

Analysis of Dynamically Scheduled Tile Algorithms for Dense Linear Algebra on Multicore Architectures

Azzam Haidar*, Hatem Ltaief*, Asim YarKhan* and Jack Dongarra*†‡§

*Department of Electrical Engineering and Computer Science,
University of Tennessee, Knoxville

†Computer Science and Mathematics Division, Oak Ridge National Laboratory,
Oak Ridge, Tennessee

‡School of Mathematics & School of Computer Science,
University of Manchester

§Research reported here was partially supported by the National Science Foundation and Microsoft Research.

Abstract—The objective of this paper is to analyze the dynamic scheduling of dense linear algebra algorithms on shared-memory, multicore architectures. Current numerical libraries, e.g., LAPACK, show clear limitations on such emerging systems mainly due to their coarse granularity tasks. Thus, many numerical algorithms need to be redesigned to better fit the architectural design of the multicore platform. The PLASMA library (Parallel Linear Algebra for Scalable Multi-core Architectures) developed at the University of Tennessee tackles this challenge by using tile algorithms to achieve a finer task granularity. These tile algorithms can then be represented by Directed Acyclic Graphs (DAGs), where nodes are the tasks and edges are the dependencies between the tasks. The paramount key to achieve high performance is to implement a runtime environment to efficiently schedule the DAG across the multicore platform. This paper studies the impact on the overall performance of some parameters, both at the level of the scheduler, e.g., window size and locality, and the algorithms, e.g., Left Looking (LL) and Right Looking (RL) variants. We conclude that some commonly accepted rules for dense linear algebra algorithms may need to be revisited.

I. INTRODUCTION

The scientific high performance computing community has recently faced dramatic hardware changes with the emergence of multicore architectures. Most of the fastest high performance computers in the world, if not all, mentioned in the last Top500 list [1] released in November 2009 are now based on multicore architectures. This confronts the scientific software community with both a daunting challenge and a unique opportunity. The challenge arises from the disturbing mismatch between the design of systems based on this new chip architecture – hundreds of thousands of nodes, a million or more cores, reduced bandwidth and memory available to cores – and the components of the traditional software stack, such as numerical libraries, on which scientific

applications have relied for their accuracy and performance. The state of the art, high performance dense linear algebra software libraries, i.e., LAPACK [4] have shown limitations on multicore architectures [3]. The performance of LAPACK relies on the use of a standard set of Basic Linear Algebra Subprograms (BLAS) [14], [20] within which nearly all of the parallelism occurs following the expensive fork-join paradigm. Moreover, its large stride memory accesses have further exacerbated the problem, and it becomes judicious to efficiently develop existing or new numerical linear algebra algorithms suitable for such hardware.

As discussed by Buttari *et al.* in [9], a combination of several parameters define the concept of *tile algorithms* and are essential to match the architecture associated with the cores: (1) Fine Granularity to reach a high level of parallelism and to fit the core small caches; (2) Asynchronicity to prevent any global barriers; (3) Block Data Layout (BDL), a high performance data representation to perform efficient memory access; and (4) Dynamic Data Driven Scheduler to ensure any queued tasks can immediately be processed as soon as all their data dependencies are satisfied.

The PLASMA library (Parallel Linear Algebra for Scalable Multi-core Architectures) [25] jointly developed by the University of Tennessee, the University of California Berkeley and the University of Denver Colorado, tackles this challenge by using tile algorithms to achieve high performance. These tile algorithms can then be represented by Directed Acyclic Graphs (DAGs), where nodes are the tasks and edges are the dependencies between the tasks. The paramount key is to implement a runtime environment to efficiently schedule the DAG across the multicore platform.

This paper studies the impact on the overall performance of some parameters, both at the level of the sched-

uler, e.g., window size and locality, and the algorithms, e.g., Left Looking (LL) and Right Looking (RL) variants. The conclusion of this study claims that some commonly accepted rules for dense linear algebra algorithms may need to be revisited.

The remainder of this paper is as follows: Section II recalls the mechanisms behind block algorithms (e.g., LAPACK) and explains the different looking variants (LL and RL). Section III describes the concept of tile algorithms. Section IV introduces the dynamic scheduler of DAGs. Section V shows some performance results. Related work in the area is mentioned in Section VI. Section VII summarizes the paper and presents future work.

II. BLOCK ALGORITHMS

In this section, we review the paradigm behind the state-of-the-art numerical software, namely the LAPACK library for shared-memory. In particular, we focus on three widely used factorizations in the scientific community, i.e., QR, LU and Cholesky, which are the first steps toward solving numerical linear systems. All the kernels mentioned below have freely available reference implementations as part of either the BLAS or LAPACK.

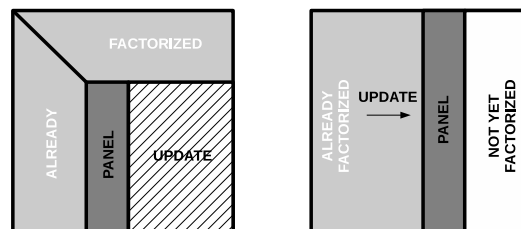
A. Description and Concept

The LAPACK library provides a broad set of linear algebra operations aimed at achieving high performance on systems equipped with memory hierarchies. The algorithms implemented in LAPACK leverage the idea of blocking to limit the amount of bus traffic in favor of a high data reuse that is present in the higher level memories that are also the fastest ones. The idea of blocking revolves around an important property of Level-3 BLAS operations (Matrix-Matrix multiplication), the so called surface-to-volume property, which states that ($\theta(n^3)$) floating point operations are performed on ($\theta(n^2)$) data. Because of this property, Level-3 BLAS operations can be implemented in such a way that data movement is limited and reuse of data in the cache is maximized. Block algorithms consist of recasting linear algebra algorithms in a way that only a negligible part of computations is done in Level-2 BLAS operations (Matrix-Vector multiplication, where no data reuse possible) while most is done in Level-3 BLAS. Most of these algorithms can be described as the repetition of two fundamental steps (see Fig. 1):

- Panel factorization: depending of the linear algebra operation that has to be performed, a number of transformations are computed for a small portion of the matrix (the so called panel). These transformations, computed by means of Level-2 BLAS operations, can be accumulated.

