

BLISlab: A Sandbox for Optimizing GEMM

Jianyu Huang

Robert A. van de Geijn

Draft
February 8, 2016

Abstract

1 Introduction

Matrix-matrix multiplication (GEMM) is frequently used as a simple example with which to raise awareness of how to optimize modern processors. A reason is that the operation is simple to describe, challenging to fully optimize, and of practical importance. In this paper, we walk the reader through the techniques that underly the currently fastest implementations for CPU architectures.

1.1 A minimal history

Need to mention BLAS3 [1] paper.

The advent of cache-based architected, high-performance implementation of GEMM necessitated careful attention to the amortization of the cost of data movement between memory layers and computation with that data [2]. To keep this manageable, it helps to realize that only a “kernel” that performs a matrix-matrix multiplication with relatively small matrices needs to be highly optimized, since computation with larger matrices can be blocked to then use such a kernel without an adverse impact on overall performance. This insight was first explicitly advocated in

Bo Kågström, Per Ling, Charles Van Loan.

GEMM-based level 3 BLAS: high-performance model implementations and performance evaluation benchmark.

ACM Transactions on Mathematical Software (TOMS).

Volume 24 Issue 3, p.268-302, Sept. 1998.

For more than a decade after that paper, the intricacies of high-performance optimization of GEMM was considered to be sufficiently complex that it should be left to the hardware vendors, yielding IBM’s ESSL, Intel’s MKL, Cray’s ???, and AMD’s ACML libraries, or auto-generated as advocated in papers on the Portable High Performance ANSI C (PHiPAC) guidelines for writing high-performance matrix-matrix multiplication in C [3] and the Automatically Tuned Linear Algebra Software (ATLAS) [4].

Around 2000, Kazushige Goto revolutionized how GEMM is implemented on current CPUs with his techniques that were first published in the paper

Kazushige Goto, Robert A. van de Geijn.

Anatomy of high-performance matrix multiplication.

ACM Transactions on Mathematical Software (TOMS).

Volume 34 Issue 3, May 2008, Article No. 12.

At the end of this note we will discuss the major insights in this paper.

1.2 The BLIS-like Library Instantiation Software (BLIS)

More recently, the BLAS-like Library Instantiation Software (BLIS) “refactored” the approach pioneered by Goto, exposing additional loops around a *micro-kernel*, as described in

Field G. Van Zee, Robert A. van de Geijn.
BLIS: A Framework for Rapidly Instantiating BLAS Functionality.
ACM Transactions on Mathematical Software (TOMS).
Volume 41 Issue 3, June 2015, Article No. 14.

One goal of the BLIS paper was to further expose the layering of Goto’s approach while simultaneously improving portability by reducing how much code must be written at a low level (e.g., in assembly code).

1.3 You too can optimize like a pro

The purpose of this note is to expose the basic techniques that underlie the best implementations of GEMM so that you too can achieve high-performance for such operations.

2 Step 1: The Basics

2.1 Simple matrix-matrix multiplication

In our discussions, we will consider the computation

$$C := AB + C$$

where A , B , and C are $m \times k$, $k \times n$, $m \times n$ matrices. respectively. Letting

$$A = \begin{pmatrix} \alpha_{0,0} & \cdots & \alpha_{0,k-1} \\ \vdots & & \vdots \\ \alpha_{m-1,0} & \cdots & \alpha_{m-1,k-1} \end{pmatrix}, B = \begin{pmatrix} \beta_{0,0} & \cdots & \beta_{0,n-1} \\ \vdots & & \vdots \\ \beta_{k-1,0} & \cdots & \beta_{k-1,n-1} \end{pmatrix}, \text{ and } C = \begin{pmatrix} \gamma_{0,0} & \cdots & \gamma_{0,n-1} \\ \vdots & & \vdots \\ \gamma_{m-1,0} & \cdots & \gamma_{m-1,n-1} \end{pmatrix}.$$

$C := AB + C$ computes

$$\gamma_{i,j} := \sum_{p=0}^{k-1} \alpha_{i,p} \beta_{p,j} + \gamma_{i,j}.$$

If A , B , and C are stored as floating point numbers in two-dimensional arrays **A**, **B**, and **C**, the following pseudocode computes $C := AB + C$:

```
for i=0:m-1
  for j=0:n-1
    for p=0:k-1
      C( i,j ) := A( i,p ) * B( p,j ) + C( i,j )
    endfor
  endfor
endfor
```

Counting a multiply and an add separately, the computation requires $2mnk$ floating point operations flops.

