

Using Hybrid CPU-GPU Platforms to Accelerate the Computation of the Matrix Sign Function

Peter Benner¹, Enrique S. Quintana-Ortí², and Alfredo Remón²

¹ Fakultät für Mathematik, Chemnitz University of Technology, D-09107 Chemnitz, Germany; benner@mathematik.tu-chemnitz.de

² Depto. de Ingeniería y Ciencia de Computadores, Universidad Jaume I, 12.071–Castellón, Spain; {quintana,remon}@icc.uji.es

Abstract. We investigate the performance of two approaches for matrix inversion based on Gaussian (LU factorization) and Gauss-Jordan eliminations. The target architecture is a current general-purpose multi-core processor connected to a graphics processor (GPU). Parallelism is extracted in both processors by linking sequential versions of the codes with multi-threaded implementations of BLAS. Our results on a system with two Intel QuadCore processors and a Tesla C1060 GPU illustrate the performance and scalability attained by the codes on this system.

Key words: Matrix sign function, hybrid platforms, GPUs, multi-core processors, linear algebra, high performance computing.

1 Introduction

Consider a matrix $A \in \mathbb{R}^{n \times n}$ with no eigenvalues on the imaginary axis, and let

$$A = T^{-1} \begin{pmatrix} J_- & 0 \\ 0 & J_+ \end{pmatrix} T, \quad (1)$$

be its Jordan decomposition, where the eigenvalues of $J_- \in \mathbb{R}^{j \times j} / J_+ \in \mathbb{R}^{(n-j) \times (n-j)}$ all have negative/positive real parts [12]. The *matrix sign function* of A is then defined as

$$\text{sign}(A) = T^{-1} \begin{pmatrix} -I_j & 0 \\ 0 & I_{n-j} \end{pmatrix} T, \quad (2)$$

where I denotes the identity matrix of the order indicated by the subscript. The matrix sign function is a useful numerical tool for the solution of control theory problems (model reduction, optimal control) [19], the bottleneck computation in many lattice quantum chromodynamics computations [10], and dense linear algebra computations (block diagonalization, eigenspectrum separation) [12, 7]. Large-scale problems as those arising, e.g., in control theory often involve matrices of dimension $n \rightarrow O(10,000 - 100,000)$ [14].

There are simple iterative schemes for the computation of the sign function. Among these, the Newton’s iteration, given by

$$\begin{aligned} A_0 &:= A, \\ A_{k+1} &:= \frac{1}{2}(A_k + A_k^{-1}), \quad k = 0, 1, 2, \dots, \end{aligned} \tag{3}$$

is specially appealing for its simplicity, efficiency, parallel performance, and asymptotic quadratic convergence [7, 3]. However, even if A is sparse, $\{A_k\}_{k=1,2,\dots}$ in general are full dense matrices and, thus, the scheme in (3) roughly requires $2n^3$ floating-point arithmetic operations (flops) per iteration.

In the past, large-scale problems have been tackled using message-passing parallel solvers based on the matrix sign function which were then executed on clusters with a moderate number of nodes/processors [2]. The result of this effort was our message-passing library PLiC [4] and subsequent libraries for model reduction (PLiCMR, see [5]) and optimal control (PLiCOC, see [18]). Using this library, 16–32 processors showed to be enough to solve problems with $n \approx 10,000$ in a few hours.

Following the recent uprise of hardware accelerators, like the graphics processors (GPUs), and the increase in the number of cores of current general-purpose processors, in this paper we evaluate an alternative approach that employs a sequential version of the codes in the PLiC library, and extracts all parallelism from tuned multi-threaded implementations of the BLAS (*Basic Linear Algebra Sub-programs*) [16, 9, 8]. The results attained in a hybrid, heterogeneous architecture composed of a general-purpose multi-core processor and a GPU demonstrate that this is a valid platform to deal with large-scale problems which, only a few years ago, would have required a distributed-memory clusters.

The rest of the paper is structured as follows. In Section 2 we elaborate on the hybrid computation of the matrix inverse on a CPU-GPU platform. This is followed by experimental results in Section 3, while concluding remarks and open questions follow in Section 4.