- Trailing submatrix update: in this step, all the transformations that have been accumulated during the panel factorization step can be applied at once to the rest of the matrix (i.e., the trailing submatrix) by means of Level-3 BLAS operations.

Although the panel factorization can be identified as a sequential execution task that represents a small fraction of the total number of FLOPS (required $\theta(n^2)$, for a total of $\theta(n^3)$), the scalability of block factorizations is limited on a multicore system. Indeed, the panel factorization is rich in Level-2 BLAS operations that cannot be efficiently parallelized on currently available shared memory machines. Moreover, the parallelism is only exploited at the level of the BLAS routines. This methodology implies a fork-join model since the execution flow of a block factorization represents a sequence of sequential operations (panel factorizations) interleaved with parallel ones (updates of the trailing submatrices).



(a) Right Looking Variant.

(b) Left Looking Variant.

Fig. 1. Block Algorithm Steps.

B. Block Cholesky Factorization

The Cholesky factorization (or Cholesky decomposition) is mainly used as a first step for the numerical solution of linear equations $Ax = b$, where A is symmetric and positive definite. Such systems arise often in physics applications, where A is positive definite due to the nature of the modeled physical phenomenon.

The Cholesky factorization of an $n \times n$ real symmetric positive definite matrix A has the form $A = LL^T$, where L is an $n \times n$ real lower triangular matrix with positive diagonal elements. In LAPACK, the double precision algorithm is implemented by the DPOTRF routine. A single step of the algorithm is implemented by a sequence of calls to the LAPACK and BLAS routines: DSYRK, DPOTF2, DGEMM, DTRSM. Due to the symmetry, the matrix can be factorized either as an upper triangular matrix or as a lower triangular matrix.

C. Block QR Factorization

Generally, a QR factorization of an $m \times n$ real matrix A is the decomposition of A as $A = QR$, where Q is an $m \times m$ real orthogonal matrix and R is an $m \times n$ real upper triangular matrix. QR factorization uses a series

of elementary Householder matrices of the general form $H = I - \tau vv^T$, where v is a column reflector and τ is a scaling factor.

Regarding the block algorithms as performed in LAPACK [4], by the DGEQRF routine, nb elementary Householder matrices are accumulated within each panel and the product is represented as $H_1 H_2 \dots H_{nb} = I - VTV^T$. Here V is an $n \times nb$ matrix in which columns are the vectors v , T is an $nb \times nb$ upper triangular matrix and nb is the block size. A single step of the algorithm is implemented by a sequence of calls to the LAPACK and BLAS routines: DGEQR2 (panel factorization kernel), DLARFT (computation of the structure T) and DLARFB (trailing submatrix update kernel).

D. Block LU Factorization

The LU factorization (or LU decomposition) with partial row pivoting of an $m \times n$ real matrix A has the form $A = PLU$, where L is an $m \times n$ real unit lower triangular matrix, U is an $n \times n$ real upper triangular matrix and P is a permutation matrix. In the block formulation of the algorithm, factorization of nb columns (the panel) is followed by the update of the remaining part of the matrix (the trailing submatrix) [13], [15]. In LAPACK the double precision algorithm is implemented by the DGETRF routine. A single step of the algorithm is implemented by a sequence of calls to the following LAPACK and BLAS routines: DGETF2, DLASWP (apply pivoting), DTRSM, DGEMM, where DGETF2 implements the panel factorization and the other routines implement the updates.

E. Block Looking Variants

Algorithmic variants exist for the factorizations explained above. The two main ones are called Left Looking (LL) and Right Looking (RL). They only differ on the location of the update applications with regards to the panel. The RL variant operates on the current panel and applies the corresponding updates to the right (see Fig. 1(a)). On the contrary, the LL variant (also called the "lazy" variant) applies all updates coming from the left up to the current panel (see Fig. 1(b)) and therefore delays subsequent updates of the remaining columns of the matrix.

III. TILE ALGORITHMS

In this section, we describe a solution that removes the fork-join overhead seen in block algorithms. Based on tile algorithms, this new model is currently used in shared memory libraries, such as PLASMA and FLAME (University of Texas Austin) [2].

A. Description and Concept

A solution to this fork-join bottleneck in block algorithms has been presented in [10], [11], [17], [19], [22]. The approach consists of breaking the panel factorization and trailing submatrix update steps into smaller tasks that operate on a block-column (i.e., a set of b contiguous columns where b is the block size). The algorithm can then be represented as a Directed Acyclic Graph (DAG) where nodes represent tasks, either panel factorization or update of a block-column, and edges represent dependencies among them.

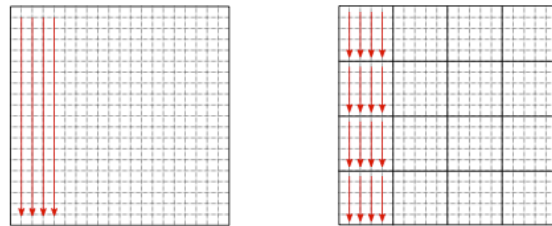


Fig. 2. Translation from LAPACK Layout to Block Data Layout

The execution of the algorithm is performed by asynchronously scheduling the tasks in a way that dependencies are not violated. This asynchronous scheduling results in an out-of-order execution where slow, sequential tasks are hidden behind parallel ones. The following sections describe the tile algorithm paradigm applied to a class of factorizations, i.e., Cholesky, LU and QR where finer granularity of the operations and higher flexibility for the scheduling can be achieved.

B. Tile Cholesky Factorization

The tile Cholesky algorithm described in Figure 3 is identical to the block Cholesky algorithm implemented in LAPACK, except for processing the matrix by tiles. Otherwise, the exact same operations are applied.

C. Tile QR Factorization

Here a derivative of the block algorithm is used called the *tile QR* factorization. The ideas behind the tile QR factorization are well known. The tile QR factorization was initially developed to produce a high-performance "out-of-memory" implementation (typically referred to as "out-of-core") [16] and, more recently, to produce a high performance implementation on "standard" (x86 and alike) multicore processors [10], [11], [12] and on the CELL processor [19].

