



Topic Notes: Parallel Programming Intro

Given an multicore/SMT or a computer with multiple processors on separate chips (a *symmetric multiprocessor (SMP)*), how can we make use of the multiple processing units?

This level of parallelism is at a much higher level than the instruction-level parallelism we looked at before. There, the compiler and/or architecture takes a single program made up of a sequential series of instructions, and executes those instructions in parallel in a way that produces the same result as a one-by-one sequential execution of the instructions.

For a computer with multiple processors, we need to provide multiple streams of instructions to be executed by the processors. A single stream of instructions will only make use of one of our processors at a time.

The easiest way to program this systems is to program them just like a regular single-processor system, but to run multiple programs at once. Each program being run will be assigned to a CPU by the operating system.

However, we would like to consider an approach where a single program can make use of these multiple CPUs.

If we are going to do this, we first need to think about how we would break down the problem to be solved into components that can be executed in parallel, then write a program to achieve it.

Consider some examples:

- Taking a census of Troy.

One person doing this would visit each house, count the people, and ask whatever questions are supposed to be asked. This person would keep running counts. At the end, this person has gathered everything.

If there are two people, they can work concurrently. Each visits some houses, and they need to “report in” along the way or at the end to combine their information. But how to split up the work?

- Each person could do what the individual was originally doing, but would check to make sure each house along the way had not yet been counted.
- Each person could start at the city hall, get an address that has not yet been visited, go visit it, then go back to the city hall to report the result and get another address to visit. Someone at city hall keeps track of the cumulative totals. This is nice because neither person will be left without work to do until the whole thing is done. This is the *master-slave* method of breaking up the work.

- The city could be split up beforehand. Each could get a randomly selected collection of addresses to visit. Maybe one person takes all houses with even street numbers and the other all houses with odd street numbers. Or perhaps one person would take everything north of Hoosick St. and the other everything south of Hoosick St. The choice of how to divide up the city may have a big effect on the total cost. There could be excessive travel if one person walks right past a house that has not yet been visited. Also, one person could finish completely while the other still has a lot of work to do. This is a *domain decomposition* approach.
- Grading a stack of exams. Suppose each has several questions. Again, assume two graders to start.
 - Each person could take half of the stack. Simple enough. But we still have the potential of one person finishing before the other.
 - Each person could take a paper from the “ungraded” stack, grade it, then put it into the “graded” stack.
 - Perhaps it makes more sense to have each person grade half of the *questions* instead of half of the exams, maybe because it would be unfair to have the same question graded by different people. Here, we could use variations on the approaches above. Each takes half the stack, grades his own questions, then they swap stacks.
 - Or we form a *pipeline*, where each exam goes from one grader to the next to the finished pile. Some time is needed to start up the pipeline and drain it out, especially if we add more graders. These models could be applied to the census example, if different census takers each went to every house to ask different questions.
 - Suppose we also add in a “grade totaler and recorder” person. Does that make any of the approaches better or worse?
- Adding two $1,000,000 \times 1,000,000$ matrices.
 - Each matrix entry in the sum can be computed independently, so we can break this up any way we like. Could use the master-slave approach, though a domain decomposition would probably make more sense. Depending on how many processes we have, we might break it down by individual entries, or maybe by rows or columns.

In each of these cases, we have taken what we might normally think of as a *sequential* process, and taken advantage of the availability of *concurrent processing* to make use of multiple workers (processing units).

Some Terminology

Sequential Program: sequence of actions that produce a result (statements + variables), called a process, task, or thread (of control). The state of the program is determined by the code, data, and a *single* program counter.

Concurrent Program: two or more processes that work together. Big difference: *multiple* program counters.

To cooperate, the processes need *communication* and *synchronization*, which can be achieved through *shared variables*, or *message passing*

How to Achieve Parallelism

- We need to determine where concurrency is possible, then break up the work accordingly
 - This is easiest if a compiler can do this for you – take your sequential program and extract the concurrency automatically. This is sometimes possible, especially with fixed-size array computations.
 - If the compiler can't do it, it is possible to give “hints” to the compiler to tell it what is safe to parallelize.
 - But often, the parallelization must be done explicitly: the programmer has to create the threads or processes, assign work to them, and manage necessary communication.
-

Finding Concurrency

We find opportunities for parallelism by looking for parts of the sequential program that can be run in any order.

