

Level-3 BLAS on a GPU: Picking the Low Hanging Fruit

Francisco D. Igual Gregorio Quintana-Ortí
Depto. de Ingeniería y Ciencia de Computadores
Universidad Jaume I
12.071–Castellón, Spain
{figual,gquintan}@icc.uji.es

Robert A. van de Geijn
Department of Computer Sciences
The University of Texas at Austin
Austin, TX 78712
rvdg@cs.utexas.edu

Technical Report DICC 2009-04-01

FLAME Working Note #37

First published: April 30, 2009

Updated: May 21, 2009

Abstract

The arrival of hardware accelerators has created a new gold rush to be the first to deliver their promise of high performance for numerical applications. Since they are relatively hard to program, with limited language and compiler support, it is generally accepted that one needs to roll up one's sleeves and tough it out, not unlike the early days of distributed memory parallel computing (or any other period after the introduction of a drastically different architecture). In this paper we remind the community that while this is a noble endeavor, there is a lot of low hanging fruit that can be harvested easily. Picking this low hanging fruit benefits the scientific computing community immediately and prototypes the approach that the further optimizations may wish to follow. We demonstrate this by focusing on a widely used set of operations, the level-3 BLAS, targeting the NVIDIA family of GPUs.

*Insanity: doing the same thing over and over again
and expecting different results.*

– Albert Einstein (1879-1955)

1 Introduction

Every time a new architecture arrives, there is a mad dash for high performance. Since compilers, languages, and tools are still rudimentary, this means that some experts roll up their sleeves and achieve high performance the old-fashioned way: they earn it. The problem is that often there are only a few with the right expertise and interest, and therefore this yields only a few routines that are highly optimized. Furthermore, it is acceptable for code that achieves high performance to be messy. When others then come into the picture, they use such implementations as their inspiration, meaning that programmability does not enter the picture until much later in the game. In this paper, we show how insights from the FLAME project, in particular the importance of having a

family of algorithms at one's disposal, allow considerable performance gains to be attained with minimal effort. We do so by focusing on the familiar and important matrix-matrix operations that are part of the Basic Linear Algebra Subprograms (BLAS) [5] and targeting the NVIDIA family of GPUs.

The arrival of NVIDIA's GPUs and IBM's Cell Broadband Engine and the recognition that they can be used for computation outside of the field of graphics has created the latest gold rush for performance. In scientific computing this has meant that considerable effort has been expended on implementing the most important kernel: matrix-matrix multiplication (GEMM). Very admirable performance has been achieved [13].

Yet, even operations that are very similar to GEMM, e.g., the other level-3 BLAS, did not achieve decent performance in the CUBLAS library for the NVIDIA GPUs when we started this study. Worse, the effort required to achieve the high performance for GEMM is daunting enough that experts like ourselves have stayed on the sideline, focusing our efforts on using the GEMM implementation for high-level operations like Cholesky factorization by using the accelerators only to compute subproblems that are matrix-matrix multiplications [10, 9, 8]. We all hoped that soon other functionality would be ported to the GPUs, but that some other poor soul would do it for us.

In this paper we once again show that as new functionality and optimizations appear, there are, for those of us who have an aversion to hard work, opportunities to quickly and easily help improve performance in the short run while simultaneously prototyping how performance can eventually be improved by those who are willing to code at a lower level.

This paper is organized as follows: In Section 2 we briefly review the three commonly encountered matrix-matrix multiplication algorithms and use this to remind the reader of the FLAME notation for presenting algorithms. In Section 3 we discuss the corresponding algorithms for various level-3 BLAS operations, where these algorithms have been modified to take advantage of special structure in the matrices. The benefits of picking the right algorithmic variant is illustrated in Section 4. Concluding remarks are found in the final section.

2 Matrix-Matrix Product

At the top level, there are three variants of matrix-matrix product, which we have come to refer to as matrix-panel product (GEMM_MP) based, panel-matrix product (GEMM_PM) based, and (outer) panel-panel product (GEMM_PP) based (also known as rank-k updating) [6]. We will discuss these briefly in this section, so that we can refer to them later as we discuss algorithms for the other matrix-matrix operations.

In Figure 1, we illustrate the GEMM_MP based algorithm. At the beginning of the iteration, C_0 and B_0 have already be updated and used, respectively. In the current iteration the next panel of matrix C is updated: $C_1 := C_1 + AB_1$. Then, the advancement for the next iteration shifts C_1 and B_1 to the next blocks of data making blocks C_0 and B_0 larger since they contain more processed data. This visual description of the algorithm motivates the algorithm, in FLAME notation, given in Figure 2. In that figure, we also give the GEMM_PM and GEMM_PP based algorithms. Although all three perform the same number of floating point operations, the final performance that is achieved can be very different depending on the matrix shapes and cache subsystems.

