

CPDS Laboratory Assignment

Laboratory Assignment 0: Parallel Execution Environment

J.R. Herrero

Fall 24



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH

Contents

1	Experimental setup	2
1.1	Node architecture and memory	3
1.2	Execution modes: interactive vs queued	3
1.3	Serial compilation and execution	4
1.4	Compilation and execution of OpenMP programs	6
1.4.1	Compiling OpenMP programs	6
1.4.2	Executing OpenMP programs	7
1.4.3	Strong vs. weak scalability	7
1.5	Compiling and Executing MPI programs	7
1.5.1	Strong scalability with processes and threads	7

Deliverable

Session 1

Experimental setup

The objective of this first assignment is to familiarise yourself with the hardware and software environment that you will be using during the semester to do all laboratory assignments in CPDS. From your PC/terminal booted with Linux¹ you will access **boada**, a multiprocessor server located at the Computer Architecture Department at UPC. From a machine within the UPC network you will be able to access **boada** directly. However, to be able to access the **boada** server from a computer outside the UPC you will need to use a UPC VPN connection. To connect to it you will have to establish a connection using the secure shell command: "`ssh -X cpds1XYZ@boada.ac.upc.edu`", being **cpds1XYZ** the username assigned to you. Option `-X` is necessary in order to forward the X11 data and be able to open remote windows in your local desktop². You should change the password for your account using "`ssh -t cpds1XYZ@boada.ac.upc.edu passwd`"³.

Once you are logged in you will find yourself in any of the interactive nodes: **boada-6** to **boada-8**, where you can execute interactive jobs and from where you can submit execution jobs to the rest of the nodes in the machine. In fact, **boada** is composed of several nodes (named **boada-1** to **boada-15**), equipped with five different processor generations, as shown in the following table:

Node name	Processor generation	Interactive	Partition
boada-1 to 4	Intel Xeon E5645	No	execution2
boada-6 to 8	Intel Xeon E5-2609 v4	Yes	interactive
boada-9	Intel Xeon E5-1620 v4 + Nvidia K40c	No	cuda9
boada-10	Intel Xeon Silver 4314 + 4 x Nvidia GeForce RTX 3080	No	cuda
boada-11 to 14	Intel Xeon Silver 4210R	No	execution
boada-15	Intel Xeon Silver 4210R + ASUS AI CRL-G116U-P3DF	No	iacard

However, in this course you are going to use only nodes **boada-6** to **boada-8** interactively, nodes **boada-11** to **boada-14** through the **execution** queue and, later on, node **boada-10** through the **cuda** queue. Submission of jobs for *batched* execution in a queue is explained in the next section. The rest of the nodes have restricted access and CPDS users are not allowed to send jobs to their corresponding queues (Partitions). Note that the all the nodes within a particular **Partition** are identical. Therefore, you do not need to worry if the execution was done in node **boada-11**, **boada-12**, or any of the other 2 nodes within that partition, when submitting to the **execution** queue.

All nodes have access to a shared NAS (*Network-Attached Storage*) disk; you can access it through `/scratch/nas/1/cpdsXXYY` (in fact this is your *home directory*, check by typing `pwd` in the command line). In addition, each node in **boada** has its own local disk which can be used to store temporary files non visible to other nodes; you can access it through `/scratch/1/cpdsXXYY`.

¹You can also access from your laptop, booted with Linux, Windows or MacOS X, if a secure shell client is installed.

²Option `-Y` if you are connecting from a MacOS X laptop with XQuartz. With Windows, if you use putty (<https://www.chiark.greenend.org.uk/~sgtatham/putty/latest.html>) then you will also need xming: (<https://wiki.centos.org/HowTos/Xming>). Alternatively, you can use MobaXterm (<https://mobaxterm.mobatek.net/download.html>) which has everything integrated.

³The `passwd` command will be executed in **boada**. After entering the old password correctly twice you will be asked for your new password also twice.

