

# Enhancing the Performance of Dense Linear Algebra Solvers on GPUs in the MAGMA Project

Marc Baboulin

University of Coimbra (Portugal)

Jim Demmel

University of California, Berkeley (USA)

Jack Dongarra

University of Tennessee, Knoxville (USA)  
Oak Ridge National Laboratory (USA)  
University of Manchester (UK)

Stanimir Tomov

University of Tennessee, Knoxville (USA)

Vasily Volkov

University of California, Berkeley (USA)

## 1 INTRODUCTION TO DENSE LINEAR ALGEBRA (DLA) FOR GPUS

DLA algorithms, due to a high ratio of floating point calculations to data required, have been of high performance. Therefore, special purpose architectures have not been able to significantly accelerate them up until recently.

\*Fathallah et al. study SDDMM (in 2004) to conclude GPUs almost always outperform CPUs for 40-1000x<sup>3</sup> produced 12 Gflops in single precision, comparable with a 3 GHz Pentium 4

\*Gallopo et al. (in 2005) had similar results on LU 6.7 Gflops in single precision on an NVIDIA 7800, compared to 3.4 GHz Pentium 4 at the time! This has changed as CPUs move to multicore architectures with an exponentially growing gap between processor speed and memory land bandwidth shared between cores, while GPUs have consistently outpaced them both in performance and memory bandwidth.

First CUDA GPU results to significantly outperform CPUs on DLA started appearing at the beginning of 2008 (illustrated also on Figure 1 for the GEMM operation)

- In January 2008 Volkov and J. Demmel [1] reported on SDDMM kernel (among others) to significantly outperform the CUBLAS library (25 Gflops vs more than 100 Gflops in their implementation) and an LU factorization running at up to 140 Gflops in single precision arithmetic (SP)
- In March S. Tomov et al. [2] presented at PASC08 Cholesky factorization running at up to 180 Gflops in SP using Volkov's SDDMM kernel (also described in [3])
- In May Volkov and J. Demmel [4] described LU, QR, and Cholesky running at up to 180 Gflops in SP
- In May, Dongarra et al. [5] reported on SP Cholesky running at 320 Gflops on a pre-release NVIDIA card

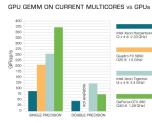


Figure 1: Performance GEMM on current multicore vs GPU. Note that in SP the GPU is 4x faster than the expensive processor. In DP the GPU is 10x faster than the expensive processor. The Intel Xeon X5405 is a 3.0 GHz Pentium 4.

## 2 GPU ALGORITHMS ≠ TRADITIONAL ALGORITHMS

### 2.1 EXTRACTING PARALLELISM

- Splitting algorithms into tasks
  - The concept of representing algorithms as Directed Acyclic Graphs (DAG) where the nodes represent the sub-tasks and the edges the dependencies.
  - Heterogeneity-aware splitting
  - Scheduling task execution
    - Crucial for performance, for example scheduling tasks on the critical path as soon as possible / trees more parallelism
- GPU triangular solvers through exploiting writing the triangular matrix
  - Significantly accelerates both TRM (up to 3 times) and TRS (order of magnitude)



### 2.2 HETEROGENEITY-AWARE ALGORITHMS

- In particular, algorithms for hybrid GPU + multicore computing should split the computation to fully exploit the power that each of the hybrid components offers.
- Small tasks of low parallelism to be executed on the CPU (for example tasks on the critical path)
  - Bigger tasks of high parallelism to be executed on the GPU
  - Proper scheduling should explore asynchronous behavior on CPU and GPU
  - Blocking strategies
    - Varying block sizes (as in QR LR)
    - Two-level blocking (as in Cholesky LR)
- Note: see [4] for other heterogeneity-related techniques and detail on 1, 4
- Work partitioning especially for hybrid GPU + Multicore IR



EXAMPLE  
Two-level blocking with writing for Single GPU + 8-core CPU host  
The matrix is split into two parts: the top part is processed by the CPU and the bottom part is processed by the GPU. The CPU part is further split into two parts: the top part is processed by the CPU and the bottom part is processed by the GPU.

### 2.3 INNOVATIVE DATA STRUCTURES

Non-traditional data layouts may be beneficial in reducing communication costs. For example, to avoid severely parallelized shared memory access in pivoting on the GPU, the matrix is laid out in the GPU memory in row-major order. This reduces the pivoting overhead from 50% of the total computation to only 1% on NVIDIA G7800 in SP

### 2.4 PRECONDITIONING FOR REDUCED PIVOTING

- Here Preconditioning is a Random Butterfly Transformation (RBT) with cost of applying negligible compared to the factorization, and meant to transform the original matrix into a "sufficiently random" matrix  $A'$  with probability close to 1 pivoting would not be needed [3]
  - We use unitary RBT (does not change the condition number) can be represented as  $2 \times 2$  block matrix where the blocks are diagonal matrices
  - To solve  $Ax = b$  using butterfly matrices  $U$  and  $V$  we solve  $U^T A' V^T = U^T b$ , followed by  $x = V^{-1} U^{-T} U^T b$
- RBT helps but in general accuracy is reduced. Two ways to improve it if needed are GSR
  - Add-linked pivoting (LP) within the block size (determined by IR) or more for example LP-LNB+IR is a LU with pivoting limited to the last  $NB-IR$  rows of the current panel
  - Add derivative refinement in the working precision

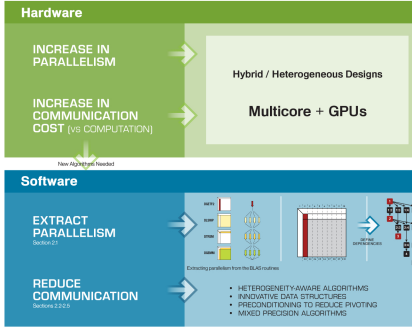
### 2.5 MIXED PRECISION ALGORITHMS

Accelerating solvers using mixed precision arithmetic:



## HARDWARE TO SOFTWARE TRENDS

Note: For the trends in Multicore architectures, we refer to [6], [7], [8] and [9].