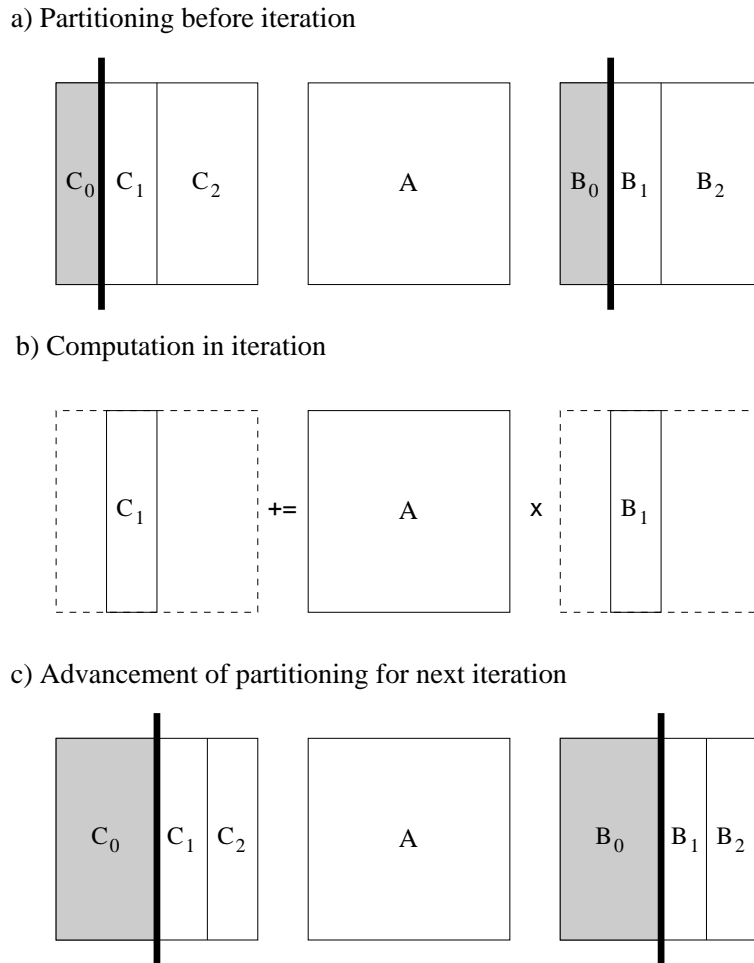


Figure 1: A visualization of the algorithm for matrix-panel variant of matrix-matrix product. Dark background means block is already processed.

3 Accelerating the CUBLAS

The level-3 BLAS operations are variations of the matrix-matrix product. We will study several: symmetric matrix-matrix multiplication (SYMM), symmetric rank-k update (SYRK), symmetric rank-2k update, (SYR2K), triangular matrix-matrix multiplication (TRMM), and triangular solve with multiple right-hand sides (TRSM).

It is well-known that for each operation there are algorithms that cast most computation in terms of matrix-matrix multiplication, as was pioneered in [7]. Moreover, as part of the FLAME project we have long advocated that it is important to have multiple algorithmic variants at our disposal so that the best algorithm can be chosen for each situation [4]. The FLAME methodology advocates systematic derivation of these variants [3, 11]. In Section 4 we will show that this is again the case for GPUs. We view our ability to rapidly develop different algorithms as a way of performing *software acceleration*, the natural (and much needed) counterpart to *hardware acceleration*. It yields a cheap (in terms of effort) boost to performance.

The algorithms presented in this section correspond naturally to the matrix-matrix multiplica-

<p>Algorithm: GEMM_MP(A, B, C)</p> <p>Partition $B \rightarrow (B_L B_R), C \rightarrow (C_L C_R)$ where B_L has 0 columns, C_L has 0 columns while $n(B_L) < n(B)$ do Determine block size b Repartition $(B_L B_R) \rightarrow (B_0 B_1 B_2),$ $(C_L C_R) \rightarrow (C_0 C_1 C_2)$ where B_1 has b columns, C_1 has b columns</p> <hr/> $C_1 := C_1 + AB_1$ <hr/> <p>Continue with $(B_L B_R) \leftarrow (B_0 B_1 B_2),$ $(C_L C_R) \leftarrow (C_0 C_1 C_2)$ endwhile</p>

<p>Algorithm: GEMM_PM(A, B, C)</p> <p>Partition $A \rightarrow \begin{pmatrix} A_T \\ A_B \end{pmatrix}, C \rightarrow \begin{pmatrix} C_T \\ C_B \end{pmatrix}$ where A_T has 0 rows, C_T has 0 rows while $m(A_T) < m(A)$ do Determine block size b Repartition $\begin{pmatrix} A_T \\ A_B \end{pmatrix} \rightarrow \begin{pmatrix} A_0 \\ A_1 \\ A_2 \end{pmatrix}, \begin{pmatrix} C_T \\ C_B \end{pmatrix} \rightarrow \begin{pmatrix} C_0 \\ C_1 \\ C_2 \end{pmatrix}$ where A_1 has b rows, C_1 has b rows</p> <hr/> $C_1 := C_1 + A_1 B$ <hr/> <p>Continue with $\begin{pmatrix} A_T \\ A_B \end{pmatrix} \leftarrow \begin{pmatrix} A_0 \\ A_1 \\ A_2 \end{pmatrix}, \begin{pmatrix} C_T \\ C_B \end{pmatrix} \leftarrow \begin{pmatrix} C_0 \\ C_1 \\ C_2 \end{pmatrix}$ endwhile</p>
--

