

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

Nicolas Richart Emmanuel Lanti

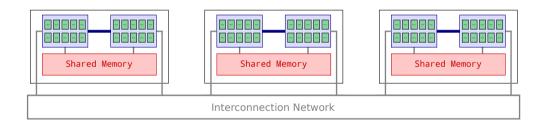
Course based on V. Keller's lecture notes

15th - 19th of November 2021

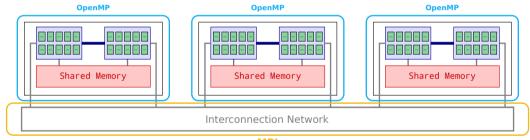


- Hybrid MPI + OpenMP programming
- Partitioned point-to-point communications
- Matching probe/receive









MPI

- Thread safety?
- Which thread/process can/will call the MPI library?
- MPI process placement in the case of multi-CPU processors?
- Data visibility? OpenMP private?
- Does my problem fits with the targeted machine?
- Levels of parallelism within my problem?



```
hybrid/hello world.cc
   #include <iostream>
 2 #include <mpi.h>
   #include <omp.h>
   int main(int argc, char *argv[]) {
     int provided, size, rank, nthreads, tid;
     MPI_Init_thread(&argc, &argv, MPI_THREAD_SINGLE, &provided);
     MPI_Comm_size(MPI_COMM_WORLD, &size);
     MPI Comm rank(MPI COMM WORLD, &rank):
10
11
12
     #pragma omp parallel default(shared) private(tid. nthreads)
13
       nthreads = omp_get_num_threads():
14
       tid = omp_get_thread_num();
15
       std::printf("Hello from thread %i out of %i from process %i out of %i\n", tid, nthreads, rank,
16
       size);
17
     MPI Finalize():
18
     return 0;
19
20 }
```



Compilation using the GNU g++ compiler:

\$> mpicxx -fopenmp hello_world.cc -o hello_world

Compilation using the Intel C++ compiler:

\$> mpiicpc -qopenmp hello_world.cc -o hello_world



Submission script the clusters

```
#!/bin/bash
#SBATCH --nodes 1
#SBATCH --ntasks 2
#SBATCH --cpus-per-task 3
export OMP_NUM_THREADS=3
srun -n 2 ./hello world
```

It will start 2 MPI processes each will spawn 3 threads

```
Hello from thread 0 out of 3 from process 0 out of 2
Hello from thread 1 out of 3 from process 0 out of 2
Hello from thread 0 out of 3 from process 1 out of 2
Hello from thread 1 out of 3 from process 1 out of 2
Hello from thread 2 out of 3 from process 0 out of 2
Hello from thread 2 out of 3 from process 1 out of 2
```



- Change your MPI initialization routine
 - ▶ MPI Init is replaced by MPI Init thread
 - ▶ MPI Init thread has two additional parameters for the level of thread support required, and for the level of thread support provided by the library implementation

```
int MPI_Init_thread(int *argc, char ***argv, int required, int *provided)
```

■ Make sure that the *provided* support matches the *required* one

```
if (provided < required)
2 MPI_Abort(MPI_COMM_WORLD, EXIT_FAILURE);</pre>
```

 Add OpenMP directives as long as you stick to the level of thread safety you specified in the call to MPI_Init_thread

SCITAS N. Richart, E. Lanti 8 / 15





- MPI THREAD SINGLE
 - ▶ Only one thread will execute (no multi-threading)
 - ► Standard MPI-only application
- MPI_THREAD_FUNNELED
 - ▶ Only the Master Thread will make calls to the MPI library
 - ▶ A thread can determine whether it is the master thread by a call to MPI_Is_thread_main
- MPI THREAD SERIALIZED
 - ▶ Only one thread at a time will make calls to the MPI library, but all threads are eligible to make such calls
- MPI THREAD MULTIPLE
 - ► Any thread may call the MPI library at any time





- MPI THREAD SINGLE
 - ▶ Only one thread will execute (no multi-threading)
 - Standard MPI-only application
- MPI THREAD FUNNELED
 - ▶ Only the Master Thread will make calls to the MPI library
 - A thread can determine whether it is the master thread by a call to MPI_Is_thread_main
- MPI THREAD SERIALIZED
 - ▶ Only one thread at a time will make calls to the MPI library, but all threads are eligible to make such calls
- MPI THREAD MULTIPLE
 - ▶ Any thread may call the MPI library at any time

