

GTC 2012:

NEW ADVANCES IN GPU LINEAR ALGEBRA

Kyle Spagnoli

EM Photonics

5/16/2012



QUICK “ABOUT US”



EM Photonics

Accelerated Computing Solutions

- » HPC/GPU Consulting Firm
- » Specializations in:
 - » Electromagnetics
 - » Image Processing
 - » Fluid Dynamics
 - » Linear Algebra



INTRODUCTION TO OUR LIBRARIES

» CULA Dense

» Linear algebra routines

» CULA Sparse

» Iterative sparse system solvers
and preconditioners

» pCULA

» Scalable solvers for multiple GPUs

» Ongoing work



CULA | dense



CULA | sparse





INTRODUCTION - COMMON POINTS

- » Easy to use
 - » No GPU programming experience necessary
 - » `dgetrf(...)` → `culaDgetrf(...)`
- » Exhaustively tested and benchmarked
 - » Accuracy & stability first!
- » Cross platform
 - » Linux, Windows, Mac OS X
- » Multiple languages
 - » C/C++, Fortran, Python, Matlab



EM Photonics

Accelerated Computing Solutions

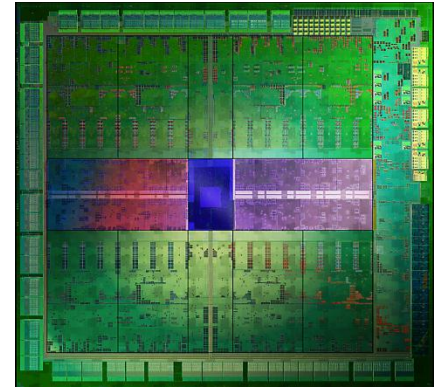
CULA DENSE



CULA DENSE – INTRODUCTION

- » First released in 2009
- » LAPACK and BLAS implementations
 - » Host or device memory
 - » Almost 300 routines
- » Upcoming release (R15)
 - » Tuned for Kepler architecture
 - » Now free for personal academic use

CULA | dense



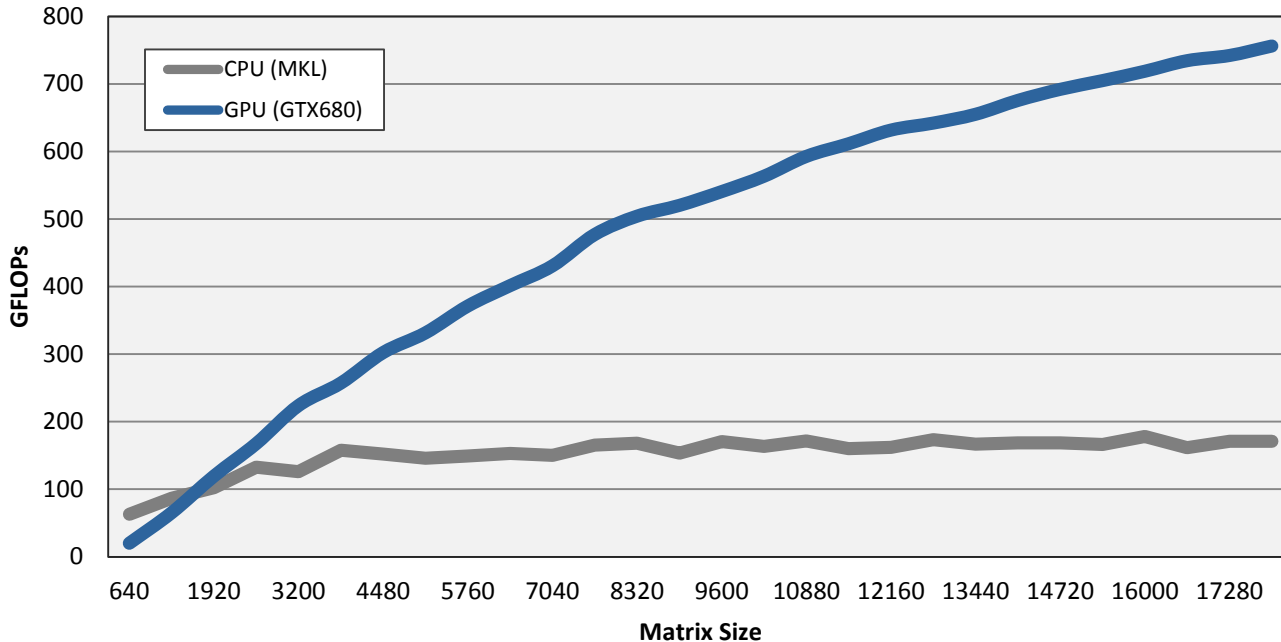


CULA DENSE - FUNCTIONALITY

LAPACK		BLAS
LU factorization	Cholesky factorization	Matrix-matrix multiply
QR decomposition	Orthogonal factorization	Matrix-vector multiply
Least squares	System solve	Rank updates
Eigenvalue routines	Matrix inversion	Conjugate
Singular value decomposition	Auxiliary routines	Transpose