<p>Algorithm: GEMM_PP(A, B, C)</p> <p>Partition $A \rightarrow (A_L A_R),$ $B \rightarrow \begin{pmatrix} B_T \\ B_B \end{pmatrix}$ where A_L has 0 columns, B_T has 0 rows while $n(A_L) < n(A)$ do Determine block size b Repartition $(A_L A_R) \rightarrow (A_0 A_1 A_2), \begin{pmatrix} B_T \\ B_B \end{pmatrix} \rightarrow \begin{pmatrix} B_0 \\ B_1 \\ B_2 \end{pmatrix}$ where A_1 has b columns, B_1 has b rows</p> <hr/> $C := C + A_1 B_1$ <hr/> <p>Continue with $(A_L A_R) \leftarrow (A_0 A_1 A_2), \begin{pmatrix} B_T \\ B_B \end{pmatrix} \leftarrow \begin{pmatrix} B_0 \\ B_1 \\ B_2 \end{pmatrix}$ endwhile</p>

Figure 2: Algorithms for computing matrix-matrix product: Top-left, matrix-panel variant; top-right, panel-matrix variant; bottom-center, panel-panel variant.

tion algorithms given in Section 2, except that they take advantage of the special structure of one of the matrices. Thus, the ...PP algorithm corresponds to the GEMM_PP algorithm, etc.

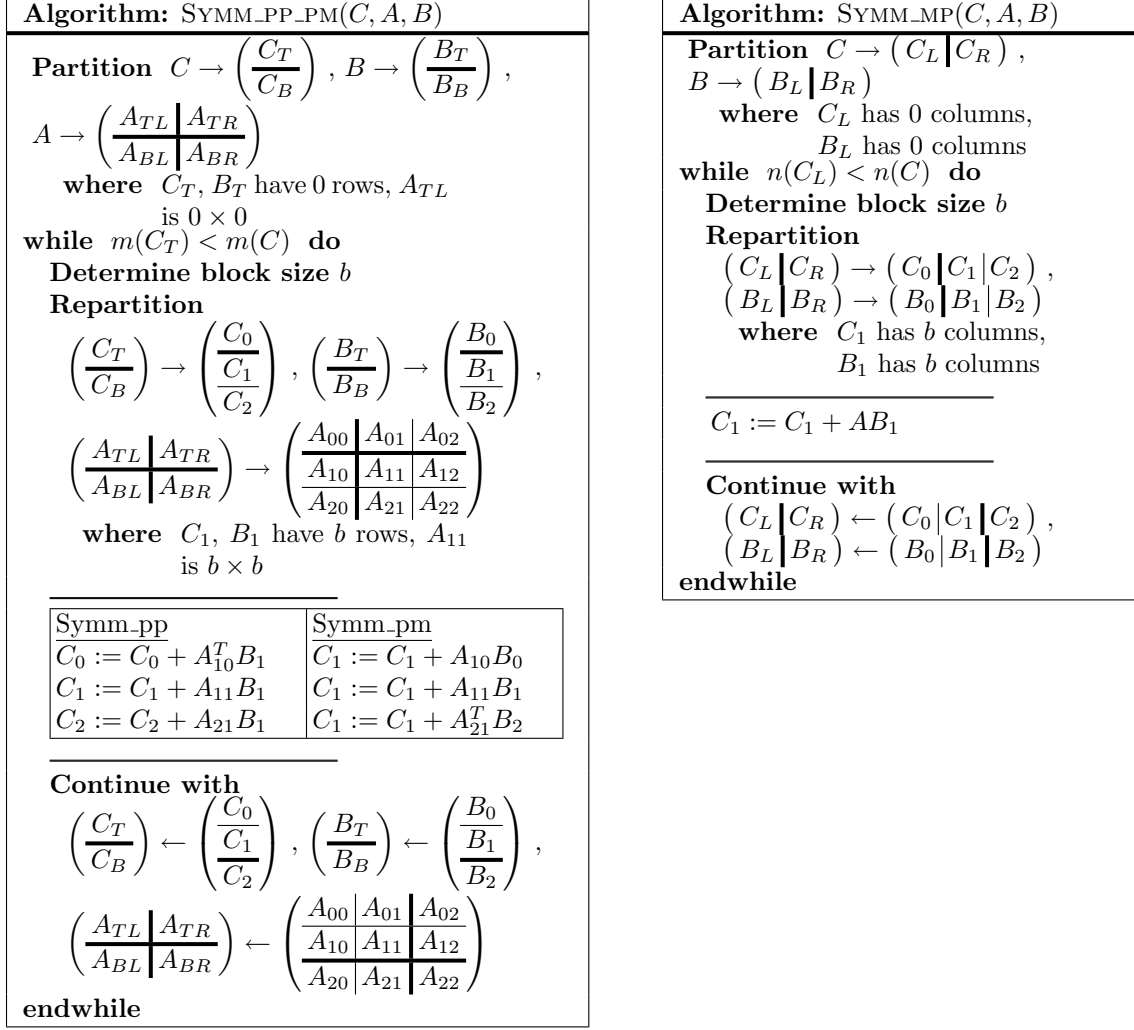


Figure 3: Algorithms for computing SYMM.

