#### Introduction to OpenMP

Chapter 5.1-5. Bryan Mills, PhD

Spring 2017

# OpenMP

- An API for shared-memory parallel programming.
- MP = multiprocessing
- Designed for systems in which each thread or process can potentially have access to all available memory.
- System is viewed as a collection of cores or CPU's, all of which have access to main memory.

#### A shared memory system



#### OpenMP

```
int main() {
    printf("The output:\n");
#pragma omp parallel num_threads(3)
    {
        printf("Hello World\n");
    }
    /* Resume Serial section*/
    printf("Done\n");
}
```

#include<omp.h>

The output: Hello World Hello World Hello World Done

# OpenMP

- API that abstracts away some of the details of pthreads
  - Simple library functions
  - Pragma directives (compiler hints)
  - Environment variables
- Compiler translates openmp code into pthreads calls
- Starts with a **master** thread, calls to openmp create a **team** of **slaves**.

# Pragmas

- Special preprocessor instructions.
- Typically added to a system to allow behaviors that aren't part of the basic C specification.
- Compilers that don't support the pragmas ignore them.
- Always starts with '#' in first column!

```
#pragma omp parallel num_threads(3)
```

# OpenMp pragmas

- # pragma omp parallel
  - Most basic parallel directive.
  - The number of threads that run the following structured block of code is determined by the run-time system.

# Pragma Clause

- Text that modifies a directive.
- The num\_threads clause can be added to a parallel directive.
- It allows the programmer to specify the number of threads that should execute the following block.

# pragma omp parallel num\_threads ( thread\_count )

# Specifying number of threads

- Multiple ways of specifying the number of threads.
  - Explicitly in pragma

#pragma omp parallel num\_threads(3)

– Environment Variable

export OMP\_NUM\_THREADS=3

– Call library function

omp\_set\_num\_threads(int number)1

#### Easy to use

```
void hello() {
  int my_rank = omp_get_thread_num();
  int total = omp get num threads();
  printf("Hello World from rank:%d total:%d\n",
      my rank,
      total);
}
int main() {
  printf("The output:\n");
# pragma omp parallel
  hello();
 printf("Done\n");
}
```

#### Independent Threads

#### Just like pthreads order is not guaranteed

Hello World from rank: 1 total: 4

- Hello World from rank: 3 total: 4
- Hello World from rank: 2 total: 4
- Hello World from rank: 0 total: 4

Or any combination: Hello World from rank: 0 total: 4 Hello World from rank: 2 total: 4 Hello World from rank: 3 total: 4 Hello World from rank: 1 total: 4

# Of note...

- There may be system-defined limitations on the number of threads that a program can start.
- The OpenMP standard doesn't guarantee that this will actually start thread\_count threads.
- Most current systems can start hundreds or even thousands of threads.
- Unless we're trying to start a lot of threads, we will almost always get the desired number of threads.

# Compiler Doesn't Support OpenMP?

int my\_rank = omp\_get\_thread\_num();

# ifdef \_OPENMP int my\_rank = omp\_get\_thread\_num(); # else int my\_rank = 0; # endif

#### The trapezoidal rule



# A First OpenMP Version

1) We identified two types of tasks:

- a) computation of the areas of individual trapezoids, andb) adding the areas of trapezoids.
- 2) There is no communication among the tasks in the first collection, but each task in the first collection communicates with task 1b.
- 3) We assumed that there would be many more trapezoids than cores.

So we aggregated tasks by assigning a contiguous block of trapezoids to each thread (and a single thread to each core).

#### Assignment of trapezoids to threads



#### **Race Condition**

Time	Thread 0	Thread 1		
0	global_result = 0 to register	finish my_result		
1	<pre>my_result = 1 to register</pre>	global_result = 0 to register		
2	add my_result to global_result	my_result = 2 to register		
3	<pre>store global_result = 1</pre>	add my_result to global_result		
4		<pre>store global_result = 2</pre>		

Unpredictable results when two (or more) threads attempt to simultaneously execute:

global\_result += my\_result ;

### Mutual exclusion

# pragma omp critical
global\_result += my\_result ;
only one thread can execute

the following structured block at a time

```
int main(int argc, char* argv[]) {
  double global result = 0.0; /* Store result in global result */
  int a = 0;
                             /* Left and right endpoints
                                                                  */
  int b = 10.0;
  int n = 100000;
                               /* Total number of trapezoids
                                                                  */
# pragma omp parallel num threads(thread count)
  Trap(a, b, n, &global result);
}
void Trap(double a, double b, int n, double* global result p) {
  double my result;
  int my rank = omp get thread num();
  int thread count = omp get num threads();
  double h = (b-a)/n;
  int local n = n/thread count;
  double local a = a + my rank*local n*h;
  double local b = local a + local n*h;
 my result = (f(local a) + f(local b))/2.0;
  for (int i = 1; i \le local n-1; i++) {
    double x = local a + i*h;
   my result += f(x);
  }
 my result = my result*h;
# pragma omp critical
  *global result p += my result;
}
```

# Scope

- In serial programming, the scope of a variable consists of those parts of a program in which the variable can be used.
- In OpenMP, the scope of a variable refers to the set of threads that can access the variable in a parallel block.

