



HPC Programming on Intel Many-Integrated-Core Hardware with MAGMA Port to Xeon Phi

DOI:

[10.1155/2015/502593](https://doi.org/10.1155/2015/502593)

Document Version

Final published version

[Link to publication record in Manchester Research Explorer](#)

Citation for published version (APA):

Dongarra, J., Gates, M., Haidar, A., Jia, Y., Kabir, K., Luszczek, P., & Tomov, S. (2015). HPC Programming on Intel Many-Integrated-Core Hardware with MAGMA Port to Xeon Phi. *Scientific Programming*, 2015, [502593]. <https://doi.org/10.1155/2015/502593>

Published in:

Scientific Programming

Citing this paper

Please note that where the full-text provided on Manchester Research Explorer is the Author Accepted Manuscript or Proof version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version.

General rights

Copyright and moral rights for the publications made accessible in the Research Explorer are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

Takedown policy

If you believe that this document breaches copyright please refer to the University of Manchester's Takedown Procedures [<http://man.ac.uk/04Y6Bo>] or contact uml.scholarlycommunications@manchester.ac.uk providing relevant details, so we can investigate your claim.



Research Article

HPC Programming on Intel Many-Integrated-Core Hardware with MAGMA Port to Xeon Phi

Jack Dongarra,^{1,2,3} Mark Gates,¹ Azzam Haidar,¹ Yulu Jia,¹ Khairul Kabir,¹ Piotr Luszczek,¹ and Stanimire Tomov¹

¹University of Tennessee, Knoxville, TN 37996, USA

²Oak Ridge National Laboratory, Oak Ridge, TN 37831, USA

³University of Manchester, Manchester M13 9PL, UK

Correspondence should be addressed to Piotr Luszczek; luszczek@eecs.utk.edu

Received 1 May 2014; Accepted 22 November 2014

Academic Editor: Roman Wyrzykowski

Copyright © 2015 Jack Dongarra et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

This paper presents the design and implementation of several fundamental dense linear algebra (DLA) algorithms for multicore with Intel Xeon Phi coprocessors. In particular, we consider algorithms for solving linear systems. Further, we give an overview of the MAGMA MIC library, an open source, high performance library, that incorporates the developments presented here and, more broadly, provides the DLA functionality equivalent to that of the popular LAPACK library while targeting heterogeneous architectures that feature a mix of multicore CPUs and coprocessors. The LAPACK-compliance simplifies the use of the MAGMA MIC library in applications, while providing them with portably performant DLA. High performance is obtained through the use of the high-performance BLAS, hardware-specific tuning, and a hybridization methodology whereby we split the algorithm into computational tasks of various granularities. Execution of those tasks is properly scheduled over the heterogeneous hardware by minimizing data movements and mapping algorithmic requirements to the architectural strengths of the various heterogeneous hardware components. Our methodology and programming techniques are incorporated into the MAGMA MIC API, which abstracts the application developer from the specifics of the Xeon Phi architecture and is therefore applicable to algorithms beyond the scope of DLA.

1. Introduction and Background

Solving linear systems of equations and eigenvalue problems is fundamental to scientific computing. The popular LAPACK library [1], and in particular its vendor optimized implementations such as Intel's MKL [2] or AMD's ACML [3], has been the software of choice to provide solver routines for dense matrices on shared memory systems. This paper considers a redesign of the LAPACK algorithms and their implementation to add efficient support for heterogeneous systems of multicore processors with Intel Xeon Phi coprocessors. This is not the first time that DLA libraries have needed a redesign to be efficient on new architectures, notable examples being the transition from LINPACK [4] to LAPACK [1] in the 1980s to make algorithms cache-friendly. Also, ScaLAPACK [5] in the 1990s added support for distributed memory systems. And at present time, the PLASMA and MAGMA

libraries [6] target efficiency on, respectively, multicore and heterogeneous architectures.

The Intel Xeon Phi coprocessor is a hardware accelerator that made its debut in the late 2012 as a platform for high-throughput technical computing. It is sometimes known under an alternative name of Many Integrated Cores (MICs). For the purposes of this paper, the common mode of operation for the device is called off-load. However, the stand-alone and reverse off-load modes are also valid possibilities. When in off-load mode, the device receives work from the host processor and reports back as soon as the computational task completes. Any such assignment of work proceeds and completes without the host device being involved. In a typical scenario, the host is an Intel x86 CPU such as Sandy Bridge, Ivy Bridge, or even more recent Haswell and Ivy Town. The CPU may monitor the activity of communication and/or computation through an event-based interface and can also

pursue its own computational activities between events. This is very similar to the operation of hardware accelerators based on throughput-oriented GPUs and compute-capable FPGAs that are specialized for certain types of workloads beyond what could be achieved on standard multicore CPUs. In fact, Xeon Phi is often considered to be an alternative to the hardware accelerators from AMD and NVIDIA despite the fact that there exist many technical differences between the three.

The development of new high-performance numerical libraries is a complex endeavor, which requires meticulous accounting for the extreme levels of parallelism, heterogeneity, and wide variety of accelerators and coprocessors available in the current architectures. Challenges vary from new algorithmic designs to choices of programming models, languages, and frameworks that ease the development, future maintenance, and portability. This paper addresses these issues while presenting our approach and algorithmic designs in the development of the MAGMA MIC [7] library. Specific differences between the GPU-based MAGMA [6] and the MIC version are elaborated upon in Section 3.

To provide a uniform portability across a variety of coprocessors/accelerators, we developed an API that abstracts the application developer from the low level specifics of the architecture. In particular, we use low level vendor libraries, like SCIF for Intel Xeon Phi (see Section 5), to define API for memory management and off-loading computations to coprocessors and/or accelerators.

To deal with the extreme level of parallelism and heterogeneity in the current architectures, MAGMA MIC uses a hybridization methodology, described in Section 6, where we split the algorithms of interest into computational tasks of various granularities and properly schedule those tasks' execution over the heterogeneous hardware. Thus, we use a Directed Acyclic Graph (DAG) approach to parallelism and scheduling that has been developed and successfully used for dense linear algebra libraries such as PLASMA and MAGMA [6], as well as in general task-based approaches to parallelism, such as runtime systems like StarPU [8] and SMPs [9].

Obtaining high performance depends on a combination of algorithmic and hardware-specific optimizations, discussed in Section 6.4. This is in addition to the use of high-performance low-level libraries, which we address in Section 5. This has implications on the resulting software: in order to maintain the performance portability across hardware, it is necessary to provide in the library a number of algorithmic variations that are tunable, for example, at installation time. This is the basic premise of autotuning—a prominent example of these kinds of advanced optimization techniques.

A performance study is presented in Section 7. Besides verifying our approach and confirming the appeal of the Intel Xeon Phi coprocessors for high-performance DLA, the results open up a number of future work opportunities discussed in Section 8 that concludes the paper.