The algorithm is based on the idea of annihilating matrix elements by square tiles instead of rectangular panels (block columns). The algorithm produces "essentially" the same R factor as the classic algorithm, e.g., the implementation in the LAPACK library (elements may differ in sign). However, a different set

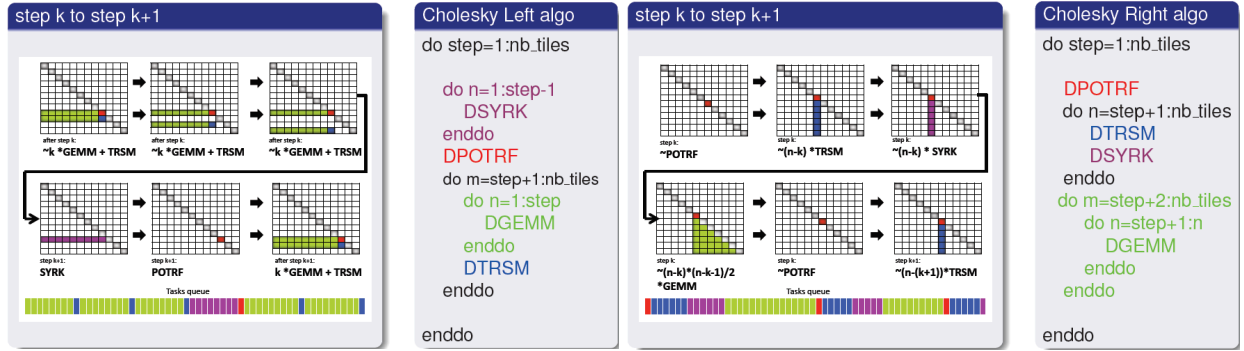


Fig. 3. Tile Cholesky factorization left v.s. right looking algorithm.

of Householder reflectors is produced and a different procedure is required to build the Q matrix. The tile QR algorithm described in Fig. 4 relies on four basic operations implemented by four computational kernels:

CORE_DGEQRT: The kernel performs the QR factorization of a diagonal tile and produces an upper triangular matrix R and a unit lower triangular matrix V containing the Householder reflectors. The kernel also produces the upper triangular matrix T as defined by the compact WY technique for accumulating Householder reflectors [6], [23]. The R factor overrides the upper triangular portion of the input and the reflectors override the lower triangular portion of the input. The T matrix is stored separately.

CORE_DTSQRT: The kernel performs the QR factorization of a matrix built by coupling the R factor, produced by **CORE_DGEQRT** or a previous call to **CORE_DTSQRT**, with a tile below the diagonal tile. The kernel produces an updated R factor, a square matrix V containing the Householder reflectors and the matrix T resulting from accumulating the reflectors V . The new R factor overrides the old R factor. The block of reflectors overrides the corresponding tile of the input matrix. The T matrix is stored separately.

CORE_DORMQR: The kernel applies the reflectors calculated by **CORE_DGEQRT** to a tile to the right of the diagonal tile, using the reflectors V along with the matrix T .

CORE_DSSMQR: The kernel applies the reflectors calculated by **CORE_DTSQRT** to two tiles to the right of the tiles factorized by **CORE_DTSQRT**, using the reflectors V and the matrix T produced by **CORE_DTSQRT**.

D. Tile LU Factorization

Here a derivative of the block algorithm is used called the *tile LU* factorization. Similarly to the tile QR

algorithm, the tile LU factorization originated as an “out-of-memory” (“out-of-core”) algorithm [22] and was recently rediscovered for the multicore architectures [11], [12].

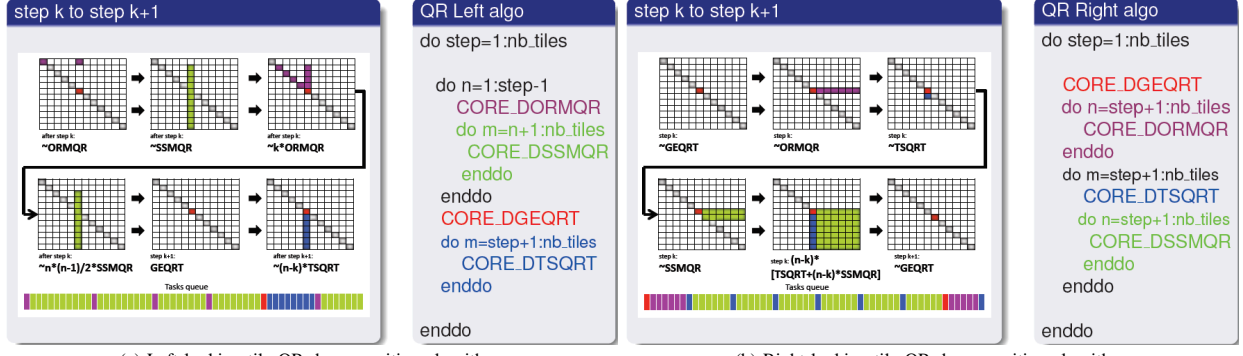
Again, the main idea here is the one of annihilating matrix elements by square tiles instead of rectangular panels. The algorithm produces different U and L factors than the block algorithm (e.g., the one implemented in the LAPACK library). In particular we note that the L matrix is not lower unit triangular anymore. Another difference is that the algorithm does not use partial pivoting but a different pivoting strategy. The tile LU algorithm relies on four basic operations implemented by four computational kernels:

CORE_DGETRF: The kernel performs the LU factorization of a diagonal tile and produces an upper triangular matrix U , a unit lower triangular matrix L and a vector of pivot indexes P . The U and L factors override the input and the pivot vector is stored separately.

CORE_DTSTRF: The kernel performs the LU factorization of a matrix built by coupling the U factor, produced by **DGETRF** or a previous call to **CORE_DTSTRF**, with a tile below the diagonal tile. The kernel produces an updated U factor and a square matrix L containing the coefficients corresponding to the off-diagonal tile. The new U factor overrides the old U factor. The new L factor overrides the corresponding off-diagonal tile. A new pivot vector P is created and stored separately. Due to pivoting, the lower triangular part of the diagonal tile is scrambled and also needs to be stored separately as L' .

