmp = get();                     // c4
23           Pthread_cond_signal(&cond);          // c5
24           Pthread_mutex_unlock(&mutex);        // c6
25           printf("%d\n", tmp);
26       }
27   }
```

Figure 30.6: **Producer/Consumer: Single CV And If Statement**

this should be obvious: only put data into the buffer when count is zero (i.e., when the buffer is empty), and only get data from the buffer when count is one (i.e., when the buffer is full). If we write the synchronization code such that a producer puts data into a full buffer, or a consumer gets data from an empty one, we have done something wrong (and in this code, an assertion will fire).

This work is going to be done by two types of threads, one set of which we'll call the **producer** threads, and the other set which we'll call **consumer** threads. Figure 30.5 shows the code for a producer that puts an integer into the shared buffer loops number of times, and a consumer that gets the data out of that shared buffer (forever), each time printing out the data item it pulled from the shared buffer.

### A Broken Solution

Now imagine that we have just a single producer and a single consumer. Obviously the put() and get() routines have critical sections within them, as put() updates the buffer, and get() reads from it. However, putting a lock around the code doesn't work; we need something more. Not surprisingly, that something more is some condition variables. In this (broken) first try (Figure 30.6), we have a single condition variable cond and associated lock mutex.

| $T_{c1}$ | State | $T_{c2}$ | State | $T_p$ | State | **Count** | *Comment* |
|---|---|---|---|---|---|---|---|
| c1 | Running | | Ready | | Ready | 0 | |
| c2 | Running | | Ready | | Ready | 0 | |
| c3 | Sleep | | Ready | | Ready | 0 | Nothing to get |
| | Sleep | | Ready | p1 | Running | 0 | |
| | Sleep | | Ready | p2 | Running | 0 | |
| | Sleep | | Ready | p4 | Running | 1 | Buffer now full |
| | Ready | | Ready | p5 | Running | 1 | $T_{c1}$ awoken |
| | Ready | | Ready | p6 | Running | 1 | |
| | Ready | | Ready | p1 | Running | 1 | |
| | Ready | | Ready | p2 | Running | 1 | |
| | Ready | | Ready | p3 | Sleep | 1 | Buffer full; sleep |
| | Ready | c1 | Running | | Sleep | 1 | $T_{c2}$ sneaks in ... |
| | Ready | c2 | Running | | Sleep | 1 | |
| | Ready | c4 | Running | | Sleep | 0 | ... and grabs data |
| | Ready | c5 | Running | | Ready | 0 | $T_p$ awoken |
| | Ready | c6 | Running | | Ready | 0 | |
| c4 | Running | | Ready | | Ready | 0 | Oh oh! No data |

Figure 30.7: **Thread Trace: Broken Solution (Version 1)**

Let's examine the signaling logic between producers and consumers. When a producer wants to fill the buffer, it waits for it to be empty (p1–p3). The consumer has the exact same logic, but waits for a different condition: fullness (c1–c3).

With just a single producer and a single consumer, the code in Figure 30.6 works. However, if we have more than one of these threads (e.g., two consumers), the solution has two critical problems. What are they?

*... (pause here to think) ...*

Let's understand the first problem, which has to do with the `if` statement before the wait. Assume there are two consumers ($T_{c1}$ and $T_{c2}$) and one producer ($T_p$). First, a consumer ($T_{c1}$) runs; it acquires the lock (c1), checks if any buffers are ready for consumption (c2), and finding that none are, waits (c3) (which releases the lock).

Then the producer ($T_p$) runs. It acquires the lock (p1), checks if all buffers are full (p2), and finding that not to be the case, goes ahead and fills the buffer (p4). The producer then signals that a buffer has been filled (p5). Critically, this moves the first consumer ($T_{c1}$) from sleeping on a condition variable to the ready queue; $T_{c1}$ is now able to run (but not yet running). The producer then continues until realizing the buffer is full, at which point it sleeps (p6, p1–p3).

Here is where the problem occurs: another consumer ($T_{c2}$) sneaks in and consumes the one existing value in the buffer (c1, c2, c4, c5, c6, skipping the wait at c3 because the buffer is full). Now assume $T_{c1}$ runs; just before returning from the wait, it re-acquires the lock and then returns. It then calls `get()` (c4), but there are no buffers to consume! An assertion triggers, and the code has not functioned as desired. Clearly, we should have somehow prevented $T_{c1}$ from trying to consume because $T_{c2}$ snuck in and consumed the one value in the buffer that had been produced. Figure 30.7 shows the action each thread takes, as well as its scheduler state (Ready, Running, or Sleeping) over time.

```
1    cond_t  cond;
2    mutex_t mutex;
3
4    void *producer(void *arg) {
5        int i;
6        for (i = 0; i < loops; i++) {
7            Pthread_mutex_lock(&mutex);         // p1
8            while (count == 1)                  // p2
9                Pthread_cond_wait(&cond, &mutex); // p3
10           put(i);                             // p4
11           Pthread_cond_signal(&cond);         // p5
12           Pthread_mutex_unlock(&mutex);       // p6
13       }
14   }
15
16   void *consumer(void *arg) {
17       int i;
18       for (i = 0; i < loops; i++) {
19           Pthread_mutex_lock(&mutex);         // c1
20           while (count == 0)                  // c2
21               Pthread_cond_wait(&cond, &mutex); // c3
22           int tmp = get();                    // c4
23           Pthread_cond_signal(&cond);         // c5
24           Pthread_mutex_unlock(&mutex);       // c6
25           printf("%d\n", tmp);
26       }
27   }
```

Figure 30.8: **Producer/Consumer: Single CV And While**

The problem arises for a simple reason: after the producer woke $T_{c1}$, but *before* $T_{c1}$ ever ran, the state of the bounded buffer changed (thanks to $T_{c2}$). Signaling a thread only wakes them up; it is thus a *hint* that the state of the world has changed (in this case, that a value has been placed in the buffer), but there is no guarantee that when the woken thread runs, the state will *still* be as desired. This interpretation of what a signal means is often referred to as **Mesa semantics**, after the first research that built a condition variable in such a manner [LR80]; the contrast, referred to as **Hoare semantics**, is harder to build but provides a stronger guarantee that the woken thread will run immediately upon being woken [H74]. Virtually every system ever built employs Mesa semantics.

### Better, But Still Broken: While, Not If

Fortunately, this fix is easy (Figure 30.8): change the if to a while. Think about why this works; now consumer $T_{c1}$ wakes up and (with the lock held) immediately re-checks the state of the shared variable (c2). If the buffer is empty at that point, the consumer simply goes back to sleep (c3). The corollary if is also changed to a while in the producer (p2).

Thanks to Mesa semantics, a simple rule to remember with condition variables is to **always use while loops**. Sometimes you don't have to recheck the condition, but it is always safe to do so; just do it and be happy.

| $T_{c1}$ | State | $T_{c2}$ | State | $T_p$ | State | Count | Comment |
|---|---|---|---|---|---|---|---|
| c1 | Running | | Ready | | Ready | 0 | |
| c2 | Running | | Ready | | Ready | 0 | |
| c3 | Sleep | | Ready | | Ready | 0 | Nothing to get |
| | Sleep | c1 | Running | | Ready | 0 | |
| | Sleep | c2 | Running | | Ready | 0 | |
| | Sleep | c3 | Sleep | | Ready | 0 | Nothing to get |
| | Sleep | | Sleep | p1 | Running | 0 | |
| | Sleep | | Sleep | p2 | Running | 0 | |
| | Sleep | | Sleep | p4 | Running | 1 | Buffer now full |
| | Ready | | Sleep | p5 | Running | 1 | $T_{c1}$ awoken |
| | Ready | | Sleep | p6 | Running | 1 | |
| | Ready | | Sleep | p1 | Running | 1 | |
| | Ready | | Sleep | p2 | Running | 1 | |
| | Ready | | Sleep | p3 | Sleep | 1 | Must sleep (full) |
| c2 | Running | | Sleep | | Sleep | 1 | Recheck condition |
| c4 | Running | | Sleep | | Sleep | 0 | $T_{c1}$ grabs data |
| c5 | Running | | Ready | | Sleep | 0 | Oops! Woke $T_{c2}$ |
| c6 | Running | | Ready | | Sleep | 0 | |
| c1 | Running | | Ready | | Sleep | 0 | |
| c2 | Running | | Ready | | Sleep | 0 | |
| c3 | Sleep | | Ready | | Sleep | 0 | Nothing to get |
| | Sleep | c2 | Running | | Sleep | 0 | |
| | Sleep | c3 | Sleep | | Sleep | 0 | Everyone asleep... |

Figure 30.9: **Thread Trace: Broken Solution (Version 2)**

However, this code still has a bug, the second of two problems mentioned above. Can you see it? It has something to do with the fact that there is only one condition variable. Try to figure out what the problem is, before reading ahead. DO IT!

*... (another pause for you to think, or close your eyes for a bit) ...*

Let's confirm you figured it out correctly, or perhaps let's confirm that you are now awake and reading this part of the book. The problem occurs when two consumers run first ($T_{c1}$ and $T_{c2}$), and both go to sleep (c3). Then, a producer runs, put a value in the buffer, wakes one of the consumers (say $T_{c1}$), and goes back to sleep. Now we have one consumer ready to run ($T_{c1}$), and two threads sleeping on a condition ($T_{c2}$ and $T_p$). And we are about to cause a problem to occur: things are getting exciting!

The consumer $T_{c1}$ then wakes by returning from `wait()` (c3), re-checks the condition (c2), and finding the buffer full, consumes the value (c4). This consumer then, critically, signals on the condition (c5), waking one thread that is sleeping. However, which thread should it wake?

Because the consumer has emptied the buffer, it clearly should wake the producer. However, if it wakes the consumer $T_{c2}$ (which is definitely possible, depending on how the wait queue is managed), we have a problem. Specifically, the consumer $T_{c2}$ will wake up and find the buffer empty (c2), and go back to sleep (c3). The producer $T_p$, which has a value to put into the buffer, is left sleeping. The other consumer thread, $T_{c1}$, also goes back to sleep. All three threads are left sleeping, a clear bug; see Figure 30.9 for the brutal step-by-step of this terrible calamity.

Signaling is clearly needed, but must be more directed. A consumer should not wake other consumers, only producers, and vice-versa.

