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 <code>join()</code>); how should such a wait be implemented? Let's look at Figure 30.1.

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

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

```
int done = 0;
1
    pthread_mutex_t m = PTHREAD_MUTEX_INITIALIZER;
2
    pthread_cond_t c = PTHREAD_COND_INITIALIZER;
3
5
   void thr_exit() {
       Pthread_mutex_lock(&m);
       done = 1;
       Pthread_cond_signal(&c);
       Pthread_mutex_unlock(&m);
9
10
11
   void *child(void *arg) {
12
      printf("child\n");
13
       thr_exit();
14
       return NULL;
15
16
17
   void thr_join() {
18
      Pthread mutex lock(&m);
19
20
       while (done == 0)
21
         Pthread_cond_wait(&c, &m);
       Pthread_mutex_unlock(&m);
22
23
24
   int main(int argc, char *argv[]) {
25
       printf("parent: begin\n");
26
       pthread_t p;
       Pthread_create(&p, NULL, child, NULL);
28
       thr_join();
29
       printf("parent: end\n");
30
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 implementations. First, you might be wondering if we need the state variable done. What if the code looked like the example below? Would this work?

```
void thr_exit() {
    Pthread_mutex_lock(&m);
    Pthread_cond_signal(&c);
    Pthread_mutex_unlock(&m);
}

void thr_join() {
    Pthread_mutex_lock(&m);
    Pthread_cond_wait(&c, &m);
    Pthread_mutex_unlock(&m);
}
```

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 <code>done</code>; 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!

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 correctness; 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 returning. 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 basic requirements of using condition variables properly. To make sure you understand, we now go through a more complicated example: the **producer/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; consumers 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 program into another, e.g., <code>grep foo file.txt | wc -1</code>. This example runs two processes concurrently; <code>grep writes</code> lines from <code>file.txt</code> with the string <code>foo</code> 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 <code>pipe</code> system call). The other end of this pipe is connected to the standard input of the process <code>wc</code>, which simply counts the number of lines in the input stream and prints out the result. Thus, the <code>grep</code> process is the producer; the <code>wc</code> process is the consumer; between them is an in-kernel bounded buffer; you, in this example, are just the happy user.

```
1
  int buffer;
  int count = 0; // initially, empty
2
   void put(int value) {
      assert (count == 0);
       count = 1;
      buffer = value;
   int get() {
10
     assert (count == 1);
11
       count = 0;
12
       return buffer;
13
```

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

```
void *producer(void *arg) {
       int i;
       int loops = (int) arg;
3
       for (i = 0; i < loops; i++) {
5
           put(i);
   void *consumer(void *arg) {
9
10
      int i;
       while (1) {
11
           int tmp = get();
           printf("%d\n", tmp);
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, lest¹ 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

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

```
cond_t cond;
1
   mutex_t mutex;
2
3
4
   void *producer(void *arg) {
5
       int i;
       for (i = 0; i < loops; i++) {
6
           Pthread_mutex_lock(&mutex); // pl
           if (count == 1)
                                                  // p2
9
                Pthread_cond_wait(&cond, &mutex); // p3
            Pthread_cond_signal(&cond); // p4
Pthread_mutov -- '...
10
11
           Pthread_mutex_unlock(&mutex);
                                                  // p6
12
        }
13
   }
14
15
   void *consumer(void *arg) {
16
17
        for (i = 0; i < loops; i++) {
18
                                                 // c1
           Pthread_mutex_lock(&mutex);
19
           if (count == 0)
                                                   // c2
               Pthread_cond_wait(&cond, &mutex); // c3
21
                              // c4
al(&cond); // c5
            int tmp = get();
22
           Pthread_cond_signal(&cond); // c5
Pthread_mutex_unlock(&mutex); // c6
23
24
           printf("%d\n", tmp);
25
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	ĺ	Ready		Ready	0	
c3	Sleep	ĺ	Ready		Ready	0	Nothing to get
	Sleep	ĺ	Ready	p1	Running	0	
	Sleep	ĺ	Ready	p2	Running	0	
	Sleep	ĺ	Ready	p4	Running	1	Buffer now full
	Ready	ĺ	Ready	p5	Running	1	T_{c1} awoken
	Ready	ĺ	Ready	p6	Running	1	
	Ready	ĺ	Ready	p1	Running	1	
	Ready	ĺ	Ready	p2	Running	1	
	Ready	ĺ	Ready	р3	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	с6	Running		Ready	0	-
c4	Running	ĺ	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.

