

## 263-0007-00: Advanced Systems Lab

Assignment 1: 100 points

Due Date: Th, March 6th, 17:00

<https://acl.inf.ethz.ch/teaching/fastcode/2025/>

Questions: fastcode@lists.inf.ethz.ch

### Academic integrity:

All homeworks in this course are single-student homeworks. The work must be all your own. Do not copy any parts of any of the homeworks from anyone including the web. Do not look at other students' code, papers, or exams. Do not make any parts of your homework available to anyone, and make sure no one can read your files. The university policies on academic integrity will be applied rigorously.

### Submission instructions (read carefully):

- (Submission)  
Homework is submitted through the [Moodle system](#) and through [Code Expert](#) for coding exercises. The enrollment link for Code Expert is <https://expert.ethz.ch/enroll/SS25/asl>.
- (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.
- (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.**
- (Code)  
When compiling the final code, ensure that you use optimization flags (e.g. for GCC use the flag “-O3”) unless indicated otherwise.
- (Neatness)  
5 points in this homework are given for neatness.

### Additional instructions:

- If you have an x86 processor, make sure to [disable frequency scaling](#) in your computer to get accurate timing measurements.

### Exercises:

#### 1. (15 pts) Get to know your machine

Determine and create a table for the following microarchitectural parameters of your computer:

- (a) Processor manufacturer, name, number and microarchitecture (e.g. Skylake, Ice-Lake, etc).
- (b) CPU base frequency.
- (c) CPU maximum frequency. Does your CPU support Frequency scaling (Turbo Boost or a similar technology)?
- (d) Phase in the Intel's development model: Tick, Tock or Optimization. (if applicable)

X86 processors offer two different floating-point instruction sets, namely x87 and SSE/SSE2, that can perform scalar floating-point operations. For example, a floating-point division can be performed using either FDIV (from x87) or DIVSD (from SSE2) assembly instructions. The x87 instruction set, however, is becoming deprecated but is still supported for backward compatibility.

Consider the following instruction:

```
1 #define MAX(a, b) ((a) > (b) ? (a) : (b))
```

Create a C/C++ file that uses this macro to compute the maximum of two *double-precision floating-point* numbers. What instruction is generated by the compiler in the following cases? Does the instruction belong to the x87 instruction set? Explain your results. *Hint*: you can output the assembly using the `-S` flag using `GCC`.

- (e) Compile using `-O0`, i.e no optimizations.
- (f) Compile using `-O3`.

Note:

- the compiler will remove the computation if you do not use its result,
- if you use values that are known at compile time, i.e if you hardcode the computation by writing something along the lines of `MAX(5.8, 7.8)`, the compiler will simplify the instruction.

For one core and **without** using SIMD vector instructions, determine the following about your machine. In (g)-(j), make sure to use the correct floating-point instruction. Be careful not to choose one from the x87 instruction set in case you have an x86 processor. Provide the reference where you found the latency and throughput information.

- (g) Maximum theoretical floating-point peak performance in flops/cycle.
- (h) Latency [cycles], throughput [ops/cycle] and instruction name for both double- and single-precision floating-point multiplication.
- (i) Latency [cycles], throughput [ops/cycle] and instruction name for both double- and single-precision square root operation.
- (j) Latency [cycles], throughput [ops/cycle] and instruction name for both double- and single-precision division operation.

Notes:

- Latency and throughput information for Intel’s and AMD’s processors can be found in [Agner Fog’s instruction tables](#), [uops](#) and [Intel’s optimization manual](#).
- Intel calls throughput what is in reality the gap = 1/throughput.
- Latency and throughput information for Apple’s M-series processors can be found in [Dougall Johnson’s documentation](#) and in the following semester project [writeup](#).
- The manufacturer’s website will contain information about the on-chip details. (e.g. [Intel 64 and IA-32 Architectures Optimization Reference Manual](#)).
- On Unix/Linux systems, typing `cat /proc/cpuinfo` in a shell will give you enough information about your CPU for you to be able to find an appropriate manual for it on the manufacturer’s website (typically AMD or Intel).
- For Windows 7/10 “Control Panel/System and Security/System” will show you your CPU, for more info “[CPU-Z](#)” will give a very detailed report on the machine configuration.
- For macOS there is `sysctl machdep.cpu.brand_string`.
- Throughout this course, we will consider the FMA instruction as two floating-point operations.

## 2. (20 pts) Symmetric Matrix Multiplication

In this exercise, we provide a [folder](#) for computing  $C = \alpha AB^T$ , with  $A$  being a  $n \times n$  symmetric matrix and  $B$  being a  $n \times n$  matrix. Note, since  $A$  is symmetric, we only need to store half of it. The content of the folder include

- A C file `symm.c` that includes the actual computation,
- A header file `tsc_x86.h` that allows reading the time stamp counter (TSC) on x86 machines,
- A header file `arm_vct.h` that allows reading the VCT registers<sup>1</sup> on ARM machines,
- A header file `kperf.h` that allows reading the Processor Monitoring Unit (PMU) on ARM machines, which can contain various metrics (cycle count, instructions issued, ...). Requires `sudo` access. More information on VCT and PMU can be found [here](#).