2. Related Work

Intel Xeon Phi [10, 11] is a family of Intel coprocessors known before under the MICs (Many Integrated Cores) moniker. Knights Corner (KNC) is the first official product

accelerator card in a series that will be followed by Knights Landing (KNL). Phi is a hardware platform based on x86 instruction set with modifications for throughput-oriented workloads. In some sense, Phi may be regarded as an alternative to NVIDIA's compute GPU cards that require CUDA programming [12] or AMD's compute GPU cards that are programmed with OpenCL [13] and the AMD's GPU libraries [14].

Phi's use for scientific applications that require solution to PDEs (Partial Differential Equations) was studied and under some scenarios revealed opportunities and advantages [15, 16].

There is a rich area of work on execution environments that begin with serial code and result in parallel execution, often using task superscalar techniques, for example, Jade [17], Cilk [18], Sequoia [19], OmpSS [20], Habanero [21], StarPU [8], or the DepSpawn [22] project.

3. Differences between GPU and MIC Versions of MAGMA

We mostly focus on the CUDA-based version of MAGMA for the comparison because it is the basis for functional interface and, in terms of the feature set, it is our aim to reproduce it on the Intel MIC coprocessor.

Fundamentally, hardware accelerators require refactoring of the existing code base to accommodate the new compute device and include it harmoniously into the mix with the CPU so that the performance gains may be fully realized. In terms of raw performance across a broad spectrum of applications, the most efficient programming language is CUDA [12]. Our experiments show that it easily outperforms portable standards-based APIs such as OpenCL [23]. While it might be tempting to include CUDA in the family of languages derived from C and C++, it is worth noting that the clear syntactic differences from the base language (mostly C++ and its 1998 standard) form an easily distinguishable delineation of the computational spaces of the CPU and the GPU. At the CPU code level the triple-chevron launch notation, for example, `<<<blocks, threadsPerBlock>>>gpu.kernel(args)`, launches GPU kernels by means of incompatible syntax that requires NVIDIA's own `nvcc` compiler. This divergent syntax has spurred over the years a number of ways to simplify the coding with the use of directive-based code and as of lately, these efforts have coalesced into the OpenACC initiative [24, 25], directive-based approach that hides some of the CUDA complexity behind compiler's pragma syntax.

The directive-based approach is what Intel MIC featured from the beginning and this is what MAGMA's port to the coprocessor used. However, the MIC port of MAGMA accommodated changes in the interfaces, feature set, and performance levels. Thus, the end user was shielded from the effects of the growth of the platform and the flux of the software ecosystem. A particular example of such an underlying change was the early use of SCIF (see Section 5) which was essential for exchanging noncontiguous memory regions between the host and the device with a very low overhead. This has been progressively phased out as the

TABLE 1: Programming models for the Intel Xeon Phi coprocessors and their current status and properties.

Programming model/API	Status	Portability	Overhead	Language support
SCIF	Mature	No	None	No
COI	Mature	Yes	Minimal	Yes
OpenMP 4.0	Early	Yes	Varies	Yes
OpenCL	Experimental	Yes	Minimal	No

need for SCIF diminished with richer functionality available through the directives and improvements in the Linux kernel drivers and runtime overheads. From the user perspective, this change was transparent for programming on Xeon Phi while the recent changes in event-driven APIs of CUDAs had to be percolated to MAGMA’s publicly visible interface.

Another departure from the CUDA-based MAGMA was the device- and software-specific tuning and optimization (described in more detail in Section 6). There is very little commonality between the targeted systems, both in terms of hardware and software. The Xeon Phi implementation has to balance the performance sensitivity of the BLAS calls in MKL, custom kernels, and their mapping onto the much different hardware substrate. Similarly, the levels and layering of parallelism nesting (software threads, hardware threads, versus BLAS threads) are anything but what is presented to the CUDA programmer. Despite the differences, however, MAGMA’s external interface remains almost indistinguishable.

4. Compiler Support for Off-Load

In this paper, we consider the off-load mode as the primary mode of operation for the Xeon Phi coprocessor. The device receives work from the host processor and reports back upon completion of the assignment without the host being involved in between these two events. This is very similar to the operation of network off-load engines, specifically, the TCP off-load engines (TOEs) that feature an optimized implementation of the TCP stack that handles the majority of the network traffic to lessen the burden of the main processor, which handles other operating system and user application tasks.

The off-load mode for the Xeon Phi devices has direct support from the compiler in that it is possible to issue requests to the device and ascertain the completion of tasks directly from the user’s C/C++ code. The support for this mode of operation is offered by the Intel compiler through Phi-specific `pragma` directives: `offload`, `offload_attribute`, `offload_transfer`, and `offload_wait` [10]. This is very closely related to the off-load directives now included in the OpenMP 4 standard. In fact, the two are syntactically and semantically equivalent, barring the difference in the “omp” prefix for the OpenMP syntax. A similar standard for GPUs is called OpenACC. A summary of various programming methods on Xeon Phi is provided in Table 1. From our rudimentary experiments we concluded that the compiler directive overhead is very close to that of the Common Offload Interface (COI) library.

5. Programming Model: Host-Device with a Server Based on LLAPI

For many scientific applications, the off-load model offered by the Intel compiler, described in Section 4, is sufficient. This is not the case for a fully equivalent port of MAGMA to the Xeon Phi because of the very rich functionality that MIC MAGMA inherits from both its CUDA and OpenCL ports. We had to use the LLAPI (low-level API) based on Symmetric Communication InterFace (SCIF) that offers, as the name suggests, a very low level interface to the host and device hardware. The use of this API is discouraged for most workloads as it tends to be error-prone and offers very little abstraction on top of the hardware interfaces. What motivated us to use it for the port of our library was (1) the asynchronous events capability that allows low-latency messaging between the host and the device to notify about completion of kernels on Xeon Phi as well as (2) the possibility of hiding the cost of data transfer between the host and the device which requires the transfer of submatrices to overlap with the computation. The direct access to the DMA (Direct Memory Access) engine allowed us to maximize the bandwidth of data transfers over the PCI Express bus. The only requirement was that the memory regions for transfer be page-aligned and pinned to guarantee their fixed location in the physical memory. Figure 1(a) shows the interaction between the host and the server running on the Xeon Phi and responding to requests that are remote invocations of numerical kernels on data that have already been transferred to the device.

6. Hybridization Methodology and Optimization Strategies

The hybridization methodology used in CUDA MAGMA [26], adopted for MIC MAGMA, is an extension of the task-based approach for parallelism and developing DLA on homogeneous multicore systems [6]. In particular,

- (i) the computation is split into BLAS-based tasks of various granularities, with their data dependencies, as shown in Figure 1(b);
- (ii) small, nonparallelizable tasks with significant control-flow are scheduled on the CPUs;
- (iii) large, parallelizable tasks are scheduled on Xeon Phi.

The difference with multicore algorithms is the task splitting, which here is of various granularities to make different tasks suitable for particular architectures and the scheduling itself.