SYMM Operation For this operation we focus on $C := AB + C$, where A is symmetric and only the lower triangular part of this matrix is stored. In Figure 3 we give three algorithmic variants.

SYRK Operation We will focus on a representative case of this operation: $C := C - AA^T$, where C is symmetric and only the lower triangular part of this matrix is stored and computed. In Figure 6 we give three algorithmic variants for this operation.

SYR2K Operation For this operation we focus on $C := C - AB^T - BA^T$, where C is symmetric and only the upper triangular part of this matrix is stored and computed. In Figure 7 we give three algorithmic variants.

TRMM Operation For this operation we focus on $C := AB + C$, where A is upper triangular. In Figure 8 we give three algorithmic variants.

TRSM Operation For this operation we focus on $XA^T = B$, where A is lower triangular and B is overwritten with the solution X . In Figure 9 we give three algorithmic variants.

Other Cases of the BLAS-3 Operations The same technique can be applied to the other cases of the BLAS-3 operations above presented. Similarly, the same technique can be applied to other BLAS-3 operations. We expect achieving similar results since the issues are the same.

4 Experimental Results

The current revision of this paper primarily differs from the original one in that it compares performance against CuBLAS 2.2, which was released immediately after our first version of this note was published. And this is where Einstein’s quote becomes relevant: we repeated the same experience, and expected different results... In addition, we have developed and evaluated other BLAS-3 operations, and we have evaluated our methods on rectangular matrices, as they are often used in LAPACK.

The target platform used in the experiments was a NVIDIA T10 GPU (a single GPU of a four GPU NVIDIA Tesla S1070) with 4 GBytes of RAM. The system is connected to a workstation with one Intel Xeon QuadCore E5405 processors (4 cores) at 2.83 GHz with 8 GBytes of DDR2 RAM. CUBLAS Release 2.2 and single precision real floating-point arithmetic were employed in the experiments. Performance is measured in terms of GFLOPS (billions of floating-point operations—flops—per second). The time to transfer data from the host to the memory of the GPU has not been included in the performance results.

Figure 4 reports the performances for the operations reported in the previous section on square matrices. Left plots show the performance in GFLOPS of both the new methods and the corresponding CUBLAS routines; right plots summarize the speedups obtained by the new operations against the corresponding routines in CUBLAS. All new codes achieve about 300 or more GFLOPS on square matrices.

Figure 5 reports the performances obtained by some operations on rectangular matrices, as they are often used inside LAPACK codes. Again, left plot shows the raw performance in GFLOPS of both new methods and CUBLAS routines; right plot summarizes the speedups obtained by the new operations against the corresponding routines in CUBLAS.

The results in both figures show the benefits of our approach. We believe them to be representative of other cases of the presented level-3 BLAS (those where matrices may have been transposed and/or stored in the other triangular part of the array) and the other level-3 BLAS.

The improvement of performances could have been even larger if we had used storage-by-blocks, a well-known modification used in more recent software. We did not employ it to keep full compatibility with NVIDIA CUBLAS.

5 Conclusion

We have demonstrated that with relatively little effort considerable performance gains can be attained when new architectures arrive. The key is to pay attention to the fact that there are many different algorithmic variants for the same operation and to program them in a productive manner. The programs we wrote for this paper required a few hours of time and could have been developed by a relative novice.