2 High-Performance Matrix Inversion

As equation (3) reveals, the application of Newton’s method to the sign function requires, at each iteration, the computation of a matrix inverse. We next review two different methods for the computation of this operation, based on the LU factorization and Gauss-Jordan transformations.

2.1 Matrix inversion via the LU factorization

The traditional approach to compute the inverse of a matrix $A \in \mathbb{R}^{n \times n}$ is based on Gaussian elimination (i.e., the LU factorization), and consists of the following three steps:

1. Compute the LU factorization $PA = LU$, where $P, L, U \in \mathbb{R}^{n \times n}$, P is a permutation matrix, and L and U are, respectively, unit lower and upper triangular factors [12].

2. Invert the triangular factor $U \rightarrow U^{-1}$.
3. Solve the system $XL = U^{-1}$ for X .
4. Undo the permutations $A^{-1} := XP$.

LAPACK [1] is a high-performance linear algebra library which provides routines that cover the functionality required in the previous steps. In particular, routine `getrf` yields the LU factorization (with partial pivoting) of a nonsingular matrix (Step 1), while routine `getri` computes the inverse matrix of A using the LU factorization obtained by `getrf` (Steps 2–4).

The computational cost of computing a matrix inverse following the previous four steps is $2n^3$ flops. The algorithm sweeps through the matrix four times (one per step) and presents a mild load imbalance, due to the work with the triangular factors.

2.2 Matrix inversion via Gauss-Jordan elimination

The Gauss-Jordan elimination algorithm [11] (GJE) for matrix inversion is, in essence, a reordering of the computation performed by matrix inversion methods based on Gaussian elimination, and hence requires the same arithmetic cost.

Figure 1 illustrates a blocked version of the GJE procedure for matrix inversion using the FLAME notation. There $m(A)$ stands for the number of rows of matrix A . We believe the rest of the notation to be intuitive; for further details, see [13, 6]. (A description of the unblocked version, called from inside the blocked one, can be found in [20]; for simplicity, we hide the application of pivoting during the factorization, but details can be found there as well.) The bulk of the computations in the procedure can be cast in terms of the matrix-matrix product, an operation with a high parallelism. Therefore, GJE is a highly appealing method for matrix inversion on emerging architectures like GPUs, where many computational units are available, and the matrix-matrix product is highly tuned.

We next introduce three implementations for the GJE method (with partial pivoting) on two parallel architectures: a multi-core CPU architecture and a GPU from NVIDIA. The following variants differ on what part of the computation is performed on the CPU (the general-purpose processor or host), and which part is off-loaded to the hardware accelerator (the GPU or device). They all try to reduce the number of communications between the memory spaces of the host and the device.

Implementation on a multi-core CPU: GJE(CPU). In this first variant all operations are performed on the CPU. Parallelism is obtained from a multi-threaded implementation of BLAS for general-purpose processors. Since most of the computations are cast in terms of products of matrices, high performance can be expected from this variant.