We will post all necessary files to do each laboratory assignment in `/scratch/nas/1/cpds0/sessions`. For the session today, you need to extract the files from file `a0.tar.gz` from that location by uncompress it **into** your home directory in `boada`.

From the the root of your home directory unpack the files with the following command line: `"tar -zxvf /scratch/nas/1/cpds0/sessions/a0.tar.gz"`.

In order to set up all environment variables you have to process the `environment.bash` file now available in your home directory with `"source ~/environment.bash"`. **Note:** since you have to do this every time you login in the account or open a new console window, we recommend that you add this command line in the `.bashrc` file in your home directory, a file that is executed every time a new session is initiated. If you have unpacked the `a0.tar.gz` files in your home directory then this file has already been added automatically to your home directory and it will be used from your next connection to `boada`. From the next time you log in to `boada` it will be automatically read, and environment variables will be properly set effortless.

In case you need to transfer files from `boada` to your local machine (laptop or desktop in laboratory room), or viceversa, you have to use the secure copy `scp` command. For example if you type the following command `"scp cpds1XYZ@boada.ac.upc.edu:a0/pi/pi_seq.c ."` in your local machine you will be copying the source file `pi_seq.c` located in directory `a0/pi` of your home directory in `boada` to the current directory, represented with the `"."`, in the local machine, with the same name.

1.1 Node architecture and memory

The first thing you will do is to investigate the architecture of the available nodes in `boada`. Run `sbatch submit-arch.sh` command in directory `a0/arch`. This command will enqueue `submit-arch.sh` script. More detail about queue execution below. This script will execute the `lscpu` and `lstopo` commands in order to obtain information about the hardware in one of the nodes of execution queue (`boada-11` to `14`). The execution of these two commands using `submit-arch.sh` script generates three files, where `number` may be `11`, `12`, `13` or `14`: 1) `lscpu-boada-number`, 2) `lstopo-boada-number`, and `map-boada-number.fig`. You can use the `xfig` command to visualise the output file generated (`map-boada-number.fig`) and export to a different format (PDF or JPG, for example) using `File → Export` in order to include it in your deliverable for this laboratory assignment⁴.

Those files will help you to figure out:

- the number of sockets, cores per socket and threads per core in a specific node;
- the amount of main memory in a specific node, and each NUMA node;
- the cache memory hierarchy (L1, L2 and L3), private or shared to each core/socket.

Fill in the table in the "Deliverable" section with the main characteristics of the node. Draw the architecture of the node based on the information generated by the tools above. The `"--of fig map.fig"` option for `lstopo` can be very useful for that purpose.

1.2 Execution modes: interactive vs queued

There are two ways to execute your programs in `boada`:

1. via a queueing system (in one of the nodes `boada-11` to `boada-14`);
2. interactively (in any of the login nodes `boada-6` to `boada-8`). In this case, the system limits the number of cores to be used in parallel executions to two.

⁴In the `boada` Linux distribution you can use `xpdf` to open pdf files and `display` to visualise graphics files. You can also use the `"fig2dev -L pdf map.fig map.pdf"` command to convert from `.fig` to `.pdf`; look for alternative output graphic languages by typing `"man fig2dev"`.

It is mandatory to use option 1 when you want to execute scripts that require several processors in a node, ensuring that your job is executed in isolation (and therefore reporting reliable performance results) and to avoid adding additional load to the interactive node accessed by all users; the execution starts as soon as a node is available. When using option 2 your execution starts immediately but will share resources with other programs and interactive jobs, not ensuring representative timing results. Usually, scripts for both options (`submit-xxxx.sh` and `run-xxxx.sh`, respectively) will be provided:

- Queueing a job for execution: `"sbatch [-p partition] ./submit-xxxx.sh"` . Additional parameters may be specified, if needed by the script, after the script name. If you do not specify the name of the partition with `"-p partition"` your script will run on the `execution` partition by default. Use `"squeue"` to ask the system about the status of your job submission. You can use `"scancel"` followed by the job identifier to remove a job from the queueing system. Note that partition names associated to each node name are shown in the last column of the table above. After the execution in an available node associated to the specified partition, in addition to the files being generated by the script, two additional files will be created. Their name will have the script name followed by an `".e"` and an `".o"` and the job identifier. They will contain the messages sent to the `standard error` and `standard output` respectively during the execution of the job. You should check them to be sure results make sense.
- Interactive execution: `./run-xxxx.sh`. Additional parameters may be specified after the script name. Jobs interactively executed have a short time limit to be executed.