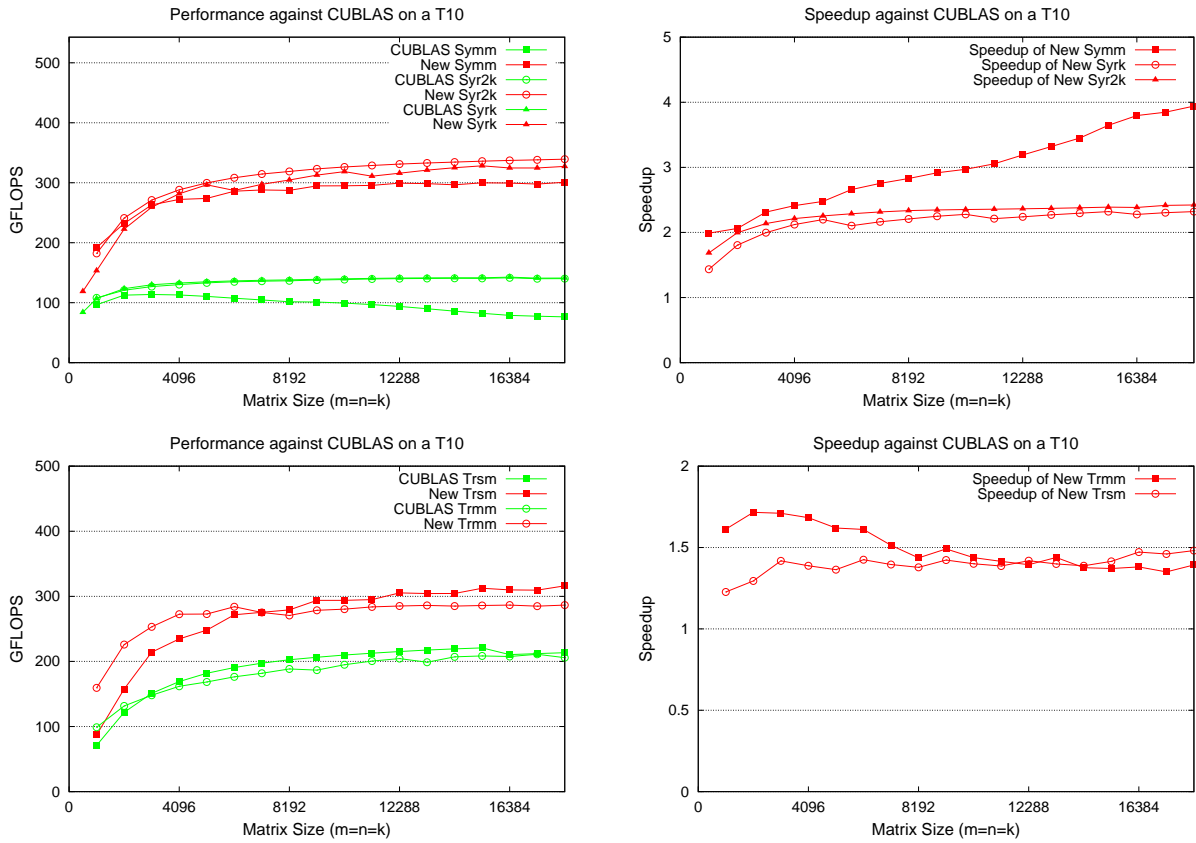


Figure 4: Performances (left) and speedups (right) of the new implementations and equivalent CUBLAS routines on square matrices.

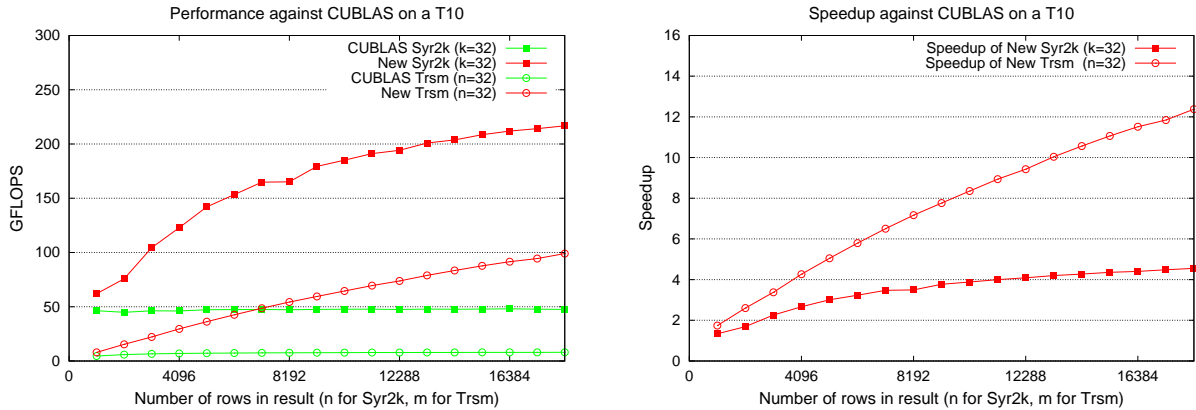


Figure 5: Performances (left) and speedups (right) of the new implementations and equivalent CUBLAS routines on rectangular matrices.

Undoubtedly, in response to this paper, there will be a flurry of activity to further improve the performance of the CUBLAS by coding at a much lower level and throwing programmability out the door. In this case, we have indirectly made a contribution to the scientific computing

community because faster libraries will then become available sooner. But we are confident that this just means that we will be able to write yet another paper on how to improve the performance of high level routines, with functionality similar to that of LAPACK [1]. And before you know it, a new shift in computer architecture will come along and the mad dash will start all over again. Thus the quote from Einstein.

We are working on an tool, FLAMES2S [12], that can automatically translate algorithms represented in code with the FLAME/C API, used to implement our `libflame` library [14], to low-level code that uses loops and indexing. This tool could easily generate the code that was created manually for the experiments in this paper. With that, we will make further progress towards overcoming the programmability problem for this class of operations and codes.

Recently Bientinesi developed a symbolic system for automatically generating families of algorithms from a high-level description of the target operation [2]. Central to this result is a methodology for deriving an algorithm from a given loop-invariant. Combining the automatic system with FLAMES2S would make it possible to derive algorithms and code from loop-invariants, simplifying dramatically the developments of programs.

Acknowledgements

The researchers at the Universidad Jaime I were supported by projects CICYT TIN2008-06570-C04-01, and P1 1B2007-19 of the *Fundación Caixa-Castellón/Bancaixa* and UJI.

We thank the other members of the FLAME team for their support.

