

[Click to show/hide PDF Layer](#)

Lab 1

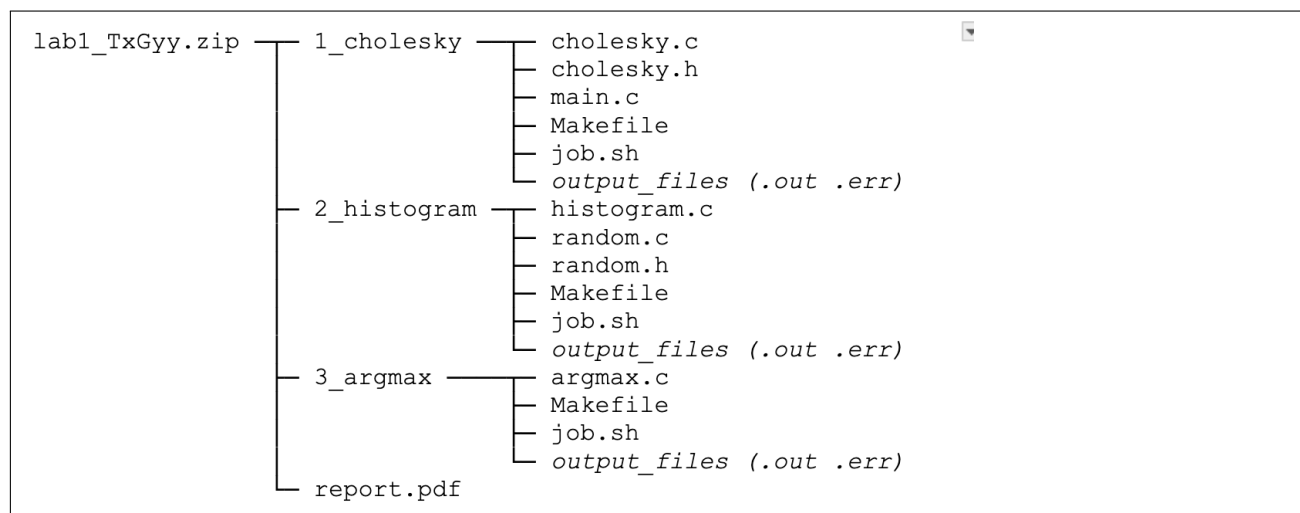
Pablo Arias^{*}, Daniel Santos-Oliván[†], Álvaro Moure[‡]

Instructions

In this lab, we will learn how to use OpenMP to parallelize code for shared-memory computers. We will explore the basics of OpenMP through its two main paradigms: *threads* and *tasks*. This session consists of three mathematical problems, each with several questions.

You must compile your answers and explanations into a report, detailing how you solved each problem and addressing all the questions. This lab assignment is worth **20% of the total course grade**, and the deadline for submission is **April 22 at 12:00h**.

Each group must submit a compressed file named `lab1-TxGyy.zip`, where TxGyy is your group identifier. A `.tar` or a `.tgz` file is also accepted (e.g., `lab1-T1G1.zip` or `lab1-T2G21.zip`). The compressed file **must** contains three folders (`1_cholesky`, `2_histogram`, and `3_argmax`) and **all** the requested files with the following structure. Additional files will not be considered and will not be penalized.



A sample file named `lab1-TxGyy.zip` containing the reference codes has been published in Aula Global. You need to create your job scripts to perform your tests according to what is asked. The Makefiles should not be modified unless the exercise says so. Focus on the code and the work requested for each exercise.

If you have any questions, please post them on the lab class forum in Aula Global. However, do not post your code in the Forum. If you have other questions regarding the assignment that you consider cannot be posted in the Forum (e.g., personal matters or code), please contact the lab responsible person.

^{*} pablo.arias@upf.edu.

[†] daniel.santos@upf.edu.

[‡] alvaro.moure@upf.edu.

Criteria

The codes will be tested and evaluated on the same cluster where you work. The maximum grade on each part will only be given to these exercises that solve in the most specific way and tackle all the functionalities and work requested. All the following criteria will be applied while reviewing your labs in the cluster.

