

T21

Condition Variables

Referência principal

Ch.30 of *Operating Systems: Three Easy Pieces* by Remzi and Andrea Arpaci-Dusseau (pages.cs.wisc.edu/~remzi/OSTEP/)

Discutido em classe em 10 de outubro de 2018

Waiting for a condition to hold

- Locks aren't the only primitives needed to build concurrent applications.
- Often a thread wishes to check whether a ***condition*** holds before continuing execution.
 - For example, a parent thread might want to wait until a child thread terminates before continuing.
 - This is frequently modelled as a `join()`.
- How should such a wait be implemented?

How would a parent thread wait for its child?

```
1. #include <stdio.h>

2. void *child(void *arg) {
3.     printf("child\n") ;
4.     // XXX how to indicate we are done?
5.     return NULL;
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. }
```

Expected output

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

Trying to use a shared variable

```
1. volatile int done = 0;

2. void *child(void *arg) {
3.     printf("child\n");
4.     done = 1;
5.     return NULL;
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.    while (done == 0)
12.        ; // spin
13.    printf("parent: end\n");
14.    return 0;
15. }
```

Does the shared variable solve the problem?

- The shared variable approach will often work.
- However, it is very inefficient as the parent spins and wastes CPU time.
- Instead, we would like to have some way to put the parent to sleep until the condition we are waiting for (e.g. the termination of the child thread) comes true.

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 – just spinning until the condition becomes true – is grossly inefficient, wastes CPU cycles, and in some cases, can be incorrect.
- Thus, **how should a thread wait for a condition?**

Condition Variables

- A condition variable is an explicit queue that threads can put themselves on when some state of execution is not as desired.
- Some other thread, when it changes that state, can then wake one (or more) of those waiting threads and thus allow them to continue.
- The approach goes back to Dijkstra's concept of "*private semaphores*" (1968).
- A similar idea was later named a "*condition variable*" by Hoare in his work on monitors (1974).

POSIX¹ Operations on Condition Variables

- The examples below refer to the wrapper functions declared in `mythreads.h`.
- Condition variables are declared by writing something like
`pthread_cond_t c;`
- A condition variable can be declared and initialized at the same time by
`pthread_cond_t cv = PTHREAD_COND_INITIALIZER;`
- A thread can wait on a condition variable by calling
`pthread_cond_wait(pthread_cond_t *c, pthread_mutex_t *m);`
- A thread can signal on a condition variable by calling
`pthread_cond_signal(pthread_cond_t *c);`

¹The Portable Operating System Interface (POSIX) is a family of standards specified by the IEEE Computer Society for maintaining compatibility between operating systems. POSIX defines the application programming interface (API), along with command line shells and utility interfaces, for software compatibility with variants of Unix and other operating systems.

The semantics of `wait()` and `signal()`

- You may have noticed that `wait()` also takes a *mutex* as a parameter.
 - It assumes that the calling thread holds this mutex at the point of call.
 - The responsibility of `wait()` is to release the lock and put the calling thread to sleep (atomically).
 - Releasing the lock is mandatory to enable some other thread to change the state, make the condition true and then `signal()` on it.
 - After some other thread's signal, `wait()` will re-acquire the lock before waking up the thread and returning to the caller.
 - This complexity is required to prevent certain race conditions from occurring when a thread is trying to put itself to sleep.
 - Let's take a look at an attempt at solving the “join problem” to understand this better.

Trying to solve the “join problem” using a condition variable

- Let us assume that the two required actions on our previous draft are implemented by
 - a `thr_join()` function that will be called by the parent
 - a `thr_exit()` function that will be called by the child

```
1. void *child(void *arg) {  
2.     printf("child\n") ;  
3.     thr_exit( );  
4.     return NULL;  
5. }  
  
6. int main(int argc, char *argv[]) {  
7.     printf("parent: begin\n");  
8.     pthread_t p;  
9.     Pthread_create(&p, NULL, child, NULL);  
10.    thr_join();  
11.    printf("parent: end\n");  
12.    return 0;  
13. }
```

Trying to solve the “join problem” using a condition variable

- Let us use a shared variable **done** to indicate that the child has finished.
- Initially, **done** will be **zero**. This will be changed to **one** by the child.
- A condition variable **cv** will be used to synchronize the parent and child threads.
 - The parent will **wait on cv** for the child to finish.
 - The child will **signal on cv** to indicate it has finished.
- Finally, to prevent race conditions, we will use a **mutex m** to protect every access to **done**.