CULA DENSE - PERFORMANCE

CULA Dense - Cholesky Factorization (SPOTRF)



Performance numbers *include* transfer time across PCI-Express (Gen2) bus

CPU

Intel Core i7 2600K

GPU

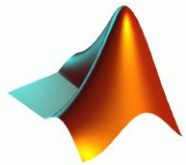
NVIDIA GTX 680 (1.5 GB)





CULA DENSE – LINK INTERFACE

- » GPU acceleration with **no** code changes!
 - » Intercepts calls to BLAS & LAPACK libraries
 - » Analyze routine, parameters, and hardware
 - » Forward to GPU if appropriate
 - » Pass-through to CPU otherwise



```
SET LAPACK_VERSION = cula_link.dll  
SET BLAS_VERSION = cula_link.dll
```



EM Photonics

Accelerated Computing Solutions

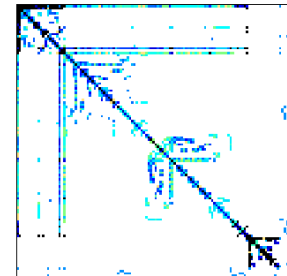
CULA SPARSE (ITERATIVE)



CULA SPARSE – INTRODUCTION

- » First released in 2011
- » Iterative solvers and preconditioners
- » Multiple matrix storage formats supported
- » Upcoming release (S3)
 - » Tuned for Kepler
 - » Free for personal academic use