```
cond_t cond;
   mutex_t mutex;
2
3
   void *producer(void *arg) {
5
      int i;
       for (i = 0; i < loops; i++) {
6
          (i = 0; 1 \ 100ps, __,
Pthread_mutex_lock(&mutex);
                                            // p1
           while (count == 1)
                                               // p2
9
              Pthread_cond_wait(&cond, &mutex); // p3
                                     // p4
10
           put(i);
                                                // p5
11
           Pthread_cond_signal(&cond);
           Pthread_mutex_unlock(&mutex);
                                               // p6
12
       }
13
14
15
   void *consumer(void *arg) {
16
17
       for (i = 0; i < loops; i++) {
18
           Pthread_mutex_lock(&mutex);
                                               // c1
19
          while (count == 0)
                                               // c2
              Pthread_cond_wait(&cond, &mutex); // c3
21
                                // c4
           int tmp = get();
22
           Pthread_cond_signal(&cond);
                                              // c5
// c6
23
           Pthread_mutex_unlock(&mutex);
          printf("%d\n", tmp);
25
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	0 0
	Sleep	c2	Running		Ready	0	
	Sleep	c3	Sleep		Ready	0	Nothing to get
	Sleep		Sleep	p1	Running	0	0 0
	Sleep		Sleep	p2	Running	0	
	Sleep		Sleep	p4	Running	1	Buffer now full
	Ready		Sleep	p5	Running	1	T_{c1} awoken
	Ready		Sleep	р6	Running	1	
	Ready		Sleep	p1	Running	1	
	Ready		Sleep	p2	Running	1	
	Ready		Sleep	р3	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}
с6	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.

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

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.

```
int buffer[MAX];
1
2 int fill_ptr = 0;
3 int use_ptr = 0;
4 int count = 0;
   void put (int value) {
    buffer[fill_ptr] = value;
      fill_ptr = (fill_ptr + 1) % MAX;
9
      count++;
10
11
  int get() {
12
  int tmp = buffer[use_ptr];
13
      use_ptr = (use_ptr + 1) % MAX;
15
     count--;
16
      return tmp;
               Figure 30.11: The Final Put And Get Routines
  cond_t empty, fill;
  mutex_t mutex;
   void *producer(void *arg) {
4
     int i;
5
       for (i = 0; i < loops; i++) {
        Pthread_mutex_lock(&mutex);
                                              // p1
         while (count == MAX)
                                             // p2
             Pthread_cond_wait(&empty, &mutex); // p3
9
          put(i);
10
11
12
13
14
  void *consumer(void *arg) {
     int i;
17
       for (i = 0; i < loops; i++) {
18
                                              // c1
          Pthread_mutex_lock(&mutex);
19
                                              // c2
20
          while (count == 0)
            Pthread_cond_wait(&fill, &mutex); // c3
         int tmp = get();
                                    // c4
22
         Pthread_mutex_unlock(&mutex); // c6
printf("%d\n", tmo):
23
24
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.

```
// how many bytes of the heap are free?
   int bytesLeft = MAX_HEAP_SIZE;
   // need lock and condition too
   cond t c;
   mutex_t m;
  void *
   allocate(int size) {
     Pthread_mutex_lock(&m);
10
       while (bytesLeft < size)
11
          Pthread_cond_wait(&c, &m);
12
      void *ptr = ...; // get mem from heap
bytesLeft -= size;
13
       Pthread_mutex_unlock(&m);
15
       return ptr;
16
17
   void free(void *ptr, int size) {
       Pthread_mutex_lock(&m);
       bytesLeft += size;
        Pthread_cond_signal(&c); // whom to signal??
       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?