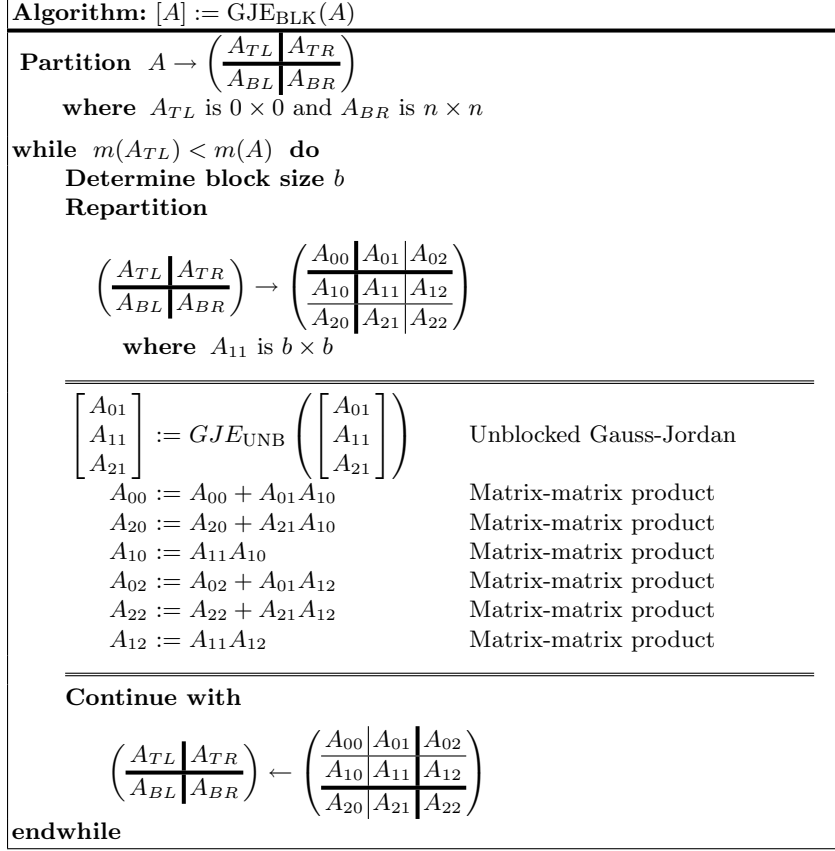


Fig. 1. Blocked algorithm for matrix inversion via GJE without pivoting.

Implementation on a many-core GPU: $\text{GJE}(\text{GPU})$. This is the GPU-analogue to the previous variant. The matrix is first transferred to the device; all computations proceed there next; and the result (the matrix inverse) is finally moved back to the host.

Hybrid implementation: $\text{GJE}(\text{Hybrid})$. While most of the operations performed in the GJE algorithm are well suited for the GPU, a few are not. This is the case for fine-grained operations, where the low computational cost and data dependencies deliver low performance on massively parallel architectures like the GPU. To solve this problem, we propose a hybrid implementation. In this new approach, operations are performed in the most convenient device, exploiting the capabilities of both architectures.

In particular, in this variant the matrix is initially transferred to the device. At the beginning of each iteration, the current column block, composed

of $[A_{01}^T, A_{11}^T, A_{21}^T]^T$ is moved to the CPU and factorized there. The result is immediately transferred back to the device, where all remaining computations (matrix-matrix products) are performed. This pattern is repeated until the full matrix inverse is computed. The inverse is finally transferred from the device memory to the host.

3 Experimental results

In this section we evaluate four parallel multi-threaded codes to compute the inverse of a matrix:

- LAPACK(CPU): The four steps of the LAPACK approach, with all computations carried out on the CPU and parallelism extracted by using a multi-threaded implementation of BLAS; see subsection 2.1)
- GJE(CPU), GJE(GPU), and GJE(Hybrid): The implementations described in subsection 2.2.

Two different implementations of the BLAS (Goto BLAS [21] (version 1.26) and Intel MKL [15] (version 10.1)) were used to execute operations on the general-purpose processor, while on the NVIDIA GPU, CUBLAS [17] (version 2.1) was the library that we used.

All experiments employed single precision and the results always include the cost of data transfers between the host and device memory spaces. The target platform consists of two Intel Xeon QuadCore processors connected to a Tesla C1060 GPU. Table 3 offers more details on the hardware.

Processors	#cores	Frequency (GHz)	L2 cache (MB)	Memory (GB)
Intel Xeon QuadCore E5405	8	2.33	12	8
Nvidia TESLA c1060	240	1.3		4

Table 1. Hardware employed in the experiments.

Figure 2 reports the GFLOPS (10^9 flops per second) rates attained by the different implementations of the inversion codes operating on matrices with sizes between 1,000 and 8,000. Several algorithmic block sizes (parameter b in Figure 1) were tested but, for simplicity, hereafter the results in the Figure correspond to those obtained with the optimal block size.