Exercises that will not be evaluated:

- A code that does not compile.
- A code giving wrong results.
- A code that does not adhere to all the input/output requests.
- `lab1_TxGyy.zip` delivered files not structured or named as described previously.

Exercises with penalty: A code with warnings in the compilation.

1. Cholesky factorization

The Cholesky factorization (or Cholesky decomposition) is a factorization of a symmetric, positive definite matrix A into two matrices: the lower matrix L and the upper matrix U , where these names refer to the lower and upper part of the diagonal of the matrix. In the Cholesky factorization, the generated matrices are the transpose of the other, $L = U^T$, and vice versa.

In this exercise, we will have to finish the sequential version of Cholesky in `cholesky()` and then parallelize it in `cholesky_omp()`. We will calculate U from A and then calculate L by doing the transpose of matrix U . **The matrix transpose operation has to be implemented to optimize the efficiency of L1 cache memory usage.** The formulas used to calculate the elements of U are the following, depending if they are diagonal elements or non-diagonal elements:

$$u_{ji} = \frac{a_{ji} - \sum_{k=0}^{i-1} u_{kj} u_{ki}}{u_{ii}}$$

off-diagonal elements

$$u_{ii} = \sqrt{a_{ii} - \sum_{k=0}^{i-1} u_{ki}^2}$$

diagonal elements

It is important to notice that we only need to compute the corresponding elements of the diagonal and the lower elements (for L) and upper elements (for U). For the sake of simplicity, L and U may be square matrices, and have the rest of the elements be zero.

Once calculated the Cholesky factorization we want to check that it has been calculated correctly. This is done by performing $B = LU$ and later verifying that $B = A$. Multiply matrices LU and calculate B . Matrix B will have the same size as A . After generating B , compare the elements of B and A and check that the difference between elements is less than 0.001% as follows:

$$Error[\%] = \left| \frac{B_{ij} - A_{ij}}{A_{ij}} \right| \times 100$$

In the file `cholesky.c` there are the two C functions that we will work with. We will need to first finish the sequential implementation and later parallelize all 5 steps for the OpenMP version. The steps of our functions are the following.

1. Matrix initialization for A, L, U and B (already done)
2. Compute Cholesky factorization for U
3. Calculate L from U^T : iterate only over the non-zero elements.
4. Calculate $B = LU$: matrix multiplication. Iterate only over non-zero elements of L .
5. Check if all elements of A and B have a difference smaller than 0.001%.

Here is a sample output of the code (you might have different results). Notice that `argv[1]` is the size n of the square matrix A .

```
1 Sequential Cholesky
2 Initialization: 0.160174
3 Cholesky: 52.790788
4 L=U.T: 0.051511
5 B=LU: 88.952372
6 Matrices are equal
7 A==B?: 0.0
8
9 OpenMP Cholesky
10 Initialization: 0.161941
11 Cholesky : 3.670381
12 L=U.T: 0.003441
13 B=LU: 5.572723
14 Matrices are equal
15 A==B?: 0
```

Report Questions 1

Cholesky (40%)

1. Expose your parallelization strategy to divide the work in the Cholesky algorithm and in the matrix multiplication. Justify the selection of the scheduler and chunk size and compare different schedulers with different chunk sizes and show the results.
2. Make two plots: one for the speedup of the Cholesky factorization and another for the matrix multiplication for $n = 3000$. Use 1, 2, 4, 8, and 16 cores for a strong scaling test. Plot the ideal speedup in the figures and use a logarithmic scale to print the results. Discuss the results.

2. Histogram

In this exercise, we provide a program that will fill an array with pseudo-random values, build a histogram of that array, and then compute statistics. This can be used as a simple test of the quality of a random number generator.

The code is sequential. The goal of this exercise is to parallelize it using four different methods: using critical, atomic, locks and a reduction.

The code requires no arguments.