2.2 Set up

To let you efficiently learn about how to efficiently compute, you start your project with much of the infrastructure in place. We have structured the subdirectory, **step1**, somewhat like a project that implements a real library might. This may be overkill for our purposes, but how to structure a software project is a useful skill to learn.

Consider Figure 3, which illustrates the directory structure for subdirectory **step1**:

```

step1
├── README
├── sourceme.sh
├── makefile
├── dgemm
│   ├── my_dgemm.c
│   ├── bl_dgemm_ref.c
│   └── bl_dgemm_util.c
├── include
│   ├── bl_dgemm.h
│   ├── bl_dgemm_ref.h
│   └── bl_config.h
├── lib
│   ├── libblislab.a
│   └── libblislab.so
├── makefile.in.files
│   ├── make.intel.inc
│   ├── make.gnu.inc
│   └── make.inc
└── test
    ├── makefile
    ├── test_bl_dgemm.c
    ├── run_bl_dgemm.sh
    ├── test_bl_dgemm.x
    └── tacc_run_bl_dgemm.sh

```

Figure 1: Structure of directory `step1`.

`README` Is a file that describes the contents of the directory and how to compile and execute the code.

`sourceme.sh` Is a file that configures the environment variables.

`BLISLAB_USE_INTEL` determines whether you use Intel compiler or GNU compiler.

`BLISLAB_USE_BLAS` determines whether your reference GEMM adopts BLAS implementation (if you have BLAS installed on your machine), or the simple triple loops implementation.

`OMP_NUM_THREADS` and `BLISLAB_IC_NT` determines your thread number for parallel version of your code. For your first step, you can just set them both to 1.

`dgemm` Is the subdirectory routines that implement GEMM can be found. In it

`bl_dgemm_ref` contains the routine `dgemm_ref` that is a simple implementation of GEMM that you will use to check the correctness of your implementations.

`my_dgemm` contains the routine `dgemm` that that initially is a simple implementation of GEMM and that you will optimize as part of the first step on your way to mastering how to optimize GEMM.

`bl_dgemm_util` contains utility routines that will come in handy later.

`include` This directory contains include files with various macro definitions and other header information.

`lib` This directory will hold libraries generated by your implemented source files (`libblislab.so` and `libblislab.a`). You can also install a reference library (e.g. OpenBLAS) in this directory to compare your performance.

`test` This directory contains “test drivers” and correctness/performance checking scripts for the various implementations.

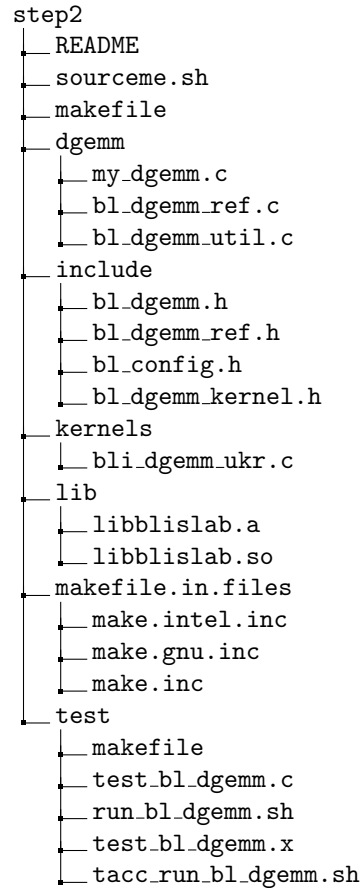


Figure 2: Structure of directory `step2`.