Before we look at the matrix-matrix multiply, we step back and look at a simpler example:

```
1: a = 10;  
2: b = a + 5;  
3: c = a - 3;  
4: b = 7;  
5: a = 3;  
6: b = c - a;  
7: print a, b, c;
```

Which statements can be run in a different order (or concurrently) but still produce the same answers at the end?

- 1 has to happen before 2 and 3, since they depend on a having a value.
- 2 and 3 can happen in either order.
- 4 has to happen after 2, but it can happen before 3.
- 5 has to happen after 2 and 3, but can happen before 4.

- 6 has to happen after 4 (so 4 doesn't clobber its value) and after 5 (because it depends on its value)
- 7 has to happen last.

This can be formalized into a set of rules called *Bernstein's conditions* to determine if a pair of tasks can be executed in parallel:

Two tasks P_1 and P_2 can execute in parallel if all three of these conditions hold:

1. $I_1 \cap O_2 = \emptyset$
2. $I_2 \cap O_1 = \emptyset$
3. $O_1 \cap O_2 = \emptyset$

where I_i and O_i are the input and output sets, respectively, for task i (Bernstein, 1966). The *input set* is the set of variables read by a task and the *output set* is the set of variables modified by a task.

Back to our example, let's see what can be done concurrently.

```

/* initialize matrices, just fill with junk */
for (i=0; i<SIZE; i++) {
    for (j=0; j<SIZE; j++) {
        a[i][j] = i+j;
        b[i][j] = i-j;
    }
}

/* matrix-matrix multiply */
for (i=0; i<SIZE; i++) { /* for each row */
    for (j=0; j<SIZE; j++) { /* for each column */
        /* initialize result to 0 */
        c[i][j] = 0;

        /* perform dot product */
        for(k=0; k<SIZE; k++) {
            c[i][j] = c[i][j] + a[i][k]*b[k][j];
        }
    }
}

sum=0;
for (i=0; i<SIZE; i++) {
    for (j=0; j<SIZE; j++) {
        sum += c[i][j];
    }
}

```

The initialization can all be done in any order – each i and j combination is independent of each other, and the assignment of $a[i][j]$ and $b[i][j]$ can be done in either order.

In the actual matrix-matrix multiply, each $c[i][j]$ must be initialized to 0 before the sum can start to be accumulated. Also, iteration k of the inner loop can only be done after row i of a and column j of b have been initialized.

Finally, the sum contribution of each $c[i][j]$ can be added as soon as that $c[i][j]$ has been computed, and after sum has been initialized to 0.

That *granularity* seems a bit cumbersome, so we might step back and just say that we can initialize a and b in any order, but that it should be completed before we start computing values in c . Then we can initialize and compute each $c[i][j]$ in any order, but we do not start accumulating sum until c is completely computed.

But all of these dependencies in this case can be determined by a relatively straightforward computation. Seems like a job for a compiler! (And in this case, it can be.)

Unfortunately, not everything can be parallelized by the compiler:

If we change the initialization code to:

```
for (i=0; i<SIZE; i++) {
  for (j=0; j<SIZE; j++) {
    if ((i == 0) || (j == 0)) {
      a[i][j] = i+j;
      b[i][j] = i-j;
    }
    else {
      a[i][j] = a[i-1][j-1] + i + j;
      b[i][j] = b[i-1][j-1] + i - j;
    }
  }
}
```

it can't be parallelized, so no matter how many processors we throw at it, we can't speed it up.

Approaches to Parallelism

Automatic parallelism is great, when it's possible. We got it for free (at least once we bought the compiler)! It does have limitations, though:

- some potential parallelization opportunities cannot be detected automatically – can add directives to help
- bigger complication – this executable cannot run on distributed-memory systems

Parallel programs can be categorized by how the cooperating processes communicate with each other:

- **Shared Memory** – some variables are accessible from multiple processes. Reading and writing these values allow the processes to communicate.
- **Message Passing** – communication requires explicit messages to be sent from one process to the other when they need to communicate.