CULA | sparse



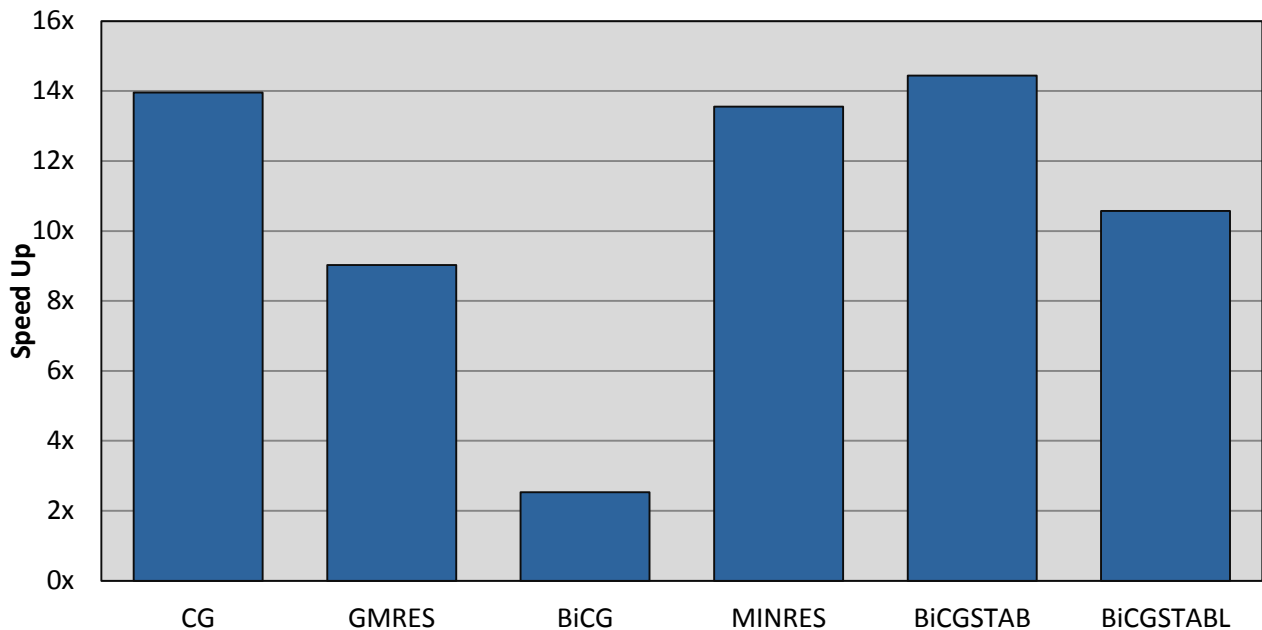


CULA SPARSE - FUNCTIONALITY

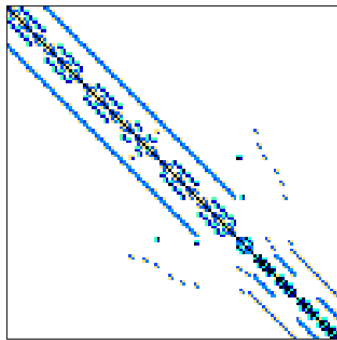
Solvers	Preconditioners	Data
CG	Jacobi	Double / Complex
BiCG	Block Jacobi	CSR / CSC / COO
BiCG-Stab / (L)	ILU0	
GMRES	Reordered ILU0	
MINRES		

CULA SPARSE - PERFORMANCE

Iterative Solver Performance

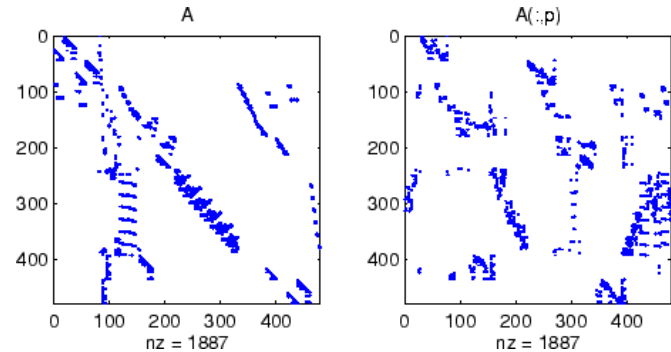


System Size = 1.5M
GPU = NVIDIA C2070
CPU = Xeon X5560 (MKL)



CULA SPARSE – PERFORMANCE FEATURES

- » Hybrid performance
 - » CPU begins working during initial transfer
 - » Preconditioner generation
 - » Initial iterations
- » Matrix reordering
 - » Can increase parallelism





EM Photonics

Accelerated Computing Solutions

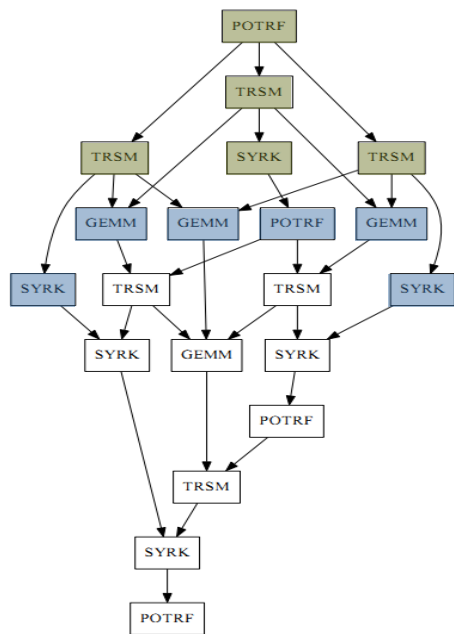
PCULA – MULTI-GPU + CPU PERFORMANCE

PCULA – INTRODUCTION

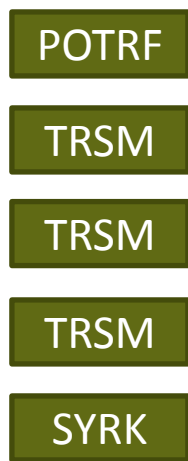
- » Scale to multiple GPUs and CPUs in a single node
- » Currently in alpha release
- » Greatly increased performance, scalability, and functionality coming soon!



