# Parallel Programming
## Single-core optimization, MPI, OpenMP, and hybrid programming

**Nicolas Richart**
**Emmanuel Lanti**

Course based on V. Keller's lecture notes

27th of October - 1st of November 2024

SCITAS

# OpenMP

- Understand the context of shared memory
- Understand more in detail the architecture of a node
- Get familiar with the OpenMP execution and memory model
- Getting some speedup with Task Level Parallelism

- October 1997: Fortran version 1.0
- Late 1998: C/C++ version 1.0
- June 2000: Fortran version 2.0
- April 2002: C/C++ version 2.0
- June 2005: Combined C/C++ and Fortran version 2.5
- May 2008: Combined C/C++ and Fortran version 3.0
- **July 2011: Combined C/C++ and Fortran version 3.1**
- July 2013: Combined C/C++ and Fortran version 4.0
- November 2015: Combined C/C++ and Fortran version 4.5
- November 2018: Combined C/C++ and Fortran version 5.0
- November 2020: Combined C/C++ and Fortran version 5.1
- November 2021: Combined C/C++ and Fortran version 5.2

- Specification:
  - ▶ Full specification
  - ▶ RefCard
- Terms:

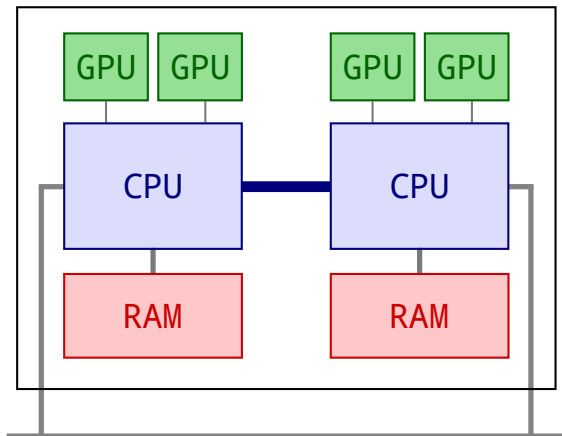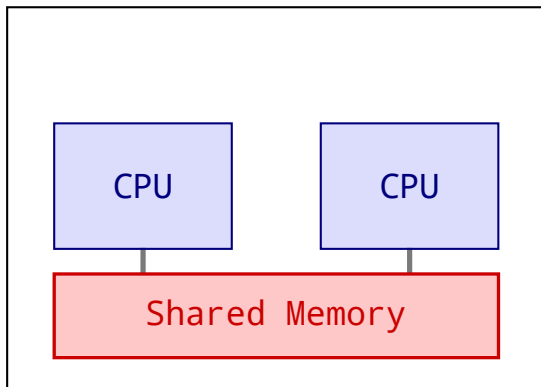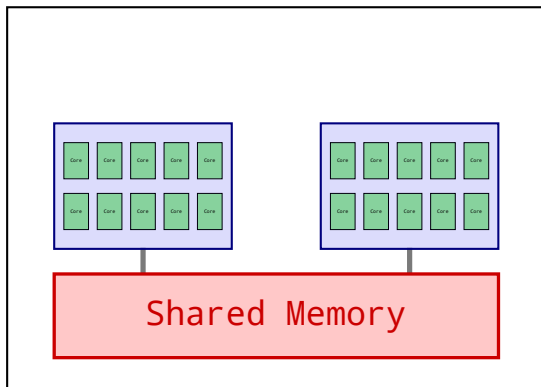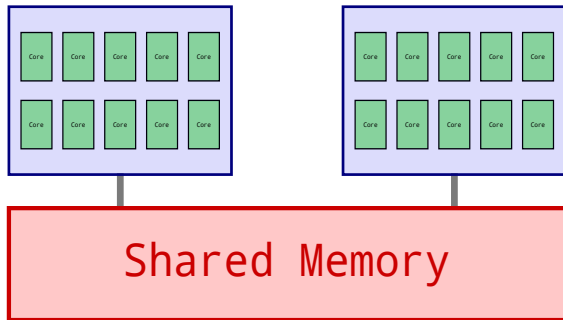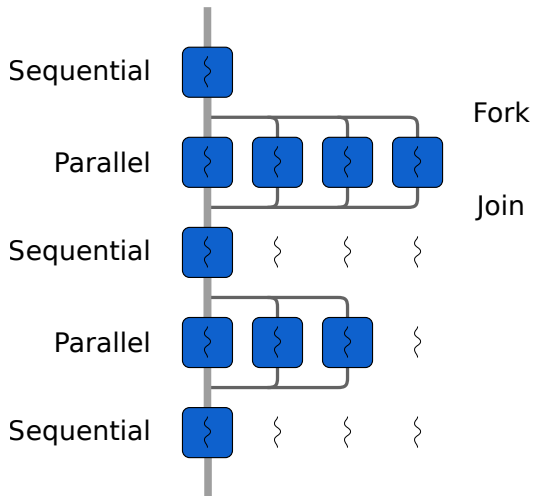|  |  |
|---|---|
| thread | an execution entity with a stack and a static memory (*threadprivate memory*) |
| OpenMP thread | a *thread* managed by the OpenMP runtime |
| processor | an hardware unit on which one or more *OpenMP thread* can execute |
| directive | a base language mechanism to specify OpenMP program behavior |
| construct | an OpenMP executable directive and the associated statement, loop nest or structured block, if any, not including the code in any called routines. That is, the lexical extent of an executable directive. |

