



# Matrix Algebra on GPU and Multicore Architectures (MAGMA)

## *PART II: Hands-on training*

**Stan Tomov**

Innovative Computing Laboratory  
Department of Computer Science  
University of Tennessee, Knoxville

**CSCS-USI Tutorial on GPU-enabled numerical libraries**  
**Lugano, Switzerland**  
September 14-15, 2013

# Outline

- **Install MAGMA**
- **Using MAGMA**
  - Documentation, naming conventions, and functionality
- **Solving a linear system of equations**  
 $Ax = b$
- **Solving an eigenvalue problem**  
 $Ax = \lambda x$
- **DGEMM example**
- **Writing a hybrid algorithm**

# Install MAGMA

- **Get MAGMA 1.4**
  - > wget <http://icl.cs.utk.edu/projectsfiles/magma/downloads/magma-1.4.0.tar.gz>
- **Unpack the library**
  - > tar zxvf magma-1.4.0.tar.gz
- **Prerequisites: LAPACK, BLAS, and CUDA toolkit**
  - > module load cudatoolkit intel
- **Modify 'make.inc' (where are CUDA & LAPACK, and for what GPU)**

See examples in make.inc.{mkl-gcc | mkl-icc | mkl-ilp64 | mkl-shared | acml | atlas | goto | shared},






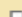




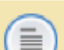
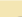


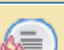
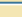



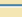
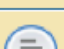

  - > cp make.inc.mkl-gcc make.inc

and modify make.inc by setting GPU\_TARGET, MKLROOT, and CUDADIR, e.g.,

```
GPU_TARGET ?= Kepler
MKLROOT      = $(INTEL_PATH)/mkl
CUDADIR      = $(CRAY_CUDATOOLKIT_DIR)
```
- **Create libmagma.a in 'lib' and testing drivers in 'testing'**
  - > make
- **For more information on installation, read file 'README'**

# Using MAGMA

- Support is provided through the MAGMA user forum  
<http://icl.cs.utk.edu/magma/forum/viewforum.php?f=2>

	<b>NaN errors with dpotrf and dpotrf_gpu</b> by fletchjp » Tue Dec 28, 2010 7:08 pm	9	3095	by fletchjp  Thu Sep 12, 2013 12:03 pm
	<b>GMRES on magma</b> by nitinb60 » Sun Aug 12, 2012 8:15 pm	4	1064	by fletchjp  Thu Sep 12, 2013 11:02 am
	<b>Help please with choice of forums</b> by fletchjp » Sat Apr 13, 2013 2:44 pm	4	1628	by fletchjp  Thu Sep 12, 2013 10:42 am
	<b>the error when I compile magma1.4.0 in vs2010</b> by Eva Joo » Mon Sep 09, 2013 4:57 am	2	62	by Eva Joo  Thu Sep 12, 2013 4:54 am
	<b>Undefined reference to cuda functions within libmagma</b> by Matt Phillips » Mon Sep 09, 2013 10:40 pm	1	63	by Matt Phillips  Mon Sep 09, 2013 11:17 pm
	<b>crash testing strsv, using OpenBLAS</b> by Matt Phillips » Thu Sep 05, 2013 9:15 pm	0	101	by Matt Phillips  Thu Sep 05, 2013 9:15 pm
	<b>MAGMA Installation: CLAPACK reference BLAS problem</b> by psrivas2 » Tue Sep 03, 2013 4:07 am	0	84	by psrivas2  Tue Sep 03, 2013 4:07 am
	<b>Running MAGMA across several GPUs on several nodes</b> by hsahasra » Tue Aug 13, 2013 1:32 pm	1	310	by mgates3  Fri Aug 30, 2013 3:48 pm
	<b>magma_init/finalize missing in fortran interface</b> by stachon » Wed Aug 28, 2013 4:02 am	1	116	by mgates3  Fri Aug 30, 2013 2:13 pm
	<b>Error: BLAS/LAPACK routine 'magma_' gave error code -7</b> by christianHEL » Thu Aug 22, 2013 2:37 pm	2	197	by christianHEL  Fri Aug 23, 2013 3:50 pm
	<b>Problems testing dsyevd with magma-1.4.0</b> by dougrabson » Thu Aug 15, 2013 5:44 am	1	173	by mgates3  Wed Aug 21, 2013 2:35 pm