1.3 Serial compilation and execution

Next you will get familiar with the compilation and execution steps for both sequential and parallel applications. For this, we will use some codes which compute digits of number Pi. As a reference, Figure 1.1 shows a python code which computes the area under one quadrant of the unit circle, i.e. a circle of radius equal to 1, defined by the curve $y = \sqrt{1 - x^2}$.

```
def midpoint_riemann_sum():
    N = 1000000000

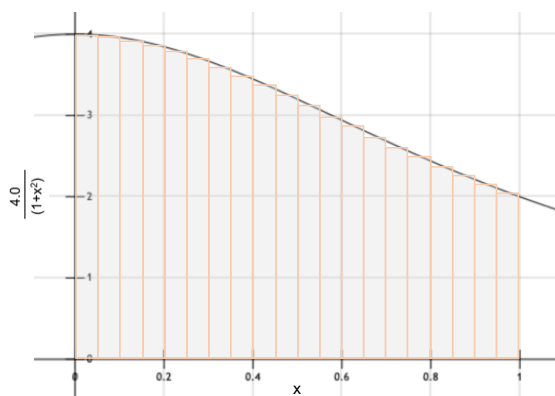
    delta_x = (1 - 0) / N
    x = 0
    pi = 0
    start = time.time()

    while x < (1-delta_x):
        x_next = x + delta_x
        x_mid = (x_next + x) / 2
        f_x = math.sqrt(1 - math.pow(x_mid, 2))
        pi += f_x * delta_x
        x += delta_x

    pi = 4 * pi
    end = time.time()
    print(end - start)
    print("Estimate: " + str(pi))
    print("Actual:   " + str(math.pi))
```

Figure 1.1: Python code for computing Pi

The execution time for this code on `boada-1` took 522.1 seconds to compute the estimate 3.14159266367805 (actual value is 3.141592653589793). As you can see, it took a long time while it only computed 7 correct digits of Pi. We aim at estimating digits of Pi much faster. For that we are going to use some C codes which we will parallelize. You are going to use a very simple code, `pi_seq.c`, which you can find inside the `a0/OpenMP/pi` directory. `pi_seq.c` performs the computation of the Pi number by computing the integral of the equation in Figure 1.2. This time the algorithm is using the derivative of



Mathematically, we know that:

$$\int_0^1 \frac{4.0}{(1+x^2)} dx = \pi$$

We can approximate the integral as a sum of rectangles:

$$\sum_{i=0}^N F(x_i) \Delta x \approx \pi$$

Where each rectangle has width Δx and height $F(x_i)$ at the middle of interval i .

Figure 1.2: One way to compute Pi

the arctangent function of x , which is equal to 1 divided by $(1 + x^2)$ which is another way to estimate Pi. Again, the equation can be solved by computing the area defined by the function, which at its turn can be approximated by dividing the area into small rectangles and adding up its area. This is one of the many ways to estimate digits of number Pi. Both algorithms above multiply by 4 because the functions compute $\pi/4$. Figure 1.3 shows a simplified version of the code you have in `pi_seq.c`. Variable `num_steps` defines the number of rectangles, whose area is computed in each iteration of the `i` loop.

```
static long num_steps = 100000;
void main () {
    double x, pi, step, sum = 0.0;

    step = 1.0/(double) num_steps;
    for (long int i=0; i<num_steps; ++i) {
        x = (i+0.5)*step;
        sum += 4.0/(1.0+x*x);
    }
    pi = step * sum;
}
```

Figure 1.3: Serial code for Pi

Figure 1.4 shows the compilation and execution flow for a sequential program. We will always compile programs to generate binary executable files through a **Makefile**, with multiple targets that specify the rules to compile each program version. In this course we are going to use `gcc` (the C front-end from the *GNU Compiler Collection*). .

