# Lecture 9 — Dependencies and Speculation

*Patrick Lam*

We talked about C++11 atomics. Someone asked about whether there was a Pthreads equivalent. Nope, not really.

gcc supports atomics via extensions:
https://gcc.gnu.org/onlinedocs/gcc/_005f_005fatomic-Builtins.html

OS X has atomics via OS calls:
https://developer.apple.com/library/mac/documentation/Cocoa/Conceptual/Multithreading/ThreadSafety/ThreadSafety.html

etc...

Reference: http://stackoverflow.com/questions/1130018/unix-portable-atomic-operations

# Dependencies

I've said that some computations appear to be "inherently sequential". We talked about a bunch of real-life analogies:

- must extract bicycle from garage before closing garage door

- must close washing machine door before starting the cycle

- must be called on before answering questions? (sort of)

- students must submit assignment before course staff can mark the assignment

Note that, in this lecture, we are going to assume that memory accesses follow the sequentially consistent memory model. For instance, if you declared all variables to be C++11 atomics, that would be fine. This reasoning is not guaranteed to work in the presence of undefined behaviour, which exists when you have data races.

**Main Idea.** A *dependency* prevents parallelization when the computation $XY$ produces a different result from the computation $YX$.

**Loop- and Memory-Carried Dependencies.** We distinguish between *loop-carried* and *memory-carried* dependencies. In a loop-carried dependency, an iteration depends on the result of the previous iteration. For instance, consider this code to compute whether a complex number $x_0 + iy_0$ belongs to the Mandelbrot set.

```
// Repeatedly square input, return number of iterations before
// absolute value exceeds 4, or 1000, whichever is smaller.
int inMandelbrot(double x0, double y0) {
  int iterations = 0;
  double x = x0, y = y0, x2 = x*x, y2 = y*y;
  while ((x2+y2 < 4) && (iterations < 1000)) {
    y = 2*x*y + y0;
    x = x2 - y2 + x0;
    x2 = x*x; y2 = y*y;
    iterations++;
```

```
    }
    return iterations;
}
```

In this case, it's impossible to parallelize loop iterations, because each iteration *depends* on the $(x, y)$ values calculated in the previous iteration. For any particular $x_0 + iy_0$, you have to run the loop iterations sequentially.

Note that you can parallelize the Mandelbrot set calculation by computing the result simultaneously over many points at once. Indeed, that is a classic "embarassingly parallel" problem, because the you can compute the result for all of the points simultaneously, with no need to communicate.

On the other hand, a memory-carried dependency is one where the result of a computation *depends* on the order in which two memory accesses occur. For instance:

```
int val = 0;

void g() { val = 1; }
void h() { val = val + 2; }
```

What are the possible outcomes after executing `g()` and `h()` in parallel threads?

## RAW, WAR, WAW and RAR

The most obvious case of a dependency is as follows:

```
int y = f(x);
int z = g(y);
```

This is a read-after-write (RAW), or "true" dependency: the first statement writes `y` and the second statement reads it. Other types of dependencies are:

|  | Read 2nd | Write 2nd |
|---|---|---|
| **Read 1st** | Read after read (RAR) | Write after read (WAR) |
|  | No dependency | Antidependency |
| **Write 1st** | Read after write (RAW) | Write after write (WAW) |
|  | True dependency | Output dependency |

The no-dependency case (RAR) is clear. Declaring data immutable in your program is a good way to ensure no dependencies.

Let's look at an antidependency (WAR) example.

```
void antiDependency(int z) {
  int y = f(x);
  x = z + 1;
}
```

```
void fixedAntiDependency(int z) {
  int x_copy = x;
  int y = f(x_copy);
  x = z + 1;
}
```

Why is there a problem?

Finally, WAWs can also inhibit parallelization:

```
void outputDependency(int x, int z) {          void fixedOutputDependency(int x, int z) {
  y = x + 1;                                     y_copy = x + 1;
  y = z + 1;                                     y = z + 1;
}                                              }
```

In both of these cases, renaming or copying data can eliminate the dependence and enable parallelization. Of course, copying data also takes time and uses cache, so it's not free. One might change the output locations of both statements and then copy in the correct output. These are usually more useful when it's not just one access, but some sort of longer computation.

## Loop-carried Dependencies

As we said last time, a loop-carried dependency is one where an iteration depends on the result of the previous iteration. Let's look at a couple of examples.

Initially, a[0] and a[1] are 1. Can we run these lines in parallel?

```
a[4] = a[0] + 1;
a[5] = a[1] + 2;
```

http://www.youtube.com/watch?v=jjXyqcx-mYY. (This one is legit! Really!)

It turns out that there are no dependencies between the two lines. But this is an atypical use of arrays. Let's look at more typical uses.

What about this? (Again, all elements initially 1.)

```
for (int i = 1; i < 12; ++i)
    a[i] = a[i-1] + 1;
```

Nope! We can unroll the first two iterations:

```
a[1] = a[0] + 1
a[2] = a[1] + 1
```

Depending on the execution order, either a[2] = 3 or a[2] = 2. In fact, no out-of-order execution here is safe—statements depend on previous loop iterations, which exemplifies the notion of a *loop-carried dependency*. You would have to play more complicated games to parallelize this.

