# Not Your Father's Math Library
# *MAGMA for Dense Matrix Problems*

*Matrix Algebra on GPU and Multicore Architectures (MAGMA)*
*Parallel Linear Algebra for Scalable Multi-core Architectures (PLASMA)*

**Jack Dongarra & Stan Tomov**

**University of Tennessee**

**Oak Ridge National Laboratory**

# Major Changes to Algorithms/Software

- **Must rethink the design of our algorithms and software**
  - **Manycore and Hybrid architectures are disruptive technology**
    - Similar to what happened with cluster computing and message passing
  - **Rethink and rewrite the applications, algorithms, and software**

  - **Data movement is expensive**
  - **Flops are cheap**

# Challenges for Software/Libraries

1. **Synchronization**
   - **Break Fork-Join model**

2. **Communication**
   - **Use methods which have lower bound on communication**

3. **Mixed precision methods**
   - **SP:DP; 2x speed of ops and 2x speed for data movement**

4. **Autotuning**
   - **Today's machines are too complicated, build "smarts" into software to adapt to the hardware**
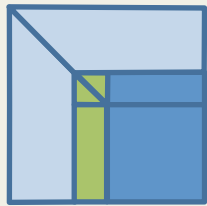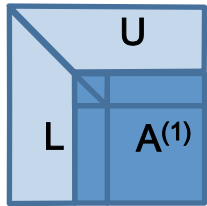
5. **Fault resilient algorithms**
   - **Implement algorithms that can recover from failures/bit flips**
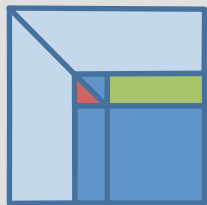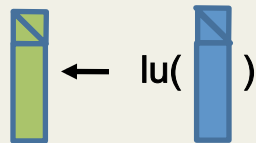
# Fork-Join Parallelization of LU and QR.

