

CS 380C: Advanced Topics in Compilers

Assignment 7: Fast matrix multiply

Due: April 26th

April 12, 2016

Late submission policy: Submission can be at the most 2 days late. There will be a 10% penalty for each day after the due date (cumulative).

In class, we described the structure of the optimized MMM code produced by ATLAS. The “computational heart” of this code is the mini-kernel that multiplies an $NB \times NB$ block of matrix A with an $NB \times NB$ block of matrix B into an $NB \times NB$ block of matrix C , where NB is chosen so that the working set of this computation fits in the cache. The mini-kernel itself is performed by repeatedly calling a micro-kernel that multiplies an $MU \times 1$ column vector of matrix A with a $1 \times NU$ row vector of matrix B into an $MU \times NU$ block of matrix C . The values of MU and NU are chosen so that the micro-kernel can be performed out of the registers. Pseudocode for the mini-kernel is shown below (note that this code assumes that NB is a multiple of MU and NU).

```
//mini-kernel
for (int j = 0; j < NB; j += MU)
  for(int i = 0; i < NB; i += NU)
    load C[i..i+MU-1, j..j+NU-1] into registers
    for (int k = 0; k < NB; k++)
      //micro-kernel
      load A[i..i+MU-1,k] into registers
      load B[k,j..j+NU-1] into registers
      multiply A's and B's and add to C's
    store C[i..i+MU-1, j..j+NU-1]
```

The data type in the matrix should be doubles.

For each optimization below:

- Compile your code in ICC with flags ‘-O3 -fp-model precise’.
- Submit your run to the job scheduler on Stampede at TACC - use the ‘serial’ queue. Since the values you obtain will depend a lot on the machine you use, you must use Stampede for the numbers you report.

- Report the performance of your code in GFLOPS. The number of floating point operations in matrix multiplication is $2 * N^3$, where $N \times N$ is the size of each matrix and so, FLOPS is given by $2 * N^3 / time$ (1 GFLOP = 10^9 FLOPS).
- Run the code to multiply matrices of various sizes (at least 5) and plot a graph of GFLOPS vs. matrix size. Explain your results briefly.

Note: Each optimization is cumulative, i.e., you implement one optimization on top of another.

As a reference, for 4096x4096 size matrices, the performance of *ijk* matrix multiplication version on Stampede is around 0.29 GFLOPS, the performance of matrix multiplication using BLAS is around 22 GFLOPS.

Register-blocking

20 points

To measure the impact of register-blocking without cache-blocking, implement register-blocking by writing a function for performing MMM, using the mini-kernel code with $NB = N$ (you should verify that this implements MMM correctly). You can use the results in the [Yotov et al. study of the ATLAS system](#) to determine good values for MU and NU, or you can experiment with different values to find good values.

Note: You will have to write some clean-up code to handle leftover pieces of the matrices when the value of N is not a multiple of your values for MU and NU.

Cache-blocking

20 points

Modify the above code to implement both register-blocking and cache-blocking. You will have to wrap three loops around the mini-kernel to get a full MMM. Use any method you wish to determine a good value for NB.

Note: You will have to write some clean-up code to handle leftover pieces of the matrices when the value of N is not a multiple of your value for NB.

Data copying

20 points

You can improve the performance of your kernel by copying blocks of data into contiguous storage as explained in the [Yotov et al. paper](#). Modify the above code to implement data copying (along with register-blocking and cache-blocking).

Vectorization

20 points

Modify the above code (along with register-blocking, cache-blocking, and data copying) to use vector registers and vector intrinsics instead of scalar registers.

Fastest

20 points

We will have a competition for the best performing matrix multiplication code. Your code should be correct for any square matrices; the performance will be measured only for 4096x4096 size matrices. You may use any of the techniques described above or from the literature but your code should not call any external library and you are not allowed to copy source code from anywhere. You will be graded on performance. You will not get any points if we are not able to reproduce your results.

Implementation notes

- Refer [PHiPAC coding guidelines](#) for writing “portable assembly language programs” in the C language.
- Stampede on TACC: Use the login node only for development - do not run or debug any executable on it. Run and debug your applications using the job scheduler. Read [this](#) to learn how to submit jobs.
- To understand the performance of your code, use **PAPI** to measure performance counters like L1 cache misses.
To program with PAPI on stampede, run:
`module load papi`
For help on using the module, run:
`module help papi`
For more information on using modules, check [this](#).
To see which papi counters are available on a host, run:
`papi_avail`
Read the PAPI [manual](#) for more information, including example code.

Deliverables

Submit (to canvas) an archive (preferably, `.tar.gz/.tgz`) of your (one) fastest code containing all optimizations (cumulative) and a report containing all plots and analysis. Briefly describe a way to verify correctness of your code. You will not be given any points if we are not able to verify correctness for randomly chosen square matrices.