- OpenMP directives are written as pragmas: `#pragma omp`
- Use the conditional compilation flag `#if defined _OPENMP` for the preprocessor

EPFL

- OpenMP directives are written as pragmas: `#pragma omp`
- Use the conditional compilation flag `#if defined _OPENMP` for the preprocessor

- Compilation using the GNU compiler:

```
$> g++ -fopenmp ex1.c -o ex1
```

- Compilation using the Intel compiler:

```
$> icpc -qopenmp ex1.c -o ex1
```

**openmp/hello.cc**

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

int main() {

# pragma omp parallel
  {
    auto mysize = omp_get_num_threads();
    auto myrank = omp_get_thread_num();
    std::printf("Hello from thread %i out of %i\n", myrank, mysize);
  }

  return 0;
}
```

**openmp/hello.cc**

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

int main() {

#pragma omp parallel
  {
    auto mysize = omp_get_num_threads();
    auto myrank = omp_get_thread_num();
    std::printf("Hello from thread %i out of %i\n", myrank, mysize);
  }

  return 0;
}
```

```
$ OMP_NUM_THREADS=4 ./openmp/hello
Hello from thread 2 out of 4
Hello from thread 1 out of 4
Hello from thread 0 out of 4
Hello from thread 3 out of 4
```

EPFL

openmp/hello_cond.cc

```cpp
 6  int main() {
 7    int mysize = 1;
 8    int myrank = 0;
 9
10  #if defined(_OPENMP)
11  #pragma omp parallel
12    {
13      mysize = omp_get_num_threads();
14      myrank = omp_get_thread_num();
15  #endif
16      std::printf("Hello from thread %i out of %i\n", myrank, mysize);
17  #if defined(_OPENMP)
18    }
19  #endif
20    return 0;
21  }
```

EPFL

- Default implementation dependent (usually max hardware thread)
- At runtime in the code

```
1  omp_set_num_threads(nthreads);
```

- With en environment variable

```
$> export OMP_NUM_THREADS=4
```

This is the mother of all constructs in OpenMP. It starts a parallel execution.

Syntax

```
#pragma omp parallel [clause[[,] clause]...]
{
   structured-block
}
```

where *clause* is one of the following:

- **if** or **num_threads** : conditional clause
- **default(private** | **firstprivate** | **shared** | **none)** : default data scoping
- **private(***list***)**, **firstprivate(***list***)**, **shared(***list***)** or **copyin(***list***)** : data scoping
- **reduction(***operator* **:** *list***)**

- In the `pi.cc` add a function call to get the number of threads.
- Compile using the porper options for OpenMP
- Test that it works by varying the number of threads **export OMP_NUM_THREADS**
- To vary the number of threads in a **sbatch** job you can set the number of threads to the number of cpus per task.

```bash
#!/bin/bash
#SBATCH -c <nthreads>

export OMP_NUM_THREADS=$SLURM_CPUS_PER_TASK
<my_openmp_executable>
```

Work-sharing constructs are possible in three "flavours" :

- **sections** construct
- **single** construct
- **workshare** construct (only in Fortran)

Syntax

```c
#pragma omp [parallel] sections [clause]
{
    #pragma omp section
    {
        code_block
    }
}
```

where *clause* is one of the following:

