

# 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)

Stanimire 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

- Fatahalian et al. study SGEMM (in 2004) to conclude CPUs almost always outperform GPUs (only ATI X800XT produced 12 Gflop/s in single precision, comparable with a 3 GHz Pentium 4)

- Galoppo et al. (in 2005) had similar results on LU (5.7 Gflop/s in single precision on an NVIDIA 7800, compared to 3.4GHz Pentium 4 at the time)

This has changed as CPUs move to multi/manycores with an exponentially growing gap between processor speed and memory (and 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 V.Volkov and J. Demmel [1] reported on SGEMM kernel (among others) to significantly outperform the CUBLAS library (125 Gflop/s vs more than 180 Gflop/s in their implementation) and an LU factorization running at up to 140 Gflop/s in single precision arithmetic (SP)

- In March S. Tomov et al. [2] presented at PPSC08 Cholesky factorization running at up to 160 Gflop/s in SP using Volkov's SGEMM kernel (also described in [3])

- In May V.Volkov and J. Demmel [4] described LU, QR, and Cholesky running at up to 180 Gflop/s in SP

- In May, Dongarra et al. [5] reported on SP Cholesky running at 325 Gflop/s on a pre-release NVIDIA card

GPU GEMM ON CURRENT MULTICORES vs GPUs

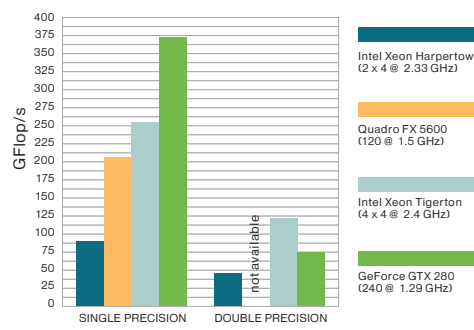


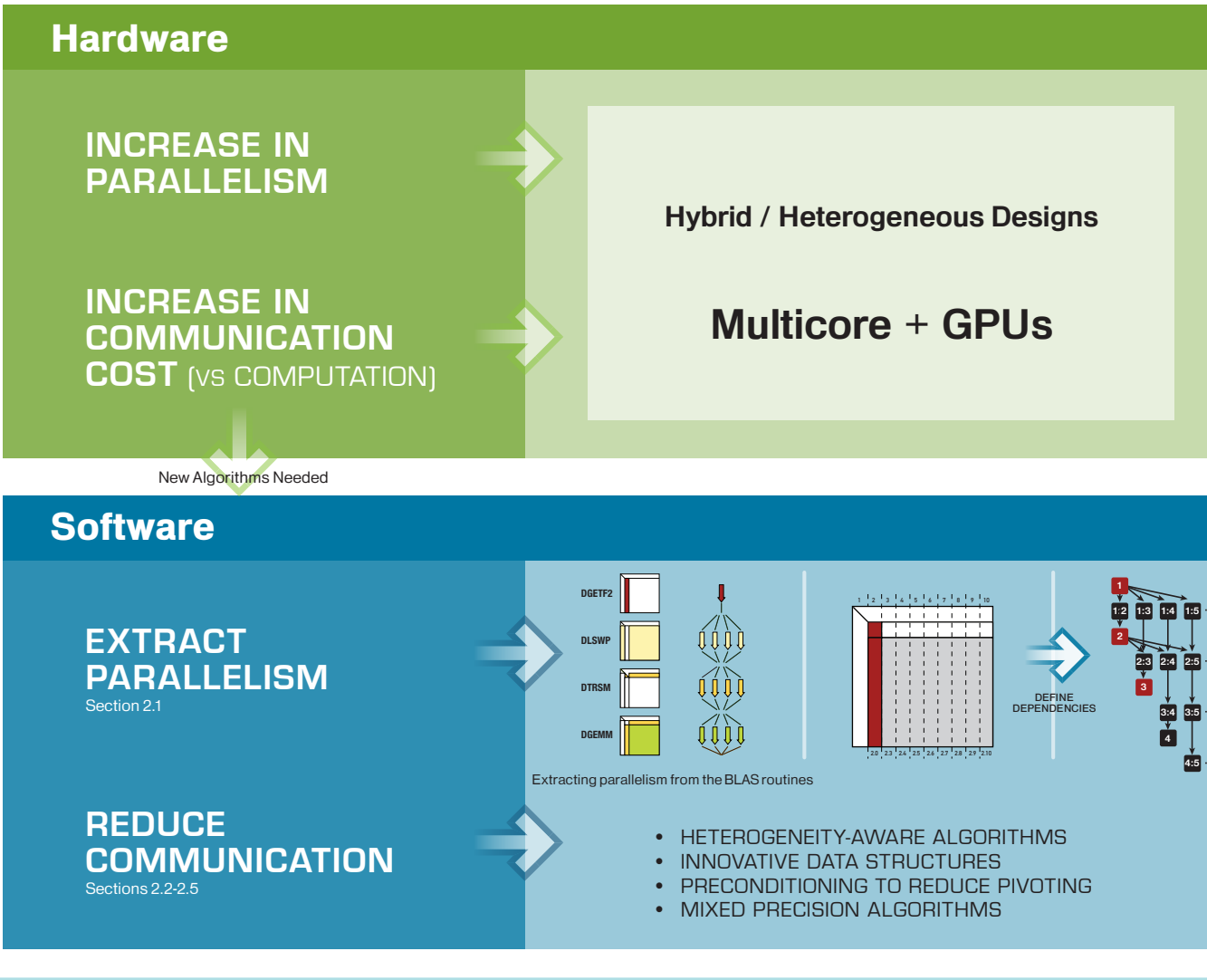