In the following steps you will compile `pi_seq.c` using the **Makefile** and execute it interactively and through the queueing system, with the appropriate timing commands to measure its execution time:

1. Open the **Makefile** file, identify the **target** you have to use to compile the sequential code. Observe how the compiler is invoked. Execute the command line `"make target_identified"` in order to generate the binary executable file. Alternatively, you can just type `"make"` which will compile all programs.
2. Interactively execute the binary generated to compute number pi sequentially by doing 1.000.000.000 iterations using the `run-seq.sh` script (`"run-seq.sh pi_seq 1000000000"`) which returns the user and system CPU time, the **elapsed** time, also called **wall** time, and the % of CPU used (using the GNU time program in `/usr/bin/time`). In addition, the program itself also reports the elapsed execution time using `gettimeofday`. Look at the source code and identify the function invocations and data structures required to measure execution time. Note that CPU time is the quantity of

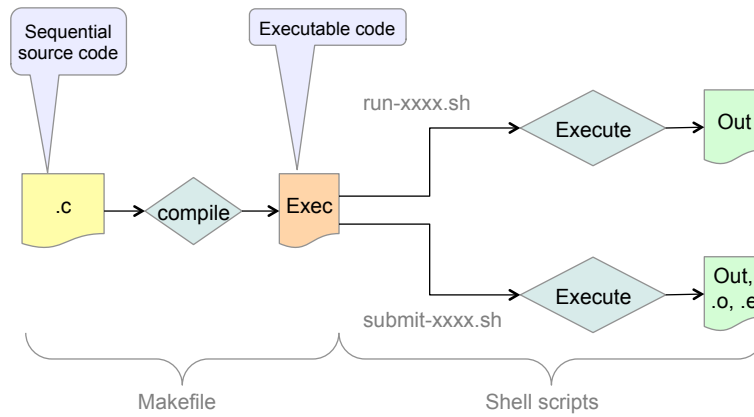


Figure 1.4: Compilation and execution flow for sequential program.

processor time taken by the process. This does not indicate duration. Elapsed time represents the total duration of the execution. If we exploit parallelism, elapsed time can be less than the total CPU time because several cores are used to do the work and it can get done faster.

3. Submit the execution to the queuing system using the "`sbatch submit-seq.sh`" command and use "`squeue`" to see that your script is running. Look at `submit-seq.sh` script and the results generated (the standard output and error of the script and the `time-pi.seq-boada{11-14}` file).

1.4 Compilation and execution of OpenMP programs

In this section we are going to use **OpenMP**, the standard for parallel programming using shared-memory, to express parallelism in the C programming language. Note that **OpenMP** will be explained in detail in the second part of the course, and the goal now is to have a first contact and observe that we can speed-up the execution of our programs. Therefore, in this section we will see how to compile and execute parallel programs in **OpenMP** and we will observe the performance improvements. Figure 1.5 shows the compilation and execution flow for an OpenMP program. The main difference with the flow shown in Figure 1.4 is that now the **Makefile** will include the appropriate compilation flag to enable **OpenMP**.

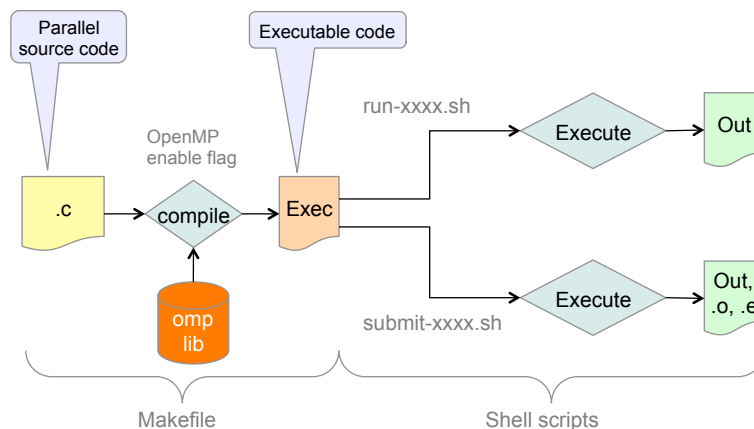


Figure 1.5: Compilation and execution flow for OpenMP.

