

Producer/Consumer (Bounded Buffer)

- one or more *producer* threads add to a buffer
- one or more *consumer* threads remove from the buffer
- the buffer has a fixed size
 - can't add when it is full
 - can't remove when it is empty

First Try

```
int buffer;
int count = 0; // initially, empty

void put(int value) {
    assert(count == 0);
    count = 1;
    buffer = value;
}

int get() {
    assert(count == 1);
    count = 0;
    return buffer;
}
```

```
void *producer(void *arg) {
    int i;
    int loops = (int) arg;
    for (i = 0; i < loops; i++) {
        put(i);
    }
}

void *consumer(void *arg) {
    while (1) {
        int tmp = get();
        printf("%d\n", tmp);
    }
}
```

- does this work?
 - no, unless the put and get calls are interleaved exactly right – two puts in a row or two gets in a row will result in a violated assertion
 - also race conditions within put and get themselves

Second Try

add a condition variable to wait for expected buffer status


```
int loops; // must initialize somewhere...
cond_t cond;
mutex_t mutex;

void *producer(void *arg) {
    int i;
    for (i = 0; i < loops; i++) {
        pthread_mutex_lock(&mutex); // p1
        if (count == 1) // p2
            pthread_cond_wait(&cond, &mutex); // p3
        put(i); // p4
        pthread_cond_signal(&cond); // p5
        pthread_mutex_unlock(&mutex); // p6
    }
}

void *consumer(void *arg) {
    int i;
    for (i = 0; i < loops; i++) {
        pthread_mutex_lock(&mutex); // c1
        if (count == 0) // c2
            pthread_cond_wait(&cond, &mutex); // c3
        int tmp = get(i); // c4
        pthread_cond_signal(&cond); // c5
        pthread_mutex_unlock(&mutex); // c6
        printf("%d\n", tmp);
    }
}
```

two consumers, one producer

T _{c1}	State	T _{c2}	State	T _p	State	Count	Comment
c1	Run			Ready	Ready	0	
c2	Run			Ready	Ready	0	
c3	Sleep			Ready	Ready	0	Nothing to get
	Sleep			p1	Run	0	
	Sleep			p2	Run	0	
	Sleep			p4	Run	1	Buffer now full
	Sleep			p5	Run	1	T _{c1} awoken
	Ready			p6	Run	1	
	Ready			p1	Run	1	
	Ready			p2	Run	1	
	Ready			p3	Sleep	1	Buffer full; sleep
	Ready	c1	Run	Sleep	1	1	T _{c2} sneaks in ...
	Ready	c2	Run	Sleep	1	1	
	Ready	c4	Run	Sleep	0	0	... and grabs data
	Ready	c5	Run	Ready	0	0	T _p awoken
	Ready	c6	Run	Ready	0	0	
c4	Run			Ready	0	0	Oh oh! No data

doesn't work  assumption is that when wait returns, the buffer is ready for us – but with more than one producer or consumer, another one could sneak in and invalidate that before we get a chance to run

Third Try

while instead of if


two consumers, one producer

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

```
int loops; // must initialize somewhere...
cond_t cond;
mutex_t mutex;

void *producer(void *arg) {
    int i;
    for (i = 0; i < loops; i++) {
        pthread_mutex_lock(&mutex); // p1
        while (count == 1) // p2
            pthread_cond_wait(&cond, &mutex); // p3
        put(i); // p4
        pthread_cond_signal(&cond); // p5
        pthread_mutex_unlock(&mutex); // p6
    }
}

void *consumer(void *arg) {
    int i;
    for (i = 0; i < loops; i++) {
        pthread_mutex_lock(&mutex); // c1
        while (count == 0) // c2
            pthread_cond_wait(&cond, &mutex); // c3
        int tmp = get(i); // c4
        pthread_cond_signal(&cond); // c5
        pthread_mutex_unlock(&mutex); // c6
        printf("%d\n", tmp);
    }
}
```

 the problem of a buffer status change between waking and access is fixed but a single condition variable means that both producers and consumers wait, and thus a thread can be woken by another thread of the same kind rather than the opposite...and if it can't do its job, it goes back to sleep and no one else gets signaled