Figure 1

Peak measured GEMM performance on current multicore (from Intel) and GPU (from NVIDIA) architectures.

Note that in SP the GTX 280 is 10 x faster than a quad-core processor (at 2.33GHz) and still 75 GFlop/s faster than an entire quad-socket quad-core Intel Xeon Tigerton system (cores running at 2.4 GHz)

## HARDWARE TO SOFTWARE TRENDS

Note: For the trends in Multicore architectures, we refer to the PLASMA project [6] and to new algorithms minimizing communication [8, 10, 11]



## 2 GPU ALGORITHMS ≠ TRADITIONAL ALGORITHMS

### 2.1 EXTRACTING PARALLELISM

1. Splitting Algorithms into tasks

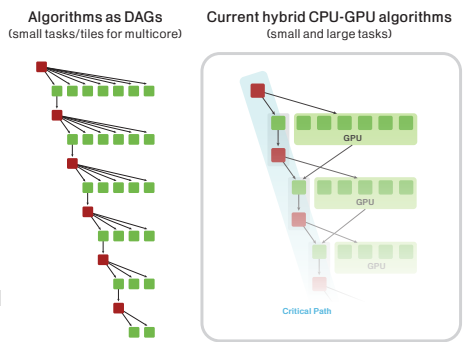
- The concept of representing algorithms as Directed Acyclic Graphs (DAGs) where the nodes represent the sub-tasks and the edges the dependencies
- Heterogeneity-aware splitting

2. Scheduling task execution

- Crucial for performance, for example scheduling tasks on the critical path 'as soon as possible' frees more parallelism

3. GPU triangular solvers through explicitly inverting the triangular matrix

- Significantly accelerates both TRSM (up to 3 times) and TRSV (order of magnitude)



### 2.2 HETEROGENEITY-AWARE ALGORITHMS

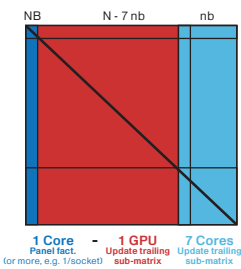
In particular, algorithms for hybrid GPU + multicore computing should split the computation to fully exploit the power that each of the hybrid components offers.

