

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

Assignment 2: 100 points

Due Date: Thu March 13 17:00

<http://www.inf.ethz.ch/personal/markusp/teaching/263-2300-ETH-spring14/course.html>

Questions: fastcode@lists.inf.ethz.ch

Submission instructions (read carefully):

- (Submission)
We set up a SVN Directory for everybody in the course. The Url of your SVN Directory is <https://svn.inf.ethz.ch/svn/pueschel/students/trunk/s14-fastcode/YOUR.NETZH.LOGIN/> You should see sub-directory for each homework.
- (Late policy)
You have 3 late days, but can use at most 2 on one homework. Late submissions have to be emailed to fastcode@lists.inf.ethz.ch.
- (Formats)
If you use programs (such as MS-Word or Latex) to create your assignment, convert them to PDF and submit to svn in the top level of the respective homework directory. Call it homework.pdf.
- (Plots)
For plots/benchmarks, be concise, but provide necessary information (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 (lecture 5).
- (Neatness)
5% of the points in a homework are given for neatness.

Exercises:

1. *Short project info (10 pts)* Go to the [list of mile stones 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. *Polynomial Evaluation (25 pts)* [Code needed](#)

The code in `poly.c` contains a function for evaluating a polynomial of degree N at a point $x_0 = 0.2$. N is given as a runtime parameter from terminal.

- (a) Inspect the function `poly` and determine its op count (double additions and multiplications only).
- (b) Assign the value computed in the previous point to the macro `OPCOUNT` in `poly.c`. Compile the code disabling vectorization and determine its runtime and performance for $N = 2^i$, $i = 7, \dots, 11$. Collect the results in a small table.

Now implement a function `horner` which uses [Horner scheme](#) for evaluating the polynomial. The function should provide exactly the same signature of the function `poly`.

- (c) What would the op count be for the new implementation?

- (d) What performance would you expect to obtain using `horner` instead of `poly`? Why? In the discussion provide all the numbers at the base of your assumption (e.g., the latency of addition and multiplication on your CPU).
- (e) In `poly.c`, change the value of the macros `FUNC` and `OPCOUNT` to `horner` and the cost from `??`, respectively. Recompile your code (always disabling vectorization) and determine its runtime and performance for $N = 2^i$, $i = \{7..11\}$. Again, collect your results in a table.
- (f) Compare the results in the two tables obtain from [2b](#) and [2e](#): Which code exhibits higher performance? Which lower runtime? Briefly discuss.
- (g) Identify key performance limitations in function `horner` and implement an optimized version of the function called `horner2`. Remember to adjust the `FUNC` macro before recompiling the code. **Hint:** A polynomial $p(x)$ of degree N can be expressed as $p(x) = \sum_{i=0}^k x^i p_i(x)$ with $k|N$. The latter formulation supports one specific optimization that helps attaining 80% of peak when compiling with `icc` or `gcc` on Sandy Bridge (i.e., disabling vectorization, 1.6 flops/cycle).

Report compiler, version, and flags. Submit your code to the SVN.

3. Optimization Blockers (40 pts) Code needed

Download, extract and inspect the code. Your task is to optimize the function called `superslow` (guess why it's called like this?) in the file `comp.c`. The function runs over an $n \times n$ matrix and performs some computation on each element. In its current implementation, `superslow` involves several optimization blockers. Your task is to optimize the code.

Edit the Makefile if needed (architecture flags specifying your processor). Running `make` and then the generated executable verifies the code and outputs the performance (the flop count is underestimated, since the trigometric functions are ignored) of `superslow`. Proceed as follows

- (a) Identify optimization blockers discussed in the lecture and remove them.
- (b) For every optimization you perform, create a new function in `comp.c` that has the same signature and register it to the timing framework through the `register_function` procedure in `comp.c`. Let it run and, if it verifies, determine the performance.
- (c) In the end, the innermost loop should be free of any procedure calls and operations other than adds and mults.
- (d) When done, rerun all code versions also with optimization flags turned off (`-O0` in the Makefile).
- (e) Create a table with the performance numbers. Two rows (optimization flags, no optimization flags) and as many columns as versions of `superslow`. Briefly discuss the table.
- (f) Submit your `comp.c` to the SVN

What speedup do you achieve?

4. Locality of a Convolution (20 pts)

Consider the following C code, which computes a form of 2D convolution. The function takes as input a matrix A of size $n \times n$, and a vector H of size k^2 ; the result is stored in a matrix B of size $n \times n$.

```

assert(N > K && K > 2);
double A[N][N], B[N][N], H[K*K];
for (int i = 0; i < N-K; i++)
    for (int j = 0; j < N-K; j++)
        for (int k = 0; k < K*K; k++)
            B[i][j] = A[i+1][j+2] * A[i+k/K][j+k%K] * H[K*K-k-1];
for (int i = N-K; i < N; i++) {
    for (int j = N-K; j < N; j++)
        for (int k = 0; k < K*K; k++)
            B[i][j] = A[i-1][j-2] * A[i-k/K][j-k%K] * H[K*K-k-1];

```

Inspecting the data accesses, where do you see

- (a) Temporal locality?
- (b) Spatial locality?