**Parallelize the update:**

- Easy and done in any reasonable software.
- This is the $2/3n^3$ term in the FLOPs count.
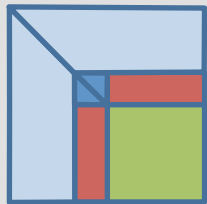- Can be done efficiently with LAPACK+multithreaded BLAS
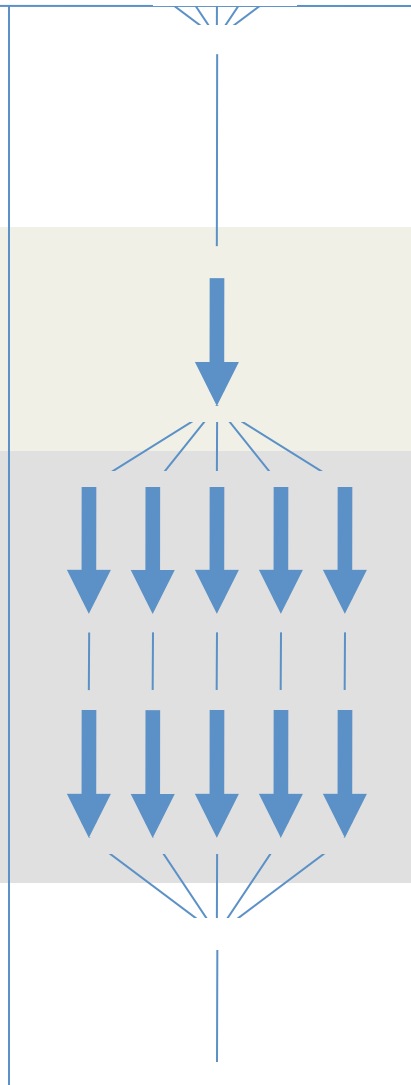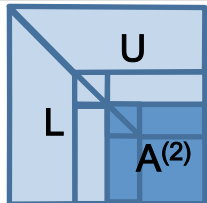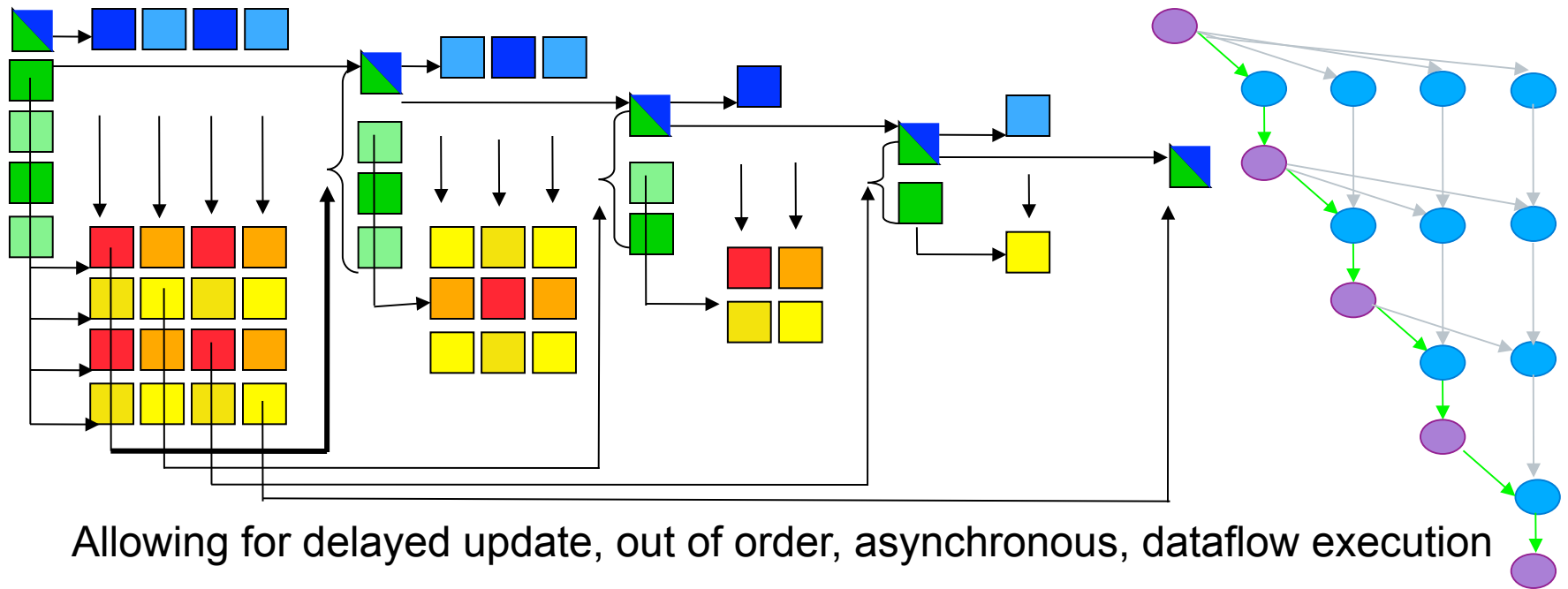
dgemm

U

L $A^{(1)}$

dgetf2

$\leftarrow$ lu( )

dtrsm (+ dswp)

$\leftarrow$ \

dgemm

$\leftarrow$ -

U

L $A^{(2)}$

Step 1 → Step 2 → Step 3 → Step 4 · · ·

Allowing for delayed update, out of order, asynchronous, dataflow execution

# PLASMA/MAGMA: Parallel Linear Algebra s/w for Multicore/Hybrid Architectures

- **Objectives**
  - **High utilization of each core**
  - **Scaling to large number of cores**
  - **Synchronization reducing algorithms**

- **Methodology**
  - **Dynamic DAG scheduling (QUARK)**
  - **Explicit parallelism**
  - **Implicit communication**
  - **Fine granularity / block data layout**

- **Arbitrary DAG with dynamic scheduling**

Cholesky
4 x 4

Fork-join parallelism

DAG scheduled parallelism

Time

# Pipelining: Cholesky Inversion
## 3 Steps: Factor, Invert L, Multiply L's

POTRF

TRTRI

LAUUM

POTRI

48 cores
POTRF, TRTRI and LAUUM.
The matrix is 4000 x 4000, tile size is 200 x 200,