- '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 asynchronicity between CPU and GPU
- Blocking strategies

- Varying block sizes (as in QR [4])
- Two-level blocking (as in Cholesky [4])

Note: see [4] for other heterogeneity-related techniques and detail on 1..4

- Work partitioning (specific) for hybrid GPU + Multicore [6]



**EXAMPLE**  
An LU factorization work splitting for Single GPU + 8 cores CPU host  
The first N - nb columns reside on the GPU and are processed by 1 GPU + 1 core, the rest resides and is processed by the remaining cores

The algorithm: (e.g. in solving  $Ax = b$ ,  $A$  is SPD)

```
1 Factorization  $A = LL^T$  (GP)
2 Solve  $Ly = b$  (GP)
3 Solve  $L^T x = y$  (GP)
do  $k = 1, 2, \dots$ 
4  $r = b - Ax$  (GP)
5 Solve  $Ly = r$  (GP)
6 Solve  $L^T z = y$  (GP)
7  $x = x + z$  (CP)
check convergence
done
```

This technique is feasible because it uses much faster SP arithmetic for the bulk of the computation and DP iterations to raise the accuracy (if needed).

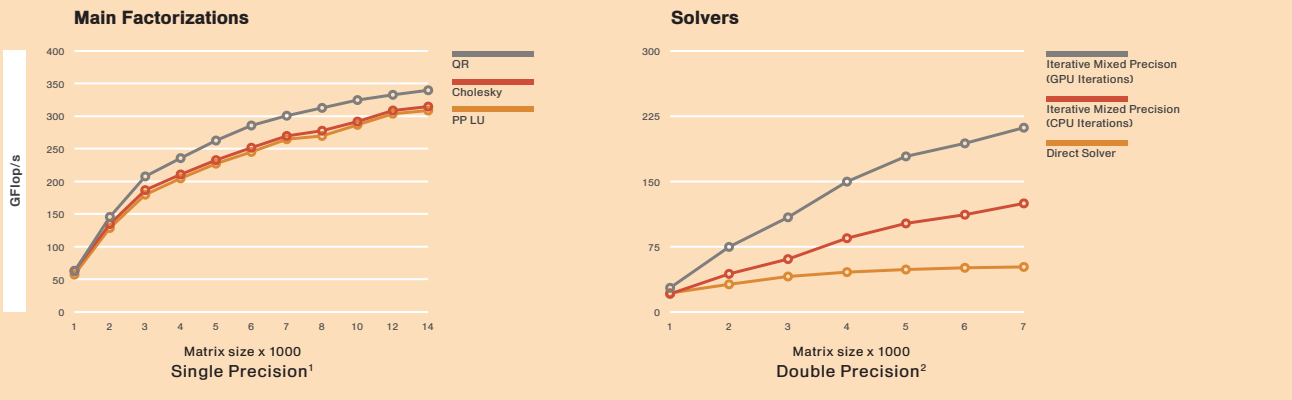
Ratios between SP and DP on the GTX 280 are much higher than those on traditional CPUs:  
Theoretical peak: ~ 10  
SGEMM/DSGEMM: ~ 5  
One sided factorizations: ~ 4 (to 5)  
often leading to speedup factors of 4!