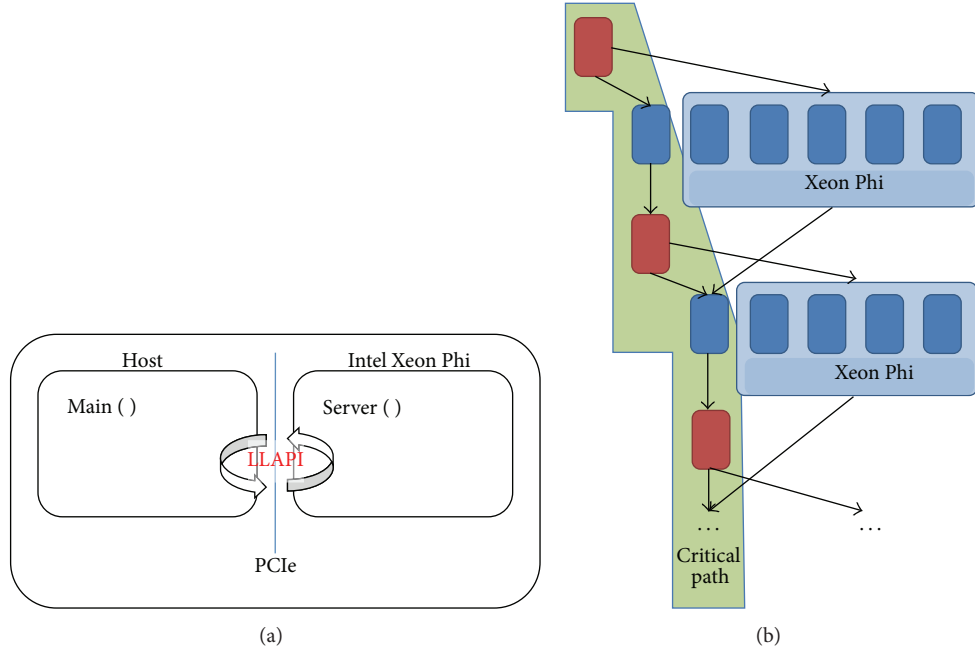


FIGURE 1: (a) MIC MAGMA programming model with a LLAPI server mediating requests between the host CPU and the Xeon Phi device. (b) DLA algorithm as a collection of BLAS-based tasks and their dependencies. The algorithm's critical path is, in general, scheduled on the CPUs and large data-parallel tasks on the Xeon Phi.

Specific algorithms using this methodology, and covering the main classes of DLA, are described in the subsections below.

6.1. Design and Functionality. The MIC MAGMA interface is similar to LAPACK. For example, compare LAPACK's LU factorization interface to MIC MAGMA's:

```
lapackf77_dgetrf(&M, &N, hA, &lda, ipiv,
                &info)
magma_dgetrf_mic(M, N, dA, 0, ldda, ipiv,
                &info, queue)
```

Here, `hA` is the typical CPU pointer (`double *`) to the matrix of interest in the CPU memory and `dA` is a pointer in the Xeon Phi memory (its type is `magmaDouble_ptr`). The last argument in every MIC MAGMA call is Xeon Phi queue, through which the computation will be streamed on the Xeon Phi device (its type is `magma_queue_t`).

To abstract the user away from knowing the low-level directives, library functions (such as BLAS), CPU-Phi data transfers, and memory allocations and deallocations are redefined in terms of MIC MAGMA data types and functions. This design allows us to more easily port the MIC MAGMA library to many devices as was the case for the GPU accelerators that either use CUDA [12] or OpenCL [13, 23] and eventually to merge them in order to maintain a single source code tree with conditional compilation options that allow seamless targeting of specific hardware. Also, the MIC MAGMA wrappers provide a complete set of functions for programming hybrid high-performance numerical libraries. Thus, not only users but also application developers can opt to use the MIC MAGMA wrappers. MIC MAGMA provides

the four standard floating-point arithmetic precisions: single and double precision real as well as single and double precision complex. It has routines for the so-called one-sided factorizations (LU, QR, and Cholesky), and recently we are developing the two-sided factorizations (Hessenberg and bi- and tridiagonal reductions), linear system and least squares solvers, matrix inversions, symmetric and nonsymmetric standard eigenvalue problems, SVD, and orthogonal transformation routines.

6.2. Task Distribution Based on Hardware Capability. Programming models that raise the level of abstraction are of great importance for reducing software development efforts. A traditional approach has been to organize algorithms in terms of BLAS calls, where hardware specific optimizations would be hidden in BLAS implementations such as Intel's MKL or AMD's ACML. This is still valid and used but has shown some drawbacks on new architectures. In particular, parallelization is achieved using a fork-join approach since each BLAS call, for example, a matrix-matrix multiplication, can be performed in parallel (fork) but a synchronization is needed before performing the next call (join). The number of synchronizations thus can become prohibitive bottlenecks for performance on highly parallel devices such as the MICs. This type of programming has been popularized under the Bulk Synchronous Processing name [27].

Instead, the algorithms (like matrix factorizations) are broken into computational tasks (e.g., panel factorizations followed by trailing submatrix updates) and pipelined for execution on the available hardware components (see below). Moreover, particular tasks are scheduled for execution on the hardware components that are best suited for them. Thus,

```

(1) PanelStartReceivingon CPU( $P_1$ );
(2) for  $P_i = P_1, P_2, \dots$  do
(3)   PanelFactorizeon CPU( $P_i$ );
(4)   PanelSendto MIC( $P_i$ );
(5)   TrailingMatrixUpdateon MIC( $P_{i+1}$ );
(6)   PanelStartReceivingon CPU( $P_{i+1}$ );
(7)   TrailingMatrixUpdateon MIC( $P_{i+2}, \dots$ );

```

ALGORITHM 1: Two-phase (first: panel, two: update) factorization of $A = [P_1, P_2, \dots]$ with lookahead of depth 1. Matrix A and the result are assumed to reside on the MIC memory.

this task distribution based on *hardware capability* allows the user for the efficient use of each hardware component. In the case of DLA factorizations, the less parallel panel tasks are scheduled for execution on multicore CPUs and the parallel updates mainly on the MICs. We illustrate this in Algorithm 1.

6.3. *LU, QR, and Cholesky Factorizations for Intel Xeon Phi.* The one-sided factorization routines implemented and currently available through MIC MAGMA are as follows:

`magma_zgetrf_mic` computes an LU factorization of a general M -by- N matrix A using partial pivoting with row interchanges;

`magma_zgeqrf_mic` computes a QR factorization of a general M -by- N matrix A ;

`magma_zpotrf_mic` computes the Cholesky factorization of a complex Hermitian positive definite matrix A .

Routines in all standard four floating-point precision arithmetics are available, following LAPACK's naming convention. Namely, the first letter of the routine name (after the prefix `magma_`) indicates the precision – **z**, **c**, **d**, or **s** for double complex, single complex, double real, or single real, respectively. The suffix `_mic` indicates that the input and the output matrices are in the Xeon Phi memory.