```
1   cond_t  empty, fill;
2   mutex_t mutex;
3
4   void *producer(void *arg) {
5       int i;
6       for (i = 0; i < loops; i++) {
7           Pthread_mutex_lock(&mutex);
8           while (count == 1)
9               Pthread_cond_wait(&empty, &mutex);
10          put(i);
11          Pthread_cond_signal(&fill);
12          Pthread_mutex_unlock(&mutex);
13      }
14  }
15
16  void *consumer(void *arg) {
17      int i;
18      for (i = 0; i < loops; i++) {
19          Pthread_mutex_lock(&mutex);
20          while (count == 0)
21              Pthread_cond_wait(&fill, &mutex);
22          int tmp = get();
23          Pthread_cond_signal(&empty);
24          Pthread_mutex_unlock(&mutex);
25          printf("%d\n", tmp);
26      }
27  }
```

Figure 30.10: **Producer/Consumer: Two CVs And While**

### The Single Buffer Producer/Consumer Solution

The solution here is once again a small one: use *two* condition variables, instead of one, in order to properly signal which type of thread should wake up when the state of the system changes. Figure 30.10 shows the resulting code.

In the code above, producer threads wait on the condition **empty**, and signals **fill**. Conversely, consumer threads wait on **fill** and signal **empty**. By doing so, the second problem above is avoided by design: a consumer can never accidentally wake a consumer, and a producer can never accidentally wake a producer.

### The Final Producer/Consumer Solution

We now have a working producer/consumer solution, albeit not a fully general one. The last change we make is to enable more concurrency and efficiency; specifically, we add more buffer slots, so that multiple values can be produced before sleeping, and similarly multiple values can be consumed before sleeping. With just a single producer and consumer, this approach is more efficient as it reduces context switches; with multiple producers or consumers (or both), it even allows concurrent producing or consuming to take place, thus increasing concurrency. Fortunately, it is a small change from our current solution.

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

Figure 30.11: **The Final Put And Get Routines**

```
1   cond_t empty, fill;
2   mutex_t mutex;
3
4   void *producer(void *arg) {
5       int i;
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;
18      for (i = 0; i < loops; i++) {
19          Pthread_mutex_lock(&mutex);            // c1
20          while (count == 0)                     // c2
21              Pthread_cond_wait(&fill, &mutex);  // c3
22          int tmp = get();                       // c4
23          Pthread_cond_signal(&empty);           // c5
24          Pthread_mutex_unlock(&mutex);          // c6
25          printf("%d\n", tmp);
26      }
27  }
```

Figure 30.12: **The Final Working Solution**

The first change for this final solution is within the buffer structure itself and the corresponding put() and get() (Figure 30.11). We also slightly change the conditions that producers and consumers check in order to determine whether to sleep or not. Figure 30.12 shows the final waiting and signaling logic. A producer only sleeps if all buffers are currently filled (p2); similarly, a consumer only sleeps if all buffers are currently empty (c2). And thus we solve the producer/consumer problem.

TIP: USE WHILE (NOT IF) FOR CONDITIONS

When checking for a condition in a multi-threaded program, using a `while` loop is always correct; using an `if` statement only might be, depending on the semantics of signaling. Thus, always use `while` and your code will behave as expected.

Using while loops around conditional checks also handles the case where **spurious wakeups** occur. In some thread packages, due to details of the implementation, it is possible that two threads get woken up though just a single signal has taken place [L11]. Spurious wakeups are further reason to re-check the condition a thread is waiting on.

## 30.3 Covering Conditions

We'll now look at one more example of how condition variables can be used. This code study is drawn from Lampson and Redell's paper on Pilot [LR80], the same group who first implemented the **Mesa semantics** described above (the language they used was Mesa, hence the name).

The problem they ran into is best shown via simple example, in this case in a simple multi-threaded memory allocation library. Figure 30.13 shows a code snippet which demonstrates the issue.

As you might see in the code, when a thread calls into the memory allocation code, it might have to wait in order for more memory to become free. Conversely, when a thread frees memory, it signals that more memory is free. However, our code above has a problem: which waiting thread (there can be more than one) should be woken up?

Consider the following scenario. Assume there are zero bytes free; thread $T_a$ calls `allocate(100)`, followed by thread $T_b$ which asks for less memory by calling `allocate(10)`. Both $T_a$ and $T_b$ thus wait on the condition and go to sleep; there aren't enough free bytes to satisfy either of these requests.

At that point, assume a third thread, $T_c$, calls `free(50)`. Unfortunately, when it calls signal to wake a waiting thread, it might not wake the correct waiting thread, $T_b$, which is waiting for only 10 bytes to be freed; $T_a$ should remain waiting, as not enough memory is yet free. Thus, the code in the figure does not work, as the thread waking other threads does not know which thread (or threads) to wake up.

The solution suggested by Lampson and Redell is straightforward: replace the `pthread_cond_signal()` call in the code above with a call to `pthread_cond_broadcast()`, which wakes up *all* waiting threads. By doing so, we guarantee that any threads that should be woken are. The downside, of course, can be a negative performance impact, as we might needlessly wake up many other waiting threads that shouldn't (yet) be awake. Those threads will simply wake up, re-check the condition, and then go immediately back to sleep.

