

# PPCES 2015: MPI Lab

16–17 March 2014

Hristo Iliev, iliev@itc.rwth-aachen.de

Portions thanks to:

Christian Iwainsky, Sandra Wienke

## Abstract

The purpose of this Lab is to familiarise you with the basic concepts of MPI. Tasks 1-3 will introduce you to the principles of basic point-to-point communication. Tasks 4-6 will practice the usage of collective communications and MPI in general. Furthermore, you will get to know Vampir, which is a performance analysis tool but is also very well suited for visualization of message passing, which we will be using throughout this workshop.

Note: the flow of exercises does not strictly match the structure of the course. It is not required that you do the exercises in the same order in which they appear in this document.

## Before you start

Before you start, log onto one of our Linux cluster frontends (**cluster-x** or **cluster-x2**) using X-Win32. Then, download the MPI lab archive from the Sharepoint server (link to the location is provided on the PPCES webpage). Extract it to a suitable location (e.g. create directory named **MPILab**, go there and extract the downloaded archive using “tar -xf ../mpilab2015.tar”). The lab archive contains stubs for the exercises described below. Intermediate solution stubs are provided where appropriate. Solutions to all problems are also provided in the **solutions** folder. We kindly advise you not to look at the solutions before you’ve tried your best to solve each exercise on your own.

The following make targets will be available for all exercises.

```
make [release | debug]           # build the program
make run [NPROCS=<#processes>]   # run the program
make runParallel [NPROCS=<#processes>] # run the program in parallel (ex. 6 only)
make clean                       # clean the directory
make vampir                      # trace the program and start Vampir
```

## 0. Hello, MPI!

The purpose of this exercise is to get you familiar with the very basics of MPI programming on the RWTH Compute cluster. Start with the minimal MPI program shown in the “Part 1: General Structure of MPI Programs” part of the tutorial slides and add a line that makes each process print its rank and the total number of processes in the MPI job. Where the line should be placed? Use the following commands to compile and run your program:

C/C++:

```
$MPICC -o hello hello.c
$MPIEXEC -n 4 ./hello
```

Fortran:

```
$MPIFC -o hello hello.f90
$MPIEXEC -n 4 ./hello
```

## 1. Ping Pong

One basic MPI program using point-to-point communication is a “ping pong” between two MPI processes. A ping-pong program stub can be found in directory **1\_pingPong**. Complete the source code parts marked with “TODO” in this program. Note: You can use Vampir to visualize the behavior of your code.

- Modify the ping-pong stub such that the first process of the MPI program transmits its input to the second process. The second process should then print the received value and send its negative value back to the first process, which should again print the received value.
- Modify the ping-pong program such that each rank sends an individually random number of elements. Hint: Note that you may have to transmit in an additional message the number of elements to the second process before sending them.
- What is the behavior of the program for NPROCS=1 and NPROCS=3? Modify the code to display an error message for too few processes and to execute properly for larger number of processes.
- Implement part b) of the assignment without explicitly sending the number of elements.
- Bonus task: Implement a loop to send/receive messages with different sizes. How does the message size influence the time being spent in MPI functions? You may use **MPI\_Wtime()** to measure process local time and set the maximum array size to  $2^{26}$  to make the impact of the data size clearly visible.

## 2. Sending and Receiving

One basic usage of MPI is to send and receive data. This however can lead to unexpected situations if not done in the correct way. A send-receive skeleton can be found in directory **2\_sendReceive**.

- Look at the given program and execute it with 2 processes. What is happening?  
Note: You can abort the program execution by hitting Ctrl-c.
- Modify the given program using **MPI\_Send()** and **MPI\_Recv()** such that it becomes a correct MPI program and completes execution.
- Can the send and receive operations be replaced with a single MPI call? Use the correct operation to replace the send and receive pair.
- Modify your code to utilize non-blocking communication primitives.
- Change the code to work with more than 2 MPI processes. In this case the messages should be sent to and received from the next higher rank.  
Hint: Will a special treatment be necessary for the last rank?