POTRF+TRTRI+LAUUM: 25 (7t-3)
Cholesky Factorization alone: 3t-2

Pipelined: 18 (3t+6)

# Hybrid Algorithms

## A methodology to use all available resources:

- **MAGMA uses HYBRIDIZATION methodology based on**
  - Representing linear algebra algorithms as collections of TASKS and DATA DEPENDENCIES among them
  - Properly SCHEDULING tasks' execution over multicore and GPU hardware components

- **Successfully applied to fundamental linear algebra algorithms**
  - One and two-sided factorizations and solvers
  - Iterative linear and eigen-solvers

- **Productivity**
  - High-level
  - Leveraging prior developments
  - Exceeding in performance homogeneous solutions

Hybrid CPU+GPU algorithms
(small tasks for multicores and large tasks for GPUs)



GPU

GPU

GPU

Critical Path

# MAGMA Performance (single GPU)



MAGMA LU in double precision on single GPU (C2050)

1,090 MFlop/W*

GPU MAGMA

55 MFlop/W*

CPU LAPACK

Gflop/s — Matrix Size: 1024, 2048, 3072, 4032, 5184, 6016, 7040, 8064, 9088, 10112

GPU Fermi C2050 (448 CUDA Cores @ 1.15 GHz)
+ Intel Q9300 (4 cores @ 2.50 GHz)
DP peak **515** + **40** GFlop/s
Power * ~**220** W

CPU AMD Istanbul
[ 8 sockets x 6 cores (48 cores) @2.8GHz ]
DP peak **538** GFlop/s
Power * ~**1,022** W

* Computation consumed power rate (total system rate minus idle rate), measured with KILL A WATT PS, Model P430

# MAGMA Performance (scaling)



MAGMA LU in double precision on multi-GPUs (Fermi C2070)

Keeneland system, using one node
3 NVIDIA GPUs (M2070 @ 1.1 GHz, 5.4 GB)
2 x 6 Intel Cores (X5660 @ 2.8 GHz, 23 GB)

10

# Cholesky Factorization (DP)

- Weak scalability on many nodes (NSF Keeneland system; node: 2-Intel 6 core & 3-Nvidia M2070)
- Input size: 34560, 46080, 69120, 92160, 138240, 184320, 276480, 460800



(node: 2-Intel 6 core & 3-Nvidia M2070)

# The Need for HP Linear Algebra

## Electronic structure calculations

- **Density functional theory**

  Many-body Schrödinger equation (exact but exponential scaling)

  $$\{-\sum_i \frac{1}{2}\nabla_i^2 + \sum_{i,j}\frac{1}{|r_i - r_j|} + \sum_{i,l}\frac{Z}{|r_i - R_l|}\}\Psi(r_1,..r_N) = E\Psi(r_1,..r_N)$$

  · Nuclei fixed, generating external potential (system dependent, non-trivial)
  · N is number of electrons

  **Kohn Sham Equation: The many body problem of interacting electrons is reduced to non-interacting electrons (single particle problem) with the same electron density and a different effective potential (cubic scaling).**

  $$\{-\frac{1}{2}\nabla^2 + \int\frac{\rho(r')}{|r - r'|}dr' + \sum_l\frac{Z}{|r - R_l|} + V_{xc}\}\psi_i(r) = E_i\psi_i(r)$$

  $$\rho(r) = \sum_i |\psi_i(r)|^2 = |\Psi(r_1,..r_N)|^2$$

  · $V_{xc}$ represents effects of the Coulomb interactions between electrons

  · $\rho$ is the density (of the original many-body system)

  $V_{xc}$ is not known except special cases ⇒ use approximation, e.g. Local Density Approximation (LDA) where $V_{xc}$ depends only on $\rho$

- A model leading to self-consistent iteration computation with need for HP LA (e.g, **diagonalization** and **orthogonalization)**