```
1    // how many bytes of the heap are free?
2    int bytesLeft = MAX_HEAP_SIZE;
3
4    // need lock and condition too
5    cond_t  c;
6    mutex_t m;
7
8    void *
9    allocate(int size) {
10       Pthread_mutex_lock(&m);
11       while (bytesLeft < size)
12           Pthread_cond_wait(&c, &m);
13       void *ptr = ...; // get mem from heap
14       bytesLeft -= size;
15       Pthread_mutex_unlock(&m);
16       return ptr;
17   }
18
19   void free(void *ptr, int size) {
20       Pthread_mutex_lock(&m);
21       bytesLeft += size;
22       Pthread_cond_signal(&c); // whom to signal??
23       Pthread_mutex_unlock(&m);
24   }
```

Figure 30.13: **Covering Conditions: An Example**

Lampson and Redell call such a condition a **covering condition**, as it covers all the cases where a thread needs to wake up (conservatively); the cost, as we've discussed, is that too many threads might be woken. The astute reader might also have noticed we could have used this approach earlier (see the producer/consumer problem with only a single condition variable). However, in that case, a better solution was available to us, and thus we used it. In general, if you find that your program only works when you change your signals to broadcasts (but you don't think it should need to), you probably have a bug; fix it! But in cases like the memory allocator above, broadcast may be the most straightforward solution available.

## 30.4   Summary

We have seen the introduction of another important synchronization primitive beyond locks: condition variables. By allowing threads to sleep when some program state is not as desired, CVs enable us to neatly solve a number of important synchronization problems, including the famous (and still important) producer/consumer problem, as well as covering conditions. A more dramatic concluding sentence would go here, such as "He loved Big Brother" [O49].

# References

[D68] "Cooperating sequential processes"
Edsger W. Dijkstra, 1968
Available: http://www.cs.utexas.edu/users/EWD/ewd01xx/EWD123.PDF
*Another classic from Dijkstra; reading his early works on concurrency will teach you much of what you need to know.*

[D72] "Information Streams Sharing a Finite Buffer"
E.W. Dijkstra
Information Processing Letters 1: 179180, 1972
Available: http://www.cs.utexas.edu/users/EWD/ewd03xx/EWD329.PDF
*The famous paper that introduced the producer/consumer problem.*

[D01] "My recollections of operating system design"
E.W. Dijkstra
April, 2001
Available: http://www.cs.utexas.edu/users/EWD/ewd13xx/EWD1303.PDF
*A fascinating read for those of you interested in how the pioneers of our field came up with some very basic and fundamental concepts, including ideas like "interrupts" and even "a stack"!*

[H74] "Monitors: An Operating System Structuring Concept"
C.A.R. Hoare
Communications of the ACM, 17:10, pages 549–557, October 1974
*Hoare did a fair amount of theoretical work in concurrency. However, he is still probably most known for his work on Quicksort, the coolest sorting algorithm in the world, at least according to these authors.*

[L11] "Pthread_cond_signal Man Page"
Available: http://linux.die.net/man/3/pthread_cond_signal
March, 2011
*The Linux man page shows a nice simple example of why a thread might get a spurious wakeup, due to race conditions within the signal/wakeup code.*

[LR80] "Experience with Processes and Monitors in Mesa"
B.W. Lampson, D.R. Redell
Communications of the ACM. 23:2, pages 105-117, February 1980
*A terrific paper about how to actually implement signaling and condition variables in a real system, leading to the term "Mesa" semantics for what it means to be woken up; the older semantics, developed by Tony Hoare [H74], then became known as "Hoare" semantics, which is hard to say out loud in class with a straight face.*

[O49] "1984"
George Orwell, 1949, Secker and Warburg
*A little heavy-handed, but of course a must read. That said, we kind of gave away the ending by quoting the last sentence. Sorry! And if the government is reading this, let us just say that we think that the government is "double plus good". Hear that, our pals at the NSA?*

## Homework

This homework lets you explore some real code that uses locks and condition variables to implement various forms of the producer/consumer queue discussed in the chapter. You'll look at the real code, run it in various configurations, and use it to learn about what works and what doesn't, as well as other intricacies.

The different versions of the code correspond to different ways to "solve" the producer/consumer problem. Most are incorrect; one is correct. Read the chapter to learn more about what the producer/consumer problem is, and what the code generally does.

The first step is to download the code and type `make` to build all the variants. You should see four:

- `main-one-cv-while.c`: The producer/consumer problem solved with a single condition variable.

- `main-two-cvs-if.c`: Same but with two condition variables and using an `if` to check whether to sleep.

- `main-two-cvs-while.c`: Same but with two condition variables and `while` to check whether to sleep. **This is the correct version.**

- `main-two-cvs-while-extra-unlock.c`: Same but releasing the lock and then reacquiring it around the fill and get routines.

It's also useful to look at `pc-header.h` which contains common code for all of these different main programs, and the `Makefile` so as to build the code properly.

See the README for details on these programs.

## Questions

1. Our first question focuses on `main-two-cvs-while.c` (the working solution). First, study the code. Do you think you have an understanding of what should happen when you run the program?

2. Now run with one producer and one consumer, and have the producer produce a few values. Start with a buffer of size 1, and then increase it. How does the behavior of the code change when the buffer is larger? (or does it?) What would you predict num_full to be with different buffer sizes (e.g., -m 10) and different numbers of produced items (e.g., -l 100), when you change the consumer sleep string from default (no sleep) to -C 0,0,0,0,0,0,1?

3. If possible, run the code on different systems (e.g., Mac OS X and Linux). Do you see different behavior across these systems?

4. Let's look at some timings of different runs. How long do you think the following execution, with one producer, three consumers, a single-entry shared buffer, and each consumer pausing at point c3 for a second, will take?

```
prompt> ./main-one-cv-while -p 1 -c 3 -m 1 -C
        0,0,0,1,0,0,0:0,0,0,1,0,0,0:0,0,0,1,0,0,0 -l 10 -v -t
```

5. Now change the size of the shared buffer to 3 (-m 3). Will this make any difference in the total time?

6. Now change the location of the sleep to c6 (this models a consumer taking something off the queue and then doing something with it for a while), again using a single-entry buffer. What time do you predict in this case?

```
prompt> ./main-one-cv-while -p 1 -c 3 -m 1 -C
        0,0,0,0,0,0,1:0,0,0,0,0,0,1:0,0,0,0,0,0,1 -l 10 -v -t
```

7. Finally, change the buffer size to 3 again (-m 3). What time do you predict now?

8. Now let's look at main-one-cv-while.c. Can you configure a sleep string, assuming a single producer, one consumer, and a buffer of size 1, to cause a problem with this code?

9. Now change the number of consumers to two. Can you construct sleep strings for the producer and the consumers so as to cause a problem in the code?

10. Now examine main-two-cvs-if.c. Can you cause a problem to happen in this code? Again consider the case where there is only one consumer, and then the case where there is more than one.

11. Finally, examine main-two-cvs-while-extra-unlock.c. What problem arises when you release the lock before doing a put or a get? Can you reliably cause such a problem to happen, given the sleep strings? What bad thing can happen?

# Semaphores

As we know now, one needs both locks and condition variables to solve a broad range of relevant and interesting concurrency problems. One of the first people to realize this years ago was **Edsger Dijkstra** (though it is hard to know the exact history [GR92]), known among other things for his famous "shortest paths" algorithm in graph theory [D59], an early polemic on structured programming entitled "Goto Statements Considered Harmful" [D68a] (what a great title!), and, in the case we will study here, the introduction of a synchronization primitive called the **semaphore** [D68b,D72]. Indeed, Dijkstra and colleagues invented the semaphore as a single primitive for all things related to synchronization; as you will see, one can use semaphores as both locks and condition variables.

THE CRUX: HOW TO USE SEMAPHORES
How can we use semaphores instead of locks and condition variables? What is the definition of a semaphore? What is a binary semaphore? Is it straightforward to build a semaphore out of locks and condition variables? To build locks and condition variables out of semaphores?

## 31.1 Semaphores: A Definition

A semaphore is an object with an integer value that we can manipulate with two routines; in the POSIX standard, these routines are `sem_wait()` and `sem_post()`[1]. Because the initial value of the semaphore determines its behavior, before calling any other routine to interact with the semaphore, we must first initialize it to some value, as the code in Figure 31.1 does.

---

[1]Historically, `sem_wait()` was called P() by Dijkstra and `sem_post()` called V(). P() comes from "prolaag", a contraction of "probeer" (Dutch for "try") and "verlaag" ("decrease"); V() comes from the Dutch word "verhoog" which means "increase" (thanks to Mart Oskamp for this information). Sometimes, people call them down and up. Use the Dutch versions to impress your friends, or confuse them, or both.

```
1    #include <semaphore.h>
2    sem_t s;
3    sem_init(&s, 0, 1);
```

Figure 31.1: **Initializing A Semaphore**

In the figure, we declare a semaphore s and initialize it to the value 1 by passing 1 in as the third argument. The second argument to `sem_init()` will be set to 0 in all of the examples we'll see; this indicates that the semaphore is shared between threads in the same process. See the man page for details on other usages of semaphores (namely, how they can be used to synchronize access across *different* processes), which require a different value for that second argument.

After a semaphore is initialized, we can call one of two functions to interact with it, `sem_wait()` or `sem_post()`. The behavior of these two functions is seen in Figure 31.2.

For now, we are not concerned with the implementation of these routines, which clearly requires some care; with multiple threads calling into `sem_wait()` and `sem_post()`, there is the obvious need for managing these critical sections. We will now focus on how to *use* these primitives; later we may discuss how they are built.

We should discuss a few salient aspects of the interfaces here. First, we can see that `sem_wait()` will either return right away (because the value of the semaphore was one or higher when we called `sem_wait()`), or it will cause the caller to suspend execution waiting for a subsequent post. Of course, multiple calling threads may call into `sem_wait()`, and thus all be queued waiting to be woken.

Second, we can see that `sem_post()` does not wait for some particular condition to hold like `sem_wait()` does. Rather, it simply increments the value of the semaphore and then, if there is a thread waiting to be woken, wakes one of them up.

Third, the value of the semaphore, when negative, is equal to the number of waiting threads [D68b]. Though the value generally isn't seen by users of the semaphores, this invariant is worth knowing and perhaps can help you remember how a semaphore functions.

Don't worry (yet) about the seeming race conditions possible within the semaphore; assume that the actions they make are performed atomically. We will soon use locks and condition variables to do just this.