## 3. Count-down Ring

Using send and receive operations, implement a round-robin communication that passes along an integer value, starting with the stub given in directory **3\_countdownRing**. Each time the value is received it should be decremented by a random number (use the function **random\_dec**). Once this value reaches zero or less, the process that is currently updating it should notify all other processes of its rank. Every process then should display the rank of the process which decremented the counter to zero. You can supply the initial countdown value like that:

```
make run N=<countdown>
```

Compare the following example output for this exercise:

```
> make run NPROCS=3 N=45
Counting down from 45
Process 1 has received the bomb (38 on the clock) and is still alive!
Process 2 has received the bomb (35 on the clock) and is still alive!
Process 0 has received the bomb (31 on the clock) and is still alive!
Process 0 has received the bomb (13 on the clock) and is still alive!
Process 1 has received the bomb (24 on the clock) and is still alive!
Process 1 has received the bomb (6 on the clock) and is still alive!
Process 2 has received the bomb (17 on the clock) and is still alive!
Process 2 has received the bomb (4 on the clock) and is still alive!
Process 0 lost
I am process 0 and 0 is the loser
I am process 2 and 0 is the loser
I am process 1 and 0 is the loser
```

## 4. Controller-Worker

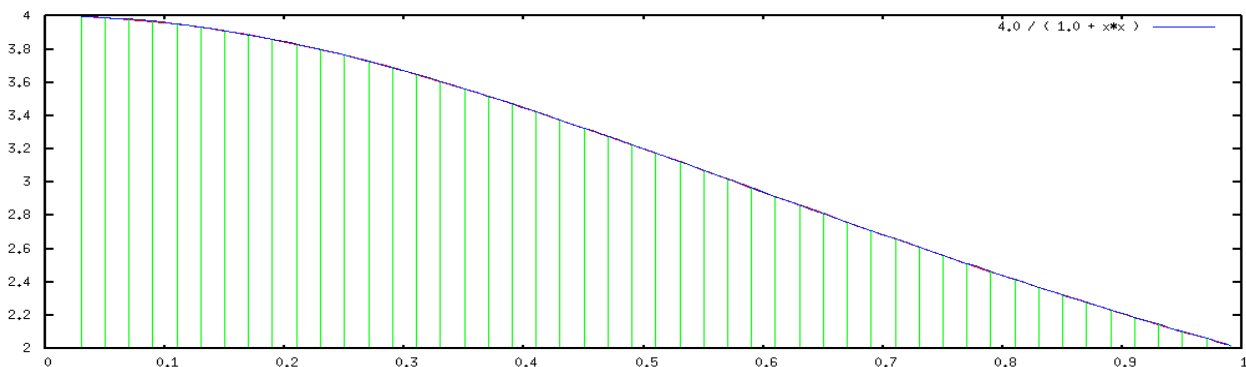
The controller-worker concept is an often-used technique in MPI applications. There one “controller” process distributes work to the other “worker” processes. In this exercise you should implement this concept using (again) point-to-point communication routines.

The trapezoid numeric function integration rule serves as a basis for this exercise. The following formula is parallelized:

$$\int_{x_0}^{x_1} f(x) dx \approx \frac{x_1 - x_0}{2N} \sum_{i=0}^{N-1} (f(x_i) + f(x_{i+1})); \quad x_i = x_0 + i * \frac{x_1 - x_0}{N}$$

As a test for this integration we use an integral representation of  $\pi$ .

$$\pi = \int_0^1 \frac{4}{1+x^2} dx$$



Code for this assignment is located in directory **4\_integration**.

- Compile and then execute the given example with different process counts, e.g. NPROCS=2, NPROCS=4. Why is no speedup observed? You can use Vampir to get an insight.
- Modify the given example so that the work is correctly done in parallel.
- Refactor (rewrite) your program so that rank 0 calls a function called **controller** and every other process calls a function **worker**. The program should still compute the integral in parallel, with the work distribution occurring in the **controller** function and the work processing in the **worker** function.
- Calling **MPI\_Send()/MPI\_Recv()** for each function evaluation is not very efficient. Modify your code to transmit ranges to be computed instead of single values. What effect does this have on the work balancing?

## 5. Your own collective communication

Often data needs to be distributed from a single process to all processes or collected from all the processes into a single rank of the MPI program. MPI provides dedicated functions to do that. In this exercise you will implement your own collective operations like broadcast, scatter, gather and all-to-all using only point-to-point MPI calls. The code stubs are located in directory **5\_myGlobals**. Note that you can compile and execute the code at any time to observe data distribution. This example is a simple toy that will be used to practice message passing by transmitting process-specific random numbers to the neighboring processes. You may specify a process rank as an argument to the executable and that process will plot its data structures so you can investigate the communication behavior.