- **private(*list*)**, **firstprivate(*list*)**, **lastprivate(*list*)**
- **reduction(*operator* : *list*)**
- Each **section** within a **sections** construct is assigned to one and only one thread

openmp/sections.cc

```
 6  # pragma omp parallel sections num_threads(4)
 7    {
 8  # pragma omp section
 9      std::printf("Thread %i handling section 1\n", omp_get_thread_num());
10  # pragma omp section
11      std::printf("Thread %i handling section 2\n", omp_get_thread_num());
12  # pragma omp section
13      std::printf("Thread %i handling section 3\n", omp_get_thread_num());
14    }
```

openmp/sections.cc

```
 6  #pragma omp parallel sections num_threads(4)
 7    {
 8  #pragma omp section
 9      std::printf("Thread %i handling section 1\n", omp_get_thread_num());
10  #pragma omp section
11      std::printf("Thread %i handling section 2\n", omp_get_thread_num());
12  #pragma omp section
13      std::printf("Thread %i handling section 3\n", omp_get_thread_num());
14    }
```

```
$ ./openmp/sections
Thread 0 handling section 1
Thread 1 handling section 2
Thread 2 handling section 3
```

Only one thread (usualy the first entering thread) executes the **single** region.

**Syntax**

```
1  # pragma omp single [clause[[,] clause] ...]
2  {
3    structured-block
4  }
```

where *clause* is one of the following:

- **private(*list*)**, **firstprivate(*list*)**
- **nowait**

Only the master thread execute the section. It can be used in any OpenMP construct

### Syntax

```
1  #pragma omp master
2  {
3      structured-block
4  }
```

Parallelization of the following loop

### Syntax

```
1  # pragma omp for [clause[[,] clause] ... ]
2  {
3    for-loop
4  }
```

where *clause* is one of the following:

- **schedule(*kind[, chunk_size]*)**
- **collapse(*n*)**
- **ordered**
- **private(*list*)**, **firstprivate(*list*)**, **lastprivate(*list*)**

openmp/for.cc

```
 6  #pragma omp parallel num_threads(2)
 7    {
 8      auto myrank = omp_get_thread_num();
 9  #pragma omp for
10      for (int i = 0; i < 6; ++i) {
11        std::printf("Thread %i handling i=%i\n", myrank, i);
12      }
13    }
```

openmp/for.cc

```cpp
6  #pragma omp parallel num_threads(2)
7    {
8      auto myrank = omp_get_thread_num();
9  #pragma omp for
10     for (int i = 0; i < 6; ++i) {
11       std::printf("Thread %i handling i=%i\n", myrank, i);
12     }
13   }
```

```
$ ./openmp/for
Thread 0 handling i=0
Thread 0 handling i=1
Thread 0 handling i=2
Thread 1 handling i=3
Thread 1 handling i=4
Thread 1 handling i=5
```

- Add a **parallel for** work sharing construct around the integral computation
- Run the code
- Run the code
- Run the code
- What can you observe on the value of $\pi$ ?

Restricts execution of the associated structured block to a single thread at a time

### Syntax

```
1  #pragma omp critical [(name) [[,] hint(hint-expression)]]
2  {
3      structured-block
4  }
```

- **name** optional to identify the construct
- **hint**-**expression** information on the expected execution
    - `omp_sync_hint_none`
    - `omp_sync_hint_uncontended`
    - `omp_sync_hint_contended`
    - `omp_sync_hint_nonspeculative`
    - `omp_sync_hint_speculative`

- To solve the data race condition from the previous exercise we can protect the computation of the sum
- Add a **critical** directive to protect the sum
- Run the code
- What can you observe on the execution time while varying the number of threads

Specifies an explicit barrier.

### Syntax

```
1  # pragma omp barrier
```

Ensures a specific storage location is accessed atomically.