Trying to solve the “join problem” using a condition variable

- Now let us try to implement `thr_exit()` and `thr_join()`, using the shared variable `done`, the condition variable `cv` and a mutex `m`.

```
1.  int done = 0;
2.  mutex_t m = PTHREAD_MUTEX_INITIALIZER;
3.  cond_t cv = PTHREAD_COND_INITIALIZER;

4.  void thr_exit() {
5.      mutex_lock(&m);
6.      done = 1;
7.      cond_signal(&cv);
8.      mutex_unlock(&m);
9.  }

10. void thr_join() {
11.     mutex_lock(&m);
12.     if (done == 0)
13.         cond_wait(&cv, &m);
14.     mutex_unlock(&m);
15. }
```

Analyzing the proposed solution to the “join problem”

- Consider that there is only one processor and examine two cases
 - a. The parent creates the child and carries on, thus calling `thr_join()`.
 - b. The child is created and runs immediately, thus printing its message and calling `thr_exit()`.
- Now, a few questions may have arisen, such as
 - Is `done` really necessary?
 - Is the `mutex` really necessary?
 - Is it the solution we wanted?
- Let us examine each one of them in turn.

Trying to get rid of `done`

- Assume that there is no `done`...

```
4. void thr_exit() {  
5.     mutex_lock(&m);  
6.     cond_signal(&c);  
7.     mutex_unlock(&m);  
8. }  
  
9. void thr_join() {  
10.    mutex_lock(&m);  
11.    cond_wait(&c, &m);  
12.    mutex_unlock(&m);  
13. }
```

```
1. int done = 0;  
2. mutex_t m = PTHREAD_MUTEX_INITIALIZER;  
3. cond_t cv = PTHREAD_COND_INITIALIZER;  
  
4. void thr_exit() {  
5.     mutex_lock(&m);  
6.     done = 1;  
7.     cond_signal(&cv);  
8.     mutex_unlock(&m);  
9. }  
  
10. void thr_join() {  
11.    mutex_lock(&m);  
12.    if (done == 0)  
13.        cond_wait(&cv, &m);  
14.    mutex_unlock(&m);  
15. }
```

- Now, re-examine case b, where the child is created and runs immediately, thus printing its message and calling `thr_exit()`.
- Does it work?

Trying to get rid of `mutex`

- Assume that there is no `mutex`...

```
4. void thr_exit() {  
5.     done = 1;  
6.     cond_signal(&c);  
7. }  
  
8. void thr_join() {  
9.     if (done == 0)  
10.        cond_wait(&c, &m);  
11. }
```

```
1. int done = 0;  
2. mutex_t m = PTHREAD_MUTEX_INITIALIZER;  
3. cond_t cv = PTHREAD_COND_INITIALIZER;  
  
4. void thr_exit() {  
5.     mutex_lock(&m);  
6.     done = 1;  
7.     cond_signal(&cv);  
8.     mutex_unlock(&m);  
9. }  
  
10. void thr_join() {  
11.     mutex_lock(&m);  
12.     if (done == 0)  
13.         cond_wait(&cv, &m);  
14.     mutex_unlock(&m);  
15. }
```

- Now, re-examine case a, where the parent creates the child and continues to run.
- Assume that after checking `(done == 0)` on line 9, the parent is interrupted just before calling `cond_wait()`
- Does it work? Or is there a race condition?

Unfortunately, there is still one bug . . .

- By now, you should have appreciated the need for **done** and the **mutex**.
- Unfortunately, our solution still has one bug.
- To make it evident, we will study a slightly more complicated example: the ***producer/consumer*** or ***bounded buffer*** problem.
- This problem was posed by Dijkstra in 1972 and led to the invention of another synchronization primitive, the ***generalized semaphore***, which we will study later.