# Scope in OpenMP

- A variable that can be accessed by all the threads in the team has shared scope.
- A variable that can only be accessed by a single thread has private scope.
- The default scope for variables declared before a parallel block is shared.

We need this more complex version to add each thread's local calculation to get *global\_result*.

void Trap(double a, double b, int n, double\* global\_result\_p)

Although we'd prefer this.

double Trap(double a, double b, int n)

global result = Trap(a, b, n);

If we use this, there's no critical section!

```
global_result += Trap(double a, double b, int n)
```

```
If we fix it like this...
```

```
global_result = 0.0;
#pragma omp parallel
{
    pragma omp critical
    global_result += Trap(a, b, n);
}
```

... we force the threads to execute sequentially.

We can avoid this problem by declaring a private variable inside the parallel block and moving the critical section after the function call.

```
global_result = 0.0;
#pragma omp parallel
{
   double my_result = Trap(a, b, n);
   pragma omp critical
   global_result += my_result;
}
```

# **Reduction operators**

- A reduction operator is a binary operation (such as addition or multiplication).
- A reduction is a computation that repeatedly applies the same reduction operator to a sequence of operands in order to get a single result.
- All of the intermediate results of the operation should be stored in the same variable: the reduction variable.

A reduction clause can be added to a parallel directive.

reduction(<operator>: <variable list>)

+, \*, -, &, |, ^, &&, ||

```
global_result = 0.0;
#pragma omp parallel reduction(+: global_result)
global_result += Trap(a, b, n);
```

# Parallel for

- Forks a team of threads to execute the following structured block.
- However, the structured block following the parallel for directive must be a for loop.
- Furthermore, with the parallel for directive the system parallelizes the for loop by dividing the iterations of the loop among the threads.

## Parallel for



# Legal forms for parallelizable for statements

	(			index++	
				++index	
			index < end	index	
			index <= end	index	
for	index = start	;	index >= end ;	index += incr	
			index > end	index -= incr	
				<pre>index = index + ind</pre>	r
				index = incr + inde	X
				index = index - ind	r,

#### Caveats

- The variable index must have integer or pointer type (e.g., it can't be a float).
- The expressions start, end, and incr must have a compatible type. For example, if index is a pointer, then incr must have integer type.
- The expressions start, end, and incr must not change during execution of the loop.
- During execution of the loop, the variable index can only be modified by the "increment expression" in the for statement.

#### Data dependencies



# What happened?

- 1. OpenMP compilers don't check for dependences among iterations in a loop that's being parallelized with a parallel for directive.
- 2. A loop in which the results of one or more iterations depend on other iterations cannot, in general, be correctly parallelized by OpenMP.

#### Estimating $\pi$

$$\pi = 4 \left[ 1 - \frac{1}{3} + \frac{1}{5} - \frac{1}{7} + \cdots \right] = 4 \sum_{k=0}^{\infty} \frac{(-1)^k}{2k+1}$$

```
double factor = 1.0;
double sum = 0.0;
for (k = 0; k < n; k++) {
   sum += factor/(2*k+1);
   factor = -factor;
}
pi_approx = 4.0 * sum;
```

# **OpenMP solution #1**

loop dependency
double factor = 1.0;
double sum = 0.0;
#pragma omp parallel for reduction(+:sum)
for (k = 0; k < n; k++) {
 sum += factor/(2\*k+1);
 factor = -factor;
}
pi approx = 4.0 \* sum;</pre>

# OpenMP solution #2

```
double factor = 1.0;
double sum = 0.0;
#pragma omp parallel for reduction(+:sum) \
    private(factor)
for (k = 0; k < n; k++) {
    sum += factor/(2*k+1);
    factor = -factor;
}
pi_approx = 4.0 * sum;
    Insures factor has
    private scope.
```

# The default clause

• Lets the programmer specify the scope of each variable in a block.

default(none)

 With this clause the compiler will require that we specify the scope of each variable we use in the block and that has been declared outside the block.

#### The default clause

```
double factor = 1.0;
double sum = 0.0;
#pragma omp parallel for reduction(+:sum) \
    default(none) private(factor)
for (k = 0; k < n; k++) {
    sum += factor/(2*k+1);
    factor = -factor;
}
pi approx = 4.0 * sum;
```