> **Syntax**
> ```
> 1  # pragma omp atomic [clause[[,] clause] ... ]
> 2  statement
> ```

where *clause* is one of the following:

- *atomic-clauses* **read**, **write**, **update**
- *memory-order-clauses* **seq_cst**, **acq_rel**, **releases**, **acquire**, **relaxed**
- or one of **capture**, **compare**, **hint(hint-expression)**, **fail(seq_cst | acquire | relaxed)**, or **weak**

- Most common source of errors
- Determine which variables are **private** to a thread, which are **shared** among all the threads
- In case of a **private** variable the variable values can be defined using:
  - ▶ **firstprivate** defines the value when entering the region
  - ▶ **lastprivate** defines the value when exiting the region (OpenMP 5.1 in C/C++)
- **default(private | firstprivate | shared | none)** can be specified
  **default(none)** means each variables should appear in a **shared** or **private** list

These attributes determines the scope (visibility) of a single or list of variables

**Syntax**

```
1  shared(list1), private(list2)
```

- The **private** clause: the data is private to each thread and non-initialized. Each thread has its own copy. `#pragma omp parallel private(i)`
- The **shared** clause: the data is shared among all the threads. It is accessible (and non-protected) by all the threads simultaneously. `#pragma omp parallel shared(array)`

These clauses determines the attributes of the variables within a parallel region:

> Syntax
>
> 1 `firstprivate(list1), lastprivate(list2)`

- The **firstprivate** super-set of **private**, variable is initialized to a copie of variable before the region
- The **lastprivate** super-set of **private** the value of the last thread exiting the region is copied

openmp/private.cc

```
8    std::printf("Thread %i sees, a, b, c: %i, %i, %g (before)\n",
9                omp_get_thread_num(), a, b, c);
10
11  #pragma omp parallel num_threads(3), private(a), firstprivate(b)
12    {
13      std::printf("Thread %i sees, a, b, c: %i, %i, %g (inside)\n",
14                  omp_get_thread_num(), a, b, c);
15      c = -1e-3;
16    }
17
18    std::printf("Thread %i sees, a, b, c: %i, %i, %g (after)\n",
19                omp_get_thread_num(), a, b, c);
```

**openmp/private.cc**

```cpp
 8    std::printf("Thread %i sees, a, b, c: %i, %i, %g (before)\n",
 9                omp_get_thread_num(), a, b, c);
10
11  #pragma omp parallel num_threads(3), private(a), firstprivate(b)
12    {
13      std::printf("Thread %i sees, a, b, c: %i, %i, %g (inside)\n",
14                  omp_get_thread_num(), a, b, c);
15      c = -1e-3;
16    }
17
18    std::printf("Thread %i sees, a, b, c: %i, %i, %g (after)\n",
19                omp_get_thread_num(), a, b, c);
```

```
$ ./openmp/private
Thread 0 sees, a, b, c: 1, 2, 3 (before)
Thread 0 sees, a, b, c: 1839769744, 2, 3 (inside)
Thread 1 sees, a, b, c: 12789424, 2, 3 (inside)
Thread 2 sees, a, b, c: 12801392, 2, -0.001 (inside)
Thread 0 sees, a, b, c: 1, 2, -0.001 (after)
```

- Create a local variable per thread
- Make each thread compute it's own sum
- After the computation of the integral use a **critical** directive to sum the local sum to a **shared** sum

### Syntax

```
1  reduction(reduction-identifier : list)
```

- *reduction-identifier*: one of the operation +, −, **\***, **&**, |, ^, **&&**, ||
- *list* item on which the reduction applies
- example: `#pragma omp for reduction(+: sum)`

- Use the **reduction** clause
- Compare the timings to the previous versions

### Syntax

```
1  schedule([modifier [, modifier] : ] kind [, chunk_size])
```

- *kind*
  - **static** iterations divided in chunks sized **chunk_size** assigned to threads in a round-robin fashion
  - **dynamic** iterations divided in chunks sized **chunk_size** assigned to threads when they request them until no chunk remains to be distributed
  - **guided** iterations divided in chunks sized **chunk_size** assigned to threads when they request them. Size of chunks is proportional to the remaining unassigned chunks.
  - **auto** The decisions is delegated to the compiler and/or the runtime system
  - **runtime** The decisions is delegated to the runtime system based on ICV