CORE_DGESSM: The kernel applies the transformations produced by the **DGETRF** kernel to a tile to the right of the diagonal tile, using the L factor and the pivot vector P .

CORE_DSSSSM: The kernel applies the transformations produced by the **CORE_DTSTRF** kernel to



(a) Left looking tile QR decomposition algorithm.

(b) Right looking tile QR decomposition algorithm.

Fig. 4. Tile QR decomposition left v.s. right looking algorithm.

the tiles to the right of the tiles factorized by `CORE_DTSTRF`, using the L' factor and the pivot vector P .

The tile LU algorithm is similar to the tile QR algorithm in Fig. 4, but with different kernels.

One topic that requires further explanation is the issue of pivoting. Since in the tile algorithm only two tiles of the panel are factorized at a time, pivoting only takes place within two tiles at a time, a scheme which could be described as *block-pairwise pivoting*. Clearly, such pivoting is not equivalent to the “standard” *partial row pivoting* in the block algorithm (e.g., LAPACK). A different pivoting pattern is produced, and also, since pivoting is limited in scope, the procedure could potentially result in a less numerically stable algorithm. More details on the numerical stability of the tile LU algorithm can be found in [11].

E. Tile Looking Variants

The algorithmic principles of the RL and the LL variants with tile algorithms are similar to block algorithms (see Section II-E). The panel and update regions of the matrix are now split into tiles. The update operations, whether for LL or RL variants, may concurrently run with the panel operations. Those variants actually highlight a trade-off between degree of parallelism (RL) and data reuse (LL) and can considerably affect the overall performance. For example, for block algorithms, Cholesky factorization is implemented with the LL variant on shared-memory (LAPACK), while for distributed memory (ScaLAPACK [7]) the RL variant is used. This paper studies whether this common rule still holds with tile algorithms on current multicore architectures.

IV. DYNAMIC DATA DRIVEN EXECUTION

A. Runtime Environment for Dynamic Task Scheduling

Restructuring linear algebra algorithms as a sequence of tasks that operate on tiles of data can remove the fork-join bottlenecks seen in block algorithms. This is accomplished by enabling out-of-order execution of tasks,

```

FOR k = 0..TILES-1
  A[k][k] <- DPOTRF(A[k][k])
  FOR m = k+1..TILES-1
    A[m][k] <- DTRSM(A[k][k], A[m][k])
  FOR n = k+1..TILES-1
    A[n][n] <- DSYRK(A[n][k], A[n][n])
    FOR m = n+1..TILES-1
      A[m][n] <- DGEMM(A[m][k], A[n][k], A[m][n])

```

Fig. 5. Pseudocode of the tile Cholesky factorization (right-looking version).

which can hide the work done by the bottleneck (sequential) tasks. We would like to schedule the sequence of tasks on shared-memory, many-core architectures in a flexible, efficient and scalable manner. In this section, we present an overview of our runtime environment for dynamic task scheduling from the perspective of an algorithm writer who is using the scheduler to create and to execute an algorithm. There are many details about the internals of the scheduler, its dependency analysis, memory management, and other performance enhancements that are not covered here. However, information about an earlier version of this scheduler can be found in [18].

As an example, we present the pseudocode for the tile Cholesky factorization in Fig. 5 as an algorithm designer might view it. Tasks in this Cholesky factorization example depend on previous tasks if they use the same tiles of data. If these dependencies are used to relate the tasks, then a directed acyclic graph (DAG) is implicitly formed by the tasks. A small DAG for a 5x5 tile matrix is shown in Fig. 6.

B. Determining Task Dependencies

a) *Description of Dependency Types:* In order for a scheduler to be able to determine dependencies between the tasks, it needs to know how each task is using its arguments. Arguments can be `VALUE`, which are copied to the task, or they can be `INPUT`, `OUTPUT`, or `INOUT`, which have the expected meanings. Given the sequential order that the tasks are added to the scheduler,

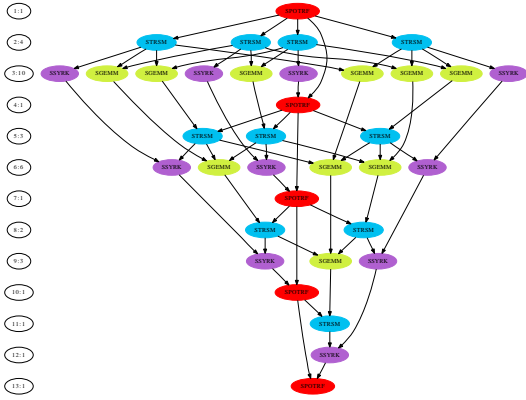


Fig. 6. DAG for a small Cholesky factorization (right looking version) with five tiles (block size 200 and matrix size 1000). The column on the left shows the depth:width of the DAG.

and the way that the arguments are used, we can infer the relationships between the tasks. A task can read a data item that is written by a previous task (read-after-write RAW dependency); or a task can write a data item that is written by previous task (write-after-write WAW dependency); a task can write a data time that is read by a previous task (write-after-read WAR dependency). The dependencies between the tasks form an implicit DAG, however this DAG is never explicitly realized in the scheduler. The structure is maintained in the way that tasks are queued on data items, waiting for the appropriate access to the data.

b) From Sequential Nested-Loop Code to Parallel Execution: Our scheduler is designed to use code very similar to the pseudocode described in Fig. 5. This is intended to make it easier for algorithm designers to experiment with algorithms and design new algorithms. In Fig. 9 we can see the final C code from the Cholesky algorithm pseudocode. Each of the calls to the core linear algebra routines is substituted by a call to a wrapper that decorates the arguments with their sizes and their usage (INPUT, OUTPUT, INOUT, VALUE). As an example, in Fig. 7 we can see how the DPOTRF call is decorated for the scheduler.