```
1    int sem_wait(sem_t *s) {
2        decrement the value of semaphore s by one
3        wait if value of semaphore s is negative
4    }
5
6    int sem_post(sem_t *s) {
7        increment the value of semaphore s by one
8        if there are one or more threads waiting, wake one
9    }
```

Figure 31.2: **Semaphore: Definitions Of Wait And Post**

```
1   sem_t m;
2   sem_init(&m, 0, X); // initialize semaphore to X; what should X be?
3
4   sem_wait(&m);
5   // critical section here
6   sem_post(&m);
```

Figure 31.3: **A Binary Semaphore (That Is, A Lock)**

## 31.2 Binary Semaphores (Locks)

We are now ready to use a semaphore. Our first use will be one with which we are already familiar: using a semaphore as a lock. See Figure 31.3 for a code snippet; therein, you'll see that we simply surround the critical section of interest with a sem_wait()/sem_post() pair. Critical to making this work, though, is the initial value of the semaphore m (initialized to X in the figure). What should X be?

*... (Try thinking about it before going on) ...*

Looking back at definition of the sem_wait() and sem_post() routines above, we can see that the initial value should be 1.

To make this clear, let's imagine a scenario with two threads. The first thread (Thread 0) calls sem_wait(); it will first decrement the value of the semaphore, changing it to 0. Then, it will wait only if the value is *not* greater than or equal to 0. Because the value is 0, sem_wait() will simply return and the calling thread will continue; Thread 0 is now free to enter the critical section. If no other thread tries to acquire the lock while Thread 0 is inside the critical section, when it calls sem_post(), it will simply restore the value of the semaphore to 1 (and not wake a waiting thread, because there are none). Figure 31.4 shows a trace of this scenario.

A more interesting case arises when Thread 0 "holds the lock" (i.e., it has called sem_wait() but not yet called sem_post()), and another thread (Thread 1) tries to enter the critical section by calling sem_wait(). In this case, Thread 1 will decrement the value of the semaphore to -1, and thus wait (putting itself to sleep and relinquishing the processor). When Thread 0 runs again, it will eventually call sem_post(), incrementing the value of the semaphore back to zero, and then wake the waiting thread (Thread 1), which will then be able to acquire the lock for itself. When Thread 1 finishes, it will again increment the value of the semaphore, restoring it to 1 again.

| Value of Semaphore | Thread 0 | Thread 1 |
|---|---|---|
| 1 | | |
| 1 | call sem_wait() | |
| 0 | sem_wait() returns | |
| 0 | (crit sect) | |
| 0 | call sem_post() | |
| 1 | sem_post() returns | |

Figure 31.4: **Thread Trace: Single Thread Using A Semaphore**

| Value | Thread 0 | State | Thread 1 | State |
|-------|----------|-------|----------|-------|
| 1 | | Running | | Ready |
| 1 | call `sem_wait()` | Running | | Ready |
| 0 | `sem_wait()` returns | Running | | Ready |
| 0 | `(crit sect: begin)` | Running | | Ready |
| 0 | *Interrupt; Switch→T1* | Ready | | Running |
| 0 | | Ready | call `sem_wait()` | Running |
| -1 | | Ready | `decrement sem` | Running |
| -1 | | Ready | `(sem<0)→sleep` | Sleeping |
| -1 | | Running | *Switch→T0* | Sleeping |
| -1 | `(crit sect: end)` | Running | | Sleeping |
| -1 | call `sem_post()` | Running | | Sleeping |
| 0 | `increment sem` | Running | | Sleeping |
| 0 | `wake(T1)` | Running | | Ready |
| 0 | `sem_post()` returns | Running | | Ready |
| 0 | *Interrupt; Switch→T1* | Ready | | Running |
| 0 | | Ready | `sem_wait()` returns | Running |
| 0 | | Ready | `(crit sect)` | Running |
| 0 | | Ready | call `sem_post()` | Running |
| 1 | | Ready | `sem_post()` returns | Running |

Figure 31.5: **Thread Trace: Two Threads Using A Semaphore**

Figure 31.5 shows a trace of this example. In addition to thread actions, the figure shows the **scheduler state** of each thread: Running, Ready (i.e., runnable but not running), and Sleeping. Note in particular that Thread 1 goes into the sleeping state when it tries to acquire the already-held lock; only when Thread 0 runs again can Thread 1 be awoken and potentially run again.

If you want to work through your own example, try a scenario where multiple threads queue up waiting for a lock. What would the value of the semaphore be during such a trace?

Thus we are able to use semaphores as locks. Because locks only have two states (held and not held), we sometimes call a semaphore used as a lock a **binary semaphore**. Note that if you are using a semaphore only in this binary fashion, it could be implemented in a simpler manner than the generalized semaphores we present here.

## 31.3 Semaphores As Condition Variables

Semaphores are also useful when a thread wants to halt its progress waiting for a condition to become true. For example, a thread may wish to wait for a list to become non-empty, so it can delete an element from it. In this pattern of usage, we often find a thread *waiting* for something to happen, and a different thread making that something happen and then *signaling* that it has happened, thus waking the waiting thread. Because the waiting thread (or threads) is waiting for some **condition** in the program to change, we are using the semaphore as a **condition variable**.

```
1    sem_t s;
2
3    void *
4    child(void *arg) {
5        printf("child\n");
6        sem_post(&s); // signal here: child is done
7        return NULL;
8    }
9
10   int
11   main(int argc, char *argv[]) {
12       sem_init(&s, 0, X); // what should X be?
13       printf("parent: begin\n");
14       pthread_t c;
15       Pthread_create(c, NULL, child, NULL);
16       sem_wait(&s); // wait here for child
17       printf("parent: end\n");
18       return 0;
19   }
```

Figure 31.6: **A Parent Waiting For Its Child**

A simple example is as follows. Imagine a thread creates another thread and then wants to wait for it to complete its execution (Figure 31.6). When this program runs, we would like to see the following:

```
parent: begin
child
parent: end
```

The question, then, is how to use a semaphore to achieve this effect; as it turns out, the answer is relatively easy to understand. As you can see in the code, the parent simply calls sem_wait() and the child sem_post() to wait for the condition of the child finishing its execution to become true. However, this raises the question: what should the initial value of this semaphore be?

*(Again, think about it here, instead of reading ahead)*

The answer, of course, is that the value of the semaphore should be set to is 0. There are two cases to consider. First, let us assume that the parent creates the child but the child has not run yet (i.e., it is sitting in a ready queue but not running). In this case (Figure 31.7, page 6), the parent will call sem_wait() before the child has called sem_post(); we'd like the parent to wait for the child to run. The only way this will happen is if the value of the semaphore is not greater than 0; hence, 0 is the initial value. The parent runs, decrements the semaphore (to -1), then waits (sleeping). When the child finally runs, it will call sem_post(), increment the value of the semaphore to 0, and wake the parent, which will then return from sem_wait() and finish the program.

The second case (Figure 31.8, page 6) occurs when the child runs to completion before the parent gets a chance to call sem_wait(). In this case, the child will first call sem_post(), thus incrementing the value of the semaphore from 0 to 1. When the parent then gets a chance to run, it will call sem_wait() and find the value of the semaphore to be 1; the parent will thus decrement the value (to 0) and return from sem_wait() without waiting, also achieving the desired effect.

| Value | Parent | State | Child | State |
|---|---|---|---|---|
| 0 | create(Child) | Running | *(Child exists; is runnable)* | Ready |
| 0 | call sem_wait() | Running | | Ready |
| -1 | decrement sem | Running | | Ready |
| -1 | (sem<0)→sleep | Sleeping | | Ready |
| -1 | *Switch→Child* | Sleeping | child runs | Running |
| -1 | | Sleeping | call sem_post() | Running |
| 0 | | Sleeping | increment sem | Running |
| 0 | | Ready | wake(Parent) | Running |
| 0 | | Ready | sem_post() returns | Running |
| 0 | | Ready | *Interrupt; Switch→Parent* | Ready |
| 0 | sem_wait() returns | Running | | Ready |

Figure 31.7: **Thread Trace: Parent Waiting For Child (Case 1)**

| Value | Parent | State | Child | State |
|---|---|---|---|---|
| 0 | create(Child) | Running | *(Child exists; is runnable)* | Ready |
| 0 | *Interrupt; Switch→Child* | Ready | child runs | Running |
| 0 | | Ready | call sem_post() | Running |
| 1 | | Ready | increment sem | Running |
| 1 | | Ready | wake(nobody) | Running |
| 1 | | Ready | sem_post() returns | Running |
| 1 | parent runs | Running | *Interrupt; Switch→Parent* | Ready |
| 1 | call sem_wait() | Running | | Ready |
| 0 | decrement sem | Running | | Ready |
| 0 | (sem≥0)→awake | Running | | Ready |
| 0 | sem_wait() returns | Running | | Ready |

Figure 31.8: **Thread Trace: Parent Waiting For Child (Case 2)**

## 31.4 The Producer/Consumer (Bounded Buffer) Problem

The next problem we will confront in this chapter is known as the **producer/consumer** problem, or sometimes as the **bounded buffer** problem [D72]. This problem is described in detail in the previous chapter on condition variables; see there for details.

### First Attempt

Our first attempt at solving the problem introduces two semaphores, empty and full, which the threads will use to indicate when a buffer entry has been emptied or filled, respectively. The code for the put and get routines is in Figure 31.9, and our attempt at solving the producer and consumer problem is in Figure 31.10.

In this example, the producer first waits for a buffer to become empty in order to put data into it, and the consumer similarly waits for a buffer to become filled before using it. Let us first imagine that MAX=1 (there is only one buffer in the array), and see if this works.

Imagine again there are two threads, a producer and a consumer. Let us examine a specific scenario on a single CPU. Assume the consumer gets to run first. Thus, the consumer will hit line c1 in the figure above, calling sem_wait(&full). Because full was initialized to the value 0,