The LAPACK code executed using all 8 cores of the two general-purpose processors yields the lowest GFLOPS rate, while the GJE algorithm using the same resources performs slightly better. Both implementations that employ the GPU outperform the ones executed only on the CPU. The Hybrid approach is the best option for small/medium matrices, while the version executed entirely on the GPU is the best for large matrices.

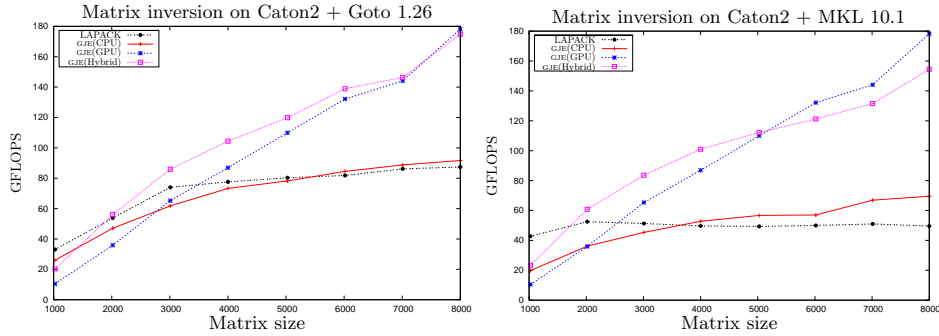


Fig. 2. Performance of the matrix inversion codes.

Figure 3 shows execution times of the Newton's iteration for the matrix sign function, using the previous matrix inversion codes. As expected, the LAPACK implementation requires the highest execution time, followed by GJE(CPU). The codes that employ the GPU are notoriously faster, specially for large matrices.

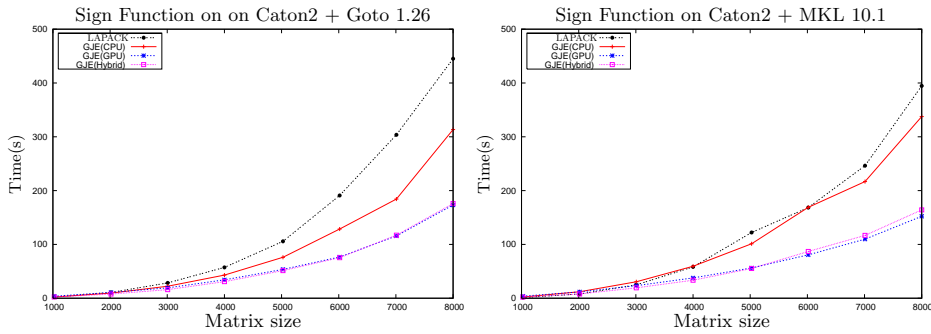


Fig. 3. Execution times of the Newton's iteration for the matrix sign function with the matrix inversion implemented using the different variants discussed.

4 Concluding remarks and Future Work

We have demonstrated the benefits of using a current GPU to off-load part of the computations in a dense linear algebra operation rich in level-3 BLAS like the matrix inversion. This operation is the basis for the computation of the matrix sign function via Newton's iteration and is also the key to the efficient solution of important problems in control theory such as model reduction or optimal control.

The evaluation of matrix inversion codes clearly identify the superior performance of the procedures based on Gauss-Jordan elimination over Gaussian elimination (the LU factorization).

Our research poses some open questions which form the basis of our ongoing and future work:

- Most applications in control theory and linear algebra require double precision but current GPUs deliver considerable lower performance when they operate with this data type. Is it possible to compute the sign function in single precision and then refine this approximation to double precision at a low cost?
- The GPU is controlled from a single thread running on a single core of the general-purpose processor(s). The calls to CUBLAS block this thread until this computation has been performed on the GPU. Can we split the computations so that a part of them is performed on the remaining general-purpose cores? Note that this requires a careful synchronization of the data transfers between the memory spaces of host and device.
- Is it possible to overlap computation on the GPU with data transfers between the CPU and the GPU memory spaces to improve the performance for the small problem sizes?