The *producer/consumer* or *bounded buffer* problem

- Imagine one or more producer threads and one or more consumer threads.
 - Producers generate data items and place them in a buffer;
 - Consumers grab items from the buffer and use them somehow.
- Because the bounded buffer is a shared resource, access to it must be synchronized, to avoid a race condition.
- Let us try to understand this problem better...
 - We will need a shared buffer, into which a producer puts data, and out of which a consumer takes data.
 - We will use a one-slot buffer, able to hold just one integer, for simplicity. We will relax this constraint later.
 - We will also use two inner routines to put a value into the shared buffer, and to get a value out of the buffer, as shown in the next slide.

The `put()` and `get()` functions

- There are two global variables:
 - The `buffer` and the `count` of numbers in it.
- `put` asserts that `buffer` is not full and puts `value` in it.
- `get` asserts that `buffer` is not empty and returns the `value` in it.
- In case an assertion fails, an exception will be raised.

```
1.  int buffer;
2.  int count = 0; // initially, empty

3.  void put(int value) {
4.      assert(count == 0);
5.      count = 1;
6.      buffer = value;
7.  }

8.  int get(void) {
9.      assert(count == 1);
10.     count = 0;
11.     return buffer;
12. }
```

The producer and consumer threads (version 1)

- There are two kinds of threads:
 - **producer**, that puts an integer in the shared buffer **loops** times;
 - **consumer**, that gets data out of the shared buffer (forever) and prints out the values it receives.
- An application is supposed to create one or more **producers** and one or more **consumers**.
- What do you think of this design? Is it correct? Does it work well?

```
1. void *producer(void *arg) {
2.     int loops = (int) arg;
3.     for (int i = 0; i < loops; i++) {
4.         put(i);
5.     }
6.     return NULL;
7. }

8. void *consumer(void *arg) {
9.     while (1) {
10.        int tmp = get();
11.        printf("%d\n", tmp);
12.    }
13.    return NULL;
14. }
```

Unfortunately, version 1 does not work...

- You certainly noticed that version 1 has flaws.
 - There is a race condition on **count** and **buffer**, which can be prevented by a **mutex**.
 - Both **producer** and **consumer** may trigger an exception by calling **put** or **get** at the “wrong” time, but their accesses can be synchronized using a **condition variable**.
- Let us create a version 2 by providing these two adjustments.

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

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

The producer and consumer threads (version 2)

```
1. void *producer(void *arg) {  
2.     int loops = *((int *)arg);  
3.     for (int i = 0; i < loops; i++) {  
4.         mutex_lock(&mutex);           // p1  
5.         if (count == 1)                // p2  
6.             cond_wait(&cond, &mutex); // p3  
7.         put(i);                        // p4  
8.         cond_signal(&cond);           // p5  
9.         mutex_unlock(&mutex);         // p6  
10.    }  
11.    return NULL;  
12. }
```

```
1. void *consumer(void *arg) {  
2.     int loops = *((int *)arg);  
3.     for (int i = 0; i < loops; i++) {  
4.         mutex_lock(&mutex);           // c1  
5.         if (count == 0)                // c2  
6.             cond_wait(&cond, &mutex); // c3  
7.         int tmp = get();                // c4  
8.         cond_signal(&cond);           // c5  
9.         mutex_unlock(&mutex);         // c6  
10.        printf("%d\n", tmp);  
11.    }  
12.    return NULL;  
13. }
```

- Let us do a thread trace to examine the behavior of this proposal.

Thread trace of version 2 with 2 consumers

t	T _{C1}	State	T _{C2}	State	T _P	State	count	Comment
1	c1	Running		Ready		Ready	0	
2	c2	Running		Ready		Ready	0	
3	c3	Asleep		Ready		Ready	0	Buffer empty; must sleep
4		Asleep		Ready	p1	Running	0	
5		Asleep		Ready	p2	Running	0	
6		Asleep		Ready	p4	Running	1	Buffer now full
7		Ready		Ready	p5	Running	1	T _{C1} awoken
8		Ready		Ready	p6	Running	1	
9		Ready		Ready	p1	Running	1	
10		Ready		Ready	p2	Running	1	
11		Ready		Ready	p3	Asleep	1	Buffer full; must sleep
12		Ready	c1	Running		Asleep	1	T _{C2} sneaks in ...
13		Ready	c2	Running		Asleep	1	
14		Ready	c4	Running		Asleep	0	... and grabs the data
15		Ready	c5	Running		Ready	0	T _P awoken
16		Ready	c6	Running		Ready	0	
17		Ready	c1	Running		Ready	0	
18		Ready	c2	Running		Ready	0	
19		Ready	c3	Asleep		Ready	0	Buffer empty; must sleep
20	c4	Running		Ready		Ready	0	Oops! No data...

