

263-0007-00: Advanced Systems Lab

Assignment 2: 80 points

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

<https://acl.inf.ethz.ch/teaching/fastcode/2022/>

Questions: fastcode@lists.inf.ethz.ch

Exercises:

1. Short project info (5 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.

Solution: This will be different for each student.

2. Optimization Blockers (30 pts)

In this exercise, we consider the following short computation that is available in Code Expert:

```
1 void slowperformance1(mat* x, mat* y, mat*z) {
2     double t1, t2, t3;
3     for (int i = 0; i < z->n1; i++) {
4         for (int j = 0; j < z->n2 - 1; j++) {
5             switch((i + j) % 2) {
6                 case 0:
7                     t1 = (C1 * mat_get(z,i,j) + mat_get(x,i,j)) / sqrt(2.0);
8                     mat_set(z,i,j,t1);
9                     break;
10                case 1:
11                    t2 = mat_get(z,i,j) / fmax(mat_get(y,0,i), mat_get(x,0,i));
12                    mat_set(z,i,j,t2);
13                    break;
14            }
15            t3 = mat_get(z,i,j+1) + C2 * mat_get(y,0,(5 * j) % 3);
16            mat_set(z,i,j+1,t3);
17        }
18    }
19 }
```

This is part of the supplied code in Code Expert:

- Read and understand the code. It enables you to register functions with the same signature, which will be timed in a microbenchmark fashion.
- Create new functions where you perform optimizations to improve the runtime. For example, loop unrolling and scalar replacement, strength reduction, inlining, and others.
- You may apply any optimization that produces the same result in exact arithmetic.
- 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, it will print the runtime in cycles.

- Implement in function *maxperformance* the implementation that achieves the best runtime. This is the one that will be autograded by Code Expert.
- For this task, the Code Expert system compiles the code using GCC 8.3.1 with the following flags: `-O3 -fno-tree-vectorize`. Note that with these flags vectorization and FMAs are disabled.
- It is not allowed to use vector intrinsics or FMAs to speedup your implementation.

Discussion:

- (a) Create a table with the runtime numbers of each new function that you created (include at least 2). Briefly discuss the table explaining the optimizations applied in each step.

Solution:

Implementation	Impl. 1	Impl. 2	Impl. 3	Impl. 4	Impl. 5
Runtime (cycles)	441.1K	114.1K	34.7K	17.4K	17.2K

The table above reports runtime in cycles for six different implementations of the above code, with optimizations turned on (`-O3 -fno-tree-vectorize`). The K stands for thousands. These numbers were recorded on a Intel(R) Xeon(R) Silver 4210 @ 2.20GHz Cascade Lake with hyper-threading disabled, and compiled with GCC 8.3.1.

Implementation 1 is the original code. Implementation 2 removes calls to `mat_set` and `mat_get` by accessing the arrays directly. Implementation 3 unrolls both loops two times. This allows to remove the switch-case. In addition, some computations are precomputed. Implementation 4 unrolls the inner loop three more times to be able to remove the expensive modulo operations (%). In addition, values that stay constant across iterations are precomputed. Finally, implementation 5 removes extra memory accesses that are repeated between iterations..

- (b) What is the speedup of function *maxperformance* compared to *slowperformance1*?

Solution: The speedup of implementation 5 is 25.6.

- (c) What is the performance in flops/cycle of your function *maxperformance*.

Solution: Implementation 5 performs around $3n_1n_2$ flops. The performance is 1.6 flops/cycles for $n_1 = 96$ and $n_2 = 97$.

3. Microbenchmarks(40 pts)

Your task is to write a program (without vector instructions, i.e., standard C) in Code Expert that benchmarks the latency and inverse throughput (also called “gap” in class) of floating point FMA instructions and the square root on doubles. For the square root we provide a C function called *sqrtsd* to access directly the instruction. In addition, the latency and gap of the function $f(a, b) = \sqrt{a^2 + 0.1 * b}$. We provide the implementation of $f(x)$ in `foo.h`. More specifically:

- Read and understand the code given in Code Expert.
- Implement the functions provided in the skeleton in file `microbenchmark.cpp`:

```
void initialize_microbenchmark_data (microbenchmark_mode_t mode);
double microbenchmark_get_fma_latency();
double microbenchmark_get_fma_gap();
double microbenchmark_get_sqrt_latency();
double microbenchmark_get_sqrt_gap();
double microbenchmark_get_foo_latency();
double microbenchmark_get_foo_gap();
```

- You can use the `initialize_microbenchmark_data` function for any kind of initialization that you may need (e.g. for initializing the input values).
- Note that the latency and gap of floating point square root can vary depending on its inputs. Thus, you are also required to find the minimum latency and gap for square root and function $f(x)$. Hint: You can try using values where performing the square root becomes trivial.
- It is not allowed to manually inline the functions in `foo.h` into the implementation of your microbenchmarks.

- Make sure that your benchmarks yield stable measurements between runs.

Additional information:

- Our Code Expert system already has Turbo Boost disabled. However, note that 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.
- For this task, our Code Expert system uses GCC 8.3.1 to compile the code with the following flags: `-O3 -fno-tree-vectorize -march=skylake`. Note that with these flags vectorization is disabled but FMAs are enabled.

Discussion:

- (a) Do the latency and gap of floating point FMA and square root match what is in the [Intel Optimization Manual](#)? If no, explain why. (You can also check [Agner's Table](#)).

Solution: Yes, the manual reports a latency and gap for FMA instructions (Skylake) of 4 and 0.5 cycles respectively which is consistent with the microbenchmarks. For the square root (`sqrtsd`), the manual reports 18 and 6 cycles for latency and gap respectively and Agner reports 15-16 and 4-6 cycles respectively. The measured latency and gap in the microbenchmarks are 18 and 6 cycles for a random value and 13 and 4.5 cycles for a trivial input value. These values are consistent with the Intel's manual. Further, the numbers by Agner are within the measured range.

- (b) Based on the dependency, latency and gap information of the floating point operations, is the measured latency and gap of function $f(x)$ close to what you would expect? Justify your answer.

Solution: Yes, function $f(x)$ consists of two multiplications that can be scheduled in parallel, an addition and a square root operation. A multiplication will be fused with the addition into an FMA instruction. This gives a theoretical latency of 4 (mult) + 4 (fma) + 18 (sqrt) = 26 cycles which is consistent with the measurements. For the theoretical gap, the square root becomes the bottleneck. Thus, the gap is 6 cycles which is also close to the measurements.

- (c) Will the latency and gap of $f(x)$ change if we compile the code with flags `-O3 -fno-tree-vectorize` (i.e., with FMAs disabled)? Justify your answer and state the expected latency and gap in case you think it will change.

Solution: No, if FMAs are disabled, the addition and multiplication in function $f(x)$ will not fuse. However, the two multiplications can execute in parallel. Thus, the latency will remain the same 4 (mul) + 4 (add) + 18 (sqrt) = 26. The gap will remain the same as the square root is still the bottleneck.