References

- [1] E. Anderson, Z. Bai, C. Bischof, L. S. Blackford, J. Demmel, Jack J. Dongarra, J. Du Croz, S. Hammarling, A. Greenbaum, A. McKenney, and D. Sorensen. *LAPACK Users' guide (third ed.)*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1999.
- [2] Paolo Bientinesi. *Mechanical Derivation and Systematic Analysis of Correct Linear Algebra Algorithms*. PhD thesis, Department of Computer Sciences, The University of Texas at Austin, August 2006.
- [3] Paolo Bientinesi, John A. Gunnels, Margaret E. Myers, Enrique S. Quintana-Ortí, and Robert A. van de Geijn. The science of deriving dense linear algebra algorithms. *ACM Trans. Math. Soft.*, 31(1):1–26, March 2005.
- [4] Paolo Bientinesi, Brian Gunter, and Robert A. Van de Geijn. Families of algorithms related to the inversion of a symmetric positive definite matrix. *ACM Trans. Math. Soft.*, 35(1).
- [5] 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.
- [6] Kazushige Goto and Robert A. van de Geijn. Anatomy of a high-performance matrix multiplication. *ACM Trans. Math. Soft.*, 34(3):12, May 2008.
- [7] B. Kågström, P. Ling, and C. Van Loan. GEMM-based level 3 BLAS: High performance model implementations and performance evaluation benchmark. *ACM Trans. Math. Soft.*, 24(3):268–302, 1998.

- [8] Mercedes Marqués, Gregorio Quintana-Ortí, Enrique S. Quintana-Ortí, and Robert van de Geijn. Solving “large” dense matrix problems on multi-core processors and gpus. In *10th IEEE International Workshop on Parallel and Distributed Scientific and Engineering Computing - PDSEC’09. Roma (Italy)*, 2009. to appear.
- [9] Mercedes Marqués, Gregorio Quintana-Ortí, Enrique S. Quintana-Ortí, and Robert van de Geijn. Using graphics processors to accelerate the solution of out-of-core linear systems. In *8th IEEE International Symposium on Parallel and Distributed Computing, Lisbon (Portugal)*, 2009. to appear.
- [10] Gregorio Quintana-Ortí, Francisco D. Igual, Enrique S. Quintana-Ortí, and Robert van de Geijn. Solving dense linear algebra problems on platforms with multiple hardware accelerators. In *ACM SIGPLAN 2009 symposium on Principles and practices of parallel programming (PPoPP’09)*, 2009.
- [11] Robert A. van de Geijn and Enrique S. Quintana-Ortí. *The Science of Programming Matrix Computations*. www.lulu.com/contents/contents/1911788/, 2008.
- [12] Richard M. Veras, Jonathan S. Monette, Enrique S. Quintana-Ort, and Robert A. van de Geijn. Transforming linear algebra libraries: From abstraction to high performance. *ACM Trans. Math. Soft.* submitted.
- [13] Vasily Volkov and James Demmel. LU, QR and Cholesky factorizations using vector capabilities of GPUs. Technical Report UCB/EECS-2008-49, EECS Department, University of California, Berkeley, May 2008.
- [14] Field G. Van Zee. *libflame: The Complete Reference*. www.lulu.com, 2009.