```
1   int buffer[MAX];
2   int fill = 0;
3   int use  = 0;
4
5   void put(int value) {
6       buffer[fill] = value;    // line f1
7       fill = (fill + 1) % MAX; // line f2
8   }
9
10  int get() {
11      int tmp = buffer[use];   // line g1
12      use = (use + 1) % MAX;   // line g2
13      return tmp;
14  }
```

Figure 31.9: **The Put And Get Routines**

```
1   sem_t empty;
2   sem_t full;
3
4   void *producer(void *arg) {
5       int i;
6       for (i = 0; i < loops; i++) {
7           sem_wait(&empty);        // line P1
8           put(i);                  // line P2
9           sem_post(&full);         // line P3
10      }
11  }
12
13  void *consumer(void *arg) {
14      int i, tmp = 0;
15      while (tmp != -1) {
16          sem_wait(&full);         // line C1
17          tmp = get();             // line C2
18          sem_post(&empty);        // line C3
19          printf("%d\n", tmp);
20      }
21  }
22
23  int main(int argc, char *argv[]) {
24      // ...
25      sem_init(&empty, 0, MAX); // MAX buffers are empty to begin with...
26      sem_init(&full, 0, 0);    // ... and 0 are full
27      // ...
28  }
```

Figure 31.10: **Adding The Full And Empty Conditions**

the call will decrement full (to -1), block the consumer, and wait for another thread to call sem_post() on full, as desired.

Assume the producer then runs. It will hit line P1, thus calling the sem_wait(&empty) routine. Unlike the consumer, the producer will continue through this line, because empty was initialized to the value MAX (in this case, 1). Thus, empty will be decremented to 0 and the producer will put a data value into the first entry of buffer (line P2). The producer will then continue on to P3 and call sem_post(&full), changing the value of the full semaphore from -1 to 0 and waking the consumer (e.g., move it from blocked to ready).

In this case, one of two things could happen. If the producer continues to run, it will loop around and hit line P1 again. This time, however, it would block, as the empty semaphore's value is 0. If the producer instead was interrupted and the consumer began to run, it would call `sem_wait(&full)` (line c1) and find that the buffer was indeed full and thus consume it. In either case, we achieve the desired behavior.

You can try this same example with more threads (e.g., multiple producers, and multiple consumers). It should still work.

Let us now imagine that MAX is greater than 1 (say MAX = 10). For this example, let us assume that there are multiple producers and multiple consumers. We now have a problem: a race condition. Do you see where it occurs? (take some time and look for it) If you can't see it, here's a hint: look more closely at the put() and get() code.

OK, let's understand the issue. Imagine two producers (Pa and Pb) both calling into put() at roughly the same time. Assume producer Pa gets to run first, and just starts to fill the first buffer entry (fill = 0 at line f1). Before Pa gets a chance to increment the fill counter to 1, it is interrupted. Producer Pb starts to run, and at line f1 it also puts its data into the 0th element of buffer, which means that the old data there is overwritten! This is a no-no; we don't want any data from the producer to be lost.

### A Solution: Adding Mutual Exclusion

As you can see, what we've forgotten here is *mutual exclusion*. The filling of a buffer and incrementing of the index into the buffer is a critical section, and thus must be guarded carefully. So let's use our friend the binary semaphore and add some locks. Figure 31.11 shows our attempt.

Now we've added some locks around the entire put()/get() parts of the code, as indicated by the NEW LINE comments. That seems like the right idea, but it also doesn't work. Why? Deadlock. Why does deadlock occur? Take a moment to consider it; try to find a case where deadlock arises. What sequence of steps must happen for the program to deadlock?

### Avoiding Deadlock

OK, now that you figured it out, here is the answer. Imagine two threads, one producer and one consumer. The consumer gets to run first. It acquires the mutex (line c0), and then calls `sem_wait()` on the full semaphore (line c1); because there is no data yet, this call causes the consumer to block and thus yield the CPU; importantly, though, the consumer still holds the lock.

A producer then runs. It has data to produce and if it were able to run, it would be able to wake the consumer thread and all would be good. Unfortunately, the first thing it does is call `sem_wait()` on the binary mutex semaphore (line p0). The lock is already held. Hence, the producer is now stuck waiting too.

```
1    sem_t empty;
2    sem_t full;
3    sem_t mutex;
4
5    void *producer(void *arg) {
6        int i;
7        for (i = 0; i < loops; i++) {
8            sem_wait(&mutex);            // line p0 (NEW LINE)
9            sem_wait(&empty);            // line p1
10           put(i);                      // line p2
11           sem_post(&full);             // line p3
12           sem_post(&mutex);            // line p4 (NEW LINE)
13       }
14   }
15
16   void *consumer(void *arg) {
17       int i;
18       for (i = 0; i < loops; i++) {
19           sem_wait(&mutex);            // line c0 (NEW LINE)
20           sem_wait(&full);             // line c1
21           int tmp = get();             // line c2
22           sem_post(&empty);            // line c3
23           sem_post(&mutex);            // line c4 (NEW LINE)
24           printf("%d\n", tmp);
25       }
26   }
27
28   int main(int argc, char *argv[]) {
29       // ...
30       sem_init(&empty, 0, MAX); // MAX buffers are empty to begin with...
31       sem_init(&full, 0, 0);    // ... and 0 are full
32       sem_init(&mutex, 0, 1);   // mutex=1 because it is a lock (NEW LINE)
33       // ...
34   }
```

Figure 31.11: **Adding Mutual Exclusion (Incorrectly)**

There is a simple cycle here. The consumer *holds* the mutex and is *waiting* for the someone to signal full. The producer could *signal* full but is *waiting* for the mutex. Thus, the producer and consumer are each stuck waiting for each other: a classic deadlock.

### Finally, A Working Solution

To solve this problem, we simply must reduce the scope of the lock. Figure 31.12 shows the final working solution. As you can see, we simply move the mutex acquire and release to be just around the critical section; the full and empty wait and signal code is left outside. The result is a simple and working bounded buffer, a commonly-used pattern in multi-threaded programs. Understand it now; use it later. You will thank us for years to come. Or at least, you will thank us when the same question is asked on the final exam.