The output must be like this. It must execute the 4 variants one after the other. This execution has been performed on a laptop. Thus, your times can be different. Remember to reinitialize the histogram after each variant.

```
1 4 threads
2 Sequential
```

```

3 histogram for 50 buckets of 1000000 values
4 ave = 20000.000000, std dev = 394.372925
5 in 0.003438 seconds
6 par with critical
7 histogram for 50 buckets of 1000000 values
8 ave = 20000.000000, std dev = 394.372925
9 in 0.077274 seconds
10 par with locks histogram for 50 buckets of 1000000 values
11 ave = 20000.000000, std dev = 394.372925
12 in 0.058147 seconds
13 par with reduction
14 histogram for 50 buckets of 1000000 values
15 ave = 20000.000000, std dev = 394.372925
16 in 0.000601 seconds

```

Report Questions 2

Histogram (30%)

1. Explain how have you solved each of the parallelizations.
2. Explain the time differences between different parallel methods if there are any.
3. Make a speedup plot for the different parallelization methods for 1, 2, 4, 8, and 16 cores. Discuss the results.

3. Argmax

The goal in this exercise is to write a function that traverses a vector v of doubles and computes the maximum value m and the index idx_m there that maximum is locate (i.e. the argmax).

You have to write your code in a file `argmax.c`. In addition to the main function your code needs to contain the following functions:

```

1 void initialize(double *v, int N) {
2     for (int i = 0; i < N; i++) {
3         v[i] = (1 - pow(0.5 - (double)i/(double)N, 2)) * cos(2*M_PI*100* (i - 0.5)/N);
4     }
5 }
6
7 // computes the argmax sequentially with a for loop
8 void argmax_seq(double *v, int N, double *m, int *idx_m) {
9 }
10
11 // computes the argmax in parallel with a for loop
12 void argmax_par(double *v, int N, double *m, int *idx_m) {
13 }
14
15 // computes the argmax recursively and sequentially
16 void argmax_recursive(double *v, int N, double *m, int *idx_m, int K) {

```

```

17 }
18
19 // computes the argmax recursively and in parallel using tasks
20 void argmax_recursive_tasks(double *v, int N, double *m, int *idx_m, int K) {
21 }

```

The recursive function `argmax_recursive` should work as follows. For input vectors with more than K elements, the function should divide the vector in two halves, compute the max and the argmax on each half by calling itself recursively, and combine the results of each half to compute the max and argmax for the whole input vector. For input vectors with less than K elements, the function should compute the max and argmax sequentially by calling `argmax_seq`.

The version with tasks, should package each call of the recursive function in a task.

The main function needs to do the following:

1. set N as $N = 4096^2$
2. receive as command line parameters (1) `nthreads` the number of threads and (2) the value K , the length at which the recursive implementations stop dividing the input,
3. set the number of threads according to `nthreads`,
4. allocate memory for v and initialize it using the provided function,
5. call the four functions above measuring their running time, storing the outputs in variables

```

seq_m      par_m      rec_m      task_m
seq_idx_m  par_idx_m  rec_idx_m  task_idx_m

```

6. print to the standard output the following string:

```

1 Running argmax with K = 2048 using 8 threads
2 sequential for      argmax: m = 1.00, idx_m=8388608, time=0.018038s
3 parallel   for      argmax: m = 1.00, idx_m=8388608, time=0.413042s
4 sequential recursive argmax: m = 1.00, idx_m=8388608, time=0.025692s
5 parallel   recursive argmax: m = 1.00, idx_m=8388608, time=0.015367s
6

```

7. free the memory allocated for v

Report Questions 3

Argmax (30%)

1. Explain the different implementations of the argmax function (sequential and recursive), and how you parallelized each of them.

2. Run the code with 2, 4 and 8 threads for a vector of size $N = 4096 * 4096$ and plot the strong speedup for both parallel implementations. For the recursive and tasks implementations plot the a strong speedup curve for the following values of K , $K = 16, 512, 2048, 4096$ and $K = 8192$. Include all curves in the same plot for better comparison. Comment on the obtained results, and provide an explanation for the behaviors of the different parallel implementations.

Hint: depending on the node load, the results might vary. Launch the job several times and keep the cases with the largest speedups. If well implemented, the tasks version should be at least two times faster using $K = 8192$ and 8 threads.

3. What is the arithmetic intensity of the argmax algorithm? Which resource (memory bandwidth or peak computing capacity) do you expect to be the bottleneck for throughput?
-