# Generalized Hermitian Eigen-Problem
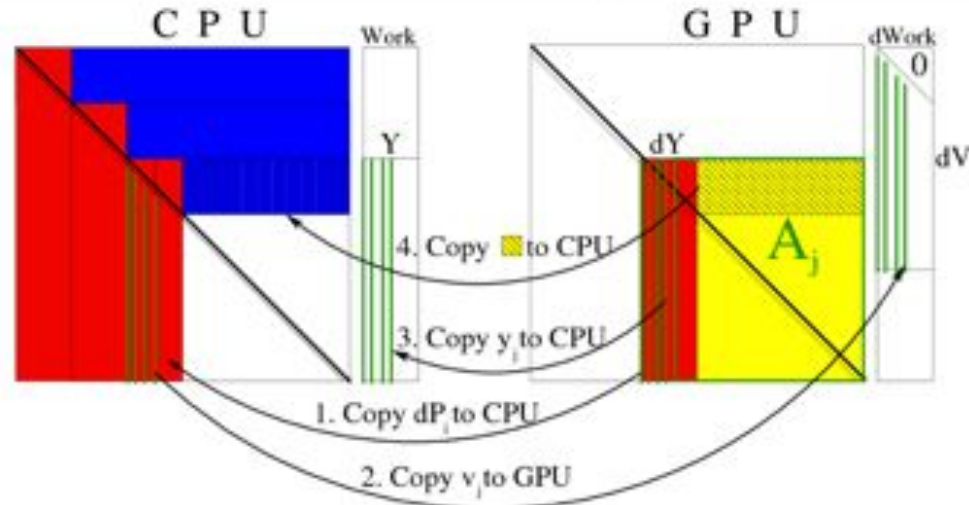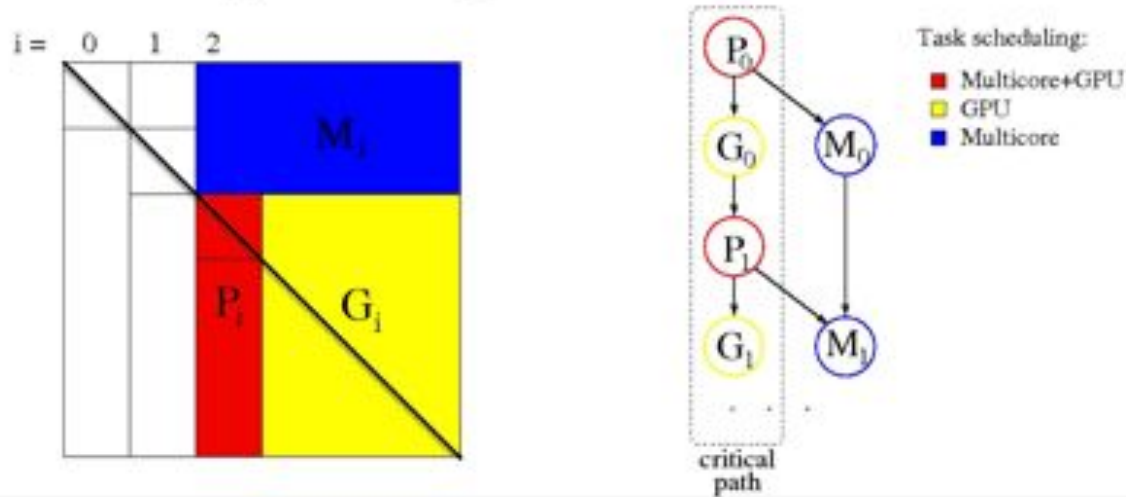
## $A x = \lambda B x$

- **Compute Cholesky factorization of**
  $B = LL^H$

- **Transform the problem to a standard eigenvalue problem**
  $\tilde{A} = L^{-1}AL^{-H}$

- **Solve Hermitian standard Eigenvalue problem**
  $\tilde{A}y = \lambda y$

- **Transform back the eigenvectors**
  $x = L^{-H} y$

Task **Splitting** & Task **Scheduling**

# Performance Comparison of Generalized Eigenproblem in Double Complex Precision



time to solution matrix size 10000, all eigenvectors

Legend:
- gen. e.p.
- cholesky dec.
- transf to st. e.p.
- st. e.p.
- tridiagonalization
- eigensolver
- 1st back transf. e.v.
- 2nd back transf. e.v.

