

MAGMA



MATRIX ALGEBRA ON GPU AND MULTICORE ARCHITECTURES (MAGMA) is a collection of next-generation linear algebra libraries for heterogeneous architectures. MAGMA is designed and implemented by the team that developed LAPACK and ScaLAPACK, incorporating the latest developments in hybrid synchronization-avoiding and communication-avoiding algorithms, as well as dynamic runtime systems. Interfaces for the current LAPACK and BLAS standards are supported to enable computational scientists to seamlessly port any linear algebra-reliant software components to heterogeneous architectures. MAGMA allows applications to fully exploit the power of current heterogeneous systems of multi/many-core CPUs and multiple GPUs to deliver the fastest possible time to accurate solution within given energy constraints.

HYBRID ALGORITHMS

MAGMA uses a hybridization methodology, where algorithms of interest are split into tasks of varying granularity, and their execution is scheduled over the available hardware components. Scheduling can be static or dynamic. In either case, small non-parallelizable tasks, often on the critical path, are scheduled on the CPU, and larger more parallelizable ones, often Level-3 BLAS, are scheduled on the GPU. When CPU-GPU communication overhead becomes significant, the entire computation might run exclusively on the GPU, leading to what is termed as "GPU-native algorithms".

PERFORMANCE & ENERGY EFFICIENCY

MAGMA LU factorization in double-precision arithmetic

CPU Intel Xeon Platinum 8460Y+ (Sapphire Rapids) 2 x 40 cores @3.7GHz

AMD Instinct 120 CUs @ 1.5GHz

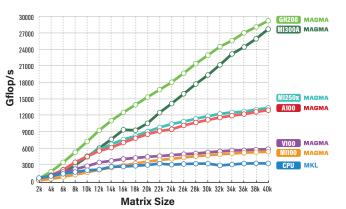
V100 **NVIDIA Volta**

NVIDIA Ampere GPU 108 SMs x 64 @ 1.41 GHz

AMD Instinct 110 CUs @ 1.7GHz

NVIDIA Grace-Hopper GH200 GPU 132 SMs @ 1.98GHz

MI300A AMD Instinct 228 CUs @ 2.1GHz





FEATURES AND SUPPORT

- MAGMA 2.9 for CUDA or HIP
- MAGMA Pre-release for SYCL/DPC++

CUDA	HIP	SYCL Pre-release	
Ø	$\mathbf{\mathscr{O}}$	\bigcirc	Linear system solvers
$\mathbf{\mathscr{O}}$	Ø	$\mathbf{\mathscr{O}}$	Eigenvalue problem solvers
Ø	Ø	Ø	Auxiliary BLAS
$\mathbf{\mathscr{O}}$	$\mathbf{\mathscr{O}}$	\odot	Batched LA
Ø	\mathbf{O}	\bigcirc	CPU/GPU Interface
$\mathbf{\mathscr{O}}$	Ø	$\mathbf{\mathscr{O}}$	Multiple precision support
Ø	\mathbf{O}		Mixed precision (including FP16)
Ø	Ø	\bigcirc	Non-GPU-resident factorizations
$\mathbf{\mathscr{O}}$	$\mathbf{\mathscr{O}}$	\odot	GPU-only factorizations
Ø	Ø		Multicore and multi-GPU support
$\mathbf{\mathscr{O}}$	Ø	$\mathbf{\mathscr{O}}$	LAPACK testing
Ø	$\mathbf{\mathscr{O}}$	\bigcirc	Linux
$\mathbf{\mathscr{O}}$	Ø		Windows
Ø	$\mathbf{\mathscr{O}}$		macOS

INDUSTRY COLLABORATION



Long-term collaboration and support on the development of MAGMA.



Intel Parallel Computing Center

The objective of the Innovative Computing Laboratory's IPCC is the development and optimization of numerical linear algebra libraries and technologies for applications, while tackling current challenges in heterogeneous Intel® Xeon Phi™ coprocessor-based High Performance Computing.



Long-term collaboration and support on the development of cIMAGMA, the OpenCL™ port of MAGMA, and hipMAGMA.

IN COLLABORATION WITH









SPONSORED BY







WITH SUPPORT FROM













MAGMA



MAGMA BATCHED

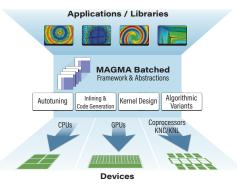
Batched factorization of a set of small matrices in parallel

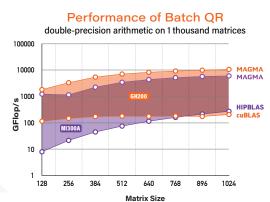
Numerous applications require factorization of many small matrices:

Sparse direct solvers Deep learning Structural mechanics High-order FEM simulations

Astrophysics

- LU, QR, and Cholesky
- Solvers and matrix inversion
- All BLAS 3 (fixed + variable)
- SYMV, GEMV (fixed + variable)





MAGMA 2.9 DRIVER ROUTINES

		MATRIX	OPERATION	ROUTINE	INTER CPU	FACES GPU
LIZED STANDARD EVP STANDARD EVP	<u>s</u>	GE	Solve using LU	{sdcz}gesv	1	1
	A N		Solve using MP	{zc,ds}gesv		/
	N S	SPD/HPD	Solve using Cholesky	{sdcz}posv	1	/
			Solve using MP	{zc,ds}posv		
	LEAST	GE GE	Solve LLS using QR	{sdcz}gels		/
	QUARE		Solve using MP	{zc,ds}geqrsv		
		GE	Compute e-values,	{sdcz}geev		
	₹		optionally e-vectors			
	ARD E	SY/HE	Computes all e-values,	{sd}syevd		/
			optionally e-vectors	{cz}heevd		
	Ā		Range (D&C)	{cz}heevdx		/
	ST		Range (B&I It.)	{cz}heevx		
			Range (MRRR)	{cz}heevr		/
	STAND.	0.5	Compute SVD,	{sdcz}gesvd		
	SVP	GE	optionally s-vectors	{sdcz}gesdd		
	EVP	SPD/HPD	Compute all e-values,	{sd}sygvd		
			optionally e-vectors	{cz}hegvd		
			Range (D&C)	{cz}hegvdx		
	Ë		Range (B&I It.)	{cz}hegvx	1	
	G		Range (MRRR)	{cz}hegvr	1	

Abbreviations

SPD/HPD

Naming Convention magma_{routine name}[_gpu]

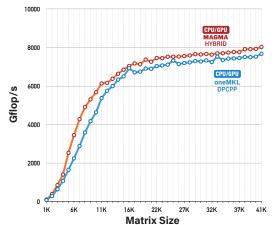
TR Triangular D&C Divide & Conquer

Symmetric/Hermitian Positive Definite

B&I It Bisection & Inverse Iteration MP Mixed-Precision Iterative Refinement

Intel GPU results on Aurora

Intel Xeon Max Series CPU + Intel Data Center Max Series GPU (codename PVC)



Results using Intel oneAPI release 2025.0.5. Aurora is a resource of the Argonne Leadership Computing Facility at Argonne National Laboratory

UPCOMING IN MAGMA 2.10

New functionality: Batch SVD

New functionality: Variable-size batch non-pivoting LU factorization

Performance Improvements for Batch Cholesky factorization and solve

Support for CUDA-13 and ROCM-7

IN COLLABORATION WITH













WITH SUPPORT FROM













SPONSORED BY