The tasks are inserted into the scheduler, which stores them to be executed when all the dependencies are satisfied. That is, a task is ready to be executed when all parent tasks have completed. The execution of ready tasks is handled by worker threads that simply wait for tasks to become ready and execute them using a combination of default tasks assignments and work stealing. The thread doing the task insertion, i.e., the thread handling the code in Fig. 9, is referred to as the master thread. Under certain circumstances, the master

```
int DSCHEDED_dpotrff( Dsched *dsched, char uplo, int n,
                    double *A, int lda, int *info )
{
    DSCHEDED_Insert_Task( dsched, TASK_core_dpotrff, 0x00,
                          sizeof(char), &uplo, VALUE,
                          sizeof(int), &n, VALUE,
                          sizeof(double)*n*n, A, INOUT | LOCALITY,
                          sizeof(int), &lda, VALUE,
                          sizeof(int), info, OUTPUT,
                          0);
}
void TASK_dpotrff(Dsched *dsched)
{
    char uplo; int n; double *A; int lda; int *info;
    dsched_unpack_args_5( dsched, uplo, n, A, lda, info );
    dpotrff_( &uplo, &n, A, &lda, info );
}
```

Fig. 7. Example of inserting and executing a task in the scheduler. The DSCHEDED_dpotrff routine inserts a task into the scheduler, passing it the sizes and pointers of arguments and their usage (INPUT, OUTPUT, INOUT, VALUE). Later, when the dependencies are satisfied and the task is ready to execute, the TASK_dpotrff routine unpacks the arguments from the scheduler and calls the actual dpotrff routine.

thread will also execute computational tasks. Fig. 8 provides an idealized overview of the architecture of the dynamic scheduler.

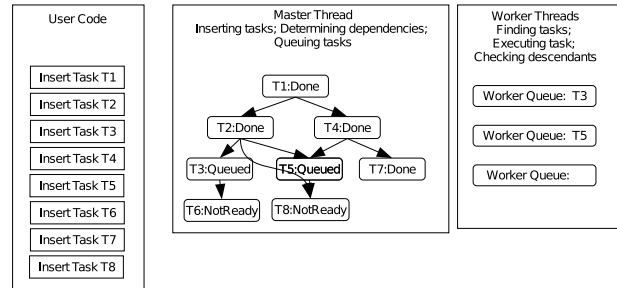


Fig. 8. Idealized architecture diagram for the dynamic scheduler. Inserted tasks go into a (implicit) DAG based on their dependencies. Tasks can be in NotReady, Queued or Done states. Workers execute queued tasks and then determine if any descendants have now become ready and can be queued.

c) Scheduling a Window of Tasks: For the applications that we are considering, the number of tasks ($\theta(n^3)$) grows very quickly with the number of TILES of data. A relatively small example of LU factorization using 20×20 tiles generates 2870 tasks (see Fig. 10), whereas using 50×50 tiles generates 42925 tasks. If we were to unfold and retain the entire DAG of tasks for a large problem, we would be able to perform some interesting analysis with respect to DAG scheduling and critical paths. However, the size of the data structures would quickly grow overwhelming. Our solution to this is to maintain a configurable window of tasks. The implicit DAG is then traversed through this sliding window, which should be large enough to ensure all cores are kept busy. When this window size is reached, the core involved in inserting tasks does not accept any

```

for ( i = 0 ; i < p ; i++ ) {
  DSCHEM_dpotrf( dsched, 'L', nb[i], A[i][i], nb[i], info );
  for ( j = i+1 ; j < p ; j++ )
    DSCHEM_dtrsm( dsched, 'R', 'L', 'T', 'N', nb[j], nb[i], 1.0, A[i][i], nb[i], A[j][i], nb[j]);
  for ( j = i+1 ; j < p ; j++ ) {
    for ( k = i+1 ; k < j ; k++ ) {
      DSCHEM_dgemm( dsched, 'N', 'T', nb[j], nb[k], nb[i], -1, A[j][i], nb[j], A[k][i], nb[k], 1, A[j][k], nb[j]);
      DSCHEM_dsyrc( dsched, 'L', 'N', nb[j], nb[i], -1.0, A[j][i], nb[j], +1.0, A[j][j], nb[j]);
    }
  }
}

```

Fig. 9. Tile Cholesky factorization that calls the scheduled core linear algebra operations.

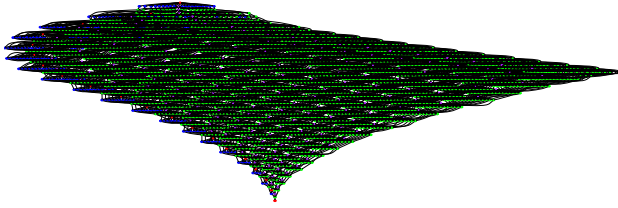


Fig. 10. DAG for a LU factorization with 20 tiles (block size 200 and matrix size 4000). The size of the DAG grows very fast with the number of tiles.

more tasks until some are completed. The usage of a window of tasks has implications in how the loops of an application are unfolded and how much look ahead is available to the scheduler. This paper discusses some of these implication in the context of dense linear algebra applications.

d) Data Locality and Cache Reuse: It has been shown in the past that the reuse of memory caches can lead to a substantial performance improvement in execution time. Since we are working with tiles of data that should fit in the local caches on each core, we have provided the algorithm designer with the ability to hint the cache locality behavior. A parameter in a call (e.g., Fig. 7) can be decorated with the LOCALITY flag in order to tell the scheduler that the data item (parameter) should be kept in cache if possible. After a computational core (worker) executes that task, the scheduler will assign by-default any future task using that data item to the same core. Note that the work stealing can disrupt the by-default assignment of tasks to cores.

The next section studies the performance impact of the locality flag and the window size on the LL and RL variants of the three tile factorizations.

V. EXPERIMENTAL RESULTS

This section describes the analysis of dynamically scheduled tile algorithms for the three factorizations (i.e., Cholesky, QR and LU) on different multicore systems. The tile sizes for these algorithm have been tuned and

are equal to $b = 200$.