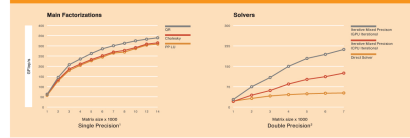


## Matrix Algebra on GPU and Multicore Architectures (MAGMA)

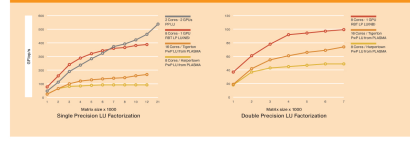
The MAGMA project, headed by the linear algebra research groups at University of Tennessee, UC Berkeley, and UC Denver, aims to develop a dense linear algebra library similar to LAPACK but for heterogeneous/hybrid architectures, starting with current Multicore+GPU systems. This transition cannot be done automatically as in many cases new algorithms that significantly differ from algorithms for conventional architectures, will be needed. Preliminary studies – on a new class of heterogeneity-aware algorithms of 'reduced communication' and 'high-parallelism, as shown in this poster – confirm that this is the case.

## 3 PERFORMANCE RESULTS

### Single Core + Single GPU



### Multicores / Multi-GPUs



## HARDWARE USED

GPU: GeForce GTX 280 (300 Cores @ 1.32 GHz)  
Host: Intel Xeon X5405 (3.0 GHz)  
Together: Intel Xeon X5405 (3.0 GHz) + 4 Cores @ 2.33 GHz  
Tightest: Intel Xeon X5405 (3.0 GHz) + 4 Cores @ 2.33 GHz + 2 GPUs  
Tallest: Intel Xeon X5405 (3.0 GHz) + 4 Cores @ 2.33 GHz

Note that GPU runs at Higher S/Mip rate than Cholesky. Tested the best fixed level of parallelism in GEMM, as a result with varying matrices.

Intel precision solvers often achieved a speedup compared to SP solvers for the speed depends on the context. Intel precision solvers often achieved a speedup compared to SP solvers for the speed depends on the context. Intel precision solvers often achieved a speedup compared to SP solvers for the speed depends on the context.

## CONCLUSIONS

- GPU computing has reached a point to significantly outperform current multicore on DLA in spite of DLA's traditionally high performance on x86 architectures.
- Architecture trends have moved towards heterogeneous (GPU + CPU) designs of increased parallelism and communication costs, and software trends have to reflect that. We addressed this with innovative heterogeneity-aware algorithms/techniques on extracting parallelism and reducing communication.
- There are significant differences between the new algorithms and those for conventional CPUs.
- The new techniques in many cases present an opportunity for trade-off between speed and accuracy.
- The need for DLA for hybrid systems will grow

## MOTIVATING OUR FUTURE WORK DIRECTIONS, AS ENVISIONED IN THE MAGMA PROJECT, TOWARDS A SELF CONTAINED DLA LIBRARY SIMILAR TO LAPACK BUT FOR HETEROGENEOUS ARCHITECTURES.

## REFERENCES

- Volkov and J. Demmel, Using GPUs to accelerate linear algebra routines, Poster at PASC08 winter meeting, January 9, 2008, <http://www.eecs.berkeley.edu/~volkov/pasc08/pasc08.pdf>
- Tomov, M. Baboulin, J. Dongarra, S. Moore, V. Nalati, G. Peterson, and D. Riche, Special purpose hardware and algorithms for accelerating dense linear algebra, Presentation at PASC, Atlanta, March 14-18, 2008, [http://www.cs.cmu.edu/~tomov/PASC\\_Tomov.pdf](http://www.cs.cmu.edu/~tomov/PASC_Tomov.pdf)
- M. Baboulin, J. Dongarra, and S. Tomov, Some issues in dense linear algebra for multicore and special purpose architectures, LAPACK Working Note 200.
- Volkov and J. Demmel, LU, QR and Cholesky factorizations using vector capabilities of GPUs, Tech. Report UOB/ECS-0008-48, ECS Department, University of California, Berkeley, May 2008.
- Dongarra, S. Moore, G. Peterson, S. Tomov, V. Nalati, and D. Riche, Exploring new architectures in accelerating CFD for Air Force applications, HPCMP User Group Conference 2008, July 14-17, 2008, <http://www.usc.edu/~dongarra/hpcmp2008.html>
- Tomov, J. Dongarra, and M. Baboulin, Towards dense linear algebra for hybrid GPU accelerated manycore systems, LAPACK Working Note 210.
- J. Langou, P. Lusignea, J. Kurzak, A. Buttari, and J. Dongarra, Exploiting the performance of 3D floating point arithmetic in obtaining 64-bit accuracy (avoiding iteration refinement) for linear systems, sci 2008, 50.
- Buttari, J. Langou, J. Kurzak, and J. Dongarra, A class of parallel fixed linear algebra algorithms for multicore architectures, LAPACK Working Note 191.
- Demmel, D. G. Li, H. Xiang, Communication-avoiding Gauss-Jordan Elimination, Supercomputing 08.
- Demmel, J. Dongarra, M. Haemmen, J. Langou, Communication Optimal Parallel and Sequential QR and LU Factorizations, SIAM J. Sci. Comp. 10 September, UC Berkeley, Tech. Report EEC-08-09.
- Demmel, J. Dongarra, M. Haemmen, J. Langou, Implementing Communication Optimal Parallel and Sequential QR and LU Factorizations, SMM J. vol 11, Comp. Submitted.