What went wrong in version 2 and how to fix it?

- The problem was caused by the assumption that on returning from **wait** the buffer would be **not full** (**producer**) or **not empty** (**consumer**).
- As we have seen this may not be the case when there are multiple producers or consumers.
- The solution is to repeat the test for fullness or emptiness of the **buffer** when returning from **wait**, i.e. let us use a **while** instead of an **if** in lines **p2** (**producer**) and **c2** (**consumer**).
- This will lead us to version 3, shown in the next slide.
- Is version 3 correct? Does it work in any situation?

The producer and consumer threads (version 3)

```
1. void *producer(void *arg) {  
2.     for (int i = 0; i < LOOPS; i++) {  
3.         mutex_lock(&mutex);           // p1  
4.         while (count == 1)           // p2  
5.             cond_wait(&cond, &mutex); // p3  
6.         put(i);                       // p4  
7.         cond_signal(&cond);          // p5  
8.         mutex_unlock(&mutex);        // p6  
9.     }  
10.    return NULL;  
11. }
```

```
1. void *consumer(void *arg) {  
2.     for (int i = 0; i < LOOPS; i++) {  
3.         mutex_lock(&mutex);           // c1  
4.         while (count == 0)           // c2  
5.             cond_wait(&cond, &mutex); // c3  
6.         int tmp = get();              // c4  
7.         cond_signal(&cond);          // c5  
8.         mutex_unlock(&mutex);        // c6  
9.         printf("%d\n", tmp);  
10.    }  
11.    return NULL;  
12. }
```

- Let us resort again to a thread trace to examine the behavior of this proposal.

Thread trace of version 3 with 2 consumers

t	T _{c1}	State	T _{c2}	State	T _p	State	count	Comment
1	c1	Running		Ready		Ready	0	
2	c2	Running		Ready		Ready	0	
3	c3	Asleep		Ready		Ready	0	Buffer empty; must sleep
4		Asleep	c1	Running		Ready	0	
5		Asleep	c2	Running		Ready	0	
6		Asleep	c3	Asleep		Ready	0	Buffer empty; must sleep
7		Asleep		Asleep	p1	Running	0	
8		Asleep		Asleep	p2	Running	0	
9		Asleep		Asleep	p4	Running	1	Buffer now full
10		Ready		Asleep	p5	Running	1	T _{c1} awoken
11		Ready		Asleep	p6	Running	1	
12		Ready		Asleep	pi	Running	1	
13		Ready		Asleep	p2	Running	1	
14		Ready		Asleep	p3	Asleep	1	Buffer full; must sleep
15	c2	Running		Asleep		Asleep	1	Recheck condition
16	c4	Running		Asleep		Asleep	0	T _{c1} grabs data
17	c5	Running		Ready		Asleep	0	Oops! Woke T _{c2} up
18	c6	Running		Ready		Asleep	0	
19	cl	Running		Ready		Asleep	0	
20	c2	Running		Ready		Asleep	0	
21	c3	Asleep		Ready		Asleep	0	Buffer empty; must sleep
22		Asleep	c2	Running		Asleep	0	
23		Asleep	c3	Asleep		Asleep	0	Everyone asleep...

What went wrong in version 3 and how to fix it?

- This time the problem was caused by waking up the “wrong” thread.
- As it may be clear now, our threads are interested in two different conditions: the **producers** need to wait for **buffer not full**, while the **consumers** need to wait for **buffer not empty**.
- Thus, we will use two different condition variables: **not_empty** and **not_full**.
- This will lead us to version 4, shown in the next slide.