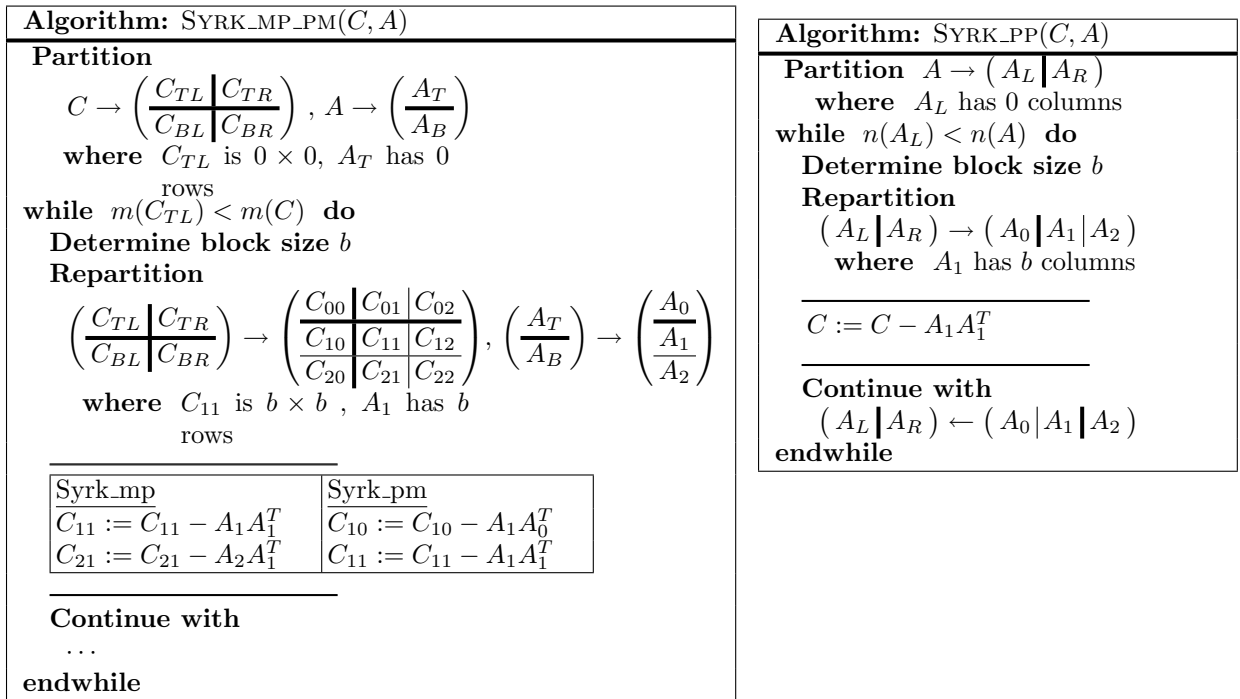


Figure 6: Algorithms for computing SYRK.

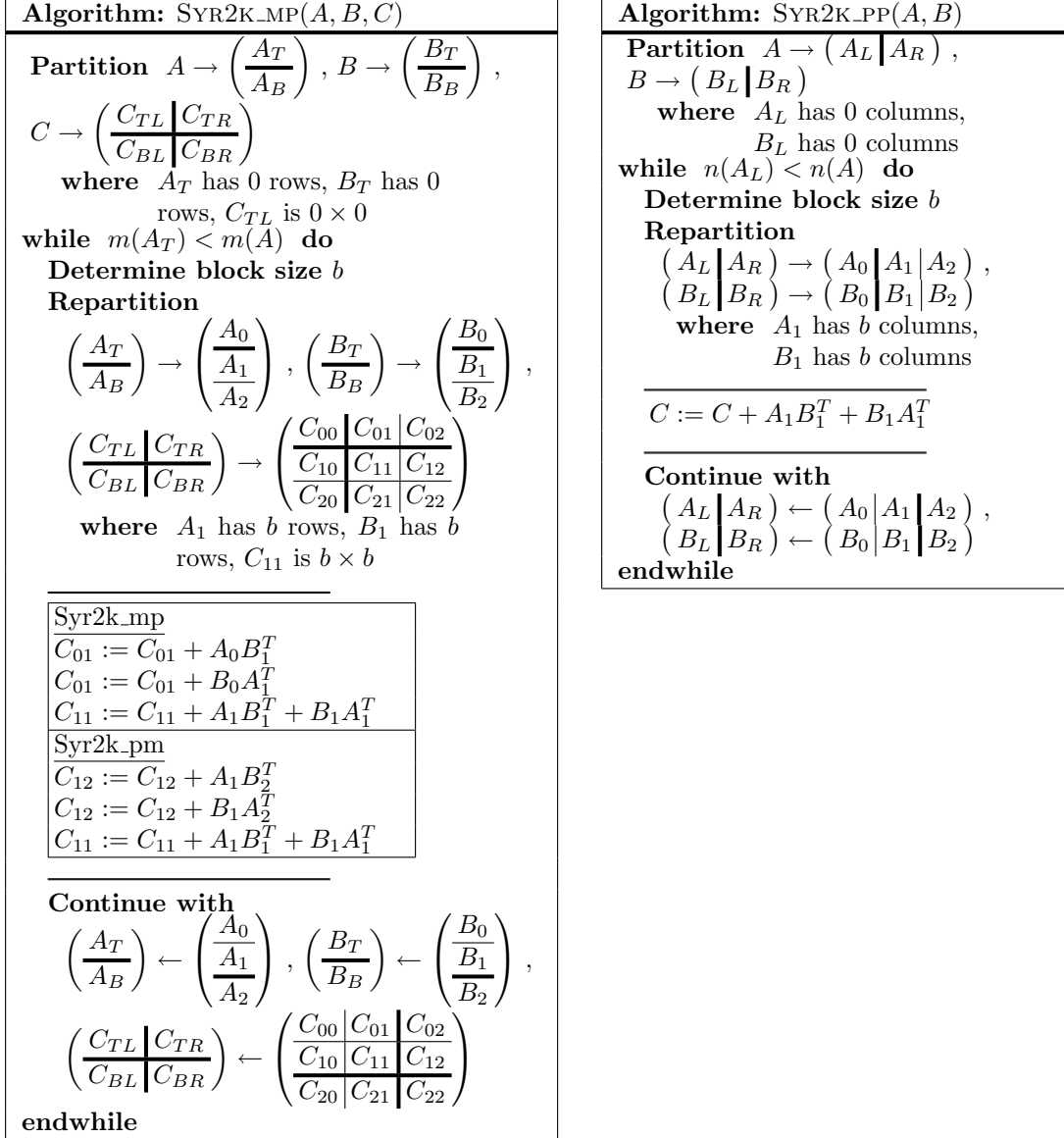


Figure 7: Algorithms for computing SYR2K.