These are functionally equivalent given appropriate operating system support. For example, one can write message-passing software using shared memory constructs, and one can simulate a shared memory by replacing accesses to non-local memory with a series of messages that access or modify the remote memory.

The automatic parallelization we have seen to this point is a shared memory parallelization, though we don't have to think about how it's done. The main implication is that we have to run the parallelized executable on a computer with multiple processors.

Our first tool for explicit parallelization will be shared memory parallelism using threads.

A Brief Intro to POSIX threads

Multithreading usually allows for the use of shared memory. Many operating systems provide support for threads, and a standard interface has been developed: *POSIX Threads* or *pthread*.

A good online tutorial is available at <https://computing.llnl.gov/computing/tutorials/pthreads/>.

You read through this and remember that it's there for reference.

A Google search for “pthread tutorial” yields many others.

Pthreads are available on the Solaris nodes in the cluster, and are standard on most modern Unix-like operating systems.

The basic idea is that we can create and destroy threads of execution in a program, on the fly, during its execution. These threads can then be executed in parallel by the operating system scheduler. If we have multiple processors, we should be able to achieve a speedup over the single-threaded equivalent.

We start with a look at a pthreads “Hello, world” program:

See: `/cs/terescoj/shared/cs2500/examples/pthreadhello`

The most basic functionality involves the creation and destruction of threads:

- `pthread_create(3THR)` – This creates a new thread. It takes 4 arguments. The first is a pointer to a variable of type `pthread_t`. Upon return, this contains a thread identifier that may be used later in a call to `pthread_join()`. The second is a pointer to a `pthread_attr_t` structure that specifies thread creation attributes. In the `pthreadhello` program, we pass in `NULL`, which will request the system default attributes. The third argument is a pointer to a function that will be called when the thread is started. This function

must take a single parameter of type `void *` and return `void *`. The fourth parameter is the pointer that will be passed as the argument to the thread function.

- `pthread_exit(3THR)` – This causes the calling thread to exit. This is called implicitly if the thread function called during the thread creation returns. Its argument is a return status value, which can be retrieved by `pthread_join()`.
- `pthread_join(3THR)` – This causes the calling thread to block (wait) until the thread with the identifier passed as the first argument to `pthread_join()` has exited. The second argument is a pointer to a location where the return status passed to `pthread_exit()` can be stored. In the `pthreadhello` program, we pass in `NULL`, and hence ignore the value.

Prototypes for `pthread` functions are in `pthread.h` and programs need to link with `libpthread.a` (use `-lpthread` at link time). When using the Sun compiler, the `-mt` flag should also be specified to indicate multithreaded code.

A slightly more interesting example:

See: `/cs/terescoj/shared/cs2500/examples/proctree_threads`

This example builds a “tree” of threads to a depth given on the command line. It includes calls to `pthread_self()`. This function returns the thread identifier of the calling thread.

Try it out and study the code to make sure you understand how it works.

A bit of extra initialization is necessary to make sure the system will allow your threads to make use of all available processors. It may, by default, allow only one thread in your program to be executing at any given time. If your program will create up to n concurrent threads, you should make the call:

```
pthread_setconcurrency(n+1);
```

somewhere before your first thread creation. The “+1” is needed to account for the original thread plus the n you plan to create.

You may also want to specify actual attributes as the second argument to `pthread_create()`. To do this, declare a variable for the attributes:

```
pthread_attr_t attr;
```

and initialize it with:

```
pthread_attr_init(&attr);
```

and set parameters on the attributes with calls such as:

```
pthread_attr_setscope(&attr, PTHREAD_SCOPE_PROCESS);
```

I recommend the above setting for threads in Solaris.

Then, you can pass in `&attr` as the second parameter to `pthread_create()`.

Any global variables in your program are accessible to all threads. Local variables are directly accessible only to the thread in which they were created, though the memory can be shared by passing a pointer as part of the last argument to `pthread_create()`.