A. Hardware Descriptions

In this study, we consider two different shared memory architectures. The first architecture (System A) is a quad-socket, quad-core machine based on an Intel Xeon EMT64 E7340 processor operating at 2.39 GHz. The theoretical peak is equal to 9.6 Gflops/s per core or 153.2 Gflops/s for the whole node, composed of 16 cores. The practical peak (measured by the performance of a GEMM) is equal to 8.5 Gflops/s per core or 136 Gflops/s for the 16 cores. The level-1 cache, local to the core, is divided into 32 kB of instruction cache and 32 kB of data cache. Each quad-core processor is actually composed of two dual-core Core2 architectures and the level-2 cache has 2×4 MB per socket (each dual-core shares 4 MB). The machine is a NUMA architecture and it provides Intel Compilers 11.0 together with the MKL 10.1 vendor library.

The second system (System B) is an 8 sockets, 6 core AMD Opteron 8439 SE Processor (48 cores total @ 2.8Ghz) with 128 Gb of main memory. Each core has a theoretical peak of 11.2 Gflops/s and the whole machine 537.6 Gflops/s. The practical peak (measured by the performance of a GEMM) is equal to 9.5 Gflops/s per core or 456 Gflops/s for the 48 cores. There are three levels of cache. The level-1 cache consist of 64 kB and the level-2 cache consist of 512 kB. Each socket is composed of 6 cores and the level-3 cache has 6 MB 48-way associative shared cache per socket. The machine is a NUMA architecture and it provides Intel Compilers 11.1 together with the MKL 10.2 vendor library.

B. Performance Discussions

In this section, we evaluate the effect of the window size and the locality feature on the LL and RL tile algorithm variants.

The nested-loops describing the tile LL variant codes are naturally ordered in a way that already promotes locality on the data tiles located on the panel. Fig. 11 shows the effect of the locality flag of the scheduler on the overall performance of the tile LL Cholesky variant. As expected, the locality flag does not really improve the performances when using small window

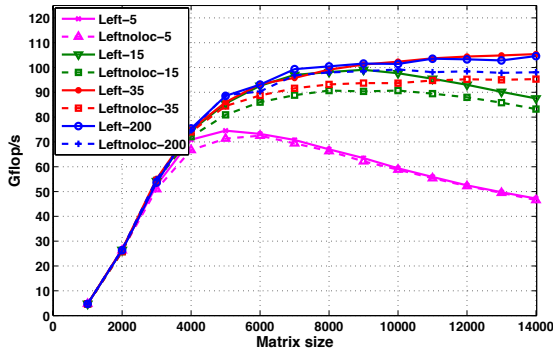


Fig. 11. Cholesky factorization (LL) on System A, 16 threads: effect of the scheduler data locality parameter.

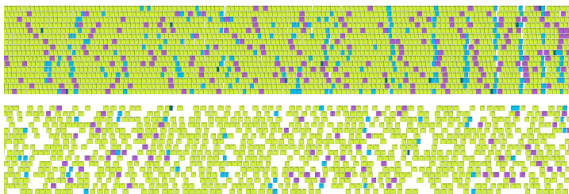


Fig. 12. The execution trace of the LL Cholesky: large (top) v.s. small (bottom) window sizes.

sizes. The scheduler is indeed not able to perform enough look-ahead to anticipate the reuse occurring in the next panels. With larger window sizes, the scheduler is now able to acquire knowledge of data reuse in the next panels and takes full advantage of it. Therefore, the locality scheduler flag permits the performance of the LL tile Cholesky variant to increase by up to 5-10% for large window sizes. The LL tile QR and LU behave in the same manner as the LL tile Cholesky. On the contrary, the RL variants for the three factorizations are not affected by this parameter. The RL variant triggers on the right side of the panel so many parallel independent tasks that the chance for any data reuse opportunities is significantly decreased.

The next parameter to be optimized for the three tile factorizations is the window size of the scheduler. Fig. 12 shows how critical the optimization of this parameter can be. The figure displays the execution trace of the tile LL Cholesky algorithm with small and large scheduler window sizes on the System A. We observe that for a small window size (bottom trace), the execution is stretched because there are not enough queued tasks to feed the cores, which makes them turn to an “idle” state (white spaces). This is further demonstrated by looking at the performance graphs from Fig. 13-15 obtained on both Systems A and B for the three factorizations. The performance is represented as a function of the matrix size. In these figures, the size of the window is defined

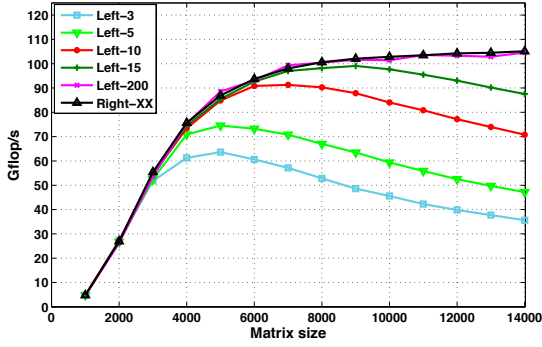
by the number indicated in the legend times the number of the available threads. By increasing the size of the scheduler window, the performances of the LL variants of the three factorizations considerably increase. For example, the overall performance of the tile LL Cholesky is multiplied by 3 on the System A and by 2 on System B when increasing the window size per core from 3 to 200. In the same way, the performance of the tile LL QR is multiplied by 4 on Systems A and B. The performance of the tile LL LU is also multiplied by 3 when increasing the window size of the scheduler on both Systems A and B. So, in other words, if the scheduler window size is small, the runtime environment system is not able to detect and anticipate the concurrent execution of independent sets of tasks. The already low performance curve encountered with small window sizes on the three factorizations starts significantly dropping even more for large matrix sizes.

On the other hand, the window size parameter has no impact on the tile RL variants. The performance curves of the RL Cholesky, QR and LU denoted by *Right-XX* outperform the LL ones regardless of the window sizes. The curves eventually meet up for large enough window sizes. Again, this is due to the high degree of parallelism offered by the RL variants, which allows many tasks to be simultaneously executed. The RL variants seem to be very attractive for shared-memory multicore systems, especially because they do not require any parameter auto-tuning like the LL variants. Therefore, the RL variants should be chosen by default in the PLASMA library distribution for those tile algorithms.