The typical hybrid computation and communication pattern for the one-sided factorizations (LU, QR, and Cholesky) is shown in Figure 2. At a given iteration, panel i is copied to the CPU and factored using LAPACK, and the result is copied back to Xeon Phi. The trailing matrix, consisting of the next panel $i + 1$ and the rest of the matrix, is updated on the Xeon Phi. After receiving panel i back from the CPU, panel $i + 1$ is updated first using panel i and the result is sent to the CPU (as being the next panel to be factored there). While the CPU starts the factorization of i , the rest of trailing matrix, panels $i + 1, i + 2, \dots$, is updated on the Xeon Phi device in parallel with the CPU factorization of panel $i + 1$. In this pattern, only data to the right of the current panel is accessed and modified, and the factorizations that use it are known as right-looking. The computation can be organized differently, to access and modify data only to the left of the panel, in which case the factorizations are known as left-looking.

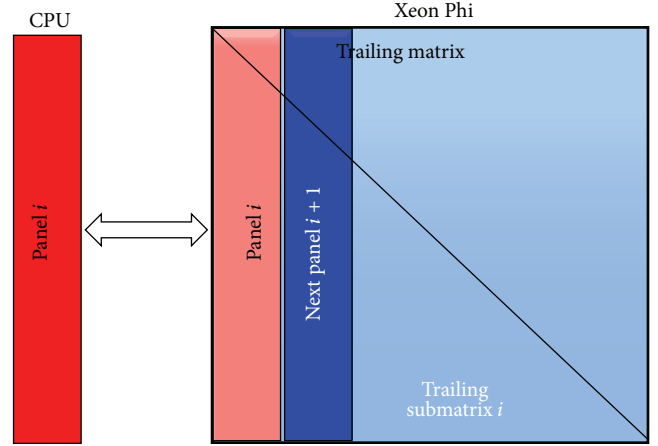


FIGURE 2: Typical computational pattern for the hybrid one-sided factorizations in MIC MAGMA.

An example of a left-looking factorization, demonstrating a hybrid implementation, is given in Algorithm 2 for the Cholesky factorization. The algorithm introduces a notion of a *blocking factor* denoted as nb , which is the algorithm-level entity that defines the number of columns in the panel and the inner dimension of the outer-product update to the trailing submatrix. Copying the panel to the CPU, in this case just a square block on the diagonal, is done on line 4. The data transfer is asynchronous, so before we factor it on the CPU (line 8), we synchronize on line 7 to enforce that the data has arrived. Note that the CPU work from line 8 is overlapped with the Xeon Phi work on line 6. This is indeed the case because line 6 is an asynchronous call/request from the CPU to Xeon Phi to start a ZGEMM operation. Thus, the control is passed to lines 7 and 8 while Xeon Phi is performing the ZGEMM. The resulting factored panel from the CPU work is sent to Xeon Phi on line 11 and used on line 14, after making sure that it has arrived through the `sync` command on line 13.

6.4. Hybrid Implementation and Optimization Techniques.

In order to explain our hybrid methodology and the optimization that we have developed, let us give a detailed analysis for the QR decomposition algorithm. While the description below only addresses the QR factorization, it is straightforward to derive with the same ideas the analysis for both the Cholesky and LU factorizations. For that we start briefly by recalling the description of the QR algorithm.

The QR factorization is a transformation that factorizes an $m \times n$ matrix A into its factors Q and R where Q is a unitary matrix of size $m \times m$ and R is an upper trapezoidal matrix of size $m \times n$. The QR algorithm can be described as a sequence of steps where, at each step, a QR of a panel is performed based on accumulating a number of Householder transformations in what is called a “*panel factorization*” which are, then, applied all at once by means of high performance Level 3 BLAS operations in what is called the “*trailing matrix update*.” Despite the fact that this approach can exploit the parallelism of the Level 3 BLAS during the trailing matrix update, it

```

(1) for  $j = 0, nb, 2\ nb, 3\ nb, \dots, n-1$  do
(2)    $jb = \min(nb, n-j)$ ;
(3)   magma_zherk_mic(MagmaUpper, MagmaConjTrans,  $jb, j, m\_one, dA(0, j), ldda, one, dA(j, j), ldda, \mathbf{queue}$ );
(4)   magma_zgetmatrix_async_mic( $jb, jb, dA(j, j), ldda, work, 0, jb, \mathbf{queue}, \&event$ );
(5)   if  $j + jb < n$  then
(6)     magma_zgemm_mic(MagmaConjTrans, MagmaNoTrans,
        $jb, n-j-jb, j, mz\_one, dA(0, j), ldd, dA(0, j+jb), ldda, \mathbf{queue}$ );
(7)   magma_event_sync_mic( $\mathbf{event}$ );
(8)   lapackf77_zpotrf(MagmaUpperStr,  $\&jb, work, \&jb, info$ );
(9)   if  $*info \neq 0$  then
(10)     $*info += j$ ;
(11)  magma_zsetmatrix_async_mic( $jb, jb, work, 0, jb, dA(j, j), ldda, \mathbf{queue}, \&event$ );
(12)  if  $j + jb < n$  then
(13)    magma_event_sync_mic( $\mathbf{event}$ );
(14)    magma_ztrsm_mic(MagmaLeft, MagmaUpper, MagmaConjTrans, MagmaNonUnit,
        $jb, n-j-jb, z\_one, dA(j, j), ldda, dA(j, j+jb), ldda, \mathbf{queue}$ );

```

ALGORITHM 2: Cholesky factorization in MIC MAGMA.

has a number of limitations when implemented on massively multithreaded system such as the Intel Xeon Phi coprocessor due to the nature of its operations. On the one hand, the panel factorization relies on Level 2 BLAS operations that cannot be efficiently parallelized on either Xeon Phi or any accelerator such as GPU-based architectures, and thus it can be considered to be close to sequential operations that limit the scalability of the algorithm. On the other hand, this algorithm is referred to as the *fork-join approach* since the execution flow will show a sequence of sequential operations (panel factorizations) interleaved with parallel ones (trailing matrix updates). In order to take advantage of the high execution rate of the massively multithreaded system, in particular, the Phi coprocessor, we redesigned the standard algorithm in a way to perform the Level 3 BLAS operations (trailing matrix update) on the Xeon Phi while performing the Level 2 BLAS operations (panel factorization) on the CPU. We also proposed an algorithmic change to remove the fork-join bottleneck and to minimize the overhead of the panel factorization by hiding its costs behind the parallel trailing matrix update. This approach can be described as the *scalable lookahead techniques* [28]. Our idea is to split the trailing matrix update into two phases: the update of the lookahead panel (panel of step $i+1$, i.e., dark blue portion of Figure 2) and the update of the remaining trailing submatrix (clear blue portion of Figure 2). Thus, during the submatrix update the CPU can receive asynchronously the panel $i+1$ and perform its factorization. As a result, our MIC MAGMA implementation of the QR factorization can be described by a sequence of the three phases described below. Consider a matrix A that can be represented as

$$A = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix}. \quad (1)$$

(i) *Phase 1: The Panel Factorization.* At a step i , this phase consists of a QR transformation of the panel $A_{i:n,i}$ as in (2).

