Parallel programming OpenMP part 1







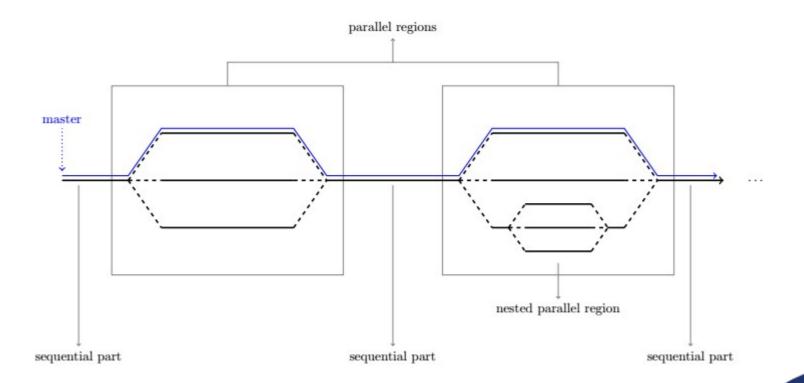
Introduction to OpenMP

- OpenMP (Open Multi-Processing) provides constructs for parallel programming in C++, C, and Fortran on Linux, MacOS, and Windows.
- A sequential code is transformed to a parallel one by adding compiler pragmas, so if a compiler does not support OpenMP, the pragmas are skipped and the output is a sequential program.
 - **OpenMP manual: 1.3 Execution model:** *more detailed (first paragraph)*
- OpenMP is widely used in software like Blender, fftw, OpenBLAS (MATLAB, Python libraries – e.g. NumPy or Scipy, R), and eigen to accelerate computations.
- It is easy to use!



Execution model

- OpenMP program starts with a single thread only (master thread).
- It is executed sequentially until it reaches a parallel region defined by OpenMP pragma.
- At the entry of parallel region, new team of threads is created. Each thread executes concurrently with the others sharing the work.





Using OpenMP

Include header file

```
#include <omp.h>
```

Cmake (multi-platform)

```
find_package(OpenMP)
if (OPENMP_FOUND)
    set(CMAKE_C_FLAGS "${CMAKE_C_FLAGS} ${OpenMP_C_FLAGS}")
    set(CMAKE_CXX_FLAGS "${CMAKE_CXX_FLAGS} ${OpenMP_CXX_FLAGS}")
    set(CMAKE_EXE_LINKER_FLAGS "${CMAKE_EXE_LINKER_FLAGS} ${OpenMP_EXE_LINKER_FLAGS}")
endif()
```

• **GCC** g++ -fopenmp ...



Hello world! In OpenMP...

lab codes/src/HelloWorld.cpp

- Demonstrated task:
 - 1) Vector multiplied by scalar value
 - 2) Sum of vector components



Runtime Library Routines

omp_get_num_procs()

Returns the number of processors that are available to the program

omp_get_num_threads()

Returns the number of threads that are currently in the team executing the parallel region from which it is called

omp_get_thread_num()

Returns the calling thread index within the current team

• ...



#pragma omp parallel

```
#pragma omp parallel

{
    cout << "This is thread" << omp_get_thread_num() << " speaking" << endl;
}

cout << "Parallel block finished" << endl;

Waits for threads to finish (barrier)
```

Output:

```
This is thread 0 speaking
This is thread 3 speaking
This is thread 2 speaking
This is thread 1 speaking
Parallel block finished
```



#pragma omp parallel

```
#pragma omp parallel num_threads(8)
{
    cout << "This is thread " << omp_get_thread_num() << " speaking" << endl;
}</pre>
```

Output:

```
This is thread 0 speaking This is thread 3 speaking This is thread 6 speaking This is thread 1 speaking This is thread 2 speaking This is thread 7 speaking This is thread 4 speaking This is thread 5 speaking This is thread 5 speaking
```



#pragma omp single

```
Block performed by single thread
#pragma omp parallel
    cout << "This is thread " << omp_get_thread_num() << " speaking" << endl;</pre>
    #pragma omp single
        cout << "The single part was done by thread " << omp_get_thread_num() << endl;</pre>
}
     Output:
    This is thread 3 speaking
    The single part was done by thread 3
    This is thread 1 speaking
    This is thread 2 speaking
    This is thread 0 speaking
```



#pragma omp sections

```
#pragma omp parallel
    Each section is performed by only one thread
        #pragma omp section
            cout << "section 1, first: " << omp get thread num() << endl;
            cout << "section 1, second: " << omp_get_thread_num() << endl;</pre>
        #pragma omp section
            cout << "section 2, first: " << omp get thread num() << endl;
            cout << "section 2, second: " << omp_get_thread_num() << endl;</pre>
                                      Waits for threads to finish (barrier).
                                      Can be changed by
                                      #pragma omp sections nowait
Output:
section 2, first: 0
section 2, second: 0
section 1, first: 1
section 1, second: 1
```



6

#pragma omp critical

```
int sum;
#pragma omp parallel
                                        Critical region, performed by all threads
                                         but not at once (mutual exclusion)
    #pragma omp critical
        cout << "Thread " << omp get thread num() << " in critical region" << endl;</pre>
        sum += omp get thread num();
cout << sum << endl;</pre>
 Output:
Thread 1 in critical region
Thread 0 in critical region
Thread 3 in critical region
Thread 2 in critical region
```