`test_bl_dgemm.c` contains the “test driver” routine `test_bl_dgemm`.

`test_bl_dgemm.x` is the executable file for `test_bl_dgemm.c`.

`run_bl_dgemm.sh` contains a bash script to collect GFLOPS result for selected problem size.

`tacc_run_bl_dgemm.sh` contains a SLURM script for you to submit the job to TACC machines to measure performance.

3 Step 2: Blocking

3.1 Set up

Figure 3 illustrates the directory structure for subdirectory `step2`. Comparing to `step1`, we have modified/added the following directories/files:

`include bl_config.h` configures the step size for different loops (MC, NC, KC, MR, NR).

`dgemm my_dgemm.c` contains the routine `dgemm` that employs the packing routines and blocking algorithms for GEMM.

`kernels` This directory contains the micro-kernel implementations for various architecture.

`bli_dgemm_ukr.c` Is a reference C implementation of micro-kernel.

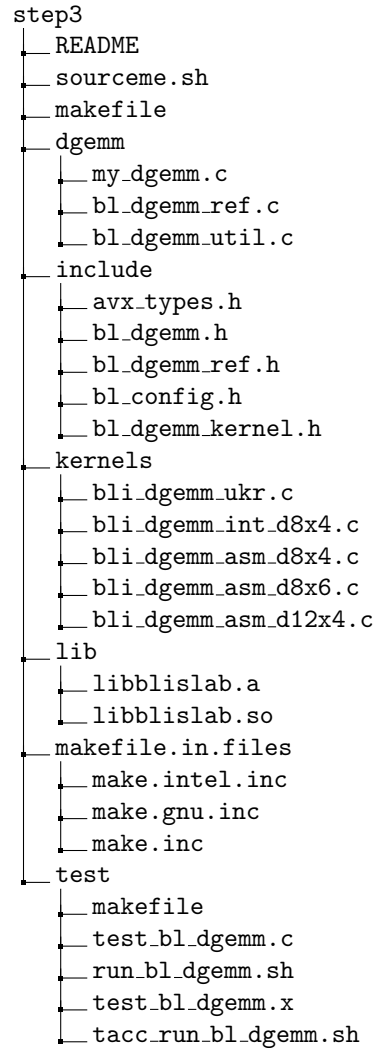


Figure 3: Structure of directory `step3`.

4 Step 3: GOTO paper: Combine Step1 (small matrix for kernel) and Step2 (blocking for framework)

4.1 Set up

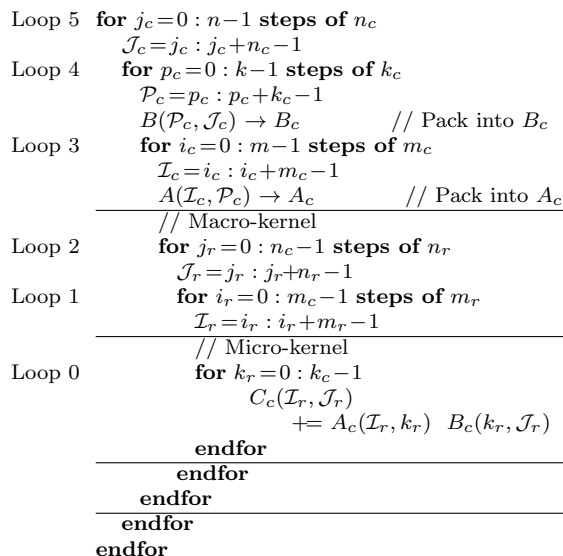
Figure 3 illustrates the directory structure for subdirectory `step3`. Comparing to `step2`, we have modified/added the following directories/files:

kernels This directory contains the micro-kernel implementations for various architecture.

`bli_dgemm_int_d8x4.c` Is a AVX intrinsics micro-kernel implementation for Sandy Bridge/Ivy Bridge micro-architecture.

`bli_dgemm_asm_d8x4.c` Is a AVX assembly micro-kernel implementation for Sandy Bridge/Ivy Bridge micro-architecture.