- Implement the **bcast\_Int** function. It should distribute an integer value from the process with rank equal to *root* to all other processes in the communicator *comm*. After the operation completes, all processes should have a copy of the data from process with rank *root*.

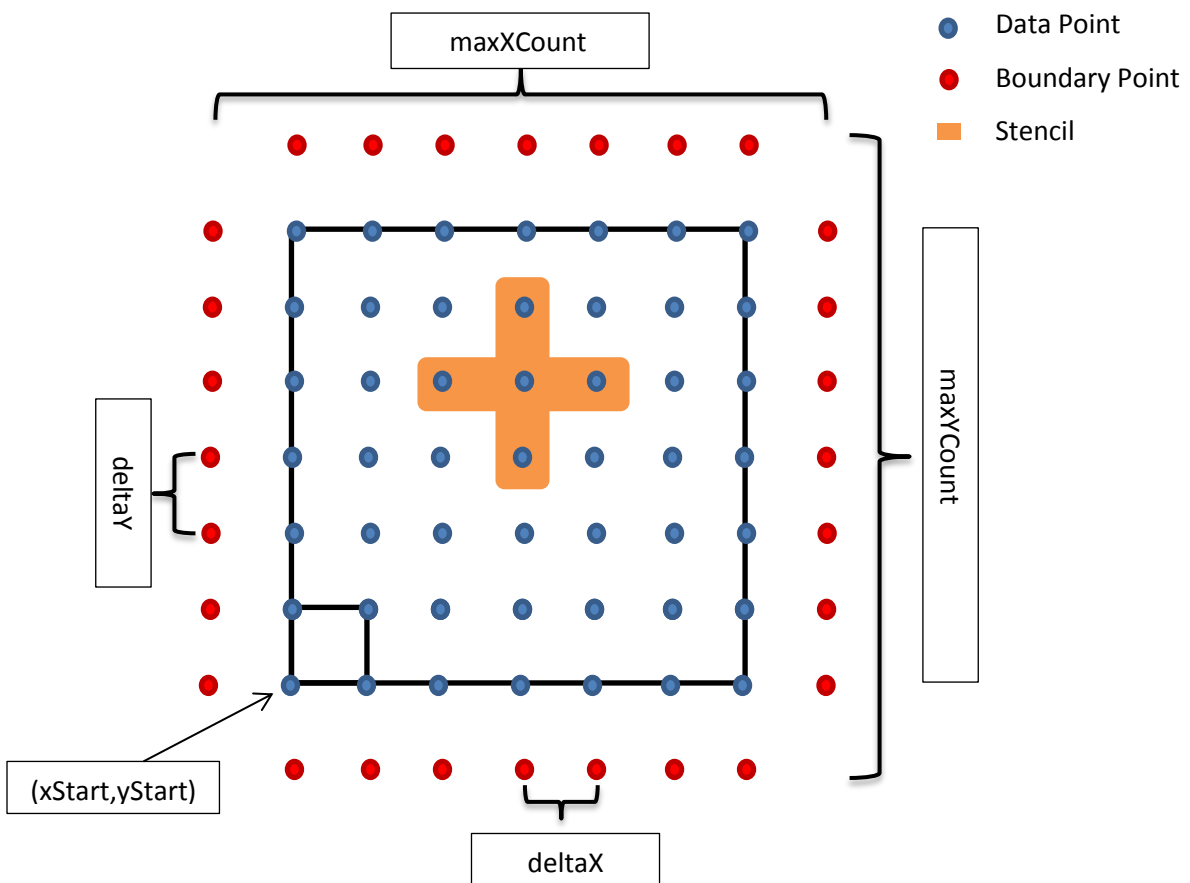
- b) Implement the **scatter\_Int** function. This function should distribute the content of the send buffer in process *root* to the participating processes. The first *sendcnt* integers of *sendbuffer* should be available in the receive buffer of rank 0, the next *sendcnt* integers at rank 1 and so forth.
  - c) Implement the **gather\_Int** function which is the reverse operation of **scatter\_Int**. It should collect *recvcnt* integers from each process in the communicator and store them in the receive buffer of the process with rank *root*.
  - d) Implement the **alltoall\_Int** function. Its operation is a combined scatter-gather. Rank 0 scatters its data to ranks 1, ..., nProcs-1, then rank 1 scatters to 0, 2, ..., nProcs-1 and so forth.
  - e) Implement the **sum\_Int** function. It should collect at process *root* a single integer value from each process in the communicator *comm* and sum all collected values. Use it to sum over all processes' ranks.  
Note:  $0 + 1 + \dots + \text{numProcs}-1 = \text{numProcs} * (\text{numProcs}-1) / 2$
  - f) Now use the corresponding collective MPI calls. What is the difference?
  - g) You may have observed that the output may have been shuffled. Consider the potential reasons for this. Fix it.
- Note: Don't forget to test your code with different process counts.