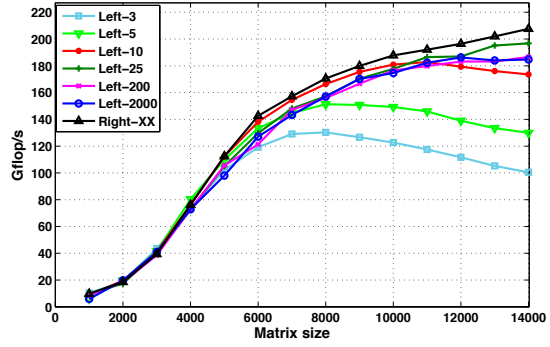
Furthermore, by closely studying some execution traces of LL and RL variants, the clear distinctions between both variants inherited from block algorithms (i.e., LAPACK) may not exist anymore in the context of data-driven, asynchronous out-of-order DAG execution. Indeed, by increasing the window sizes of the scheduler, the tile LL variants are able to perform look-ahead techniques by initiating the work on the next panel while the one on the current panel is still pending. Likewise, the tile RL variants are permitted to start processing the next panel while the updates of the current panel are still ongoing. One type of algorithm can thus morph into another type of algorithm. The definitions of LL and RL variants may then become obsolete in this particular context.

VI. RELATED WORK

The FLAME project [2] developed by the University of Texas, Austin follows the same algorithmic principle by splitting the matrix into finer blocks. However, the runtime environment (SuperMatrix) requires the explicit construction of the DAG before the actual parallel execution of the tasks. As depicted in Fig. 10, the size of a

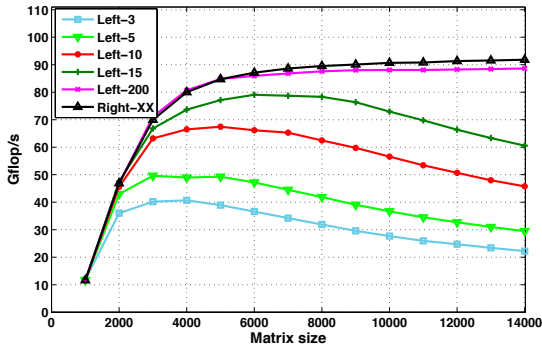


(a) System A: 16 threads

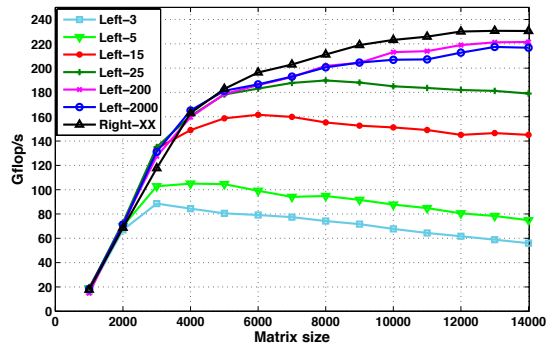


(b) System B: 48 threads

Fig. 13. Cholesky factorization (left v.s. right looking version) with different task window size.



(a) System A: 16 threads



(b) System B: 48 threads

Fig. 14. QR factorization (left v.s. right looking version) with different task window size.

DAG considerably increases with respect to the number of tiles. This may necessitate a large amount of memory to allocate the data structure which is not an option for scheduling large matrix sizes, especially when dealing with other algorithms (e.g., two-sided transformations) that generate a more complex and larger DAG compared to QR, LU and Cholesky.

Furthermore, there are many projects that are designed to provide high performance, near-transparent computing environments for shared-memory machines. Here we discuss two projects that have been considered during the design of our runtime environment.

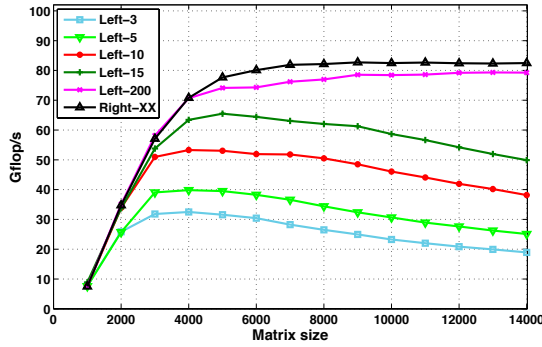
The SMP superscalar (SMPSs) project[5][21] from the Barcelona Supercomputing Center is a programming environment for shared memory, multi-core architectures focused on the ease of programming, portability and flexibility. A standard C or Fortran code can be marked up using preprocessor pragma directives to enable task level parallelism. The parameters to the functions are marked as input, output, or inout, and the data dependencies between tasks are inferred in order to determine a task DAG. A source-to-source compiler and a supporting runtime library are used to generate native code for the platform. The SMPSs project shares many similarities

with the dynamic runtime environment presented here. One difference is that our implementation uses an API to express parallelism rather than compiler pragmas, thus eliminating an intermediate step. Another difference is that our runtime environment allows for specialized flags, such as the LOCALITY flag, which enable tuning for linear algebra algorithms.

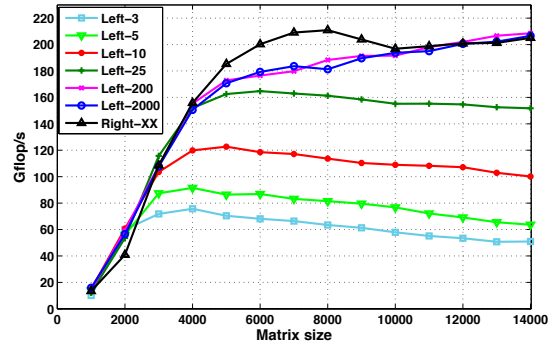
The Cilk project [8] [24] from the MIT Laboratory for Computer Science is a compiler based extension of the C language that gives the programmer a set of keywords to express task level parallelism (*cilk*, *spawn*, *sync*, *inlet*, *abort*). Cilk is well suited to algorithms that can be expressed recursively and implements a fork-join model of multi-threaded computation. Since the parallelism in our algorithms is expressed in a DAG obtained through a data dependency analysis, Cilk is not well suited to our problems.