Time shown, so lower is better

X-axis: MAGMA 1 gpu, MAGMA 2 gpus, ELPA 1 st. 1 node, ELPA 1 st. 2 nodes, ELPA 1 st. 3 nodes, ELPA 1 st. 4 nodes

Y-axis: time (s)

- **Test system:**
  **2 x Intel X5650 (6 core), 2 x Nvidia M2090**

Thomas Schulthess
Raffaele Solcà

ETH
Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

ELPA solver from:

Parallel solution of partial symmetric eigenvalue problems from electronic structure calculations ☆

T. Auckenthaler[a], V. Blum[b], H.-J. Bungartz[a], T. Huckle[a], R. Johanni[c], L. Krämer[d], B. Lang[d,*], H. Lederer[c], P.R. Willems[d]

[a] Fakultät für Informatik, Technische Universität München, D-85748 Garching, Germany
[b] Fritz-Haber-Institut der Max-Planck-Gesellschaft, D-14195 Berlin, Germany
[c] Rechenzentrum Garching der Max-Planck-Gesellschaft am Max-Planck-Institut für Plasmaphysik, D-85748 Garching, Germany
[d] Fachbereich C, Bergische Universität Wuppertal, D-42097 Wuppertal, Germany

# Two-Stage Approach to Tridiagonal Form (Communication Reducing)



- ## Reduction to band
  - **On multicore + GPUs**
  - **Performance as in the one-sided factorizations [derived from fast Level 3 BLAS]**

- ## Band to tridiagonal
  - **Leads to "irregular" (bulge chasing) computation**
    - **Done very efficiently on multicore !**
  - **GPUs are used to assemble the orthogonal Q from the transformations [needed to find the eigenvectors]**

# Performance results

**Tridiagonalization in double precision on Fermi**



- ◆ **Communication reducing**

- ◆ **Developed routines for multiGPUs obtaining scalable performance**

- ◆ **The new algorithm (2 stages approach) on a Keeneneland node bring a speedup of ~ 10 X**

**Keeneland system, using one node**
3 NVIDIA GPUs (M2070@ 1.1 GHz, 5.4 GB)
2 x 6 Intel Cores (X5660 @ 2.8 GHz, 23 GB)

# Mixed Precision Methods

- **Mixed precision, use the lowest precision required to achieve a given accuracy outcome**
  - **Improves runtime, reduce power consumption, lower data movement**
  - **Reformulate to find correction to solution, rather than solution [ Δx rather than x ].**

# Mixed Precision Solvers

MAGMA LU-based solvers on Fermi (C2050)



**FERMI**     Tesla C2050: 448 CUDA cores @ 1.15GHz
SP/DP peak is 1030 / 515 GFlop/s

- **Direct solvers**
  - Factor and solve in working precision
- **Mixed Precision Iterative Refinement**
  - Factor in single (i.e. the bulk of the computation in fast arithmetic) and use it as preconditioner in simple double precision iteration, e.g.