## 6. First parallelisation on your own

The program in directory **6\_jacobi** solves numerically a finite difference version of the screened Poisson equation:

$$(\nabla^2 - \alpha)u = \frac{d^2}{dx^2}u + \frac{d^2}{dy^2}u - \alpha u = f$$

using the iterative Jacobi method with over-relaxation.



In this exercise you must apply everything you have learnt today in order to parallelise the given serial version of the Jacobi solver, as this is what you will usually have to do on your own after this course. Change to directory **6\_jacobi**.

In general the algorithm works as follows:

- a. Initialises *u\_old* to zero

- b. Computes an approximation at every point of  $u$  (a discretization point of your solution vector) based on the former approximation  $u_{old}$  (implemented in function **one\_jacobi\_iteration**)
- c. Swaps  $u$  and  $u_{old}$
- d. Repeats from b) on until there is no significant change anymore or the maximum number of iterations is reached

**Note:** The domain is extended with two data points along each dimension to accommodate the boundary.

Depending on how confident you feel about programming with MPI, you can either start with the already partially parallelised **jacobi.c** / **jacobi.f90** and complete the TODOs in it or you can try to do the full parallelisation from scratch. We have provided the original sequential version of the solver as **jacobi\_serial.c** / **jacobi\_serial.f90**. You may find the following guiding questions helpful:

- What are the data dependencies for the computation of a single iteration?
- What communication method is suitable for computing the error?
- How would you partition the domain? Is there any difference between FORTRAN and C?
- Can you handle the case where the domain size is not divisible by the number of MPI processes?
- Can you partition the domain across the other dimension instead?

## 7. Using the Batch System

Long-running computations and benchmarks should be submitted as jobs to our batch system Platform LSF. Jobs are represented by batch scripts – basically a shell script containing all commands which must be executed together with a list of resource requirements and options for the batch system. Batch scripts are submitted using the **bsub** command. It is also possible to provide options as command-line arguments to **bsub**. In directory **7\_batch** you will find an example script **submit.sh** together with the sample MPI program **sendRecv.c**. Note that the batch script can (and should!) be tested interactively before submission. The interactive test lets you more easily spot eventual errors without the implied waiting time of the batch system. To run the batch script interactively first use the following commands:

```
chmod 755 submit.sh
./submit.sh
```

**Note:** The MPI\_EXEC line will execute with two processes only as the value of the environment variable **FLAGS\_MPI\_BATCH** depends on the execution environment. For interactive testing you have to set this variable:

```
FLAGS_MPI_BATCH="-np 4" ./submit.sh
```

After you have tested your script you can submit it to the batch system:

```
bsub < submit.sh
```

Notice the input redirection symbol. It is very important: if you omit it the job will be queued but none of the LSF options specified inside it will be respected. The batch system then responds to the submission with the job ID of the queued job:

```
Job <xxxxxx> is submitted to default queue <qqqqqq>.
```

The status of all your jobs can be monitored by the **bjobs** command. Only pending (PEND) and running (RUN) jobs as well as jobs in error state are shown; finished jobs are not listed.

You can cancel a job with the **bkill** command:

```
bkill xxxxxx
```

where xxxxxx is the job ID from the output of the **bsub** command.

Compile and run this program to check your batch script. Submit your script to the batch system and wait it to complete.

Note 1: More information about batch jobs can be found in the IT Center documentation Wiki:

<https://doc.itc.rwth-aachen.de/display/CC/Using+the+batch+system>

Note 2: It might take (quite) some time before the job starts.