The code uses different timers available to time the matrix multiplication. Inspect and understand the code and do the following:

- Using your computer, compile and run the code. Compile with the highest level of optimization provided by your compiler (with GCC, compile with the flag `-O3`). A modern compiler will automatically vectorize this routine. Ensure you get consistent timings between timers and for at least two consecutive executions. Don't forget to disable Turbo Boost (if you can). (No need to answer anything here)
- Inspect the `compute()` function in `symm.c` and answer the following:
  - Determine the exact number of floating-point additions and multiplications that it performs.
  - Determine an upper bound on its operational intensity.
- For all square matrices of sizes  $n$  between 100 and 1500, in increments of 100, create a performance plot with  $n$  on the x-axis and performance (flops/cycle) on the y-axis. Create three series such that:
  - The first series has all optimizations disabled: use flag `-O0`.
  - The second series has the major optimizations except for vectorization: use flags `-O3` and `-fno-tree-vectorize`. If you are using the clang compiler, also add `-fno-slp-vectorize` flag to disable vectorization.
  - The third series has all major optimizations enabled including vectorization: use flags `-O3`, `-ffast-math` and `-march=native`. If you are using an Apple M processor and your compiler doesn't support `-march=native` you can use `-mcpu=apple-mx` instead, where  $x$  is the number of your cpu.

**Note:** it is good practice to set long benchmarks such that they run at night; keep this in mind for your project. You can automate with a simple bash script or Makefile.
- Discuss performance variations of your plots and report the highest performance that you achieved.

### 3. (20 pts) Performance analysis and bounds

Assume that vectors  $u, w, x, y$  and  $z$  of length  $n$  are implemented using double precision floating-point and combined as follows:

$$z_i = z_i + \text{ceil}(u_i/w_i) \cdot x_i + x_i \cdot y_i.$$

We consider a Core i7 CPU with a Skylake microarchitecture. As seen in the lecture, it offers FMA instructions (as part of AVX2). Recall that we consider cost of the FMA instruction as two floating-point operations (an addition and a multiplication). Assume the bandwidths that are given in the additional material from the lecture: [Abstracted Microarchitecture](#). Assume that no optimization is performed that simplifies floating-point arithmetic (i.e. `-ffast-math` flag is not used) and that the `ceil` function is translated to a `roundsd` instruction by the compiler. Answer the following and justify your answers.

- Define a suitable detailed floating-point cost measure  $C(n)$ .
- Compute the cost  $C(n)$  of the computation.

<sup>1</sup>The ARM documentation does not state what VCT stands for, we think it should be virtual cycle timer. For more information read here: <https://developer.arm.com/documentation/102379/0101/The-processor-timers>.

- (c) Consider only one core without using vector instructions (i.e. using flag `-fno-tree-vectorize`) and determine a hard lower bound (not asymptotic) on the runtime (measured in cycles), based on each of the following cases:
- The throughput of the floating-point operations. Assume that no FMA instructions are used. Be aware that the lower bound is also affected by the available ports offered for the computation (see lecture slides).
  - The throughput of the floating-point operations where FMAs are used to fuse an addition and a multiplication (i.e. `-mfma` flag is enabled).
  - The throughput of data reads, for the following two cases: All floating-point data is L1-resident, and all floating-point data is RAM-resident. Consider the best case scenario (peak bandwidth and ignore latency). Note that arrays that are only written are also read and this read should be included.
- (d) Determine an upper bound on the operational intensity. Assume empty caches and consider only reads but note: arrays that are only written are also read and this read should be included.

#### 4. (25 pts) Basic optimization

Consider the following function that computes the square euclidian norm  $\|x-y\|_2^2$ , where  $x, y$  are vectors of doubles of length  $n$ .

```

1 void comp(double *x, double *y, int n) {
2     double s = 0.0;
3     for (int i = 0; i < n; i++) {
4         double m = x[i] - y[i];
5         s += m * m;
6     }
7     x[0] = s;
8 }

```

- Create a benchmarking infrastructure based on the timing function that produces the most consistent results in Exercise 2 for you and for all two-power sizes  $n = 2^6, \dots, 2^{23}$  create a performance plot for the function `comp` with  $n$  on the x-axis (choose logarithmic scale) and performance (in flops/cycle) on the y-axis. Randomly initialize all arrays. For all  $n$  repeat your measurements 30 times reporting the median in your plot. Compile your code with flags `-O3 -mfma2 -fno-tree-vectorize`. If you are using clang, add also the `-fno-slp-vectorize` and `-ffp-contract=fast` flags.
- Considering the latency and throughput information of floating-point operations in your machine, and the dependencies in `comp`, derive an upper bound on the performance (flops/cycles) of `comp` when using the specified flags in (a), i.e., when FMA instructions are enabled (`-mfma`) but vectorization is disabled (`-fno-tree-vectorize`).
- Perform optimizations that increase the ILP of function `comp` to improve its runtime. It is not allowed to use vector instructions. Add the performance to the previous plot (so one plot with two series in total for (a) and (c)). Compile your code with the same flags as before.
- Discuss performance variations of your plot and report the highest performance that you achieved. Also discuss the optimizations that you performed to increase the ILP.
- Enroll and submit the code of your optimized function in [Code Expert](#). Carefully read and follow the instructions given in Code Expert to submit your code.