Acknowledgments

This work was supported by the Spanish Ministry of Science and Innovation/FEDER (contracts no. TIN2005-09037-C02-02 and TIN2008-06570-C04-01) and by the Fundación Caixa-Castelló/Bancaixa (contracts no. P1B-2007-19 and P1B-2007-32).

References

1. E. Anderson, Z. Bai, C. Bischof, J. Demmel, J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen. *LAPACK Users' Guide*. SIAM, Philadelphia, PA, third edition, 1999.
2. P. Benner, J.M. Claver, and E.S. Quintana-Ortí. Parallel distributed solvers for large stable generalized Lyapunov equations. *Parallel Processing Letters*, 9(1):147–158, 1999.
3. P. Benner and E.S. Quintana-Ortí. Solving stable generalized Lyapunov equations with the matrix sign function. *Numer. Algorithms*, 20(1):75–100, 1999.
4. P. Benner, E.S. Quintana-Ortí, and G. Quintana-Ortí. A portable subroutine library for solving linear control problems on distributed memory computers. In G. Cooperman, E. Jessen, and G.O. Michler, editors, *Workshop on Wide Area Networks and High Performance Computing, Essen (Germany), September 1998*, Lecture Notes in Control and Information, pages 61–88. Springer-Verlag, Berlin/Heidelberg, Germany, 1999.
5. P. Benner, E.S. Quintana-Ortí, and G. Quintana-Ortí. State-space truncation methods for parallel model reduction of large-scale systems. *Parallel Comput.*, 29:1701–1722, 2003.

6. 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 Transactions on Mathematical Software*, 31(1):1–26, March 2005.
7. R. Byers. Solving the algebraic Riccati equation with the matrix sign function. *Linear Algebra Appl.*, 85:267–279, 1987.
8. 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.
9. Jack J. Dongarra, Jeremy Du Croz, Sven Hammarling, and Richard J. Hanson. An extended set of FORTRAN basic linear algebra subprograms. *ACM Trans. Math. Soft.*, 14(1):1–17, March 1988.
10. A. Frommer, T. Lippert, B. Medeke, and K. Schilling, editors. *Numerical Challenges in Lattice Quantum Chromodynamics*, volume 15 of *Lecture Notes in Computational Science and Engineering*. Springer-Verlag, Berlin/Heidelberg, 2000.
11. Alexandros V. Gerbessiotis and Wolfson Building. Algorithmic and Practical Considerations for Dense Matrix Computations on the BSP Model, 1997.
12. Gene H. Golub and Charles F. Van Loan. *Matrix Computations*. The Johns Hopkins University Press, Baltimore, 3rd edition, 1996.
13. John A. Gunnels, Fred G. Gustavson, Greg M. Henry, and Robert A. van de Geijn. FLAME: Formal linear algebra methods environment. *ACM Transactions on Mathematical Software*, 27(4):422–455, December 2001.
14. IMTEK, <http://www.imtek.de/simulation/benchmark/>. *Oberwolfach model reduction benchmark collection*.
15. Intel Corporation., <http://www.intel.com/>.
16. C. L. Lawson, R. J. Hanson, D. R. Kincaid, and F. T. Krogh. Basic linear algebra subprograms for Fortran usage. *ACM Trans. Math. Soft.*, 5(3):308–323, Sept. 1979.
17. Nvidia Corporation, <http://www.nvidia.com/cuda/>.
18. G. Quintana P. Benner, E. S. Quintana. Solving linear-quadratic optimal control problems on parallel computers. *Optimization Methods & Software*, 23(6):879–909, 2008.
19. P.H. Petkov, N.D. Christov, and M.M. Konstantinov. *Computational Methods for Linear Control Systems*. Prentice-Hall, Hertfordshire, UK, 1991.
20. E.S. Quintana-Ortí, G. Quintana-Ortí, X. Sun, and R.A. van de Geijn. A note on parallel matrix inversion. *SIAM J. Sci. Comput.*, 22:1762–1771, 2001.
21. Texas Advanced Computing Center, <http://www.tacc.utexas.edu/~kgoto/>.