

LU Factorization of Small Matrices: Accelerating Batched DGETRF on the GPU

Tingxing Dong*, Azzam Haidar*, Piotr Luszczek*, James Austin Harris* Stanimire Tomov*, and Jack Dongarra*^{†‡}

* University of Tennessee, Knoxville, TN 37996

[†] Oak Ridge National Laboratory, USA

[‡] University of Manchester M13 9PL, UK

{tdong, haidar, luszczek, tomov, dongarra}@utk.edu

Abstract—Gaussian Elimination is commonly used to solve dense linear systems in scientific models. In a large number of applications, a need arises to solve many small size problems, instead of few large linear systems. The size of each of these small linear systems depends on the number of the ordinary differential equations (ODEs) used in the model, and can be on the order of hundreds of unknowns. To efficiently exploit the computing power of modern accelerator hardware, these linear systems are processed in batches. To improve the numerical stability, at least partial pivoting is required, most often accomplished with row pivoting. However, row pivoting can result in a severe performance penalty on GPUs because it brings in thread divergence and non-coalesced memory accesses. In this paper, we propose a batched LU factorization for GPUs by using a multi-level blocked right looking algorithm that preserves the data layout but minimizes the penalty of partial pivoting. Our batched LU achieves up to 2.5-fold speedup when compared to the alternative CUBLAS solution on a K40c GPU.

I. INTRODUCTION

Various scientific applications use Gaussian elimination to solve dense linear systems. An important class of problems is when many small size systems, instead of few large ones, must be solved. Typically, the order of the linear systems is up to a few hundred, and their number is from a few thousand to millions. For example, subsurface transportation simulations have a number of reaction systems to solve. Each system involves computing a Jacobian matrix and iteratively applying the Gaussian elimination until an outer solver converges. The system size is typically around 100.

The one-sided factorizations such as the Cholesky, LU, and QR factorizations are based on block outer-product updates of the trailing matrix. Algorithmically, this corresponds to a sequence of two distinct phases: the *panel factorization* and the *trailing matrix update*. Implementation of these two phases can be expressed as a straightforward loop shown in Algorithm 1.

Algorithm 1 Two-phase implementation of a one-sided factorization.

```

for  $P_i \in \{P_1, P_2, \dots, P_n\}$  do
  PanelFactorize( $P_i$ )
  TrailingMatrixUpdate( $C^{(i)}$ )
end for

```

The panel factorization is latency and memory-bound due

to its predominant reliance on the Level 2 BLAS operations, MAGMA performs the panel factorization on the CPU and only uses the GPU to update the trailing matrix. A data transfer of the factorized panel from the CPU to the GPU is required at each step of the loop in Algorithm 1 to perform the trailing matrix update.

In the batched LU implementation, however, we cannot afford such a memory transfer at any step, since the trailing matrix is small and the amount of computation is not sufficient to overlap it in time with the panel factorization. Many small data transfers will take away any performance advantage enjoyed by the GPU, especially due to the fact that the data for transfer are not continuous in the memory but instead are stored with a stride called a leading dimension. Another challenge to achieving good performance is the pivoting, which is a source of thread divergence and non-coalescent memory accesses. This is the result of consecutive threads accessing the matrix elements with a stride of one column instead of one element stride when the matrix is stored in column-major format.

II. ALGORITHMIC VARIANTS

The LU factorization (also called decomposition) is the first step in solving a dense linear system of equations $Ax = b$, where $A \in \mathbb{R}^{m \times n}$. The LU factorization of A with partial pivoting has the form $PA = LU$, where $L \in \mathbb{R}^{m \times n}$ is a lower triangular with unit diagonal elements (lower trapezoidal if $m > n$), $U \in \mathbb{R}^{n \times n}$ is an upper triangular matrix (upper trapezoidal if $m < n$), and $P \in \{0, 1\}^{m \times m}$ is the row permutation matrix.

A. The Blocked Right-Looking Algorithm

The blocked right-looking variant is shown in Algorithm 2 and its patterns of access to matrix elements is depicted in Figure 1. The factorization of the m by n matrix A proceeds in $\lceil n/nb \rceil$ steps of size nb except for the last one. The computation of the above steps in the LAPACK routine `dgetrf`, involves four operations: `dgetf2`, `dtrsm`, `dgemm`, and `dlaswp`.