<sup>2</sup>For Apple M1/M2 processors, the flag `-mfma` may not be supported. If this is the case, use instead `-mcpu=apple-m1` or `-march=native`.

5. (15 pts) ILP analysis

Consider the following computations:

```
1 double comp(double a, double b, double c, double d) {
2     double t0 = a * b;
3     double t1 = b * c;
4     double t2 = c * d;
5     double t3 = t0 + t1;
6     double t4 = t3 + t2;
7     double t5 = t2 / a;
8     return t4 + t5;
9 }
```

Make the same assumptions as in Exercise 3, i.e., consider a Skylake processor, only one core without using vector instructions (using flag `-fno-tree-vectorize`), and assume that no optimization is performed that simplifies floating-point arithmetic (i.e. `-ffast-math` flag is not used).

- (a) State the latency, throughput and port usage of the double-precision floating-point addition, multiplication and division instructions.

Determine hard lower bounds (not asymptotic) on the runtime (measured in cycles) for the following cases. Base your analysis on the *latency*, *throughput*, and *dependencies* of the floating-point operations. Be aware that the lower bound is also affected by the *available ports* offered for the computation (see lecture slides). It may be useful to draw the dependency graph of the computation. Justify your answers.

- (b) Assume that the operations are issued in the order they are written, i.e., `t1` is computed after (or simultaneously with) `t0`, `t2` after (or simultaneously with) `t1` and `t0`, and so on. Determine a hard lower bound on the runtime for `comp`.
- (c) Now assume that operations can be issued out-of-order, meaning that `t1 = b * c` can be issued before `t0 = a * b`. However, data dependencies must still be respected. Determine a hard lower bound on the runtime for `comp`.

## How to disable Frequency scaling

Frequency scaling is a technology that enables a processor to run above its base operating frequency via dynamic control of the processor's clock rate. It is activated when the operating system requests the highest performance state of the processor.

### BIOS

Frequency scaling is typically enabled by default. You can only disable and enable the technology through a switch in the BIOS. No other user controllable settings are available. Once enabled, frequency scaling works automatically under operating system control. When access to BIOS is not available, few workarounds are possible:

### Linux

Linux does not provide an interface to disable frequency scaling. One alternative, that works, is disabling frequency scaling by writing into MSR registers. Assuming 2 cores, the following should work:

```
wrmsr -p0 0x1a0 0x4000850089
wrmsr -p1 0x1a0 0x4000850089
```

To enable it:

```
wrmsr -p0 0x1a0 0x850089
wrmsr -p1 0x1a0 0x850089
```

Using `wrmsr -a` should set all the registers. Check the manual entry for `wrmsr`.

This method has been criticized [here](#) and [here](#) stating that the OS can circumvent the MSR value, using an opportunistic strategy. But so far in our tests, we have observed that Linux conforms to the MSR value. An alternative method would be to use `cpupower`, as explained in the [ArchLinux Wiki](#), as well as the driver `intel_pstate`. Unfortunately, we cannot confirm the deterministic behavior between different kernel versions with this method.

In case the above methods do not work, this [guide](#) suggests using these instructions to disable frequency scaling on Intel and AMD CPUs:

```
# Intel
echo 1 > /sys/devices/system/cpu/intel_pstate/no_turbo
# AMD
echo 0 > /sys/devices/system/cpu/cpufreq/boost
```

### *Mac OS*

For Intel Macs, disabling Turbo Boost in OS X can be easily done with the [Turbo Boost Switcher for OS X](#). Note that the change is not persistent after restart. The method also writes to the MSR register, and shares the same weaknesses as the Linux approach.

Arm Macs do not have a way of disabling frequency scaling.

### *Windows*

Windows does not provide any functionality to disable Intel Turbo Boost. The only effective way to disable it is using the BIOS. On some Intel machines, however, it is possible to fix the CPU multiplier such that the resulting frequency corresponds to the nominal frequency of the CPU. [ThrottleStop](#) provides this functionality with a convenient GUI. “Disable Turbo” will effectively fix the frequency such that it corresponds to the behavior of a CPU with disabled Turbo Boost.

## Benchmarking Tips

The following [guide](#) list various steps that help with benchmarking. In particular, [setting the scaling governor to performance](#) and [using the taskset utility](#) to run a piece of code on a specific CPU. Consider using these for your project.