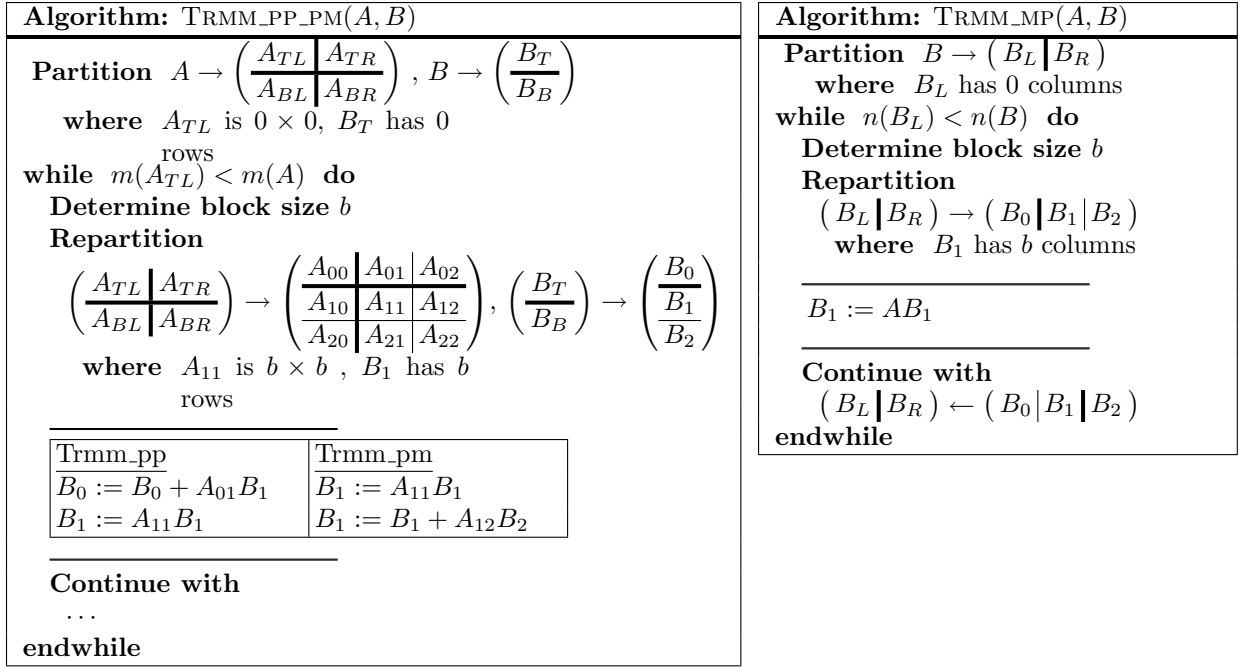


Figure 8: Algorithms for computing TRMM.

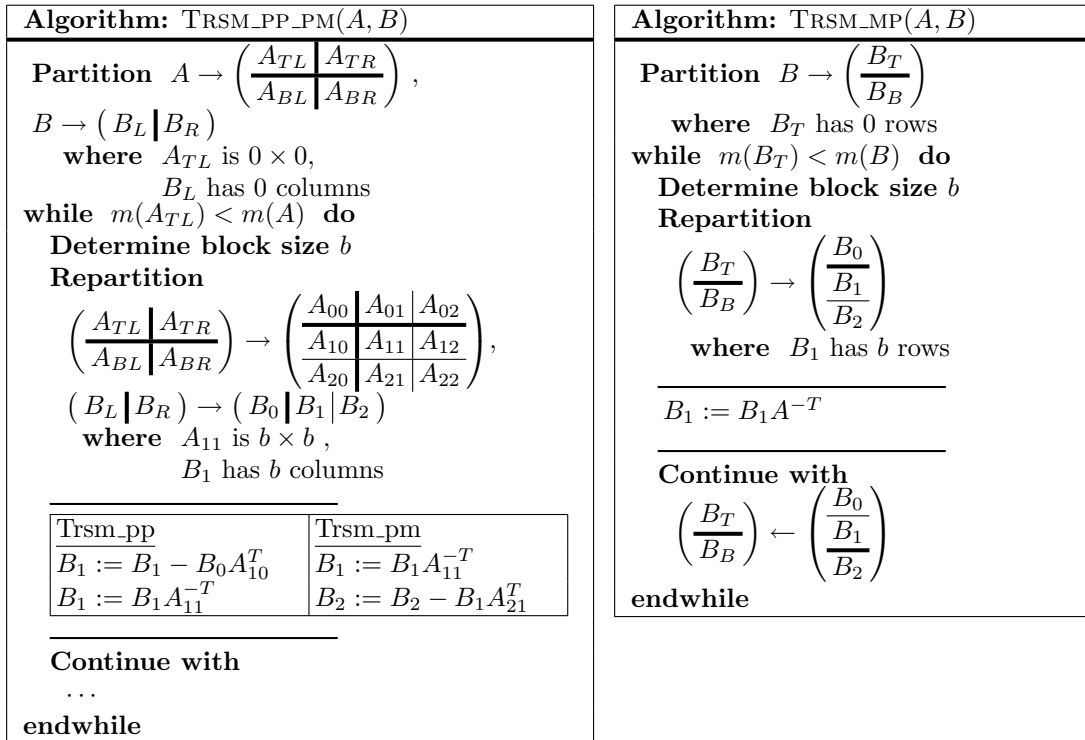


Figure 9: Algorithms for computing TRSM.