This operation consists of calling two routines: the DGEQR2 that factorizes the panel and produces nb Householder reflectors (V_{*i}) and an upper triangular matrix R_{ii} of size $nb \times nb$, which is a portion of the final R factor, and the DLARFT that generates the triangular matrix T_{ii} of size $nb \times nb$ used for the trailing matrix update. This phase is performed on the CPU:

$$\begin{bmatrix} A_{11} \\ A_{21} \\ A_{31} \end{bmatrix} \Rightarrow \begin{bmatrix} V_{11} \\ V_{21} \\ V_{31} \end{bmatrix}, [R_{1,1}], [T_{1,1}]. \quad (2)$$

(ii) *Phase 2: The Look Ahead Panel Update.* The transformation that was computed in the panel factorization needs to be applied to the rest of the matrix (trailing matrix, i.e., the blue portion of Figure 2). This phase consists in updating only the next panel (dark blue portion of Figure 2) in order to let the CPU start its factorization as soon as possible while the update of the remaining portion of the matrix is performed in phase 3. The idea is to hide the cost of the panel factorization. This operation, presented in (3), is performed on the Phi coprocessor and involves the DLARFB routine which has been redesigned as a sequence of DGEMMs to better take advantage of the Level 3 BLAS operations:

$$\begin{bmatrix} R_{12} \\ \tilde{A}_{22} \\ \tilde{A}_{32} \end{bmatrix} = [I - V_{*i} T_{ii}^T V_{*i}^T] \begin{bmatrix} A_{12} \\ A_{22} \\ A_{32} \end{bmatrix}. \quad (3)$$

(iii) *Phase 3: The Trailing Matrix Update.* Similarly to phase 2, this phase consists of applying the Householder reflectors generated during the panel factorization of step i , according to (3), to the remaining portion of the matrix (the trailing submatrix, i.e., the clear blue portion of Figure 2). This operation is also performed on the Phi coprocessor, while, in parallel to it, the CPU performs the factorization of the panel $i+1$ that has been computed in phase 2.

This hybrid technique of distribution of tasks in CPU-Phi allows us to hide the memory bound operations that occurred during the panel factorization (phase 1) by performing such operation on the CPU in parallel with the trailing submatrix update (phase 3) on the Phi coprocessor. However, one of the key parameters to performance tuning is the blocking size as the performance and the overlap between the CPU-Phi will be solely guided by it. Figure 3 illustrates the effect of the blocking factor on the performance. It is obvious that a small nb will reduce the cost of the panel factorization phase 1, but it decreases the efficiency of the Level 3 BLAS kernel of phase 2 and phase 3, thus resulting in a bad performance. On the contrary, a large nb will dramatically affect the panel factorization phase 1 which becomes slow and thus the CPU-Phi computation cannot be overlapped, providing a deterioration in the performance as shown in Figure 3. As a consequence, the challenging problem is the following: on the one hand, the blocking size nb needs to be large enough to extract high performance from Level 3 BLAS phase 3 and, on the other hand, it has to be small enough to extract efficiency (thanks to the cache speedup) from the Level 2 BLAS phase 1 and overlap CPU-Phi computation. Figure 3 shows the performance obtained for different blocking sizes and we can see a trade-off between small and large nb 's. Either $nb = 480$ or $nb = 960$ can be considered as a good choice because MKL Phi BLAS is optimized for multiples of 240. Moreover, to extract the maximum performance and allow the maximum overlap between both the CPU and the Xeon Phi coprocessor, we developed a new variant that can use a variable nb during the steps of the algorithm. The flexibility of our implementation allows an efficient task execution overlap between the CPU host and the Phi coprocessor which enables the implementation to scale almost linearly with the number of cores on the Phi coprocessor, as we can see (below) from the very good performance that is close to the practical peak obtained on such a system from matrix-matrix multiply and related dense linear algebra operations, which achieve over 70% of the theoretical peak performance. Our tuned variable implementation is represented by the red curve of Figure 3 where we can easily observe its advantages over the other variants.

The Phi-specific techniques had to be employed in order to reap the benefits of the above design in the presence of particular constraints and opportunities present on the Intel hardware. One opportunity is to choose the best one out of a number of interfaces for transferring data between the CPU and the coprocessor; refer to Table 1 for details. The Phi implementation of MAGMA seeks to minimize the latency and maximize the bandwidth of the PCIe transfers while maintaining a good computational load of both the host and the device. If the proper API for the right size of data transfer is chosen, the DMA hardware can take over and offload the transfer logistics so that the compute components can remain busy computing on matrix elements and not polluting their cache hierarchy with spurious messaging data. In particular, SCIF offers the lowest latency but the large data transfers create complexity burden of dealing with many smaller transfer requests. Higher level mechanisms, such as COI and virtual shared memory regions, carry a larger

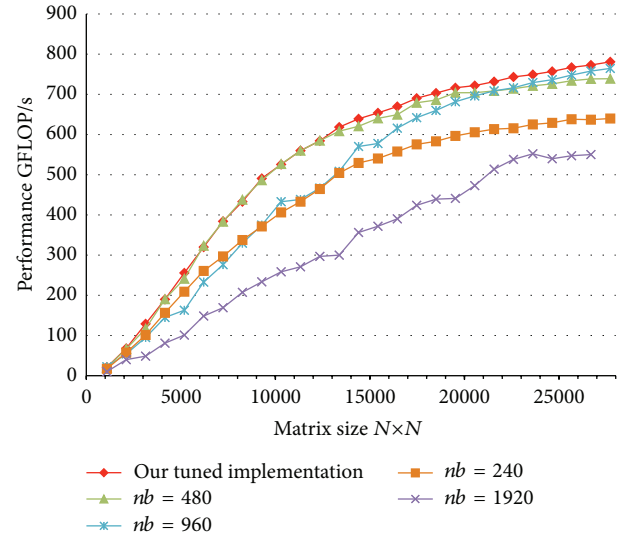


FIGURE 3: Effect of the blocking factor on performance of MAGMA MIC factorizations.

overhead but allow the handling of large volumes of data in a much more automated fashion. The switching between these interfaces occurs seamlessly behind the familiar MIC MAGMA functions.

6.5. Task-Based Runtime Model. The scheduling of tasks for execution can be static or dynamic. In either case, the small and not easy to parallelize tasks from the critical path (e.g., panel factorizations) are executed on CPUs and the large and highly parallel task (like the matrix updates) mostly on the MICs.

The use of multiple coprocessors complicates the development using static scheduling. Instead, the use of a light-weight runtime system is preferred as it can keep scheduling overhead low, while enabling the expression of parallelism through sequential-like code. The runtime system relieves the developer from keeping track of the computational activities that, in the case of heterogeneous systems, are further exacerbated by the separation between the address spaces of the main memory of the CPU and the MICs. Our runtime model is built on the QUARK [29] superscalar execution environment that has been originally used with great success for linear algebra software on just multicore platforms [30]. The conceptual work though could be replicated within other models such as StarPU [8], OmpSS [20], Cilk [18], and Jade [17], to just mention a few.