For $A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}$, where $A_{11} \in \mathbb{R}^{nb \times nb}$, the $\begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix}$ submatrix is called a panel matrix. The panels are

factorized by the `dgetf2` routine:

$$P \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} = \begin{bmatrix} L_{11} \\ L_{21} \end{bmatrix} U_{11}. \quad (1)$$

The L_{11} and U_{11} submatrices overwrite A_{11} . The transformations in this panel factorization, along with the pivoting P must be applied to the trailing matrix before the factorization proceeds to the next step. Related to the pivoting, the rows of A_{12} are permuted with the selected from the factorization pivot rows (from A_{12} and A_{22}). This is done by the `dlaswp` routine. The pivoting information is stored in a vector generated by `dgetf2`.

After the permutation, A_{12} is updated by a lower-triangular solve $A'_{12} \leftarrow L_{11}^{-1} A_{12}$ (`dtrsm`), and A_{22} is updated by the so called Schur complement: $A'_{22} = A_{22} - A'_{21} A'_{12}$ (`dgemm`). The trailing matrix A'_{22} is now considered as the new matrix to be factored in the next iteration of the loop. This algorithm keeps updating the right hand side – the trailing matrix – and hence it is called right-looking.

The `dtrsm` and the `dgemm` routines are known as Level 3 BLAS – they allow for cache-friendly implementations that scale well with computational load without overly taxing the main memory bus. Due to the use of Level 3 BLAS, the blocked implementations perform very well and reach high flops-per-second, and in particular much higher than a non-blocking implementation that relies on memory-bound operations such as the Level 2 BLAS [1].

Algorithm 2 The blocked right looking LU factorization.

```

for  $i \in \{1, 2, 3, \dots, n/nb\}$  do
  Panel Factorize  $A_{ii} = L_{ii} U_{ii}$ 
  Compute  $A_{ij} = L_{ii}^{-1} A_{ij}$ 
  Permutation  $P$ 
  Trailing Matrix Update  $A_{jj} = A_{jj} - A_{ji} A_{ij}$  where  $A_{ij} =$ 
     $a(i \times nb : n, j \times nb : n)$ 
end for

```

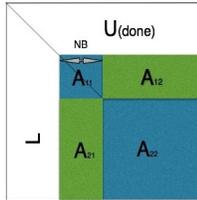


Fig. 1. Access patterns to matrix regions for the blocked right-looking LU factorization algorithm.

B. Multi-level blocked algorithm

The multi-level blocked algorithm is a variant of the blocked algorithm. The main difference is in the update of the trailing matrix. The right-looking variant operates on a current panel and updates all the way to the right. The multi-level blocked variant only applies the update to the next panel, but postpones

the update of the rest of the trailing matrix after the “k-levels” of panels are factorized.

III. BATCHED IMPLEMENTATION

We target matrices of size less than or equal to 512, since most application candidates for batched execution are of this size [2], [3].

A. Batched routines implementation

In a batched problem, each matrix is a separate problem that is solved independently. All of the routines discussed are batched and denoted by the corresponding LAPACK routine name. We have implemented the routines in the four standard precision arithmetics. For convenience, we use double precision routine name throughout the paper.

1) `dgetf2`: `dgetf2` is used to factorize a panel of size $m \times nb$ at each step of the LU factorization. It consists of three Level 1 BLAS calls (`idamax`, `dswap` and `dscal`) and one Level 2 BLAS call (`dger`). Note that a natural way of implementing `dgetf2` could be to load the panel to the GPU’s shared memory and then do the entire computation before writing the result back to the main memory. However, this direction cannot be easily implemented and cannot provide good performance for two main reason. First, the size of the shared memory is limited currently to only 48KB per streaming multiprocessor (SMX), which limits the panels that can fit at once in it. Second, saturating the shared memory per SMX can decrease performance, since only one thread-block will be mapped to a SMX at a time. Indeed, the number of threads used in the thread block will be limited, resulting in low occupancy, and subsequently poor core utilization. In our implementation of `dgetf2`, to perform the Gaussian elimination for the i^{th} column of the panel, we load only the column i to the shared memory. We found that such an implementation allows many thread-blocks to be executed by the same SMX in parallel, and thus taking a better advantage of its resources.