EPFL

---

Syntax

```
1  collapse(n)
```

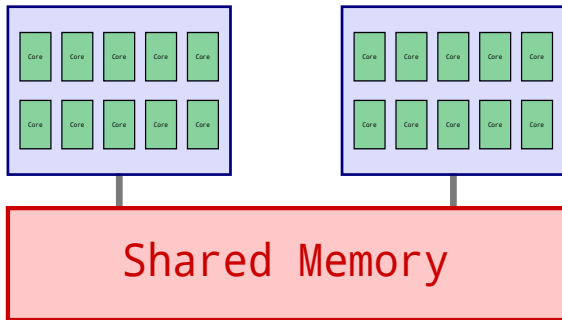Specifies how many loop are combine into a logical space
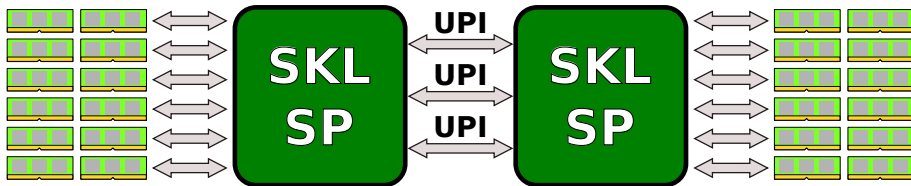
openmp/dgemm.cc

```
34  #pragma omp parallel for collapse(1) schedule(static, N / nthreads)
35    for (int i = 0; i < N; ++i)
36      for (int j = 0; j < N; ++j)
37        for (int k = 0; k < N; ++k)
38          C[i * N + j] += A[i * N + k] * B[k * N + j];
```

**openmp/dgemm.cc**

```
34 #pragma omp parallel for collapse(1) schedule(static, N / nthreads)
35   for (int i = 0; i < N; ++i)
36     for (int j = 0; j < N; ++j)
37       for (int k = 0; k < N; ++k)
38         C[i * N + j] += A[i * N + k] * B[k * N + j];
```
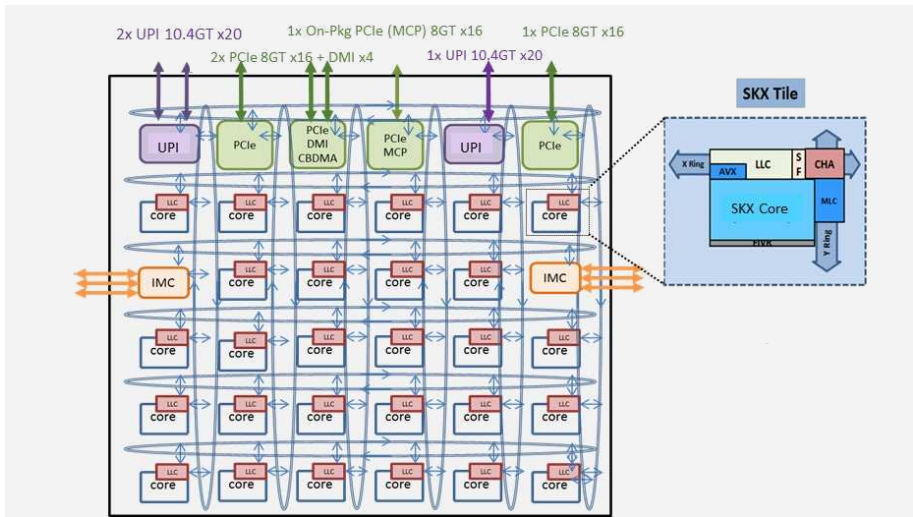
```
$ OMP_NUM_THREADS=1 ../build/openmp/dgemm
DGEMM with 1 threads, collapse(1): 21.1209 GFLOP/s (verif 2)
$ OMP_NUM_THREADS=2 ../build/openmp/dgemm
DGEMM with 2 threads, collapse(1): 40.2308 GFLOP/s (verif 2)
$ OMP_NUM_THREADS=4 ../build/openmp/dgemm
DGEMM with 4 threads, collapse(1): 72.7659 GFLOP/s (verif 2)
$ OMP_NUM_THREADS=1 ../build/openmp/dgemm
```
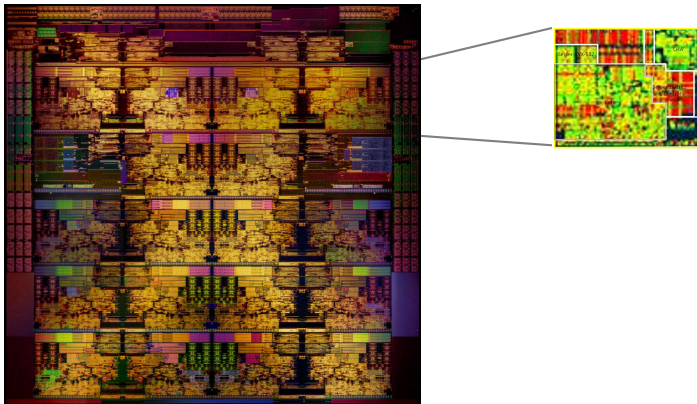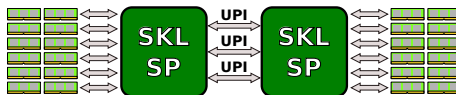
```
DGEMM with 1 threads, collapse(2): 20.358 GFLOP/s (verif 2)
$ OMP_NUM_THREADS=2 ../build/openmp/dgemm
DGEMM with 2 threads, collapse(2): 40.0818 GFLOP/s (verif 2)
$ OMP_NUM_THREADS=4 ../build/openmp/dgemm
DGEMM with 4 threads, collapse(2): 72.4462 GFLOP/s (verif 2)
```
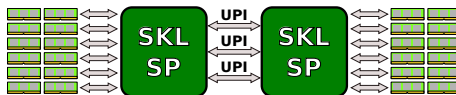
- One thread can only saturate 1 channel
- On memory bound code bandwidth saturate when # of threads $\sim$ # of channels
- If memory allocated on the other processor memory, data go through CPU interconnect (UPI $3 \times 10.4\,\text{GT/s}$)