Dynamic runtime scheduling plays an important role in translating dependences annotated at the source code level and discovered at runtime when the execution traverses the Direct Acyclic Graph of computational tasks. For example, one of the symbolic dependences of tasks in Algorithm 1 could be

$$\text{PanelFactorize}_{\text{CPU}}(P_i) \longrightarrow \text{TrailingMatrixUpdate}_{\text{MIC}}(P_{i+1}). \quad (4)$$

At runtime, this dependence formula is repeatedly applied to form a sequence of tasks:

$$\begin{aligned} \text{PanelFactorize}_{\text{CPU}}(P_1) &\longrightarrow \text{TrailingMatrixUpdate}_{\text{MIC}}(P_2) \\ \text{PanelFactorize}_{\text{CPU}}(P_2) &\longrightarrow \text{TrailingMatrixUpdate}_{\text{MIC}}(P_3) \\ &\vdots \end{aligned} \quad (5)$$

The runtime environment for scheduling maintains the current set of tasks and the future set of tasks. The completed tasks enable execution of their dependent tasks and are discarded from the system.

6.6. Improving Off-Load Mode Communication. It is well known that the off-load transfer mode copies only continuous chunks of data from and to the coprocessors. However most of the scientific application algorithms require exchanging data with 2D or 3D storage and thus this may create an issue when using the off-load transfer mode. In particular, the one-sided factorizations (Cholesky, LU, and QR) require sending the panel to the CPU and then receiving it later after being factorized by the CPU. A simple implementation loops over one direction and calls the off-load section to send and receive a contiguous vector. Such an implementation behaves poorly and as a result the communication will become expensive and will slow down the algorithm. Indeed, another alternative is to copy the 2D panel to a contiguous temporary space on the MIC and then to send it and vice versa. Hence, there are two points that need to be taken into consideration. Firstly, the copy needs to be implemented as a multithreaded operation in order to hide its cost. For that, we implemented a parallel copy that uses all of the 240 hardware threads of the MIC to perform the copy. This might be against the common wisdom that multithreading is of little help for bandwidth-limited operations such as a memory copy. This is not the experience on the MIC, where the clock frequency of the compute cores is twice as low as that of the memory, the exact opposite of which is the case in Intel x86 multicore processors. In addition to the low frequency, the current MIC hardware is to a large degree an in-order architecture with dual-pipeline execution and single-issue fetch/decode units [11] which poses constraints on the amount of bandwidth that can be utilized by a single core. These can be overcome in multiple ways, including the use of streaming loads and having the multiple threads request data. Secondly, when the MIC copies data to or from the temporary space, it should be the only kernel running; otherwise, it will run simultaneously with another executing kernel and this may slow down both of the kernels. To that end, we represented the copy kernel as a task with high priority and the scheduler is responsible for executing it as soon as possible and handling the dependencies so that no other kernel will be running at the same time. Xeon Phi requires multiple cores driving a single FPU, which is similar to Hyperthreading in the recent Intel x86 processors. In fact, the core-to-FPU ratio must be two-to-one to satisfy the data rate that a single FPU can sustain. If the ratio is lower, the FPU goes largely

underutilized because the data request rate from memory is too low.

Experiments showed that when using these optimizations the performance of the off-load communication mode is comparable to both the SCIF and the COI mode with a variance of less than 5%.

6.7. Trading Extra Computation for Higher Execution Rate. The optimization discussed here is MIC-specific but is often valid for any hardware architecture with multilayered memory hierarchy. The `dlarfb` routine used by the QR decomposition consists of two `dgemms` and one `dtrmm`. Since coprocessors are better at handling compute-bound tasks, for computational efficiency, we replace the `dtrmm` by `dgemm`, yielding 5–10% performance improvement. For the Cholesky factorization, the trailing matrix update requires a `dsyrk`. Due to uneven storage, the multidevice `dsyrk` cannot be assembled purely from regular `dsyrk` calls on each device. Instead, each block column must be processed individually. The diagonal blocks require special attention. One can use a `dsyrk` to update each diagonal block and a `dgemm` to update the remainder of each block column below the diagonal block. The small `dsyrk` operations have little parallelism and therefore their execution is inefficient on MICs. This can be improved to some degree by using `pragma` to run several `dsyrks` simultaneously. Nevertheless, because we have copied the data to the device, we can consider the space above the diagonal to be a scratch workspace. Thus, we update the entire block column, including the diagonal block, writing extra data into the upper triangle of the diagonal block, which is subsequently ignored. We do extra computation for the diagonal block but gain efficiency overall by launching fewer BLAS kernels on the device and using the more efficient `dgemm` kernels, instead of small `dsyrk` kernels.

The per-kernel improvement in performance exceeds 20% and for the entire factorization 5–10% improvement levels may be observed.

7. Performance Results

This section presents the performance results obtained by our hybrid CPU-Xeon Phi implementation in the context of the development of the state-of-the-art numerical linear algebra libraries.

7.1. Experimental Environment. Our experiments were performed on a system equipped with Intel Xeon Phi formerly known as Knights Corner. It is representative of a vast class of servers and workstations commonly used for computationally intensive workloads. We benchmarked all implementations on an Intel multicore system with dual-socket, 8-core Intel Xeon E5-2670 (Sandy Bridge) processors, each running at 2.6 GHz. Each socket has a 24 MB shared L3 cache, and each core has a private 256 KB L2 and 64 KB L1. The system is equipped with 52 Gbytes of memory. The theoretical peak for this architecture in double precision is 20.8 Gflop/s per core, giving 332 Gflops in total. The system is also equipped with Intel Xeon Phi cards with 7.7 Gbytes