2) `dlaswp`: To improve the numerical stability, pivoting is required. However, pivoting can be a performance killer for matrices stored in column major format. Indeed, a factorization directly in column-major format can be two times slower (depending on hardware and problem sizes) than implementations that transpose the matrix in order to internally use a row-major storage format [4]. Yet, experiments show that this conversion is too expensive for batched problems. In the LAPACK’s `dlaswp`, the row swapping operations are serial, that is row by row. This limits the parallelism and is one of the factors for slow `dlaswp` for matrices in the column-major format. To minimize this penalty, we proposed a parallel swapping, detailed in Section IV-A.

3) `dtrsm`: After the panel factorization (1) and the row swapping, we compute the inverse of L_{11} , L_{11}^{-1} , with the `dtrtri` routine. Then, the A'_{12} update is accomplished by a `dgemm`, $A'_{12} = L_{11}^{-1} A_{12}$. Generally, computing the inverse of a matrix may suffer from numerical stability, but since A_{11} results from the numerically stable LU with partial pivoting and its size is just $nb \times nb$, or in our case 32×32 , we do not have this problem [5].

4) *dgemm*: The goal of our batched LU is to reach the performance of the batched *dgemm*. Because of its importance, a lot of previous efforts have been focused on optimizing the *dgemm* routine. In particular for our case, *dgemm* is not only used in the trailing matrix updates but also in the implementation of the triangular matrix solvers (*dtrsm*). Since NVIDIA CUBLAS *dgemm* is written in assembly language and highly optimized on Kepler architecture, we call CUBLAS routines.

IV. VARIOUS FACTOR IMPACTS ON THE PERFORMANCE

A. Parallel swapping

We analyzed and evaluated the implementation as described above to find that more than 60% of the factorization time is spent in the swapping routine. Figure 2 shows the execution trace of 2,000 batched LU factorization of matrices of size 512. We can observe on the top trace that the classical *dlaswp* kernel is the most time consuming part of the algorithm. The swapping consists of nb successive interchanges of two rows of the matrices. The main reason that this kernel is the most time consuming is because the nb row interchanges are performed in a sequential order, and that the data of a row is not coalescent, thus the thread warps do not read/write it in parallel. CPUs for example alleviate the effect of the long latency operations and bandwidth limitations by using hierarchical caches. Accelerators on the other hand, in addition to hierarchical memories, uses thread level parallelism (TLP) where threads are grouped into warps (e.g., of 32 threads) and multiple warps assigned for execution on the same SMX unit. In order to overcome the bottleneck of swapping, we proposed to modify the kernel in order to apply all nb row swaps in parallel. This modification will also allow the coalescent write of the first nb rows of the matrix. So we changed the algorithm to generate two pivot vectors, where the first vector gives the final destination row indices for the first nb rows of the panel, and the second gives the row indices of the nb rows that must become the first nb rows of the panel. Figure 2 depicts the execution trace (bottom) when using our parallel *dlaswp* kernel. The experiment shows that this reduces the time spent in the kernel from 60% to around 10%. As a result, the gain obtained in terms of performance is around 50%, as shown in Figure 5.

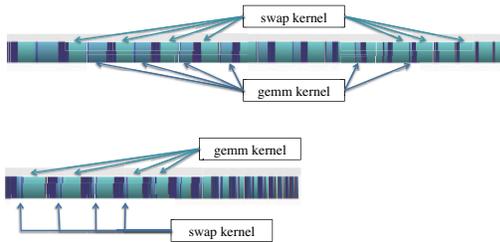


Fig. 2. Execution trace of the batched LU factorization using either classical swap (top) or our new parallel swap (bottom).