```
1    sem_t empty;
2    sem_t full;
3    sem_t mutex;
4
5    void *producer(void *arg) {
6        int i;
7        for (i = 0; i < loops; i++) {
8            sem_wait(&empty);          // line p1
9            sem_wait(&mutex);          // line p1.5 (MOVED MUTEX HERE...)
10           put(i);                    // line p2
11           sem_post(&mutex);          // line p2.5 (... AND HERE)
12           sem_post(&full);           // line p3
13       }
14   }
15
16   void *consumer(void *arg) {
17       int i;
18       for (i = 0; i < loops; i++) {
19           sem_wait(&full);           // line c1
20           sem_wait(&mutex);          // line c1.5 (MOVED MUTEX HERE...)
21           int tmp = get();           // line c2
22           sem_post(&mutex);          // line c2.5 (... AND HERE)
23           sem_post(&empty);          // line c3
24           printf("%d\n", tmp);
25       }
26   }
27
28   int main(int argc, char *argv[]) {
29       // ...
30       sem_init(&empty, 0, MAX); // MAX buffers are empty to begin with...
31       sem_init(&full, 0, 0);    // ... and 0 are full
32       sem_init(&mutex, 0, 1);   // mutex=1 because it is a lock
33       // ...
34   }
```

Figure 31.12: **Adding Mutual Exclusion (Correctly)**

## 31.5 Reader-Writer Locks

Another classic problem stems from the desire for a more flexible locking primitive that admits that different data structure accesses might require different kinds of locking. For example, imagine a number of concurrent list operations, including inserts and simple lookups. While inserts change the state of the list (and thus a traditional critical section makes sense), lookups simply *read* the data structure; as long as we can guarantee that no insert is on-going, we can allow many lookups to proceed concurrently. The special type of lock we will now develop to support this type of operation is known as a **reader-writer lock** [CHP71]. The code for such a lock is available in Figure 31.13.

The code is pretty simple. If some thread wants to update the data structure in question, it should call the new pair of synchronization operations: rwlock_acquire_writelock(), to acquire a write lock, and rwlock_release_writelock(), to release it. Internally, these simply use the writelock semaphore to ensure that only a single writer can ac-

```
1   typedef struct _rwlock_t {
2     sem_t lock;      // binary semaphore (basic lock)
3     sem_t writelock; // used to allow ONE writer or MANY readers
4     int   readers;   // count of readers reading in critical section
5   } rwlock_t;
6
7   void rwlock_init(rwlock_t *rw) {
8     rw->readers = 0;
9     sem_init(&rw->lock, 0, 1);
10    sem_init(&rw->writelock, 0, 1);
11  }
12
13  void rwlock_acquire_readlock(rwlock_t *rw) {
14    sem_wait(&rw->lock);
15    rw->readers++;
16    if (rw->readers == 1)
17      sem_wait(&rw->writelock); // first reader acquires writelock
18    sem_post(&rw->lock);
19  }
20
21  void rwlock_release_readlock(rwlock_t *rw) {
22    sem_wait(&rw->lock);
23    rw->readers--;
24    if (rw->readers == 0)
25      sem_post(&rw->writelock); // last reader releases writelock
26    sem_post(&rw->lock);
27  }
28
29  void rwlock_acquire_writelock(rwlock_t *rw) {
30    sem_wait(&rw->writelock);
31  }
32
33  void rwlock_release_writelock(rwlock_t *rw) {
34    sem_post(&rw->writelock);
35  }
```

Figure 31.13: **A Simple Reader-Writer Lock**

quire the lock and thus enter the critical section to update the data structure in question.

More interesting is the pair of routines to acquire and release read locks. When acquiring a read lock, the reader first acquires lock and then increments the readers variable to track how many readers are currently inside the data structure. The important step then taken within rwlock_acquire_readlock() occurs when the first reader acquires the lock; in that case, the reader also acquires the write lock by calling sem_wait() on the writelock semaphore, and then finally releasing the lock by calling sem_post().

Thus, once a reader has acquired a read lock, more readers will be allowed to acquire the read lock too; however, any thread that wishes to acquire the write lock will have to wait until *all* readers are finished; the last one to exit the critical section calls sem_post() on "writelock" and thus enables a waiting writer to acquire the lock.

This approach works (as desired), but does have some negatives, espe-

TIP: SIMPLE AND DUMB CAN BE BETTER (HILL'S LAW)
You should never underestimate the notion that the simple and dumb
approach can be the best one. With locking, sometimes a simple spin lock
works best, because it is easy to implement and fast. Although something
like reader/writer locks sounds cool, they are complex, and complex can
mean slow. Thus, always try the simple and dumb approach first.

This idea, of appealing to simplicity, is found in many places. One early
source is Mark Hill's dissertation [H87], which studied how to design
caches for CPUs. Hill found that simple direct-mapped caches worked
better than fancy set-associative designs (one reason is that in caching,
simpler designs enable faster lookups). As Hill succinctly summarized
his work: "Big and dumb is better." And thus we call this similar advice
**Hill's Law**.

cially when it comes to fairness. In particular, it would be relatively easy
for readers to starve writers. More sophisticated solutions to this prob-
lem exist; perhaps you can think of a better implementation? Hint: think
about what you would need to do to prevent more readers from entering
the lock once a writer is waiting.

Finally, it should be noted that reader-writer locks should be used
with some caution. They often add more overhead (especially with more
sophisticated implementations), and thus do not end up speeding up
performance as compared to just using simple and fast locking primi-
tives [CB08]. Either way, they showcase once again how we can use
semaphores in an interesting and useful way.

## 31.6 The Dining Philosophers

One of the most famous concurrency problems posed, and solved, by
Dijkstra, is known as the **dining philosopher's problem** [D71]. The prob-
lem is famous because it is fun and somewhat intellectually interesting;
however, its practical utility is low. However, its fame forces its inclu-
sion here; indeed, you might be asked about it on some interview, and
you'd really hate your OS professor if you miss that question and don't
get the job. Conversely, if you get the job, please feel free to send your OS
professor a nice note, or some stock options.

The basic setup for the problem is this (as shown in Figure 31.14): as-
sume there are five "philosophers" sitting around a table. Between each
pair of philosophers is a single fork (and thus, five total). The philoso-
phers each have times where they think, and don't need any forks, and
times where they eat. In order to eat, a philosopher needs two forks, both
the one on their left and the one on their right. The contention for these
forks, and the synchronization problems that ensue, are what makes this
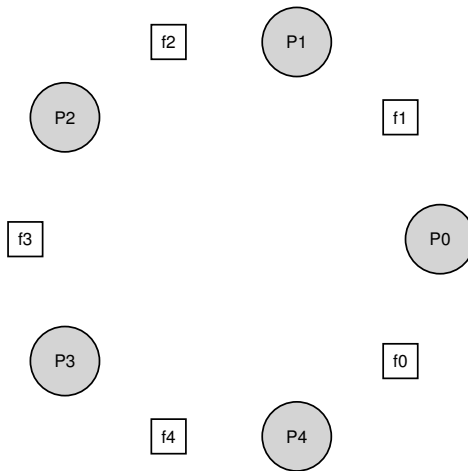a problem we study in concurrent programming.

Figure 31.14: **The Dining Philosophers**

Here is the basic loop of each philosopher:

```
while (1) {
  think();
  getforks();
  eat();
  putforks();
}
```

The key challenge, then, is to write the routines getforks() and putforks() such that there is no deadlock, no philosopher starves and never gets to eat, and concurrency is high (i.e., as many philosophers can eat at the same time as possible).

Following Downey's solutions [D08], we'll use a few helper functions to get us towards a solution. They are:

```
int left(int p)  { return p; }
int right(int p) { return (p + 1) % 5; }
```

When philosopher p wishes to refer to the fork on their left, they simply call left(p). Similarly, the fork on the right of a philosopher p is referred to by calling right(p); the modulo operator therein handles the one case where the last philosopher (p=4) tries to grab the fork on their right, which is fork 0.

We'll also need some semaphores to solve this problem. Let us assume we have five, one for each fork: sem_t forks[5].

```
1   void getforks() {
2     sem_wait(forks[left(p)]);
3     sem_wait(forks[right(p)]);
4   }
5
6   void putforks() {
7     sem_post(forks[left(p)]);
8     sem_post(forks[right(p)]);
9   }
```

Figure 31.15: **The `getforks()` And `putforks()` Routines**

### Broken Solution

We attempt our first solution to the problem. Assume we initialize each semaphore (in the `forks` array) to a value of 1. Assume also that each philosopher knows its own number (`p`). We can thus write the `getforks()` and `putforks()` routine as shown in Figure 31.15.

The intuition behind this (broken) solution is as follows. To acquire the forks, we simply grab a "lock" on each one: first the one on the left, and then the one on the right. When we are done eating, we release them. Simple, no? Unfortunately, in this case, simple means broken. Can you see the problem that arises? Think about it.

The problem is **deadlock**. If each philosopher happens to grab the fork on their left before any philosopher can grab the fork on their right, each will be stuck holding one fork and waiting for another, forever. Specifically, philosopher 0 grabs fork 0, philosopher 1 grabs fork 1, philosopher 2 grabs fork 2, philosopher 3 grabs fork 3, and philosopher 4 grabs fork 4; all the forks are acquired, and all the philosophers are stuck waiting for a fork that another philosopher possesses. We'll study deadlock in more detail soon; for now, it is safe to say that this is not a working solution.

### A Solution: Breaking The Dependency

The simplest way to attack this problem is to change how forks are acquired by at least one of the philosophers; indeed, this is how Dijkstra himself solved the problem. Specifically, let's assume that philosopher 4 (the highest numbered one) acquires the forks in a *different* order. The code to do so is as follows:

```
1   void getforks() {
2     if (p == 4) {
3       sem_wait(forks[right(p)]);
4       sem_wait(forks[left(p)]);
5     } else {
6       sem_wait(forks[left(p)]);
7       sem_wait(forks[right(p)]);
8     }
9   }
```

Because the last philosopher tries to grab right before left, there is no situation where each philosopher grabs one fork and is stuck waiting for another; the cycle of waiting is broken. Think through the ramifications of this solution, and convince yourself that it works.