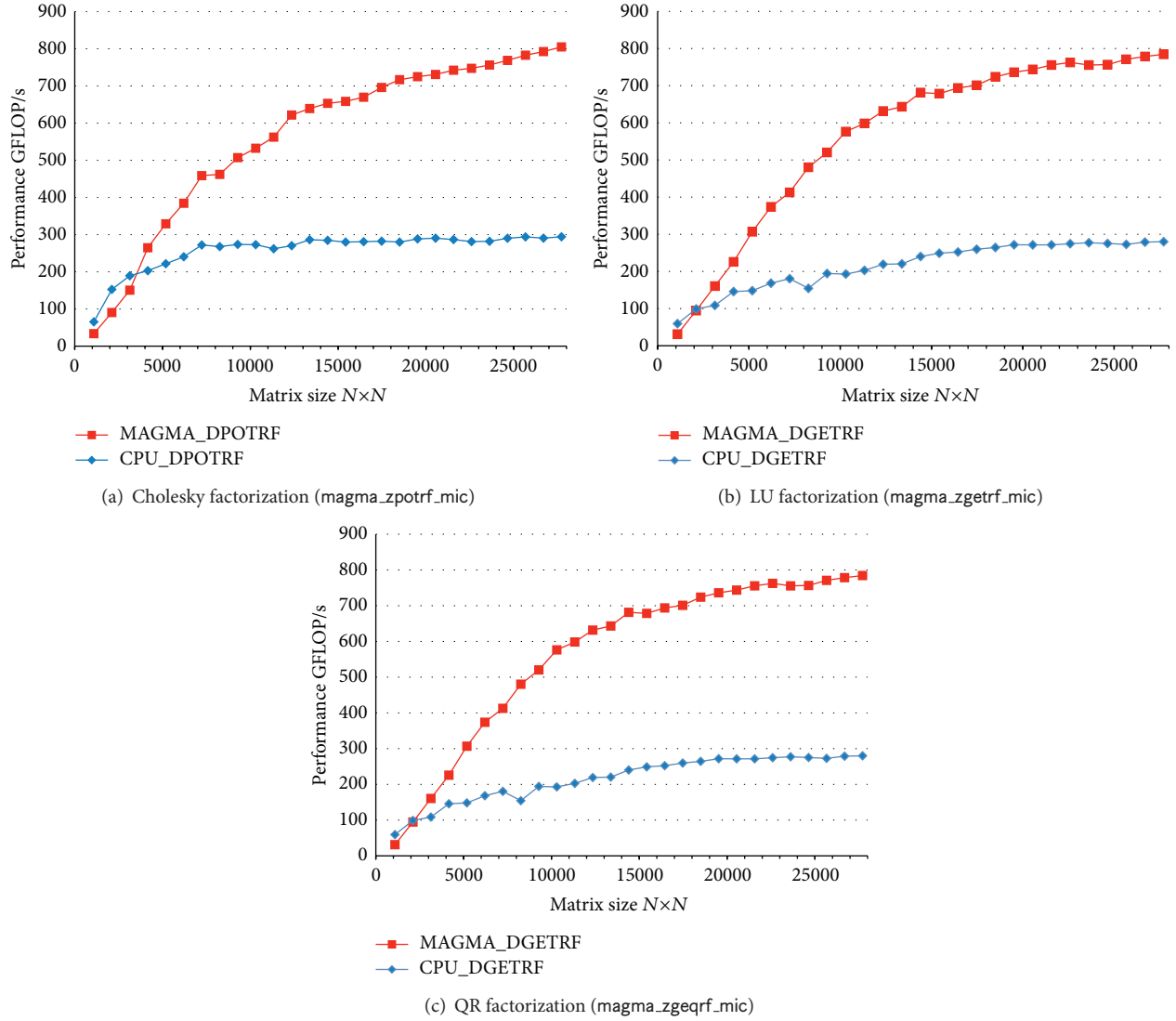


FIGURE 4: Comparison of the performance versus the optimized CPU version of the MKL libraries for the three one-sided factorizations.

per card running at 1.09 GHz and giving a double precision theoretical peak of 1046 Gflops.

There are a number of software packages available. On the CPU side we used the MKL (Math Kernel Library) [2] which is a commercial software package from Intel that is a highly optimized numerical library. On the Intel Xeon side, we used the MPSS 2.1.5889-16 as the software stack, icc 13.1.1 20130313, which comes with the composer_xe.2013.3.163 suite as the compiler and Level 3 BLAS routine GEMM from MKL 11.00.03.

7.2. Experimental Results. Figure 4 reports the performance of the three linear algebra factorization operations, the Cholesky, QR, and LU factorizations, with our hybrid implementation and compares it to the performance of the CPU implementation of the MKL libraries. For our implementation, the blocking factor has been chosen to be flexible in order to achieve the best performance. A detailed description

of how to choose this factor is included in Section 6.4 and in the results presented in this section we choose the factor to be in the range between 480 and 960. As a general rule, we use smaller blocking factors for smaller matrices and larger ones for the larger matrices. The graphs show the performance measured using all the cores available on the system (i.e., 60 for the Intel Phi and 16 for the CPU) with respect to the problem size. In order to reflect the time to completion, for each algorithm the operation count is assumed to be the same as that of the LAPACK algorithm, that is, $(1/3)N^3$, $(2/3)N^3$, and $(4/3)N^3$ for the Cholesky factorization, the LU factorization, and the QR decomposition, respectively.

Figures 4(a), 4(b), and 4(c) provide the common type of information that is characteristic of dense linear algebra computations. Clearly, our algorithms from the MIC MAGMA library, which employ hybrid techniques, deliver higher execution rates than their CPU counterparts optimized by the vendor. This is in correspondence with the difference

of the peak performance rates between the two hardware components. It should be obvious from the graphs that the combination of a CPU and a Phi coprocessor with a tuned implementation provides substantial performance benefits as opposed to a CPU-only implementation. The figures show that the MIC MAGMA hybrid algorithms are capable of completing any of the three factorization algorithms as much as twice as fast as the CPU optimized version for a matrix of size larger than 10000 and more than three times faster when the matrix size is large enough (larger than 20000). The actual curves of Figure 4 illustrate the efficiency of our hybrid techniques where we note that the performance obtained by our implementation achieves a very close level to the practical peak of the Intel Xeon Phi coprocessor computed by running the GEMM routine (which is around 850 Gflop/s). This gain is mostly obtained by two improvements. First, the nature of the operations involved on the Phi side which are mostly BLAS Level 3 operations was redesigned and implemented as a combination of vendor's DGEMM calls. For more details we denote below the routines executed on the Xeon Phi coprocessor:

- (i) the DSYRK operations for the Cholesky factorization where the DSYRK has been redesigned as a combination of DGEMM's routines,
- (ii) the DGEMM for the LU factorization,
- (iii) the DLARFB for the QR decomposition where also its has been redesigned as a combination of DGEMMs.

Second, all of the Level 2 BLAS routines that are memory bound and that represent a limit for the performance (i.e., DPOTF2, DGETF2, and DGEQR2 for Cholesky, LU, and QR factorization, resp.) are executed on the CPU side while being overlapped with the Phi coprocessor execution as described in Section 6.4.

An important remark has to be made here for the Cholesky factorization: the *left-looking* algorithm as implemented in LAPACK is considered as well optimized for memory reuse but at the price of less parallelism and thus is not suitable for massively multicore machines. This variant delivers poor performance when compared to the *right-looking* variant that allows more parallelism and thus runs at higher speed.

8. Conclusions and Future Work

In this paper, we have shown how to extend our hybridization methodology from existing systems to a new hardware platform. The challenge of the porting effort stemmed from the fact that the new coprocessor from Intel, the Xeon Phi, featured programming models and relative execution overheads that were markedly different from what we have been targeting on GPU-based accelerators. Nevertheless, we believe that the techniques used in this paper adequately adapt our hybrid algorithm to best take advantage of the new heterogeneous hardware. We have derived an implementation schema of the dense linear algebra kernels that also can be applied either to the two-sided factorization used for solving the eigenproblem and the SVD or to the sparse linear

algebra algorithms. We plan to further study the implementation of multi-Xeon Phi algorithms in a distributed computing environment. We think that the techniques presented will become more popular and will be integrated into dynamic runtime system technologies. The ultimate goal is that this integration will help to tremendously decrease development time while retaining high performance.

