

263-2300-00: How To Write Fast Numerical Code

Assignment 2: 80 points

Due Date: Th, March 14th, 17:00

<https://acl.inf.ethz.ch/teaching/fastcode/2019/>

Questions: fastcode@lists.inf.ethz.ch

Submission instructions (read carefully):

- (Submission)
Homework is submitted through the Moodle system <https://moodle-app2.let.ethz.ch/course/view.php?id=10968>.
- (Late policy)
You have 3 late days, but can use at most 2 on one homework, meaning submit latest 48 hours after the due time. For example, submitting 1 hour late costs 1 late day. Note that each homework will be available for submission on the Moodle system 2 days after the deadline. However, if the accumulated time of the previous homework submissions exceeds 3 days, the homework will not count.
- (Formats)
If you use programs (such as MS-Word or Latex) to create your assignment, convert it to PDF and name it homework.pdf. When submitting more than one file, make sure you create a zip archive that contains all related files, and does not exceed 10 MB. Handwritten parts can be scanned and included or brought (in time) to Alen's, Tyler's or Gagandeep's office.
- (Plots)
For plots/benchmarks, **provide (concise) necessary information for the experimental setup (e.g., compiler and flags) and always briefly discuss the plot and draw conclusions**. Follow (at least to a reasonable extent) the small guide to making plots from the lecture.
- (Code)
When compiling the final code, ensure that you use optimization flags (e.g. for GCC use the flag "-O3"). **Disable SSE/AVX for this exercise when compiling**. Under Visual Studio you will find it under Config / Code Generator / Enable Enhanced Instructions (should be off). With gcc you can use the flag `-fno-tree-vectorize`.
- (Neatness)
5% of the points in a homework are given for neatness.

Exercises:

1. *Short project info (10 pts)* Go to the [list of milestones for the projects](#). If you have not done that yet, please register your project there. Read through the different points and fill in the first two with the following about your project (be brief):

Point 1) An exact (as much as possible) but also short, problem specification.

For example for MMM, it could be like this:

Our goal is to implement matrix-matrix multiplication specified as follows:

Input: Two real matrices A, B of compatible size, $A \in \mathbb{R}^{n \times k}$ and $B \in \mathbb{R}^{k \times m}$. We may impose divisibility conditions on n, k, m depending on the actual implementation.

Output: The matrix product $C = AB \in \mathbb{R}^{n \times m}$.

Give the name of the algorithm you plan to consider for the problem and a precise reference (e.g., a link to a publication plus the page number) that explains it.

Point 2) A very short explanation of what kind of code already exists and in which language it is written.

2. Optimization Blockers (20 pts) Code needed

In this exercise, we consider the following short computation:

```
1 void slowperformance1(double *w, double *x, double *y, double *z, int n) {
2   for (int i = 0; i < n; i++) {
3     for (int k = 0; k < n; k++) {
4       z[i] += x[k] * y[k] / sqrt(w[i]);
5     }
6   }
7 }
```

This is part of the supplied code.

- Read and understand the supplied code. It enables you to register functions with the same signature, which will be timed in a microbenchmark fashion.
- Count a `sqrt` operation as a single floating point operation.
- Create new functions where you perform some loop unrolling and scalar replacement as discussed in the lecture to increase the performance. Explore at least three possible choices in this space, as different as possible.
- You may apply any optimization that produces the same result in exact arithmetic.
- The array `w` is guaranteed to contain only positive values.
- For every optimization you perform, create a new function in `comp.cpp` that has the same signature as `slowperformance1` and register it to the timing framework through the `register_function` function in `comp.cpp`. Let it run and, if it verifies, determine the performance (flops per cycle).
- When done, rerun all code versions also with optimization flags turned off (`-O0` in the Makefile).
- Create a table with the performance numbers. Two rows (optimization flags, no optimization flags) and as many columns as versions of `slowperformance1`. Briefly discuss the table.
- Submit your `comp.cpp` to Moodle.

What speedup do you achieve?

3. Microbenchmarks(45 pts) Code needed

Write a program (without vector instructions, i.e., standard C, and compiled with autovectorization disabled) that benchmarks the latency and throughput of a floating point multiplication and division instructions on doubles. Use the skeleton available provided and:

- Implement the functions provided in the skeleton:

```
void      microbenchmark_mode (microbenchmark_mode_t mode);
double    microbenchmark_get_mul_latency   ();
throughput_t microbenchmark_get_mul_throughput ();
double    microbenchmark_get_div_latency   ();
throughput_t microbenchmark_get_div_throughput ();
double    microbenchmark_get_sqrt_latency   ();
throughput_t microbenchmark_get_sqrt_throughput ();
```

- Use the `initialize_microbenchmark_data` function below for initialization and test modes.
- Report your CPU, operating system and compiler.
- Report any optimization flags that you used.
- Report the results for latency in cycles and throughput in operations per cycle.
- Generate `microbenchmark.s` - the assembly version of the benchmark.
- Submit only `microbenchmark.c` and `microbenchmark.s`.
- To accurately measure the number of CPU cycles, you must disable turboboost. In addition, to save power, CPUs may throttle their frequency below the nominal frequency. To ensure that the CPU is not throttled down when running the experiments, one can **warm up** the CPU before timing them.

Discussion:

- (a) Can you reach the theoretical latency / throughput of the instructions?
- (b) Does the `sqrt` instruction have the same latency and throughput for every input? If no, report the latency and throughput for as many classes of inputs as you can find. You should be able to find at least 3.
- (c) Do the throughput and latency of the `sqrt` instruction match what is in the manual? If no, why not? Hint: try looking at `microbenchmark.s`.

Note that you might have to run the same instruction many times in order to get precise benchmark results. The code skeleton uses [CMake](#) to compile and uses `RDTC` to measure performance. To compile it, you can use the following:

```
cd microbenchmark
mkdir build
cd build
cmake ..
cd ..
cmake --build build --config Release
```

Then run it, using:

```
./build/ubenchmark
```

The `zip` bundle contains detailed informations on the compilation steps in Linux / Windows / Mac OS X, available in `README.md`, as well as in the `doc` folder.