Brief Intro to Critical Sections

As you may have been shown in other contexts, concurrent access to shared variables can be dangerous.

Consider this example:

See: `/cs/terescoj/shared/cs2500/examples/pthread_danger`

Run it with one thread, and we get 100000. What if we run it with 2 threads? On a multiprocessor, it is going to give the wrong answer! Why?

The answer is that we have concurrent access to the shared variable `counter`. Suppose that two threads are each about to execute `counter++`, what can go wrong?

`counter++` really requires three machine instructions: (i) load a register with the value of `counter`'s memory location, (ii) increment the register, and (iii) store the register value back in `counter`'s memory location. Even on a single processor, the operating system could switch the process out in the middle of this. With multiple processors, the statements really could be happening concurrently.

Consider two threads running the statements that modify `counter`:

Thread A	Thread B
A_1 <code>R0 = counter;</code>	B_1 <code>R1 = counter;</code>
A_2 <code>R0 = R0 + 1;</code>	B_2 <code>R1 = R1 + 1;</code>
A_3 <code>counter = R0;</code>	B_3 <code>counter = R1;</code>

Consider one possible ordering: $A_1 A_2 B_1 A_3 B_2 B_3$, where `counter=17` before starting. Uh oh.

What we have here is a *race condition* that can lead to *interference* of the actions of one thread with another. We need to make sure that when one process starts modifying `counter`, that it finishes before the other can try to modify it. This requires *synchronization* of the processes.

When we run it on a single-processor system, the problem is unlikely to show itself - we almost certainly get the correct sum when we run it. However, there is no guarantee that this would be the case. The operating system could switch threads in the middle of the load-increment-store, resulting in a race condition and an incorrect result.

We need to make those statements that increment `counter` *atomic*. We say that the modification of `counter` is a *critical section*.

There are many solutions to the critical section problem and this is a major topic in an operating systems course. But for our purposes, at least for now, it is sufficient to recognize the problem, and use available tools to deal with it.

The pthread library provides a construct called a *mutex* (short for the *mutual exclusion* that we want to enforce for the access of the counter variable) allows us to ensure that only one thread at a time is executing a particular block of code. We can use it to fix our “danger” program:

See: `/cs/terescoj/shared/cs2500/examples/pthread_nodanger`

We declare a mutex like any other shared variable. It is of type `pthread_mutex_t`. Four functions are used:

- `pthread_mutex_init(3THR)` – initialize the mutex and set it to the unlocked state.
- `pthread_mutex_lock(3THR)` – request the lock on the mutex. If the mutex is unlocked, the calling thread acquires the lock. Otherwise, the thread is blocked until the thread that previously locked the mutex unlocks it.
- `pthread_mutex_unlock(3THR)` – unlock the mutex.
- `pthread_mutex_destroy(3THR)` – destroy the mutex (clean up memory).

A few things to consider about this:

Why isn’t the access to the mutex a problem? Isn’t it just a shared variable itself? – Yes, it’s a shared variable, but access to it is only through the pthread API. Techniques that are discussed in detail in an operating systems course (and that we may discuss more here) are used to ensure that access to the mutex itself does not cause a race condition.

Doesn’t that lock/unlock have a significant cost? – Let’s see. We can time the programs we’ve been looking at:

See: `/cs/terescoj/shared/cs2500/examples/pthread_danger_timed`

See: `/cs/terescoj/shared/cs2500/examples/pthread_nodanger_timed`

Perhaps the cost is too much if we’re going to lock and unlock that much. Maybe we shouldn’t do so much locking and unlocking. In this case, we’re pretty much just going to lock again as soon as we can jump back around through the `for` loop again.

Here’s an alternative:

See: `/cs/terescoj/shared/cs2500/examples/pthread_nodanger_coarse`

In this case, the coarse-grained locking (one thread gets and holds the lock for a long time) should improve the performance significantly. But at what cost? We’ve completely serialized the computation! Only one thread can actually be doing something at a time, so we can’t take advantage of multiple processors. If the “computation” was something more significant, we would need to be more careful about the granularity of the locking.