1.4.1 Compiling OpenMP programs

1. In the same `a0/OpenMP/pi` directory you will find an **OpenMP** version of the code for doing the computation of pi in parallel (`pi_omp.c`). Compile the **OpenMP** code using the appropriate target

in the `Makefile`). What is the compiler telling you? Is the compiler issuing a warning or an error message? Is the compiler generating an executable file? The compilation run smoothly because we used flag `-fopenmp`, but if you remove the `-fopenmp` flag and force the recompilation you will see the difference.

1.4.2 Executing OpenMP programs

1. Interactively execute the OpenMP code with 8 threads (processors) and same number of iterations (1.000.000.000) using the `run-omp.sh` script ("`run-omp.sh pi_omp 1000000000 8`"). What is the `time` command telling you about the user and system CPU time, the elapsed time, and the % of CPU used? Take a look at the script to discover how do we specify the number of threads to use in OpenMP. Note that, since any login node `boada-[678]` is shared by many users and is reserved for interactive usage, it has been configured so that our processes only have 2 cores available for their execution. Thus, we will see no improvement due to the parallelization.
2. Use `submit-omp.sh` script to queue the execution of the OpenMP code ("`sbatch submit-omp.sh`") and measure the CPU time, elapsed time and % of CPU when executing the OpenMP program using 8 threads in isolation. Do you see a major difference between interactive and queued execution?

1.4.3 Strong vs. weak scalability

Finally in this section you are going to explore the scalability of the `pi_omp.c` code when varying the number of threads used to execute the parallel code. To evaluate the scalability the ratio between the sequential and the parallel execution times will be computed. Two different scenarios will be considered: *strong* and *weak* scalability.

- In *strong* scalability the number of threads is changed with a fixed problem size. In this case parallelism is used to reduce the execution time of your program.
- In *weak* scalability the problem size is proportional to the number of threads. In this case parallelism is used to increase the problem size for which your program is executed.

We provide you with two scripts, `submit-strong-omp.sh` and `submit-weak-omp.sh`, which should be submitted to the queueing system with the `sbatch` command. The scripts execute the parallel code using from 1 (`np_NMIN`) to 20 (`np_NMAX`) threads. The problem size for strong scalability is 1.000.000.000 iterations; for weak scalability, the initial problem size is 100.000.000 which grows proportionally with the number of threads. The execution **will take some time** because several executions are done for each test (in order to get a minimum time), please be patient! And, very important, do NOT execute interactively but, instead, **submit them for execution in the execution queue**. As a result each script generates several files, including a plot (a file with a name ending in `.ps` in Postscript format) showing the resulting parallel execution time and speed-up or parallel efficiency. Visualise⁵ the plots generated and reason about how the speed-up changes with the number of threads in the two scenarios.

Change the value for `np_NMAX` in the `submit-strong-omp.sh` from 20 to 40 and execute again. Can you explain the behaviour observed in the scalability plot? Even when **Hyper-Threading** is enabled and the system can take instructions quickly from 40 threads, there are only 20 cores available!

1.5 Compiling and Executing MPI programs

1.5.1 Strong scalability with processes and threads

Next, we will test ourselves how to compile and execute parallel programs in MPI, or a combination of both MPI and OpenMP. To do so, follow similar same steps as for OpenMP in the previous section, but working with the files in directory `a0/MPI/pi`. The code calculates number Pi with a different algorithm based on the so called Monte Carlo method. Monte Carlo methods are a broad class of computational algorithms that rely on repeated random sampling to obtain numerical results. The idea is to simulate

⁵In the boada Linux distribution you can use the ghostscript `gs` command to visualise Postscript files or convert the files to PDF using the `ps2pdf` command and use `xpdf` to visualise PDF files.

random (x, y) points in a 2-D plane with domain as a square of side 1 unit. Imagine a circle inside the same domain with same diameter and inscribed into the square. We then calculate the ratio of number of points that lied inside the circle and total number of generated points. To higher the number of samples, the better the precision obtained.