VII. CONCLUSION AND FUTURE WORK

This paper presents an analysis of dynamically scheduled tile algorithms for dense linear algebra on shared-memory, multicore systems. In particular, the paper highlights the significant impact of the locality feature and the window sizes of the scheduler on the LL and



(a) System A: 16 threads



(b) System B: 48 threads

Fig. 15. LU factorization (left v.s. right looking version) with different task window size on 48 threads.

RL variants of the tile Cholesky, QR and LU. The tile RL variants of the three factorizations outperform the LL ones regardless of the scheduler locality and window size optimizations. It is also common rule-of-thumb that LL variants of linear algebra algorithms (e.g., Cholesky) are well suited to shared-memory architectures and RL variants are better suited to distributed memory architectures. However, these rules are no longer valid in the context of data-driven, asynchronous out-of-order DAG execution environments, because adjusting the size of the task window (for LL variants) provides sufficient look-ahead opportunities for one algorithmic variant can morph into another algorithmic variant.

REFERENCES

- [1] TOP500 Supercomputing Sites. <http://www.top500.org/>.
- [2] The FLAME project. <http://z.cs.utexas.edu/wiki/flame/wiki/FrontPage>, April 2010.
- [3] E. Agullo, B. Hadri, H. Ltaief, and J. Dongarra. Comparative study of one-sided factorizations with multiple software packages on multi-core hardware. In *SC '09: Proceedings of the Conference on High Performance Computing Networking, Storage and Analysis*, pages 1–12, New York, NY, USA, 2009. ACM.
- [4] E. Anderson, Z. Bai, C. Bischof, L. S. Blackford, J. W. Demmel, J. J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. Sorensen. *LAPACK Users' Guide*. SIAM, Philadelphia, PA, 1992. <http://www.netlib.org/lapack/ug/>.
- [5] Barcelona Supercomputing Center. *SMP Superscalar (SMPSs) User's Manual, Version 2.0*, 2008.
- [6] C. Bischof and C. van Loan. The WY representation for products of Householder matrices. *J. Sci. Stat. Comput.*, 8:2–13, 1987.
- [7] L. S. Blackford, J. Choi, A. Cleary, E. D'Azevedo, J. Demmel, I. Dhillon, J. J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R. C. Whaley. *ScaLAPACK Users' Guide*. SIAM, Philadelphia, PA, 1997. <http://www.netlib.org/scalapack/ug/>.
- [8] R. D. Blumofe, C. F. Joerg, B. C. Kuszmaul, C. E. Leiserson, K. H. Randall, and Y. Zhou. Cilk: An efficient multithreaded runtime system. In *Proceedings of the Fifth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPOPP)*, pages 207–216, Santa Barbara, California, July 1995.
- [9] A. Buttari, J. Langou, J. Kurzak, and J. Dongarra. Parallel tiled QR factorization for multicore architectures. *Concurr. Comput. : Pract. Exper.*, 20(13):1573–1590, 2008.
- [10] A. Buttari, J. Langou, J. Kurzak, and J. J. Dongarra. Parallel tiled QR factorization for multicore architectures. *Concurrency Computat.: Pract. Exper.*, 20(13):1573–1590, 2008.
- [11] A. Buttari, J. Langou, J. Kurzak, and J. J. Dongarra. A class of parallel tiled linear algebra algorithms for multicore architectures. *Parallel Comput. Syst. Appl.*, 35:38–53, 2009.
- [12] E. Chan, E. S. Quintana-Orti, G. Gregorio Quintana-Orti, and R. van de Geijn.
- [13] J. W. Demmel. *Applied Numerical Linear Algebra*. SIAM, 1997.
- [14] J. J. Dongarra, J. D. Croz, I. S. Duff, , and S. Hammarling. A set of level 3 basic linear algebra subprograms. *ACM Trans. Math. Soft.*, 16:1–17, 1990.
- [15] J. J. Dongarra, I. S. Duff, D. C. Sorensen, and H. A. van der Vorst. *Numerical Linear Algebra for High-Performance Computers*. SIAM, 1998.
- [16] B. C. Gunter and R. A. van de Geijn. Parallel out-of-core computation and updating the QR factorization. *ACM Transactions on Mathematical Software*, 31(1):60–78, 2005.
- [17] J. Kurzak, A. Buttari, and J. J. Dongarra. Solving systems of linear equation on the CELL processor using Cholesky factorization. *Trans. Parallel Distrib. Syst.*, 19(9):1175–1186, 2008.
- [18] J. Kurzak and J. Dongarra. Fully dynamic scheduler for numerical scheduling on multicore processors. Technical Report LAWN (LAPACK Working Note) 220, UT-CS-09-643, Innovative Computing Lab, University of Tennessee, 2009.
- [19] J. Kurzak and J. J. Dongarra. QR factorization for the CELL processor. *Scientific Programming*, 17(1-2):31–42, 2009.
- [20] C. L. Lawson, R. J. Hanson, D. Kincaid, and F. T. Krogh. Basic linear algebra subprograms for FORTRAN usage. *ACM Trans. Math. Soft.*, 5:308–323, 1979.
- [21] J. M. Perez, R. M. Badia, and J. Labarta. A Dependency-Aware Task-Based Programming Environment for Multi-Core Architectures. In *CLUSTER'08*, pages 142–151, 2008.
- [22] E. S. Quintana-Orti and R. A. van de Geijn. Updating an LU factorization with pivoting. *ACM Trans. Math. Softw.*, 35(2):11, 2008.
- [23] R. Schreiber and C. van Loan. A storage-efficient WY representation for products of Householder transformations. *J. Sci. Stat. Comput.*, 10:53–57, 1991.
- [24] Supercomputing Technologies Group, Massachusetts Institute of Technology Laboratory for Computer Science. *Cilk 5.4.6 Reference Manual*, Nov. 2001.
- [25] University of Tennessee. *PLASMA Users' Guide, Parallel Linear Algebra Software for Multicore Architectures, Version 2.0*, November 2009.