```
1   typedef struct __Zem_t {
2       int value;
3       pthread_cond_t cond;
4       pthread_mutex_t lock;
5   } Zem_t;
6
7   // only one thread can call this
8   void Zem_init(Zem_t *s, int value) {
9       s->value = value;
10      Cond_init(&s->cond);
11      Mutex_init(&s->lock);
12  }
13
14  void Zem_wait(Zem_t *s) {
15      Mutex_lock(&s->lock);
16      while (s->value <= 0)
17          Cond_wait(&s->cond, &s->lock);
18      s->value--;
19      Mutex_unlock(&s->lock);
20  }
21
22  void Zem_post(Zem_t *s) {
23      Mutex_lock(&s->lock);
24      s->value++;
25      Cond_signal(&s->cond);
26      Mutex_unlock(&s->lock);
27  }
```

Figure 31.16: **Implementing Zemaphores With Locks And CVs**

There are other "famous" problems like this one, e.g., the **cigarette smoker's problem** or the **sleeping barber problem**. Most of them are just excuses to think about concurrency; some of them have fascinating names. Look them up if you are interested in learning more, or just getting more practice thinking in a concurrent manner [D08].

## 31.7  How To Implement Semaphores

Finally, let's use our low-level synchronization primitives, locks and condition variables, to build our own version of semaphores called ... *(drum roll here)* ... **Zemaphores**. This task is fairly straightforward, as you can see in Figure 31.16.

As you can see from the figure, we use just one lock and one condition variable, plus a state variable to track the value of the semaphore. Study the code for yourself until you really understand it. Do it!

One subtle difference between our Zemaphore and pure semaphores as defined by Dijkstra is that we don't maintain the invariant that the value of the semaphore, when negative, reflects the number of waiting threads; indeed, the value will never be lower than zero. This behavior is easier to implement and matches the current Linux implementation.

Curiously, building locks and condition variables out of semaphores

TIP: BE CAREFUL WITH GENERALIZATION
The abstract technique of generalization can thus be quite useful in systems design, where one good idea can be made slightly broader and thus solve a larger class of problems. However, be careful when generalizing; as Lampson warns us "Don't generalize; generalizations are generally wrong" [L83].

One could view semaphores as a generalization of locks and condition variables; however, is such a generalization needed? And, given the difficulty of realizing a condition variable on top of a semaphore, perhaps this generalization is not as general as you might think.

is a much trickier proposition. Some highly experienced concurrent programmers tried to do this in the Windows environment, and many different bugs ensued [B04]. Try it yourself, and see if you can figure out why building condition variables out of semaphores is more challenging than it might appear.

## 31.8  Summary

Semaphores are a powerful and flexible primitive for writing concurrent programs. Some programmers use them exclusively, shunning locks and condition variables, due to their simplicity and utility.

In this chapter, we have presented just a few classic problems and solutions. If you are interested in finding out more, there are many other materials you can reference. One great (and free reference) is Allen Downey's book on concurrency and programming with semaphores [D08]. This book has lots of puzzles you can work on to improve your understanding of both semaphores in specific and concurrency in general. Becoming a real concurrency expert takes years of effort; going beyond what you learn in this class is undoubtedly the key to mastering such a topic.

# References

[B04] "Implementing Condition Variables with Semaphores"
Andrew Birrell
December 2004
*An interesting read on how difficult implementing CVs on top of semaphores really is, and the mistakes the author and co-workers made along the way. Particularly relevant because the group had done a ton of concurrent programming; Birrell, for example, is known for (among other things) writing various thread-programming guides.*

[CB08] "Real-world Concurrency"
Bryan Cantrill and Jeff Bonwick
ACM Queue. Volume 6, No. 5. September 2008
*A nice article by some kernel hackers from a company formerly known as Sun on the real problems faced in concurrent code.*

[CHP71] "Concurrent Control with Readers and Writers"
P.J. Courtois, F. Heymans, D.L. Parnas
Communications of the ACM, 14:10, October 1971
*The introduction of the reader-writer problem, and a simple solution. Later work introduced more complex solutions, skipped here because, well, they are pretty complex.*

[D59] "A Note on Two Problems in Connexion with Graphs"
E. W. Dijkstra
Numerische Mathematik 1, 269271, 1959
Available: http://www-m3.ma.tum.de/twiki/pub/MN0506/WebHome/dijkstra.pdf
*Can you believe people worked on algorithms in 1959? We can't. Even before computers were any fun to use, these people had a sense that they would transform the world...*

[D68a] "Go-to Statement Considered Harmful"
E.W. Dijkstra
Communications of the ACM, volume 11(3): pages 147148, March 1968
Available: http://www.cs.utexas.edu/users/EWD/ewd02xx/EWD215.PDF
*Sometimes thought as the beginning of the field of software engineering.*

[D68b] "The Structure of the THE Multiprogramming System"
E.W. Dijkstra
Communications of the ACM, volume 11(5), pages 341346, 1968
*One of the earliest papers to point out that systems work in computer science is an engaging intellectual endeavor. Also argues strongly for modularity in the form of layered systems.*

[D72] "Information Streams Sharing a Finite Buffer"
E.W. Dijkstra
Information Processing Letters 1: 179180, 1972
Available: http://www.cs.utexas.edu/users/EWD/ewd03xx/EWD329.PDF
*Did Dijkstra invent everything? No, but maybe close. He certainly was the first to clearly write down what the problems were in concurrent code. However, it is true that practitioners in operating system design knew of many of the problems described by Dijkstra, so perhaps giving him too much credit would be a misrepresentation of history.*

[D08] "The Little Book of Semaphores"
A.B. Downey
Available: http://greenteapress.com/semaphores/
*A nice (and free!) book about semaphores. Lots of fun problems to solve, if you like that sort of thing.*

[D71] "Hierarchical ordering of sequential processes"
E.W. Dijkstra
Available: http://www.cs.utexas.edu/users/EWD/ewd03xx/EWD310.PDF
*Presents numerous concurrency problems, including the Dining Philosophers. The wikipedia page about this problem is also quite informative.*

[GR92] "Transaction Processing: Concepts and Techniques"
Jim Gray and Andreas Reuter
Morgan Kaufmann, September 1992
*The exact quote that we find particularly humorous is found on page 485, at the top of Section 8.8: "The first multiprocessors, circa 1960, had test and set instructions ... presumably the OS implementors worked out the appropriate algorithms, although Dijkstra is generally credited with inventing semaphores many years later."*

[H87] "Aspects of Cache Memory and Instruction Buffer Performance"
Mark D. Hill
Ph.D. Dissertation, U.C. Berkeley, 1987
*Hill's dissertation work, for those obsessed with caching in early systems. A great example of a quantitative dissertation.*

[L83] "Hints for Computer Systems Design"
Butler Lampson
ACM Operating Systems Review, 15:5, October 1983
*Lampson, a famous systems researcher, loved using hints in the design of computer systems. A hint is something that is often correct but can be wrong; in this use, a signal() is telling a waiting thread that it changed the condition that the waiter was waiting on, but not to trust that the condition will be in the desired state when the waiting thread wakes up. In this paper about hints for designing systems, one of Lampson's general hints is that you should use hints. It is not as confusing as it sounds.*

NAME I SYNOPSIS I DESCRIPTION I RETURN VALUE I ERRORS I CONFORMING TO I
NOTES I SEE ALSO I COLOPHON

Search online pages

**FLOCK(2)**                    **Linux Programmer's Manual**                    **FLOCK(2)**


## NAME         top

        flock – apply or remove an advisory lock on an open file


## SYNOPSIS         top

        **#include <sys/file.h>**

        **int flock(int** *fd*, **int** *operation*);


## DESCRIPTION         top

        Apply or remove an advisory lock on the open file specified by *fd*.
        The argument *operation* is one of the following:

            **LOCK_SH**   Place a shared lock.  More than one process may hold a
                        shared lock for a given file at a given time.

            **LOCK_EX**   Place an exclusive lock.  Only one process may hold an
                        exclusive lock for a given file at a given time.

            **LOCK_UN**   Remove an existing lock held by this process.

        A call to **flock**() may block if an incompatible lock is held by
        another process.  To make a nonblocking request, include **LOCK_NB** (by
        ORing) with any of the above operations.

        A single file may not simultaneously have both shared and exclusive
        locks.

        Locks created by **flock**() are associated with an open file description
        (see open(2)).  This means that duplicate file descriptors (created
        by, for example, fork(2) or dup(2)) refer to the same lock, and this
        lock may be modified or released using any of these file descriptors.
        Furthermore, the lock is released either by an explicit **LOCK_UN**
        operation on any of these duplicate file descriptors, or when all
        such file descriptors have been closed.

        If a process uses open(2) (or similar) to obtain more than one file
        descriptor for the same file, these file descriptors are treated
        independently by **flock**().  An attempt to lock the file using one of
        these file descriptors may be denied by a lock that the calling
        process has already placed via another file descriptor.

A process may hold only one type of lock (shared or exclusive) on a
file.  Subsequent **flock**() calls on an already locked file will
convert an existing lock to the new lock mode.

Locks created by **flock**() are preserved across an execve(2).

A shared or exclusive lock can be placed on a file regardless of the
mode in which the file was opened.

## RETURN VALUE       top

On success, zero is returned.  On error, -1 is returned, and *errno* is
set appropriately.

## ERRORS       top

**EBADF**   *fd* is not an open file descriptor.

**EINTR**   While waiting to acquire a lock, the call was interrupted by
        delivery of a signal caught by a handler; see signal(7).

**EINVAL**  *operation* is invalid.

**ENOLCK**  The kernel ran out of memory for allocating lock records.

**EWOULDBLOCK**
        The file is locked and the **LOCK_NB** flag was selected.

## CONFORMING TO       top

4.4BSD (the **flock**() call first appeared in 4.2BSD).  A version of
**flock**(), possibly implemented in terms of fcntl(2), appears on most
UNIX systems.

## NOTES       top

Since kernel 2.0, **flock**() is implemented as a system call in its own
right rather than being emulated in the GNU C library as a call to
fcntl(2).  With this implementation, there is no interaction between
the types of lock placed by **flock**() and fcntl(2), and **flock**() does
not detect deadlock.  (Note, however, that on some systems, such as
the modern BSDs, **flock**() and fcntl(2) locks *do* interact with one
another.)

In Linux kernels up to 2.6.11, **flock**() does not lock files over NFS
(i.e., the scope of locks was limited to the local system).  Instead,
one could use fcntl(2) byte-range locking, which does work over NFS,
given a sufficiently recent version of Linux and a server which
supports locking.  Since Linux 2.6.12, NFS clients support **flock**()

locks by emulating them as byte-range locks on the entire file.  This
means that fcntl(2) and **flock**() locks *do* interact with one another
over NFS.  Since Linux 2.6.37, the kernel supports a compatibility
mode that allows **flock**() locks (and also fcntl(2) byte region locks)
to be treated as local; see the discussion of the *local_lock* option
in nfs(5).

**flock**() places advisory locks only; given suitable permissions on a
file, a process is free to ignore the use of **flock**() and perform I/O
on the file.

**flock**() and fcntl(2) locks have different semantics with respect to
forked processes and dup(2).  On systems that implement **flock**() using
fcntl(2), the semantics of **flock**() will be different from those
described in this manual page.

Converting a lock (shared to exclusive, or vice versa) is not
guaranteed to be atomic: the existing lock is first removed, and then
a new lock is established.  Between these two steps, a pending lock
request by another process may be granted, with the result that the
conversion either blocks, or fails if **LOCK_NB** was specified.  (This
is the original BSD behavior, and occurs on many other
implementations.)

## SEE ALSO          top

flock(1), close(2), dup(2), execve(2), fcntl(2), fork(2), open(2),
lockf(3), lslocks(8)

*Documentation/filesystems/locks.txt* in the Linux kernel source tree
(*Documentation/locks.txt* in older kernels)

## COLOPHON          top

This page is part of release 4.08 of the Linux *man-pages* project.  A
description of the project, information about reporting bugs, and the
latest version of this page, can be found at
https://www.kernel.org/doc/man-pages/.

**Linux**                              **2016-03-15**                            **FLOCK(2)**
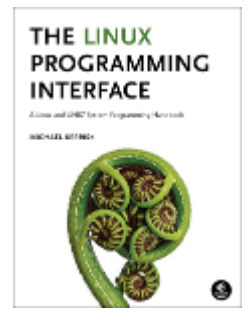
Copyright and license for this manual page

HTML rendering created 2016-10-08 by Michael Kerrisk, author of *The Linux
Programming Interface*, maintainer of the Linux *man-pages* project.

For details of in-depth **Linux/UNIX system programming training courses** that I
teach, look here.

Hosting by jambit GmbH.

2

man7.org > Linux > man-pages                    **Linux/UNIX system programming training**

Search online pages