# Using MAGMA

- Doxygen documentation  
<http://icl.cs.utk.edu/magma/docs/>



## MAGMA magma-1.4.0

Matrix Algebra on GPU and Multicore Architectures

Main Page	Modules	Namespaces	Data Structures	Files
<div>File List</div> <div><ul style="list-style-type: none"><li>▶ dgeqrn_gpu.cpp</li><li>▶ dgeqrf_mgpu.cpp</li><li>▶ dgeqrf_ooc.cpp</li><li>▶ dgeqrs3_gpu.cpp</li><li>▶ dgeqrs_gpu.cpp</li><li>▶ dgesm_gpu.cpp</li><li>▶ dgesv.cpp</li><li>▶ dgesv_gpu.cpp</li><li>▶ dgesvd.cpp</li><li>▶ dgetf2_nopiv.cpp</li><li>▶ dgetrf.cpp</li><li>▶ dgetrf2_mgpu.cpp</li><li>▶ <b>dgetrf_gpu.cpp</b></li><li>▶ dgetrf_incpiv_gpu.cpp</li><li>▶ dgetrf_m.cpp</li></ul></div>				

```
23
24 Purpose
25 =====
26 DGETRF computes an LU factorization of a general M-by-N matrix A
27 using partial pivoting with row interchanges.
28
29 The factorization has the form
30   A = P * L * U
31 where P is a permutation matrix, L is lower triangular with unit
32 diagonal elements (lower trapezoidal if m > n), and U is upper
33 triangular (upper trapezoidal if m < n).
34
35 This is the right-looking Level 3 BLAS version of the algorithm.
36 If the current stream is NULL, this version replaces it with user defined
37 stream to overlap computation with communication.
38
39 Arguments
40 =====
41 M      (input) INTEGER
42        The number of rows of the matrix A.  M >= 0.
43
44 N      (input) INTEGER
45        The number of columns of the matrix A.  N >= 0.
46
47 A      (input/output) DOUBLE_PRECISION array on the GPU, dimension
```

# Using MAGMA

- **Naming conventions (e.g., magma\_dgesv\_gpu )**
  - Prefix **magma\_** or **magmablas\_**
  - Followed by precision  
**s** single, **d** double, **c** single complex, or **z** double complex  
**ds** mixed double-single, or **zc** double complex-single complex
  - Matrix type  
**general**      **symmetric**      **hermetian**      **positive definite**  
**orthogonal**   **unitary**          **triangular**
  - Operation  
**sv**    solve                      **ev**    eigenvalue problem  
**trf**    triangular factorization   **gv**    generalized eigenvalue
  - Suffix **\_gpu** if input and output are on the GPU memory

# Using MAGMA

- MAGMA functionality**

<http://icl.cs.utk.edu/graphics/posters/files/SC12-MAGMA.pdf>

**GE** – General  
**SPD/HPD** – Symmetric/Hermitian Positive Definite  
**TR** – Triangular  
**D & C** – Divide & Conquer  
**B & I IT** – Bisection & Inverse Iteration  
**MP** – Mixed-precision Iterative Refinement  
**Naming Convention:** magma\_{routine name}\_{gpu}

## DRIVER ROUTINES IN MAGMA 1.3

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

## COMPUTATIONAL ROUTINES IN MAGMA 1.3

		MATRIX	OPERATION	ROUTINE	INTERFACES	
					CPU	GPU
LINEAR EQUATIONS	GE		LU	{sdcz}getrf	✓	✓
			Solve	{sdcz}getrs		✓
			Invert	{sdcz}getri		✓
	SPD/HPD		Cholesky	{sdcz}potrf	✓	✓
			Solve	{sdcz}potrs		✓
			Invert	{sdcz}potri		✓
ORTHOGONAL FACTORIZATIONS	TR		Invert	{sdcz}trtri	✓	
			QR	{sdcz}geqrf	✓	
	GE		QR w/ pivoting	{sdcz}geqp3		✓
			Generate Q	{sd}orgqr	✓	✓
				{cz}ungqr	✓	✓
			Multiply matrix by Q	{sd}ormqr	✓	✓
				{cz}unmqr	✓	✓
			LQ factorization	{sdcz}gelqf	✓	✓
			QL factorization	{sdcz}geqlf	✓	
			Multiply matrix by Q	{sd}ormql	✓	✓
STANDARD EVP	GE			{cz}unmql	✓	✓
			Hessenberg reduction	{sdcz}gehrd	✓	
			Generate Q	{sd}orghr	✓	
				{cz}unghr	✓	
			Tridiagonalization	{sd}sytrd	✓	
				{cz}hetrd	✓	
			Generate Q	{sd}orgtr		✓
	SY/HE			{cz}ungtr		✓
			Multiply by Q	{sd}ormtr	✓	✓
				{cz}unmtr	✓	✓
SVD	GE		Bidiagonalization	{sdcz}gebdd	✓	
			Reduction to standard form	{sd}sygst {cz}hegst	✓	✓

# Examples / Exercises

- Solving a linear system of equations  
 $Ax = b$

```
#include "magma.h"
#include "magma_lapack.h"
...
```

```
magma_init();
cublasInit();
```

```
...
double *hA, *hB;
magma_malloc_pinned((void **)&hA, lda * N * sizeof(double));
magma_malloc_cpu( (void **)&hB, nrhs* M *sizeof(double));
```

```
...
init_matrix(M, N, hA, lda);
```

```
...
magma_dgetrf( M, N, hA, lda, ipiv, &info);
lapackf77_dgetrs( MagmaNoTransStr, &N, &nrhs,
                  hA, &lda, ipiv, hB, &ldb, &info );
```

```
// A is a typical array on the CPU
// A can be allocated in pinned memory
```

```
// LU factorization (using CPU+GPU)
// Solve on CPU with LAPACK
```