In most cases MPI_THREAD_FUNNELED provides the best choice for hybrid programs

The 4 options for thread support



- MPI THREAD SINGLE
 - ▶ Only one thread will execute (no multi-threading)
 - ► Standard MPI-only application
- MPI THREAD FUNNELED
 - ▶ Only the Master Thread will make calls to the MPI library
 - A thread can determine whether it is the master thread by a call to MPI_Is_thread_main
- MPI THREAD SERIALIZED
 - ▶ Only one thread at a time will make calls to the MPI library, but all threads are eligible to make such calls
- MPI THREAD MULTIPLE
 - ▶ Any thread may call the MPI library at any time

In most cases MPI_THREAD_FUNNELED provides the best choice for hybrid programs

```
int MPI_Query_thread(int * thread_level_provided);
```

Returns the level of thread support provided by the MPI library

The 4 options for thread support



- Thread support values are monotonic, i.e.

 MPI THREAD SINGLE < MPI THREAD FUNNELED < MPI THREAD SERIALIZED < MPI THREAD MULTIPLE
- Different processes in MPI_COMM_WORLD can have different thread safety
- The level(s) of provided thread support depends on the implementation
- The rules for thread support attribution are done in the following order:
 - return provided = required
 - return the least supported level such that provided > required
 - return the highest supported level



- New feature from MPI 4.0 standard (June 2021!)
- We have already talked about persistent point-to-point communications
- Partitioned comms are just persistent comms where the message is constructed in partitions
- Typical case: multi-threading with each thread building a portion of the message



MPI partitioned communications

■ Remember the typical cycle for persistent point-to-point communications

```
Init (Start Test/Wait)* Free
```

where * means zero or more

Partitioned are very similar

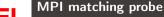
```
PInit (PStart PReady)* Free
```

```
1 MPI_Psend_init(msg, parts, count, MPI_INT, dest, tag, info, MPI_COMM_WORLD, &request);
2 MPI_Start(&request);
3 #pragma omp parallel for shared(request)
4 for (int i = 0; i < parts; ++i) {</pre>
   /* compute and fill partition #i, then mark ready: */
    MPI_Pready(i, request);
7 }
8 while(!flag) {
   /* Do useful work */
    MPI_Test(&request, &flag, MPI_STATUS_IGNORE);
    /* Do useful work */
12 }
13 MPI_Request_free(&request);
```



MPI matching probe

- We have already talked about MPI_Probe to obtain information about a message waiting to be received
- This is typically used when the size of the message is unknown (probe, allocate, receive)



EPFL

- We have already talked about MPI_Probe to obtain information about a message waiting to be received
- This is typically used when the size of the message is unknown (probe, allocate, receive)
- Care must be taken because it is a stateful method: A subsequent receive [...] will receive the message that was matched by the probe, if no other intervening receive occurs after the probe [...]

- We have already talked about MPI_Probe to obtain information about a message waiting to be received
- This is typically used when the size of the message is unknown (probe, allocate, receive)
- Care must be taken because it is a stateful method: A subsequent receive [...] will receive the message that was matched by the probe, if no other intervening receive occurs after the probe [...]
- Problem with multi-threading!
- Imagine two threads A and B that must do a Probe, Allocation, and Receive

$$A_P \longrightarrow A_A \longrightarrow A_R \longrightarrow B_P \longrightarrow B_A \longrightarrow B_R$$

but may also be

$$A_P \longrightarrow B_P \longrightarrow B_A \longrightarrow B_R \longrightarrow A_A \longrightarrow A_R$$

Thread B stole thread A's message!





- The solution of this problem is the matching probe
- MPI provides two versions, MPI_Improbe and MPI_Mprobe
- \blacksquare It allows to receive only a message matching a specific probe



- The solution of this problem is the matching probe
- MPI provides two versions, MPI_Improbe and MPI_Mprobe
- It allows to receive only a message matching a specific probe

- Counter part operation are the matching receive MPI_Imrecv and MPI_Mrecv
- They are used to receive messages that have been previously matched by a matching probe





- Always keep in mind that you are mixing threads and processes
- You will need to test your code performance on every machine
- There are no magic rule on the best configuration to use
- Often 1 MPI task per NUMA region seems to give the best performance