#pragma omp barrier

```
#pragma omp parallel

cout << "Before barrier thread " << omp_get_thread_num() << endl;

#pragma omp barrier

cout << "After barrier thread " << omp_get_thread_num() << endl;
}</pre>
```

Output:

```
Before barrier thread 0
Before barrier thread 3
Before barrier thread 1
Before barrier thread 2
After barrier thread 1
After barrier thread 2
After barrier thread 0
After barrier thread 3
```



Example

- Write a function for computing vector normalization. Split the vector into two halves, each one is processed a separate section.
 - You may use skeletonlab_codes/src/VectorNormalization.cpp

To remind normalization:

$$||\vec{\mathbf{a}}|| = \sqrt{a_1^2 + a_2^2 + \dots + a_n^2}; \ \vec{\mathbf{a}} \in \mathbb{R}^n$$

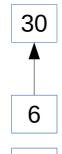
$$\vec{\mathbf{n}} = \frac{\vec{\mathbf{a}}}{||\vec{\mathbf{a}}||} = \left(\frac{a_1}{||\vec{\mathbf{a}}||}; \frac{a_2}{||\vec{\mathbf{a}}||}; \dots; \frac{a_n}{||\vec{\mathbf{a}}||}\right)$$



Sequential summing of matrix rows

```
3
                                                                     2
vector<vector<double>> matrix;
                                                                     3
                                                               5
                                                         4
                                                                     5
                                                         1
                                                               6
                                                                     3
                                                                                 2
vector<double> rowSums(matrix.size(), 0);
                                                        6
                                                                           1
for (int i = 0; i < matrix.size(); i++) {
    for (int j = 0; j < matrix[i].size(); j++) {</pre>
                                                        12
                                                                           5
                                                                                 3
                                                                     4
        rowSums[i] += matrix[i][j];
                                                        12
                                                                                 5
                                                                     1
                                                                           6
```

```
double sum = 0.0;
for (int i = 0; i < matrix.size(); i++) {
    sum += rowSums[i];
}</pre>
```







Parallel summing of matrix rows

```
Each iteration of for loop performed
#pragma omp parallel
                                       by a thread (in parallel) from the team
    #pragma omp for
    for (int i = 0; i < matrix.size(); i++) {
         for (int j = 0; j < matrix[i].size(); j++) {</pre>
             rowSums[i] += matrix[i][j];
A shorter code...
#pragma omp parallel for
for (int i = 0; i < matrix.size(); i++) {</pre>
   \Deltafor (int j = 0; j < matrix[i].size(); j++) {
        rowSums[i] += matrix[i][j];
```

Question: what happens if you write

the pragma on the inner loop?

15 / 20



If clause

```
#pragma omp parallel for if(matrix.size() >= 10)
for (int i = 0; i < matrix.size(); i++) {
    for (int j = 0; j < matrix[i].size(); j++) {
        rowSums[i] += matrix[i][j];
    }
}</pre>
```

Threads are only created for large matrices. Small matrices are summed sequentially since it does not pays off to create threads.



Reduction

Parallel aggregation of an expression, e.g., a sum

Operators:

```
sum = rowSums[0] + rowSums[1] + rowSums[2] + ... + rowSums[matrix.size() - 1];
```

```
double sum = 0.0;
```

```
double sum = 0.0;
#pragma omp parallel for reduction(+:sum)
for (int i = 0; i < matrix.size(); i++) {
    sum += rowSums[i];</pre>
```

List of variables: var₁, var₂, ..., var_n

Useful for doing multiple reductions at once



Collapse

Collapse for loops into one for distribution of the work among threads

```
int numRows = matrix.size();
int numCols = matrix[0].size();

double sum = 0.0;
#pragma omp parallel for collapse(2) reduction(+:sum)
for (int i = 0; i < numRows; i++) {
    for (int j = 0; j < numCols; j++) {
        sum += matrix[i][j];
    }
}</pre>
```



Data sharing

```
int a = 10; 
int b = 100; 
#pragma omp parallel for
for (int i = 0; i < 10; i++) {
   int c = a * b + i;
}</pre>
Thread private
```

- The sharing can be stated explicitly as a clause
 - #pragma omp parallel for private(a, b)
 - Variables a and b are private to each thread (without global initialization)
 - #pragma omp parallel for firstprivate(a, b)
 - Variables a and b are private to each thread (with global initialization)
 - #pragma omp parallel for shared(a, b)
 - Variables a and b are shared among threads
- The default policy can be set to
 - #pragma omp parallel for default(shared)
 - By default, all the variables outside of the parallel section are shared
 - #pragma omp parallel for default(none)
 - The programmer must explicitly state the sharing policy of the variables



Example

- Vector normalization using parallel for (reduction, critical section ...)
- Computation of pi using estimating the value using Monte Carlo

Samples from uniform distribution

 $\frac{4 \cdot numSamplesInCircle}{totalNumSamples} = \pi$

Derivation:

$$\frac{\text{area of the circle}}{\text{area of the square}} = \frac{\pi r^2}{4r^2} = \frac{\pi}{4}$$

$$\frac{\pi}{4} = \frac{\text{no. of points generated inside the circle}}{\text{total no. of points generated or no. of points generated inside the square}}$$

$$\pi = 4 * \frac{\text{no. of points generated inside the circle}}{\text{no. of points generated inside the square}}$$