Now consider this example—is it parallelizable? (Again, all elements initially 1.)

```
for (int i = 4; i < 12; ++i)
    a[i] = a[i-4] + 1;
```

Yes, to a degree. We can execute 4 statements in parallel at a time:

- a[4] = a[0] + 1, a[8] = a[4] + 1

- a[5] = a[1] + 1, a[9] = a[5] + 1

3

- a[6] = a[2] + 1, a[10] = a[6] + 1

- a[7] = a[3] + 1, a[11] = a[7] + 1

We can say that the array accesses have stride 4—there are no dependencies between adjacent array elements. In general, consider dependencies between iterations.

**Larger loop-carried dependency example.**    Now consider the following function.

```
// Repeatedly square input, return number of iterations before
// absolute value exceeds 4, or 1000, whichever is smaller.
int inMandelbrot(double x0, double y0) {
  int iterations = 0;
  double x = x0, y = y0, x2 = x*x, y2 = y*y;
  while ((x2+y2 < 4) && (iterations < 1000)) {
    y = 2*x*y + y0;
    x = x2 - y2 + x0;
    x2 = x*x; y2 = y*y;
    iterations++;
  }
  return iterations;
}
```

How do we parallelize this?

Well, that's a trick question. There's not much that you can do with that function. What you can do is to run this function sequentially for each point, and parallelize along the different points.

As mentioned in class, but one potential problem with that approach is that one point may take disproportionately long. The safe thing to do is to parcel out the work at a finer granularity. There are (unsafe!) techniques for dealing with that too. We'll talk about that later.

TODO L09: Refactor / Live Code Example?

# Breaking Dependencies with Speculation

Recall that computer architects often use speculation to predict branch targets: the direction of the branch depends on the condition codes when executing the branch code. To get around having to wait, the processor speculatively executes one of the branch targets, and cleans up if it has to.

We can also use speculation at a coarser-grained level and speculatively parallelize code. We discuss two ways of doing so: one which we'll call speculative execution, the other value speculation.

**Speculative Execution for Threads.**

The idea here is to start up a thread to compute a result that you may or may not need. Consider the following code:

```
void doWork(int x, int y) {
  int value = longCalculation(x, y);
  if (value > threshold) {
    return value + secondLongCalculation(x, y);
  }
  else {
```

```
      return value;
  }
}
```

Without more information, you don't know whether you'll have to execute `secondLongCalculation` or not; it depends on the return value of `longCalculation`.

Fortunately, the arguments to `secondLongCalculation` do not depend on `longCalculation`, so we can call it at any point. Here's one way to speculatively thread the work:

```
void doWork(int x, int y) {
  thread_t t1, t2;
  point p(x,y);
  int v1, v2;
  thread_create(&t1, NULL, &longCalculation, &p);
  thread_create(&t2, NULL, &secondLongCalculation, &p);
  thread_join(t1, &v1);
  thread_join(t2, &v2);
  if (v1 > threshold) {
    return v1 + v2;
  } else {
    return v1;
  }
}
```

We now execute both of the calculations in parallel and return the same result as before.

Intuitively: when is this code faster? When is it slower? How could you improve the use of threads?

We can model the above code by estimating the probability $p$ that the second calculation needs to run, the time $T_1$ that it takes to run `longCalculation`, the time $T_2$ that it takes to run `secondLongCalculation`, and synchronization overhead $S$. Then the original code takes time

$$T = T_1 + pT_2,$$

while the speculative code takes time

$$T_s = \max(T_1, T_2) + S.$$

**Exercise.** Symbolically compute when it's profitable to do the speculation as shown above. There are two cases: $T_1 > T_2$ and $T_1 < T_2$. (You can ignore $T_1 = T_2$.)

## Value Speculation

The other kind of speculation is value speculation. In this case, there is a (true) dependency between the result of a computation and its successor:

```
void doWork(int x, int y) {
  int value = longCalculation(x, y);
  return secondLongCalculation(value);
}
```

If the result of `value` is predictable, then we can speculatively execute `secondLongCalculation` based on the predicted value. (Most values in programs are indeed predictable).

```
void doWork(int x, int y) {
    thread_t t1, t2;
    point p(x,y);
    int v1, v2, last_value;
    thread_create(&t1, NULL, &longCalculation, &p);
    thread_create(&t2, NULL, &secondLongCalculation,
                  &last_value);
    thread_join(t1, &v1);
    thread_join(t2, &v2);
    if (v1 == last_value) {
      return v2;
    } else {
      last_value = v1;
      return secondLongCalculation(v1);
    }
}
```

Note that this is somewhat similar to memoization, except with parallelization thrown in. In this case, the original running time is
$$T = T_1 + T_2,$$
while the speculatively parallelized code takes time
$$T_s = \max(T_1, T_2) + S + pT_2,$$
where $S$ is still the synchronization overhead, and $p$ is the probability that `v1 != last_value`.

**Exercise.** Do the same computation as for speculative execution.

## When can we speculate?

Speculation isn't always safe. We need the following conditions:

- `longCalculation` and `secondLongCalculation` must not call each other.

- `secondLongCalculation` must not depend on any values set or modified by `longCalculation`.

- The return value of `longCalculation` must be deterministic.

As a general warning: Consider the *side effects* of function calls.

TODO L09: Bibliography