We can get a plot of the execution time and speed-up that is obtained when varying the number of MPI processes (strong scalability) by submitting a job to the `execution` queue with `sbatch submit-strong-mpi.sh`). In addition, we want to observe the combination of MPI and OpenMP. For that execute `sbatch submit-mpi2-omp.sh`) and study the results presented in the output file created after the execution finishes. Note that since the test uses 2 MPI processes the execution implies using an even number of threads between 2 and 40 even when the output presents the results for 1 to 20 threads: there are between 1 to 20 threads per MPI process; with 2 processes this results in having twice as many threads.

We hope that, with this experiments, you realize that we can speed up the code by parallelizing a code with OpenMP, with MPI, or a combination of both to use several nodes and several cores per node!

Deliverable

After completing this laboratory assignment, and before starting the next one, you will have to deliver a report in PDF format (other formats will not be accepted) containing the answers to the questions stated at the end of this document;

Your professor will open the assignment at the moodle Atenea website and set the appropriate delivery dates for the delivery. Note that, unlike the subsequent assignments which can be done in groups of 2 students, **this initial assignment must be done independently by each student, requiring Individual submission** through the moodle Atenea website.

Important: In the front cover of the document, please clearly state your name and surname, the identifier of the group (username `cpds1XYZ`), title of the assignment, date, academic course/semester, ... and any other information you consider necessary.

As part of the document, you can include any code fragment, figure or plot you need to support your explanations. In case you need to transfer files from boada to your local machine (laptop or desktop in laboratory room), or viceversa, you can use the secure copy `scp` command. For example "`scp cpds1XYZ@boada.ac.upc.edu:a0/foobar.txt local/directory/.`" executed in your local machine to copy file `foo.txt` inside directory `a0` in your home directory of boada to directory `local/directory/` with the same name.

Node architecture and memory

1. Complete the following table with the relevant architectural characteristics of the different node types available in boada:

	boada-11 to boada-14
Number of sockets per node	
Number of cores per socket	
Number of threads per core	
Maximum core frequency	
L1-I cache size (per-core)	
L1-D cache size (per-core)	
L2 cache size (per-core)	
Last-level cache size (per-socket)	
Main memory size (per socket)	
Main memory size (per node)	

2. Include in the document the architectural diagram for one of the nodes `boada-11` to `boada-14`.

Timing sequential and parallel executions

For each of the following sections show the plots or tables required. But, in addition, please **provide some reasoning about the results obtained**. For instance, comment whether the scalability is good or not, and why.

3. Plot the execution time and speed-up that is obtained when varying the number of threads (strong scalability) by submitting the jobs to the `execution` queue (section 1.4.3). Show the parallel efficiency obtained when running the weak scaling test. Explain what strong and weak scalability refer to, exemplifying your explanation with the plots that you present. Explain the reasons which cause the behaviour observed when more than 20 threads are used.
4. Plot the execution time and speed-up that is obtained when varying the number of MPI processes from 1 to 20 (strong scalability) by submitting the jobs to the `execution` queue (section 1.5.1). In addition, show in a table the elapsed execution time when executed with 2 MPI processes when varying the number of threads from 1 to 20⁶. Which was the original sequential time?; the time with 20 MPI processes?; and the time with 2 MPI processes each using 20 OpenMP threads? You can retrieve such information for 1 and 20 MPI processes from file `elapsed.txt`; and for 2 MPI processes each with 20 threads within the output file created after the execution of `sbatch submit-mpi2-omp.sh`.

⁶Note that for 2 MPI processes the scalability test implies using an even number of threads between 2 and 40.