$$x_{i+1} = x_i + (LU_{SP})^{-1} P (b - A x_i)$$

- Similar results for Cholesky & QR

## COMPUTATIONAL ROUTINES IN MAGMA 1.1

| MATRIX | | OPERATION | ROUTINE | CPU | GPU |
|---|---|---|---|---|---|
| **LINEAR EQUATIONS** | GE | LU | {sdcz}getrf | ✓ | ✓ |
| | GE | Solve | {sdcz}getrs | | ✓ |
| | GE | Invert | {sdcz}getri | | ✓ |
| | SPD/HPD | Cholesky | {sdcz}potrf | ✓ | ✓ |
| | SPD/HPD | Solve | {sdcz}potrs | | ✓ |
| | SPD/HPD | Invert | {sdcz}potri | | ✓ |
| | TR | Invert | {sdcz}trtri | ✓ | |
| **ORTHOGONAL FACTORIZATIONS** | GE | QR | {sdcz}geqrf | ✓ | |
| | GE | Generate Q | {sd}orgqr | ✓ | ✓ |
| | GE | | {cz}ungqr | ✓ | ✓ |
| | GE | Multiply matrix by Q | {sd}ormqr | ✓ | ✓ |
| | GE | | {cz}unmqr | ✓ | ✓ |
| | GE | LQ factorization | {sdcz}gelqf | ✓ | ✓ |
| | GE | QL factorization | {sdcz}geqlf | ✓ | |
| | GE | Multiply matrix by Q | {sd}ormql | ✓ | |
| | GE | | {cz}unmql | ✓ | ✓ |
| **STANDARD EVP** | GE | Hessenberg reduction | {sdcz}gehrd | ✓ | |
| | GE | Generate Q | {sd}orghr | ✓ | |
| | GE | | {cz}unghr | ✓ | |
| | SY/HE | Tridiagonalization | {sd}sytrd | ✓ | |
| | SY/HE | | {cz}hetrd | ✓ | |
| | SY/HE | Generate Q | {sd}orgtr | | ✓ |
| | SY/HE | | {cz}ungtr | | ✓ |
| | SY/HE | Multiply by Q | {sd}ormtr | ✓ | ✓ |
| | SY/HE | | {cz}unmtr | ✓ | ✓ |
| **SVD** | GE | Bidiagonalization | {sdcz}gebrd | ✓ | |
| **GENERALIZED EVP** | SPD/HPD | Reduction to standard form | {sd}sygst | ✓ | ✓ |
| | SPD/HPD | | {cz}hegst | ✓ | ✓ |

## DRIVER ROUTINES IN MAGMA 1.1

| MATRIX | | OPERATION | ROUTINE | CPU | GPU |
|---|---|---|---|---|---|
| **LINEAR EQUATIONS** | GE | Solve using LU | {sdcz}gesv | ✓ | ✓ |
| | GE | Solve using MP | {zc,ds}gesv | | ✓ |
| | SPD/HPD | Solve using Cholesky | {sdcz}posv | ✓ | ✓ |
| | SPD/HPD | Solve using MP | {zc,ds}posv | | ✓ |
| **LLS** | GE | Solve LLS using QR | {sdcz}geqrs | | ✓ |
| | GE | Solve using MP | {zc,ds}geqrsv | | ✓ |
| **STANDARD EVP** | GE | Compute e-values, optionally e-vectors | {sdcz}geev | ✓ | |
| | SY/HE | Computes all e-values, optionally e-vectors | {sd}syevd | ✓ | ✓ |
| | SY/HE | | {cz}heevd | ✓ | ✓ |
| | SY/HE | Range ( D&C ) | {cz}heevdx | | ✓ |
| | SY/HE | Range ( B & I lt. ) | {cz}heevx | ✓ | ✓ |
| | SY/HE | Range ( MRRR ) | {cz}heevr | ✓ | ✓ |
| **STAND. SVP** | GE | Compute SVD, optionally s-vectors | {sdcz}gesvd | ✓ | ✓ |
| **GENERALIZED EVP** | SPD/HPD | Compute all e-values, optionally e-vectors | {sd}sygvd | ✓ | |
| | SPD/HPD | | {cz}hegvd | ✓ | |
| | SPD/HPD | Range ( D&C ) | {cz}hegvdx | ✓ | |
| | SPD/HPD | Range ( B & I lt. ) | {cz}hegvx | ✓ | |
| | SPD/HPD | Range ( MRRR ) | {cz}hegvr | ✓ | |

# Collaborators / Support

- **MAGMA team**
  **http://icl.cs.utk.edu/magma/**

- **PLASMA team**
  **http://icl.cs.utk.edu/plasma**

- **DAGuE team**
  **http://icl.cs.utk.edu/dague/**

- **Collaborating partners**

  **University of Tennessee, Knoxville**
  **University of California, Berkeley**
  **University of Colorado, Denver**

  **INRIA, France**
  **KAUST, Saudi Arabia**