`bli_dgemm_asm_d8x6.c` Is a AVX2 assembly micro-kernel implementation for Haswell micro-architecture.



```

Users
├── rvdg
│   ├── blislab ..... The main directory.
│   │   ├── docs ..... The directory where this document and its LATEX source exist.
│   │   └── src ..... Directory where the source for the programming exercises exists.
│   │       ├── step1 ..... Source for step one and two of the exercise.
│   │       ├── step3 ..... Source for step three of the exercise.
│   │       └── step_final ..... Source for final step of the exercise.

```

`bli_dgemm_asm_d12x4.c` Is an alternative AVX2 assembly micro-kernel implementation for Haswell micro-architecture.

$$\begin{bmatrix} 2 \\ 6 \\ 5 \\ 3 \\ 4 \end{bmatrix}$$

6 Organization of the Project

One goal of this exercise is to teach the reader how a software library project is often organized into directories. We acknowledge that the structure is probably overkill for this relatively simple situation, but hope that it has value nonetheless.

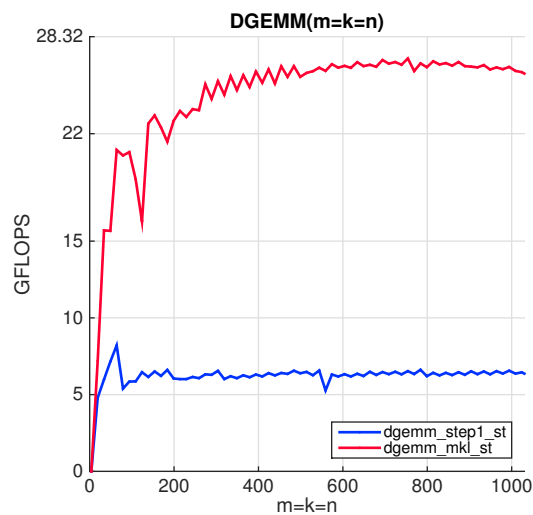


Figure 6: Step1 performance (triple loops)

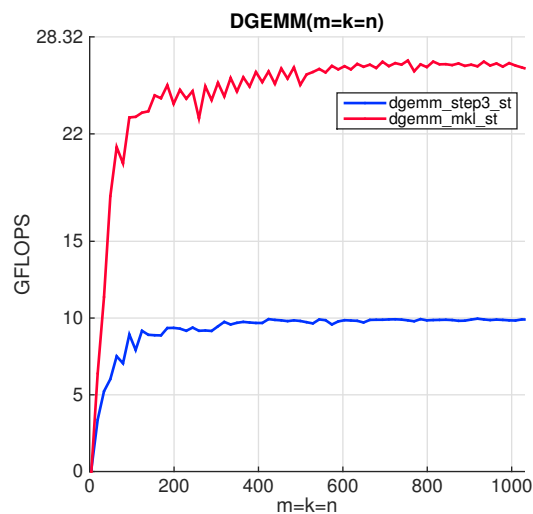


Figure 7: Step3 performance (6 loops with blocking)

7 Conclusion

Conclusion.

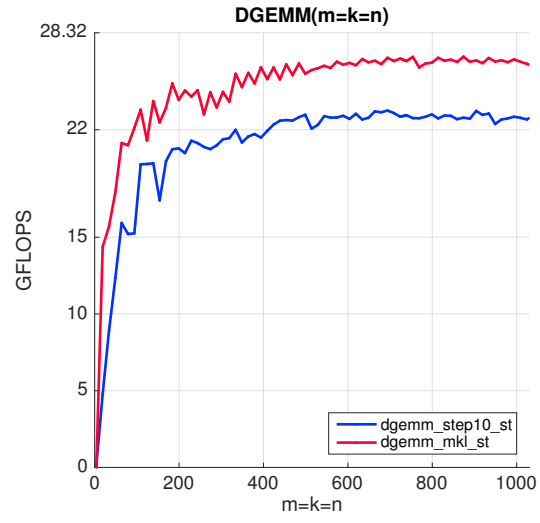


Figure 8: Step10 performance (blocking+Intrinsics)

