

Total Points 100, Due Sunday 4/11/21 11:59 pm

Q1 [10 pts]: The following is the `Producer(.)` function in a `BoundedBuffer` implementation. What is the purpose of the mutex in the following? Can we do without the mutex? In what circumstances?

```
Producer(item) {  
    emptySlots.P();  
    mutex.P();  
    Enqueue(item);  
    mutex.V();  
    fullSlots.V();  
}
```

Answer: Make sure the enqueue procedure not been interrupted by other threads. We can not do without the mutex unless we can finish the procedure in an atomic operation in the assembly code level.

Q2 [5 pts]: In the above producer code, can we change the order of the first 2 lines (i.e., the 2 lines with `P(.)`)? If so, why? If not, why not?

Answer: We can not change the order of the first 2 lines. When the queue is full, we need to wait an empty slot in the queue. If we put the “`mutex.P()`” in front of the “`emptySlots.P()`”, other threads can not consume the items in the queue. The queue will always be full. The deadlock situation will happen.

Q3 [5 pts]: How would you change the `Producer` from Q1 into an unbounded buffer, where there is no limit to its capacity. However, make sure that the `Consumer(s)` must still wait for an item to be present in the buffer before they can consume. In other words, the producers never stall, only the consumers stall in the absence of items.

Answer:

```
Producer(item) {  
    mutex.P();  
    Enqueue(item);  
    mutex.V();  
    fullSlots.V();  
}
```

Q4 [10 pts]: Using only condition variables and mutexes, write a `BoundedBuffer` that does not let a consumer consume when the buffer is below 10% of its capacity. You only need to write the `pop()` function.

Answer:

```
vector<char> BoundedBuffer::pop(){  
    unique_lock<mutex> l (mtx);  
    // keep waiting as long as q.size() == 0  
    data_avail.wait (l, [this]{return q.size() >= 0.1*cap;});
```

```

    // pop from the queue
    vector<char> data = q.front();
    q.pop();
    // notify any potential producer(s)
    slot_avail.notify_one ();
    // unlock the mutex so that others can go in
    l.unlock();
    return data;
}

```

Q5 [15 pts]: The following is a thread function that is being run from 5 threads.

```

int x=0;
void ThreadFunc(int &x){
    x += 2;
}

```

Can you have $x=8$? If yes, show the schedule of 5 threads using assembly language that would lead to $x=8$.

You must write functional assembly code as shown in class lectures - pseudo code is not OK

Answer:

```

load r1, x
    load r2, x
    add r2, r2, 2
    store x, r2
add r1, r1, 2
store x, r1
load r1, x
add r1, r1, 2
store x, r1
load r1, x
add r1, r1, 2
store x, r1
load r1, x
add r1, r1, 2
store x, r1

```

Q6 [15 pts]: If we run 5 instances of ThreadA() and 1 instance of ThreadB(), what can be the maximum number of threads active simultaneously in the Critical Section? The mutex is initially unlocked. Note that ThreadB() is buggy and mistakenly unlocks the mutex first instead of locking first. Explain your answer.

<pre> ThreadA(){ mutex.P() /* Start Critical Section */ /* End Critical Section */ mutex.V(); } </pre>	<pre> ThreadB(){ mutex.V() /* Start Critical Section */ /* End Critical Section */ mutex.P(); } </pre>
--	--

Answer: The maximum number is 3.

ThreadA() comes first, locks and mutex=0.

ThreadB() comes, unlocks and mutex=1.

ThreadA() comes, locks and mutex=0.

Other ThreadA() can not come in and there is no other ThreadB() left. So the maximum number active simultaneously in the Critical Section is 3.

Q7 [20 pts]: Consider a multithreaded web crawling and indexing program, which needs to first download a web page and then parse the HTML of that page to extract links and other useful information from it. The problem is both downloading a page and parsing it can be very slow depending on the content. Your goal is to make both these components as fast as possible. First, to speed up downloading, you delegate the task to **m** downloader threads, each with only a portion of the page to download. (Note that this is quite common in real life and a typical web browser does this all the time as long as the server supports this feature. Usually it is done through opening multiple TCP connections with the server and downloading equal sized chunks through each connection). The **M** chunks are downloaded to a single page buffer. Once all the chunks are downloaded into the buffer, you can then start parsing it. However, since you want to speed up parsing as well, you now use **n** parsing threads who again can parse the page independently, and together they take much less time.

By now, you probably see that **M** download threads are acting as Producers and **N** parser threads as Consumers. Additionally, note that the both downloader and parser threads come from a pool of **M** Producer threads and **N** Consumer threads where $M > m$ and $N > n$. Out of many of these, you have to let exactly **m** Producer threads carry out the download and then exactly **n** consumer threads parse, and then the whole cycle will repeat. IOW, in each cycle, you are employing **m** out of **M** Producer worker threads (who are all eagerly waiting) to download the page simultaneously, and then **n** out of **N** Consumers are concurrently parsing the downloaded page. You cannot assume $m=M$ or $n=N$. Assume that you can do a function call `download(URL)` to download the page and `parse(chunk)` to parse a chunk of the page. No need to be any more specific/concrete than that. The main thing of interest is the Producer-Consumer relation.

Look at the given program [1PNC.cpp](#) that works for 1 Producer and n Consumer threads. Be sure to run the program first to see how it behaves. You need to extend the program such that it works for m-producers instead of just 1. Add necessary semaphores to the program. However, you will lose points if you add unnecessary Semaphores or Mutexes. Use [Semaphore.h](#) as the Semaphore class definition. To keep things simple, declare the mutexes as instances of the Semaphore classes well. Test your program to make sure that it is correct. In your submission directory, include a file called **Q7.cpp** that contains the correct program.

Q8 [20 pts]: There are 3 sets of threads A, B, C. First 1 instance of A has to run, then 2 instances of B and then 1 instance of C, then the cycle repeats. This emulates a chain of producer-consumer relationship that we learned in class, but between multiple pairs of threads. Write code to run these sets of threads.

Assumptions and Instructions: There are 100s of A, B, C threads trying to run. Write only the thread functions with proper wait and signal operation in terms of semaphores. You can use the necessary number of semaphores as long as you declare them in global and initialize them properly with correct values. The actual operations done by A, B and C does not really matter. Submit a separate C++ file called **Q8.cpp** that includes the solution.