PCULA – TASK SCHEDULING



Completed
Tasks



Pending
Valid Tasks



Hardware



Busy

(3 tasks)



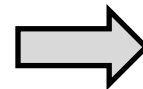
Free

(1 task)



Free

(0 tasks)



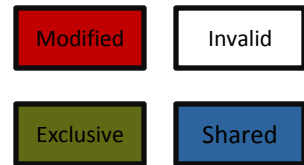
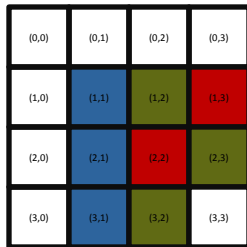


pCULA – HETEROGENEOUS TASK SCHEDULING

- » Data locality is critical
- » Hardware performance
 - » Persistent “live tuning” performance database
- » Task queue depth
 - » Too long → idle hardware if not perfect
 - » Too short → worker starvation

PCULA – OUT OF (GPU) CORE

- » Solve problems larger than GPU memory
 - » Natural extension of tiled data partitioning
 - » MESI memory coherence protocol
 - » Least recently used replacement strategy



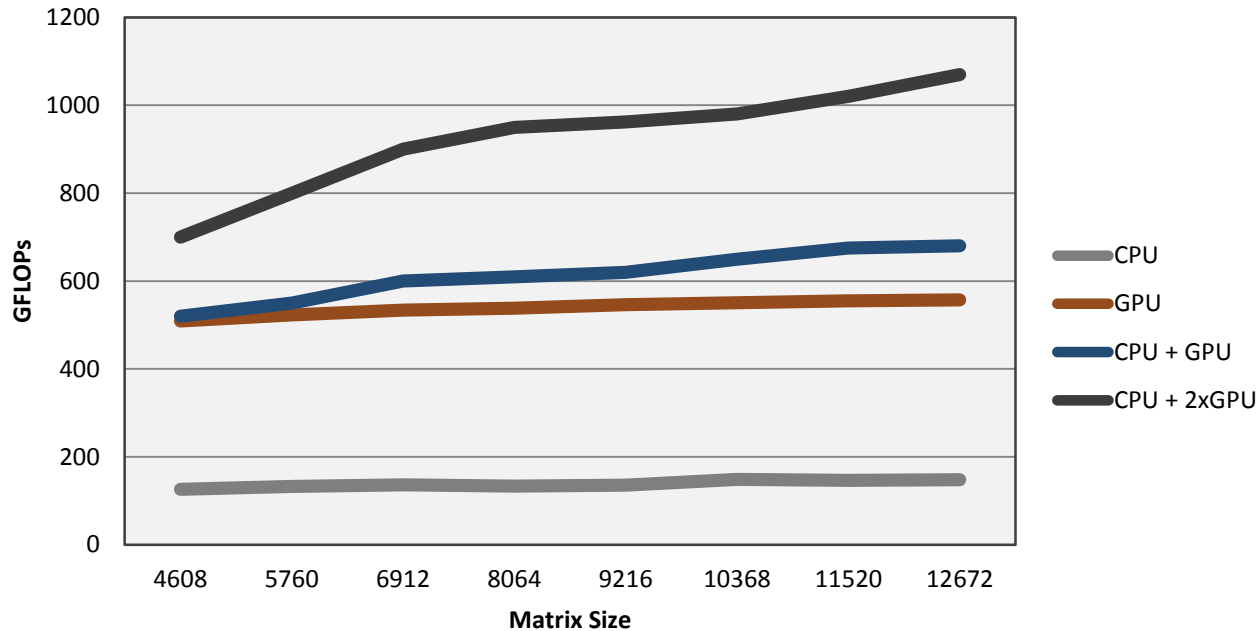


PCULA – FUNCTION LIST

- » Currently supports
 - » BLAS Routines (GEMM, TRSM, GEMV)
 - » LU Factorization & Solve (GETRF + GESV)
 - » Cholesky Factorization & Solve (POTRF + POSV)
 - » QR Factorization & Solve (GEQRF + GEQRS)
- » Eigenvalue and SVD routines in future release

pCULA - PERFORMANCE

pCULA - DGEMM Performance



Performance numbers
include transfer time across
PCI-Express (Gen2) bus

CPU
Intel Xeon 5560

GPU
2x NVIDIA C2050



EM Photonics

Accelerated Computing Solutions

ONGOING WORK



ONGOING WORK - CULA

- » CULA Dense
 - » More routines/tuning
- » CULA Sparse
 - » Direct solvers
 - » Algebraic Multi-Grid (AMG)
- » pCULA
 - » Multi-node cluster support
 - » NUMA optimizations





ONGOING WORK – C++ AMP

- » Microsoft's C++ AMP library
 - » “ampblas” development project
 - » Linear algebra to C++ AMP ecosystem
 - » Multiple talks today and tomorrow
 - » C++ AMP Lounge





CULA PARTNERS & INTEGRATORS

» Here at GTC 2012....



PGI[®]





THANKS!

Thanks!

Questions?

- » Convention hall @ booth #20
- » More information @ www.culatools.com