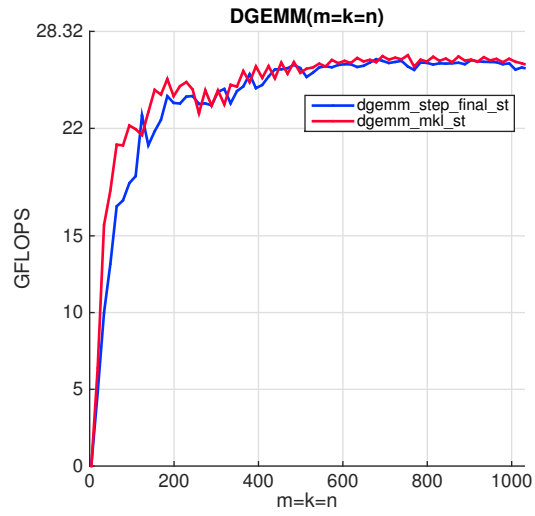


Figure 9: Step Final performance (blocking+AVX assembly)

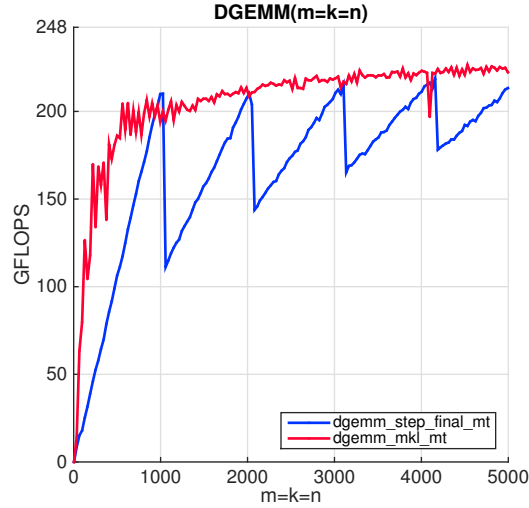


Figure 10: Step Final performance (multi-thread), sawtooth due to workload imbalance in OMP threads

Additional information

For additional information on FLAME visit

<http://www.cs.utexas.edu/users/flame/>.

Acknowledgements

We thank the other members of the FLAME team for their support. This research was partially sponsored by NSF grant CCF-***.

Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation (NSF).

References

- [1] Jack J. Dongarra, Jeremy Du Croz, Sven Hammarling, and Iain Duff. A set of level 3 basic linear algebra subprograms. *ACM Trans. Math. Soft.*, 16(1):1–17, March 1990.
- [2] Kazushige Goto and Robert A. van de Geijn. Anatomy of a high-performance matrix multiplication. *ACM Trans. Math. Soft.*, 34(3):12, May 2008. Article 12, 25 pages.
- [3] Tyler M. Smith, Robert van de Geijn, Mikhail Smelyanskiy, Jeff R. Hammond, , and Field G. Van Zee. Anatomy of high-performance many-threaded matrix multiplication. In *International Parallel and Distributed Processing Symposium 2014*, 2014.
- [4] Tze Meng Low, Francisco D. Igual, Tyler M. Smith, and Enrique S. Quintana-Ortí. Analytical modeling is enough for high performance blis. *ACM Transactions on Mathematical Software*. in review.
- [5] Field G. Van Zee, Tyler Smith, Bryan Marker, Tze Meng Low, Robert A. van de Geijn, Francisco D. Igual, Mikhail Smelyanskiy, Xianyi Zhang, Michael Kistler, Vernon Austel, John Gunnels, and Lee Killough. The blis framework: Experiments in portability. *ACM Transactions on Mathematical Software*. to appear.
- [6] Field G. Van Zee and Robert A. van de Geijn. BLIS: A framework for rapidly instantiating blas functionality (replicated computational results certified). *ACM Trans. Math. Soft.*, 41(3):14:1–14:33, June 2015.