The producer and consumer threads (version 4)

```
1.  cond_t not_empty, not_full;
2.  mutex_t mutex;
```

```
1.  void *producer(void *arg) {
2.      for (int i = 0; i < LOOPS; i++) {
3.          mutex_lock(&mutex);           // p1
4.          while (count == 1)           // p2
5.              cond_wait(&not_full, &mutex); // p3
6.          put(i);                       // p4
7.          cond_signal(&not_empty);      // p5
8.          mutex_unlock(&mutex);         // p6
9.      }
10.     return NULL;
11. }
```

```
1.  void *consumer(void *arg) {
2.      for (int i = 0; i < LOOPS; i++) {
3.          mutex_lock(&mutex);           // c1
4.          while (count == 0)           // c2
5.              cond_wait(&not_empty, &mutex); // c3
6.          int tmp = get();              // c4
7.          cond_signal(&not_full);      // c5
8.          mutex_unlock(&mutex);         // c6
9.          printf("%d\n", tmp);
10.     }
11.     return NULL;
12. }
```

- Study the code and then answer...
 - Is it correct? Does it work in any situation?
 - Why two condition variables but just one mutex?

The Final Producer/Consumer Solution

- Our last version works but it is not a general solution to the problem.
- To increase concurrency and efficiency, let us enlarge the buffer, turning it into an array of integers.
- This asks for only a couple of simple changes to our current design:
 - The **buffer** structure must be reflected on both **put()** and **get()**.
 - The **while** statement of the producer (line **p2**), must cater for the new **buffer** (shown in the next slide).

```
1.  #define MAX 100
2.
3.  int buffer[MAX];
4.  int fill_ptr = 0;
5.  int use_ptr = 0;
6.  int count = 0;
7.
8.  void put(int value) {
9.      buffer[fill_ptr] = value;
10.     fill_ptr = (fill_ptr + 1) % MAX;
11.     count++;
12. }
13.
14. int get(void) {
15.     int tmp = buffer[use_ptr];
16.     use_ptr = (use_ptr + 1) % MAX;
17.     count--;
18.     return tmp;
19. }
```

The producer and consumer threads (final version)

```
1.  cond_t not_empty, not_full;
2.  mutex_t mutex;
```

```
1.  void *producer(void *arg) {
2.      for (int i = 0; i < loops; i++) {
3.          mutex_lock(&mutex);           // p1
4.          while (count == MAX)         // p2
5.              cond_wait(&not_full, &mutex); // p3
6.          put(i);                       // p4
7.          cond_signal(&not_empty);     // p5
8.          mutex_unlock(&mutex);        // p6
9.      }
10.     return NULL;
11. }
```

```
1.  void *consumer(void *arg) {
2.      for (int i = 0; i < loops; i++) {
3.          mutex_lock(&mutex);           // c1
4.          while (count == 0)           // c2
5.              cond_wait(&not_empty, &mutex); // c3
6.          int tmp = get();              // c4
7.          cond_signal(&not_full);      // c5
8.          mutex_unlock(&mutex);        // c6
9.          printf("%d\n", tmp);
10.     }
11.     return NULL;
12. }
```

Covering Conditions

- Consider a multi-threaded memory allocation library, from which this code snippet was taken.
- Assume that there are no `bytesLeft` and that two threads call `allocate(100)` and `allocate(10)`, respectively.
 - As a result, both will wait on `c` (line 7).
- Now assume that a third thread calls `free(50)`.
 - If the first thread catches the signal, everyone will remain asleep.
- If we turn `cond_signal()` (line 16) into `code_broadcast()`, all sleeping threads could have a chance to awake and continue.

```
1.  int bytesLeft = MAX_HEAP_SIZE;
2.
3.  cond_t c;
4.  mutex_t m;
5.
6.  void *allocate(int size) {
7.      mutex_lock(&m);
8.      while (bytesLeft < size)
9.          cond_wait(&c, &m);
10.     void *ptr = ...;    // get mem from heap
11.     bytesLeft -= size;
12.     mutex_unlock(&m);
13.     return ptr;
14. }
15.
16. void free(void *ptr, int size) {
17.     mutex_lock(&m);
18.     bytesLeft += size;
19.     cond_signal(&c);    // whom to signal??
20.     mutex_unlock(&m);
21. }
```