In addition, we see an opportunity in fully automating the tuning process of various algorithmic parameters of our implementation including the blocking factor nb and the number of threads used in various computational kernels. This will become even more important as the number of linear algebra operations included in MIC MAGMA grows.

Conflict of Interests

The authors declare that there is no conflict of interests regarding the publication of this paper.

Acknowledgments

This work was funded in part by the financial support of the Russian Scientific Fund Agreement N14-11-00190. The authors would like to thank the National Science Foundation for supporting this work under Grant no. ACI-1339822 and the Department of Energy and ISTC for Big Data for supporting this research effort.

References

- [1] E. Anderson, Z. Bai, C. Bischof et al., *LAPACK User's Guide*, SIAM, Philadelphia, Pa, USA, 3rd edition, 1999.
- [2] Intel, Math Kernel Library, <https://software.intel.com/en-us/intel-mkl>.
- [3] AMD, AMD Core Math Library (ACML), <http://developer.amd.com/tools/>.
- [4] J. Dongarra, J. Bunch, C. Moler, and G. W. Stewart, *LINPACK Users' Guide*, SIAM, Philadelphia, Pa, USA, 1979.
- [5] S. L. Blackford, J. Choi, A. Cleary et al., *ScaLAPACK Users' Guide*, SIAM, Philadelphia, Pa, USA, 1997, <http://www.netlib.org/scalapack/slug/>.
- [6] E. Agullo, J. Demmel, J. Dongarra et al., "Numerical linear algebra on emerging architectures: the PLASMA and MAGMA projects," *Journal of Physics: Conference Series*, vol. 180, no. 1, Article ID 012037, 2009.
- [7] *Software Distribution of MAGMA MIC Version 1.0*, 2013, <http://icl.cs.utk.edu/magma/software/>.
- [8] C. Augonnet, S. Thibault, R. Namyst, and P.-A. Wacrenier, "StarPU: a unified platform for task scheduling on heterogeneous multicore architectures," *Concurrency Computation Practice and Experience*, vol. 23, no. 2, pp. 187–198, 2011.
- [9] Barcelona Supercomputing Center, *SMP Superscalar (SMPSs) User's Manual, Version 2.0*, 2008, <http://www.bsc.es/media/1002.pdf>.
- [10] Intel, Intel Xeon Phi coprocessor system software developers guide, <http://software.intel.com/en-us/articles/>.
- [11] J. Jeffers and J. Reinders, *Intel Xeon Phi Coprocessor High-Performance Programming*, Morgan Kaufmann Publishers, 2013.

- [12] NVIDIA Corporation, *NVIDIA CUDA C Programming Guide*, 2014, <http://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>.
- [13] Khronos OpenCL Working Group, *The OpenCL Specification, Version: 1.0 Document Revision: 48*, Khronos OpenCL Working Group, 2009.
- [14] AMD, clMath libraries: clBLAS 2.0, August 2013, <https://github.com/clMathLibraries>.
- [15] M. M. Wolf, M. A. Heroux, and E. G. Boman, "Factors impacting performance of multithreaded sparse triangular solve," in *High Performance Computing for Computational Science—VECPAR 2010: 9th International conference, Berkeley, Calif, USA, June 22–25, 2010, Revised Selected Papers*, vol. 6449 of *Lecture Notes in Computer Science*, pp. 32–44, Springer, Berlin, Germany, 2011.
- [16] S. Williams, D. D. Kalamkar, A. Singh et al., "Optimization of geometric multigrid for emerging multi- and manycore processors," in *Proceedings of the 24th International Conference for High Performance Computing, Networking, Storage and Analysis (SC '12)*, pp. 96:1–96:11, IEEE, Los Alamitos, Calif, USA, November 2012.
- [17] M. C. Rinard, D. J. Scales, and M. S. Lam, "Jade: a high-level, machine-independent language for parallel programming," *Computer*, vol. 26, no. 6, pp. 28–38, 1993.
- [18] R. D. Blumofe, C. F. Joerg, B. C. Kuszmaul, C. E. Leiserson, K. H. Randall, and Y. Zhou, "Cilk, an efficient multithreaded runtime system," *SIGPLAN Notices*, vol. 30, pp. 207–216, 1995.
- [19] K. Fatahalian, D. R. Horn, T. J. Knight et al., "Sequoia: programming the memory hierarchy," in *Proceedings of the ACM/IEEE Conference on Supercomputing (SC '06)*, ACM, New York, NY, USA, 2006.
- [20] J. M. Pérez, R. M. Badia, and J. Labarta, "A dependency-aware task-based programming environment for multi-core architectures," in *Proceedings of the IEEE International Conference on Cluster Computing (CCGRID '08)*, pp. 142–151, IEEE, Tsukuba, Japan, October 2008.
- [21] R. Barik, Z. Budimčić, V. Cavè et al., "The habanero multicore software research project," in *Proceedings of the 24th ACM SIGPLAN Conference Companion on Object Oriented Programming Systems Languages and Applications (OOPSLA '09)*, pp. 735–736, ACM, New York, NY, USA, October 2009.
- [22] C. H. González and B. B. Fraguera, "A framework for argument-based task synchronization with automatic detection of dependencies," *Parallel Computing*, vol. 39, no. 9, pp. 475–489, 2013.
- [23] P. Du, R. Weber, P. Luszczek, S. Tomov, G. Peterson, and J. Dongarra, "From CUDA to OpenCL: towards a performance-portable solution for multi-platform GPU programming," *Parallel Computing*, vol. 38, no. 8, pp. 391–407, 2012.
- [24] The OpenACC application programming interface version 1.0, 2011.
- [25] "Proposed additions for OpenACC 2.0, OpenACC application programming interface," 2012.
- [26] S. Tomov and J. Dongarra, "Dense linear algebra for hybrid GPU-based systems," in *Scientific Computing with Multicore and Accelerators*, J. Kurzak, D. A. Bader, and J. Dongarra, Eds., Chapman & Hall, CRC Press, 2010.
- [27] L. G. Valiant, "Bulk-synchronous parallel computers," in *Parallel Processing and Artificial Intelligence*, M. Reeve, Ed., pp. 15–22, JohnWiley & Sons, 1989.
- [28] P. E. Strazdins, "Lookahead and algorithmic blocking techniques compared for parallel matrix factorization," in *Proceedings of the 10th International Conference on Parallel and Distributed Computing and Systems (IASTED '98)*, Las Vegas, Nev, USA, October 1998.
- [29] A. YarKhan, *Dynamic task execution on shared and distributed memory architectures [Ph.D. thesis]*, University of Tennessee, 2012.
- [30] J. Kurzak, P. Luszczek, A. YarKhan et al., "Multithreading in the PLASMA library," in *Handbook of Multi and Many-Core Processing: Architecture, Algorithms, Programming, and Applications*, Computer and Information Science Series, Chapman & Hall/CRC, 2013.