## 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



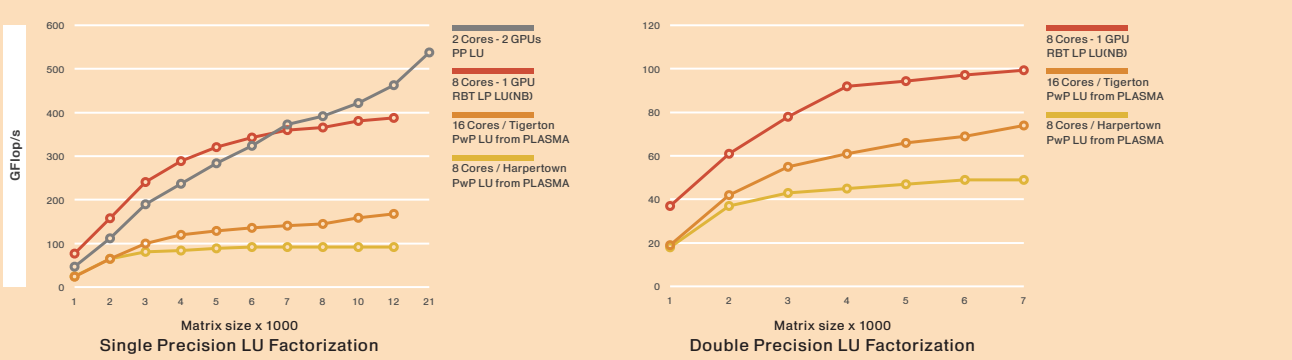
### HARDWARE USED

GPU: GeForce GTX 280 (240 Cores @ 1.30 GHz)  
Host: Intel Xeon (2 x 4 Cores @ 2.33 GHz)  
Tigerton: Intel Xeon (4 x 4 Cores @ 2.4 GHz)  
Harpertown: Intel Xeon (2 x 4 Cores @ 2.33 GHz)

<sup>1</sup> Note that QR runs at a higher MFlop rate than Cholesky: Cholesky has less thread-level parallelism in GEMM, as it deals with triangular matrices.

<sup>2</sup> Mixed precision solvers often achieve 4 x speedup compared to DP solvers but the speed depends on the conditioning of the matrix. In these performance results, we considered 3 steps of iterative refinement (on symmetric and positive definite matrices using Cholesky).

### Multicores / Multi-GPUs



## CONCLUSIONS

- GPU computing has reached a point to significantly outperform current multicores 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

- V. Volkov and J. Demmel, Using GPUs to accelerate linear algebra routines, Poster at PAR lab winter retreat, January 9, 2008, [http://www.eecs.berkeley.edu/~volkov/volkov08\\_parlab.pdf](http://www.eecs.berkeley.edu/~volkov/volkov08_parlab.pdf).
- S. Tomov, M. Baboulin, J. Dongarra, S. Moore, V. Natoli, G. Peterson, and D. Richie, Special-purpose hardware and algorithms for accelerating dense linear algebra, Presentation at PPSC, Atlanta, March 12-14, 2008, [http://www.cs.utk.edu/~tomov/PP08\\_Tomov.pdf](http://www.cs.utk.edu/~tomov/PP08_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.
- V. Volkov and J. Demmel, LU, QR and Cholesky factorizations using vector capabilities of GPUs, Tech. Report UCB/EECS-2008-49, EECS Department, University of California, Berkeley, May 2008.
- J. Dongarra, S. Moore, G. Peterson, S. Tomov, J. Allred, V. Natoli, and D. Richie, Exploring new architectures in accelerating CFD for Air Force applications, HPCMP User Group Conference 2008 (July 14-17, 2008), [http://www.cs.utk.edu/~tomov/ugc2008\\_final.pdf](http://www.cs.utk.edu/~tomov/ugc2008_final.pdf)
- S. Tomov, J. Dongarra, and M. Baboulin, Towards dense linear algebra for hybrid GPU accelerated manycore systems, LAPACK Working Note 191.
- J. Langou, J. Langou, P. Luszczyk, J. Kurzak, A. Buttari, and J. Dongarra, Exploiting the performance of 32 bit floating point arithmetic in obtaining 64 bit accuracy (revisiting iterative refinement for linear systems), sc 0 (2006), 50.
- A. Buttari, J. Langou, J. Kurzak, and J. Dongarra, A class of parallel tiled linear algebra algorithms for multicore architectures, LAPACK Working Note 191.
- J. Demmel, L. Grigori, H. Xiang, Communication-avoiding Gaussian Elimination, Supercomputing '08.
- J. Demmel, L. Grigori, M. Hoemmen, J. Langou, Communication Optimal Parallel and Sequential QR and LU Factorizations, SIAM J. Sci. Comp (submitted), UC Berkeley Tech Report EECS-2008-89.
- J. Demmel, L. Grigori, M. Hoemmen, J. Langou, Implementing Communication Optimal Parallel and Sequential QR and LU Factorizations, SIAM J. Sci Comp (submitted).