B. Nested blocking

The panel factorization as described in III-A1 goes over the nb columns and factorizes them one after another, similarly to the LAPACK algorithm. At each of the nb steps, a rank-1 update is required to update the vectors at the right hand side of the factorized column i (this operation is done by the *dger* kernel). Since we cannot load the entire panel into the shared memory of the GPU, the right hand side vectors are loaded back and forth from the main memory at every step. Thus, one can expect that the rank-1 operation is the most time consuming of the panel factorization. A detailed analysis using the profiler reveals that the *dger* kernel consists of more than 80% of the panel factorization time, and around 40% of the total LU factorization time. Similarly to the swapping kernel described above the main bottleneck here is the memory access. For that, we propose to improve the efficiency of this kernel and to reduce the memory access by using a recursive level of blocking techniques. In principle, the panel can be blocked recursively until a single element. Yet, in practice, 2-3 blocked levels are sufficient to achieve high performance. The above routines must be optimized for each blocked level, which complicates the implementation. The boost in performance obtained by this optimization is around 25%, as demonstrated in Figure 5.

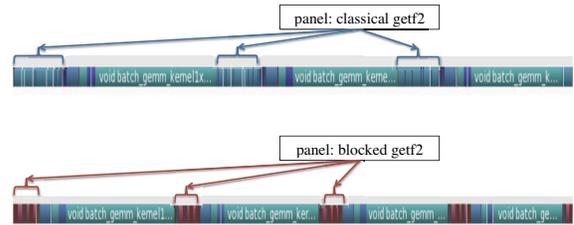


Fig. 3. Execution trace of the batched LU factorization using either classical *getf2* (top) or our recursive *getf2* (bottom).

C. Streamed dgemm

Our main goal is to achieve higher performance and to accomplish this we performed deep analysis of every kernel of the algorithm. We found that 70% of the time is spent in the batched *dgemm* kernel. An evaluation of the performance of the *dgemm* kernel using either batched or streamed *dgemm* is illustrated in Figure 4. The curves let us conclude that the streamed *dgemm* was performing better than the batched one for some cases, e.g., for $k = 32$ when the matrix size is of order of $m > 200$ and $n > 200$. We note that the performance of the batched *dgemm* is stable and does not dependent on k , in the sense that the difference in performance between $k = 32$ and $k = 128$ is minor. However it is bound by 300 Gflop/s. For that we proposed to use the streamed *dgemm* whenever it is faster, and to roll back to the batched one otherwise. The use of the streamed *dgemm* (when the size allows it) can speed up the factorization by about 20% and this is confirmed by the performance curve plotted in Figure 5.

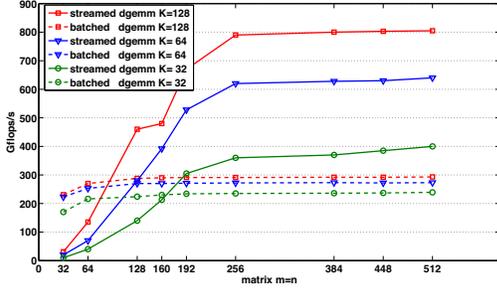


Fig. 4. Performance comparison between the streamed and the batched `dgemm` kernel for different value of K and different matrix sizes where $m = n$.

D. Multi-level blocking of the update

The performance of the streamed `dgemm` kernel as shown in Figure 4 is highly dependent on the size of the matrices. In particular, this affects the trailing matrix updates in the LU factorization which consist of rank- k operations. The performance of the streamed `dgemm` kernel is around twice higher for $k = 128$ than for $k = 32$. Since our panel size is limited to 32, the performance of the trailing matrix update is limited by the performance of the `dgemm` for $k = 32$. However, in order to achieve higher performance, we use multi-level of blocking of the trailing matrix update. The idea here is to use multi-level of blocking during the trailing matrix update. This means that at step i we only update the next panel and delay the subsequent portion of the update till step $i + l$ where we reach a value of k that is acceptable to perform the whole update of the delayed portion and then start over again.

We can observe that for this range of small matrices, increasing the value of acceptable “ k ” for example to 128 gives us the advantage of performing `dgemm` at higher speed but it reduce the number of such `dgemm` operations. The performance observed is similar for both $k = 64$ and $k = 128$ for matrices of size 512 while $k = 64$ is always outperforming $k = 128$ for smaller sizes. As a result a trade-off value of k need to be chosen depending on the matrix size. The improvement obtained by this technique is around 15%, as shown in Figure 5.