**Alternatively, there is a direct MAGMA function to solve:**  
see **testing\_dgesv.cpp** and **testing\_dgetrf.cpp**



# Examples / Exercises

- Solving an eigenvalue problem  
 $Ax = \lambda x$

```
#include "magma.h"
#include "magma_lapack.h"
...

magma_init();
cublasInit();
...
double *hA;
magma_malloc_pinned((void **)&hA, lda * N * sizeof(double)); // A is a typical array on the CPU
                                                                // A can be allocated in pinned memory
...
init_matrix(M, N, hA, lda);
...
magma_dsyeval( opts.jobz, opts.uplo,
               N, hA, lda, w1,
               h_work, lwork,
               iwork, liwork, &info );
```

see `testing_dsyeval.cpp`

# Examples / Exercises

- **DGEMM example**

```
salloc -N1 --gres=gpu:1 --time=00:30:00
module load cudatoolkit intel
export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:/apps/dom/intel/composer_xe_2013/mkl/lib/intel64/
export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:/apps/opcode/CUDA-5.0/lib64
export LD_LIBRARY_PATH=$LD_LIBRARY_PATH:/apps/dom/intel/composer_xe_2013.1.117/compiler/lib/intel64/
```

```
aprun -n1 -N1 -e OMP_NUM_THREADS=16 -d16 ./testing_dgemm -l -c
```

MAGMA 1.4.0 , capability 3.0

device 0: Tesla K20X, 732.0 MHz clock, 5759.6 MB memory, capability 3.5

Usage: ./testing\_dgemm [options] [-h|--help]

transA = N, transB = N

M	N	K	MAGMA Gflop/s (ms)	CUBLAS Gflop/s (ms)	CPU Gflop/s (ms)	MAGMA error	CUBLAS error
1088	1088	1088	612.88 ( 4.20)	969.12 ( 2.66)	3.13 ( 823.44)	6.92e-15	6.92e-15
2112	2112	2112	667.85 ( 28.21)	1100.09 ( 17.13)	6.79 (2775.89)	8.12e-15	7.95e-15
3136	3136	3136	675.16 ( 91.36)	1144.59 ( 53.89)	11.14 (5538.38)	1.13e-14	1.13e-14
...							

see file **testing\_dgemm.cpp**, **sgemm.pdf**, **dgemm\_fermi.cu**, and directory **testing**

# Examples / Exercises

- **Writing a hybrid algorithm**

Develop a hybrid CPU-GPU algorithm for this Matlab code

```
function [Q,R] = chol_qr_it(A)
    i=0;
    cn = 200;
    Q = A;
    G = Q'*Q;
    n = size(A,2);
    R = eye(n);

    while cn > 100, i = i + 1
        [u,s,v]=svd(G);
        [q,r]=qr(sqrt(s)*v');
        R = r * R;
        cn = sqrt(cond(s));
        Q = Q * inv(r);
        if cn>100
            G = Q'*Q;
        end;
    end;
    return
```

G =

Q'

Q

Computations on small  
data to be done on the  
CPU

Computations on large  
data to be done on the  
GPU

# Examples / Exercises

- **Writing a hybrid algorithm**

Develop a hybrid CPU-GPU algorithm for this Matlab code

```
function [Q,R] = chol_qr_it(A)
    i=0;
    cn = 200;
    Q = A;
    G = Q'*Q;
    n = size(A,2);
    R = eye(n);

    while cn > 100, i = i + 1
        [u,s,v]=svd(G);
        [q,r]=qr(sqrt(s)*v');
        R = r * R;
        cn = sqrt(cond(s));
        Q = Q * inv(r);
        if cn>100
            G = Q'*Q;
        end;
    end;
    return
```

**cublasSgemm( ... );      // G = Q' \* Q (GPU computation)**

**cublasGetVector( ... );      // G -> work (send the result G to the CPU)**  
**sgesvd\_( ... );      // svd G on the CPU**  
**sgeqrf\_( ... );      // QR on the CPU**

**cublasSetVector( ... );      // send r back to the the GPU)**  
**cublasStrsm( ... );      // Triangular matrix solve on the GPU**

see **testing\_dgeqqr\_gpu.cpp** and **dgeqqr\_gpu.cpp**

# Contributions

The MAGMA program style follows the general guidelines for Sca/LAPACK in terms of interfaces, copyrights and licensing, citing the authors of the software, and documentation:

<http://www.netlib.org/lapack-dev/lapack-coding/program-style.html>

- Routine Naming and Design
- Language, source formatting, timing, and testing
- See file '**ContributorsGuige.txt**'

# Collaborators and Support



## MAGMA team

<http://icl.cs.utk.edu/magma>

## PLASMA team

<http://icl.cs.utk.edu/plasma>



## Collaborating partners

University of Tennessee, Knoxville

University of California, Berkeley

University of Colorado, Denver

INRIA, France (StarPU team)

KAUST, Saudi Arabia



U.S. DEPARTMENT OF  
**ENERGY**