Fourth Try

separate condition variables for empty buffer and full buffer

```
int loops; // must initialize somewhere...
cond_t empty, full;
mutex_t mutex;

void *producer(void *arg) {
    int i;
    for (i = 0; i < loops; i++) {
        Pthread_mutex_lock(&mutex);
        while (count == 1)
            Pthread_cond_wait(&empty, &mutex);
        put(i);
        Pthread_cond_signal(&full);
        Pthread_mutex_unlock(&mutex);
    }
}

void *consumer(void *arg) {
    int i;
    for (i = 0; i < loops; i++) {
        Pthread_mutex_lock(&mutex);
        while (count == 0)
            Pthread_cond_wait(&full, &mutex);
        int tmp = get();
        Pthread_cond_signal(&empty);
        Pthread_mutex_unlock(&mutex);
        printf("%d\n", tmp);
    }
}
```

```
int buffer;
int count = 0; // initially, empty

void put(int value) {
    assert(count == 0);
    count = 1;
    buffer = value;
}

int get() {
    assert(count == 1);
    count = 0;
    return buffer;
}
```

- does this work?

😊 yes!



the buffer itself is not thread safe – why is this not a problem?



separate condition variables is a better solution than broadcast – avoids waking up threads that will just end up going back to sleep right away

Full Producer/Consumer Solution

```
int buffer[MAX];
int fill_ptr = 0;
int use_ptr = 0;
int count = 0;

void put(int value) {
    buffer[fill_ptr] = value;
    fill_ptr = (fill_ptr + 1) % MAX;
    count++;
}

int get() {
    int tmp = buffer[use_ptr];
    use_ptr = (use_ptr + 1) % MAX;
    count--;
    return tmp;
}
```



add more slots to the buffer

use a circular array implementation of a queue – use_ptr at the head, fill_ptr at the tail – to avoid shifting and achieve O(1) operations

```
int loops; // must initialize somewhere...
cond_t empty, full;
mutex_t mutex;

void *producer(void *arg) {
    int i;
    for (i = 0; i < loops; i++) {
        Pthread_mutex_lock(&mutex);
        while (count == MAX)
            Pthread_cond_wait(&empty, &mutex);
        put(i);
        Pthread_cond_signal(&full);
        Pthread_mutex_unlock(&mutex);
    }
}

void *consumer(void *arg) {
    int i;
    for (i = 0; i < loops; i++) {
        Pthread_mutex_lock(&mutex);
        while (count == 0)
            Pthread_cond_wait(&full, &mutex);
        int tmp = get();
        Pthread_cond_signal(&empty);
        Pthread_mutex_unlock(&mutex);
        printf("%d\n", tmp);
    }
}
```

wait/signal Best Practices

- it is required to hold a lock when calling wait
- it is a good idea to hold the (same) lock when calling signal
 - not always required for correctness, but avoids bugs in cases where it is needed
- use a loop instead of if to check conditions
 - protects against state being changed between ready and running
 - avoids problems from spurious wakeups (signal waking more than one thread)

When using condition variables there is always a boolean predicate involving shared variables associated with each condition wait that is true if the thread should proceed. Spurious wakeups from the `pthread_cond_wait()` or `pthread_cond_timedwait()` functions may occur. Since the return from `pthread_cond_wait()` or `pthread_cond_timedwait()` does not imply anything about the value of this predicate, the predicate should be re-evaluated upon such return.

Covering Conditions

```
// 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);
    while (bytesLeft < size)
        Pthread_cond_wait(&c, &m);
    void *ptr = ...; // get mem from heap
    bytesLeft -= size;
    Pthread_mutex_unlock(&m);
    return ptr;
}

void free(void *ptr, int size) {
    Pthread_mutex_lock(&m);
    // return some memory to the heap
    bytesLeft += size;
    Pthread_cond_signal(&c);
    Pthread_mutex_unlock(&m);
}
```

if there isn't enough memory available, wait until some is free

more memory is available – wake a thread that is waiting

...but which one? waking up a thread that needs more memory than is available will prevent one that doesn't from running

Covering Conditions

```
// 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);
    while (bytesLeft < size)
        pthread_cond_wait(&c, &m);
    void *ptr = ...; // get mem from heap
    bytesLeft -= size;
    pthread_mutex_unlock(&m);
    return ptr;
}

void free(void *ptr, int size) {
    pthread_mutex_lock(&m);
    // return some memory to the heap
    bytesLeft += size;
    pthread_cond_broadcast(&c);
    pthread_mutex_unlock(&m);
}
```

instead wake up all threads
this is a *covering condition* – waking up all threads when memory is freed covers all of the specific cases where a particular thread needs to wake up

wait/signal Best Practices

- if things only work with broadcast instead of signal, fix the bug instead
 - waking up threads unnecessarily introduces additional overhead
 - reserve broadcast for when it is truly necessary