V. PERFORMANCE RESULTS

We conducted our experiments on a NVIDIA K40c card with 11.6 GB of GDDR memory per card running at 825 MHz. The cards were connected to the host via two PCIe I/O hubs with 6 GB/s bandwidth.

CUBLAS version 5.5 features a `dgetrfBatched` routine. By comparison, our batched LU is up to $2.5\times$ faster than the CUBLAS routine as shown in Figure 5. The slowest code in the figure has performance below 60 Gflop/s and is marked as “classic” – it corresponds to the performance of the MAGMA library, which was optimized for large matrices. The *classic* implementation is improved upon by CUBLAS’ `dgetrfBatched` version (marked as “CUBLAS” in Figure 5) and the performance exceeds 70 Gflop/s. To go beyond 100 Gflop/s,

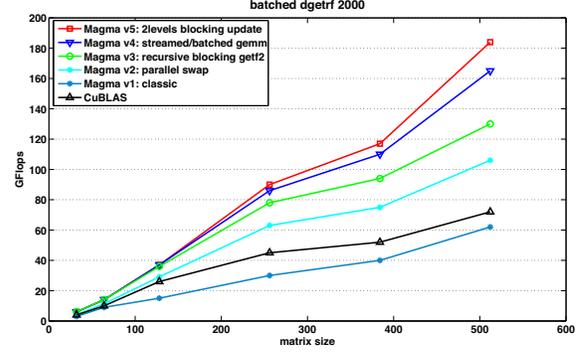


Fig. 5. Performance in Gflops/s of our different version of the batched LU factorization compared to the CUBLAS implementation.

we used the code that optimizes pivoting with *parallel swap*. Next step in performance improvement is the use of variable blocking (also called *recursive blocking getf2*), which enables performance levels that go slightly above 130 Gflop/s. The final two improvements are *streamed/batched gemm*, which moves the performance level beyond 160 Gflop/s, and finally, *2-levels blocking update* completes the set of optimizations and takes the performance beyond the 180 Gflop/s mark. Each of these optimizations is described in detail in Section III.

VI. CONCLUSIONS

The need to solve large number of small linear systems often arises in scientific computing applications. In contrast to large linear system which can expose data parallelism and can be efficiently implemented on either GPUs or CPUs, solving small linear systems is memory bound. This is due to the fact that the ratio of the computation to the data needed is very small compared to the one for large matrices. We demonstrated that GPU architectures can be used efficiently for solving many small size problems. In particular, we developed different algorithm variants and optimization techniques for the batched LU factorization on GPUs and analyzed their impacts on performance. These techniques can be used by other high level linear algebra solvers, for example, QR, Cholesky, as well. Our performance exceeded the CUBLAS `dgetrfBatched` by up to $2.5\times$.

REFERENCES

- [1] K. Gallivan, W. Jalby, and U. Meier, “The use of BLAS3 in linear algebra on a parallel processor with a hierarchical memory,” *SIAM J. Sci. Stat. Comp.*, vol. 8, 1987, 10791084.
- [2] O. Messer, J. Harris, S. Parete-Koon, and M. Chertkow, “Multicore and accelerator development for a leadership-class stellar astrophysics code,” in *Proceedings of "PARA 2012: State-of-the-Art in Scientific and Parallel Computing."*, 2012.
- [3] V. Oreste, N. A. Gawande, and A. Tumeo, “Accelerating subsurface transport simulation on heterogeneous clusters,” in *IEEE International Conference on Cluster Computing (CLUSTER 2013)*, Indianapolis, Indiana, September, 23-27 2013.
- [4] V. Volkov and J. W. Demmel, “LU, QR and Cholesky factorizations using vector capabilities of GPUs,” Tech. Rep. LAPACK Working Note 202.
- [5] D. Croz, J. J. Dongarra, and N. J. Higham, “Stability of methods for matrix inversion,” *IMA J. Numer. Anal.*, vol. 12, no. 119, 1992.