```
LOCKF(3)                    Linux Programmer's Manual                    LOCKF(3)
```

## NAME        top

       lockf - apply, test or remove a POSIX lock on an open file

## SYNOPSIS        top

       **#include <unistd.h>**

       **int lockf(int** *fd*, **int** *cmd*, **off_t** *len*)**;**

    Feature Test Macro Requirements for glibc (see feature_test_macros(7)):

       **lockf**():
           _XOPEN_SOURCE >= 500
               || /* Glibc since 2.19: */ _DEFAULT_SOURCE
               || /* Glibc versions <= 2.19: */ _BSD_SOURCE || _SVID_SOURCE

## DESCRIPTION        top

       Apply, test or remove a POSIX lock on a section of an open file.  The
       file is specified by *fd*, a file descriptor open for writing, the
       action by *cmd*, and the section consists of byte positions
       *pos*..*pos*+*len*-1 if *len* is positive, and *pos*-*len*..*pos*-1 if *len* is
       negative, where *pos* is the current file position, and if *len* is zero,
       the section extends from the current file position to infinity,
       encompassing the present and future end-of-file positions.  In all
       cases, the section may extend past current end-of-file.

       On Linux, **lockf**() is just an interface on top of fcntl(2) locking.
       Many other systems implement **lockf**() in this way, but note that
       POSIX.1 leaves the relationship between **lockf**() and fcntl(2) locks
       unspecified.  A portable application should probably avoid mixing
       calls to these interfaces.

       Valid operations are given below:

       **F_LOCK** Set an exclusive lock on the specified section of the file.
              If (part of) this section is already locked, the call blocks
              until the previous lock is released.  If this section overlaps
              an earlier locked section, both are merged.  File locks are
              released as soon as the process holding the locks closes some
              file descriptor for the file.  A child process does not
              inherit these locks.

**F_TLOCK**
Same as **F_LOCK** but the call never blocks and returns an error instead if the file is already locked.

**F_ULOCK**
Unlock the indicated section of the file.  This may cause a locked section to be split into two locked sections.

**F_TEST** Test the lock: return 0 if the specified section is unlocked or locked by this process; return -1, set *errno* to **EAGAIN** (**EACCES** on some other systems), if another process holds a lock.

## RETURN VALUE        top

On success, zero is returned.  On error, -1 is returned, and *errno* is set appropriately.

## ERRORS        top

**EACCES** or **EAGAIN**
The file is locked and **F_TLOCK** or **F_TEST** was specified, or the operation is prohibited because the file has been memory-mapped by another process.

**EBADF**  *fd* is not an open file descriptor; or *cmd* is **F_LOCK** or **F_TLOCK** and *fd* is not a writable file descriptor.

**EDEADLK**
The command was **F_LOCK** and this lock operation would cause a deadlock.

**EINVAL** An invalid operation was specified in *cmd*.

**ENOLCK** Too many segment locks open, lock table is full.

## ATTRIBUTES        top

For an explanation of the terms used in this section, see attributes(7).

| Interface | Attribute | Value |
|-----------|-----------|-------|
| **lockf**() | Thread safety | MT-Safe |

## CONFORMING TO        top

POSIX.1–2001, POSIX.1–2008, SVr4.


## SEE ALSO          top

fcntl(2), flock(2)

*locks.txt* and *mandatory-locking.txt* in the Linux kernel source
directory *Documentation/filesystems* (on older kernels, these files
are directly under the *Documentation* directory, and *mandatory-
locking.txt* is called *mandatory.txt*)


## COLOPHON          top

This page is part of release 4.08 of the Linux *man-pages* project.  A
description of the project, information about reporting bugs, and the
latest version of this page, can be found at
https://www.kernel.org/doc/man-pages/.

**GNU**                                  **2016–03–15**                                  **LOCKF(3)**

Copyright and license for this manual page

HTML rendering created 2016-10-08 by Michael Kerrisk, author of *The Linux
Programming Interface*, maintainer of the Linux *man-pages* project.

For details of in-depth **Linux/UNIX system programming training courses** that I
teach, look here.

Hosting by jambit GmbH.