

A Standard for **Batched BLAS Routines**

ABSTRACT

We propose an API for Batched Basic Linear Algebra Subprograms (Batched BLAS). We focus on multiple independent BLAS operations on small matrices that are grouped together as a single routine. We aim to provide a more efficient and portable library for multi/manycore HPC systems. We achieve 2x speedups and 3x better energy efficiency compared to vendor implementations. We also demonstrate the need for Batched BLAS and its potential impact in multiple application domains.

Batched BLAS: multiple independent BLAS operations on small DEFINITION matrices grouped together as a single routine





APPLICATIONS



Umeå

University

Batched BLAS:

Astrophysics

PROPOSED SPECIFICATION The function arguments are reminiscent of the BLAS standard.

Batched Level 3	dgemm_batch(enum	*transa,
BLAS DGEMM		enum	*transb,
Calling Sequence		integer	*m,
		integer	*n,
		integer	*k,
		double	*alpha,
		double	<pre>**arrayA,</pre>
		integer	*lda,
		double	<pre>**arrayB,</pre>
		integer	*ldb,
		double	*beta,
		double	<pre>**arrayC,</pre>
		integer	*ldc,
		integer	<pre>batch_count,</pre>
		enum	<pre>batch_opts,</pre>
		integer	<pre>*info);</pre>

Batched Level 2	dgemv_batch(enum	*trans,
BLAS DGEMV		integer	*m,
Calling Sequence		integer	*n,
		double	*alpha,
		double	**arrayA,
		integer	*lda,
		double	**X,
		integer	<pre>*incx,</pre>
		double	*beta,
		double	**y,
		integer	<pre>*incy,</pre>
		integer	<pre>batch_count,</pre>
		enum	<pre>batch_opts,</pre>
		integer	<pre>*info);</pre>

Batched Level 1	<pre>daxpy_batch(</pre>	integer	*n,
BLAS DAXPY		double	*alpha,
Calling Sequence		double	**X,
		integer	*incx,
		double	**y,
		integer	<pre>*incy,</pre>
		integer	<pre>batch_count,</pre>
		enum	<pre>batch_opts,</pre>
		integer	<pre>*info);</pre>

BATCHED BLAS PERFORMANCE

Batched Level 3 BLAS DGEMM Example

Batched Level 2 BLAS DGEMV Example

Batched DAXPY Example



APPLICATIONS OF BATCHED BLAS

Batched LAPACK DGETRF Example



Nuclear network simulation (XNet benchmark)

• 150 x 150 matrices

- batch_count = 100+
- Titan Supercomputer at ORNL
- 3x faster than MKL

SPEEDUP OF THE SOLVER AMD Opteron 6274 16-core CPUs NVIDIA Tesla K20X GPU FOR MATRIX SIZE 150



Astrophysical thermonuclear networks coupled to hydrodynamical simulations in explosive burning scenarios

• 7x faster using Batched BLAS

• batch_count = 10 to 800+



• 2x faster than MA48 factorization (HSL)



TECHNOLOGIES

Some of the technologies we may wish to utilize include:

OpenMP	CUDA	OpenCL
 Multicore 	 Fused Kernels 	
 Accelerators 	 Multiple Streams 	

ADVANTAGES

- More efficient and portable implementations
- HPC Numerical library for modern architectures
- Better hardware utilization and energy efficiency
- Encourages, as well as simplifies, community efforts to build higher-level algorithms on top of Batched BLAS

The specification is open for community discussion; we would welcome your comments: www.nlafet.eu

ACKNOWLEDGMENTS

This material is based upon work supported in part by the European Union's Horizon 2020 research and innovation programme under the NLAFET grant agreement No 671633, the U.S. National Science Foundation under Grants No. CSR 1514286 and ACI-1339822, NVIDIA, and the U.S. Department of Energy.

REFERENCES

- [1] A. Haidar, T. Dong, S. Tomov, P. Luszczek, and J. Dongarra, Framework for Batched and GPU- resident Factorization Algorithms Applied to Block Householder Transformations, ISC HPC, Springer LNCS, Frankfurt, Germany, July 12-16, 2015.
- [2] A. Haidar, S. Tomov, P. Luszczek, and J. Dongarra, MAGMA Embedded: Towards a Dense Linear Algebra Library for Energy Efficient Extreme Computing, 19th IEEE High Performance Extreme Computing Conference (HPEC 2015), Best Paper Award, IEEE, Waltham, MA, September, 2015.
- [3] A. Haidar, T. Dong, P. Luszczek, S. Tomov, and J. Dongarra. Batched matrix computations on hardware accelerators based on GPUs. International Journal of High Performance Computing Applications, first published on February 9, 2015 as doi:10.1177/1094342014567546.
- [4] J. Dongarra, I. Duff, M. Gates, A. Haidar, S. Hammarling, N. Higham, J. Hogg, P. Lara, M. Zounon, S. Relton, and S. Tomov, A Proposed API for Batched Basic Linear Algebra Subprograms, UTK Computer Science Technical Report, (available at https://bit.ly/batched-blas), March 2016.