- One thread can only saturate 1 channel
- On memory bound code bandwidth saturate when # of threads $\sim$ # of channels
- If memory allocated on the other processor memory, data go through CPU interconnect (UPI $3 \times 10.4\,\text{GT/s}$)
- How to mitigate this effects ?
  - ▶ Loop schedule
  - ▶ Memory first touch
  - ▶ Thread placements

- The variable `OMP_PLACES` describes these places in terms of the available hardware.
- The variable `OMP_PROC_BIND` describes how threads are bound to OpenMP places
- The variable `OMP_DISPLAY_AFFINITY` helps to debug the affinity

```
$ OMP_NUM_THREADS=4 OMP_DISPLAY_AFFINITY=true ./openmp/hello
OMP: pid 2115280 tid 2115280 thread 0 bound to OS proc set {0-31}
OMP: pid 2115280 tid 2115285 thread 3 bound to OS proc set {0-31}
OMP: pid 2115280 tid 2115284 thread 2 bound to OS proc set {0-31}
Hello from thread 0 out of 4
Hello from thread 3 out of 4
OMP: pid 2115280 tid 2115283 thread 1 bound to OS proc set {0-31}
Hello from thread 1 out of 4
Hello from thread 2 out of 4
```

Possible values for `OMP_PLACES` where each place corresponds to:

threads a single hardware thread on the device.

cores a single core (having one or more hardware threads) on the device.

ll_caches a set of cores that share the last level cache on the device.

numa_domains a set of cores for which their closest memory on the device is:
- the same memory; and
- at a similar distance from the cores.

sockets a single socket (consisting of one or more cores) on the device.

Possible values for OMP_PROC_BIND:

false   threads not bonded

true   threads are bonded (implementation dependant)

primary   collocate threads with the primary thread

close   place threads close to the master in the places list

spread   spread out threads as much as possible

- Memory is organized in pages
- When allocating data "nothing" happens
- Pages are allocated on the memory associated to the first thread initializing it

- Memory is organized in pages
- When allocating data "nothing" happens
- Pages are allocated on the memory associated to the first thread initializing it

- To mitigate the problem, initialize the arrays in same order they are accessed

- Data race:
  - ▶ Data accessed by multiple threads without protection
  - ▶ Lead to undetermined results

- Data race:
  - ▶ Data accessed by multiple threads without protection
  - ▶ Lead to undetermined results

- False sharing
  - ▶ Data smaller than cache-line size
  - ▶ Multiple threads accessing data in the same cache line will poison each other caches

Sub set of the routines in OpenMP

- `omp_get_num_threads()` : number of threads in the current region

- `omp_get_thread_num()` : id of the current thread

- `omp_get_max_threads()` : upper bound to the number of threads that could be used

- `omp_get_wtime()` : wall clock time in seconds

- `omp_get_wtick()` : seconds between successive clock ticks

- Now you can apply what you learn to the `poisson` code.
- Remember that 90% of the time is spend in the dumpers. So make sure you dump only once at the end of the simulation to get a validation image.