



HAL
open science

Towards a High-Performance Tensor Algebra Package for Accelerators

Marc Baboulin, Veselin Dobrev, Jack Dongarra, Christopher Earl, Joel Falcou, Azzam Haidar, Ian Karlin, Tzanio Kolev, Ian Masliah, Stanimire Tomov

► **To cite this version:**

Marc Baboulin, Veselin Dobrev, Jack Dongarra, Christopher Earl, Joel Falcou, et al.. Towards a High-Performance Tensor Algebra Package for Accelerators. Smoky Mountains Computational Sciences and Engineering Conference (SMC 2015), Aug 2015, Gatlinburg, United States. , 2015. <hal-01231234>

HAL Id: hal-01231234

<https://hal.science/hal-01231234v1>

Submitted on 19 Nov 2015

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization

Towards a High-Performance Tensor Algebra Package for Accelerators

M. Baboulin, V. Dobrev, J. Dongarra, C. Earl, J. Falcou, A. Haidar, I. Karlin, T. Kolev, I. Masliah, and S. Tomov



Abstract

Numerous important applications, e.g., high-order FEM simulations, can be expressed through tensors. Examples are computation of FE matrices and SpMV products expressed as generalized tensor contractions. Contractions by the first index can often be represented as tensor index reordering plus gemm, which is a key factor to achieve high-performance. We present ongoing work on the design of a high-performance package in MAGMA for Tensor algebra that includes techniques to organize tensor contractions, data storage, and parametrization related to batched execution of large number of small tensor contractions. We apply auto-tuning and code generation techniques to provide an architecture-aware, user-friendly interface.

Motivation

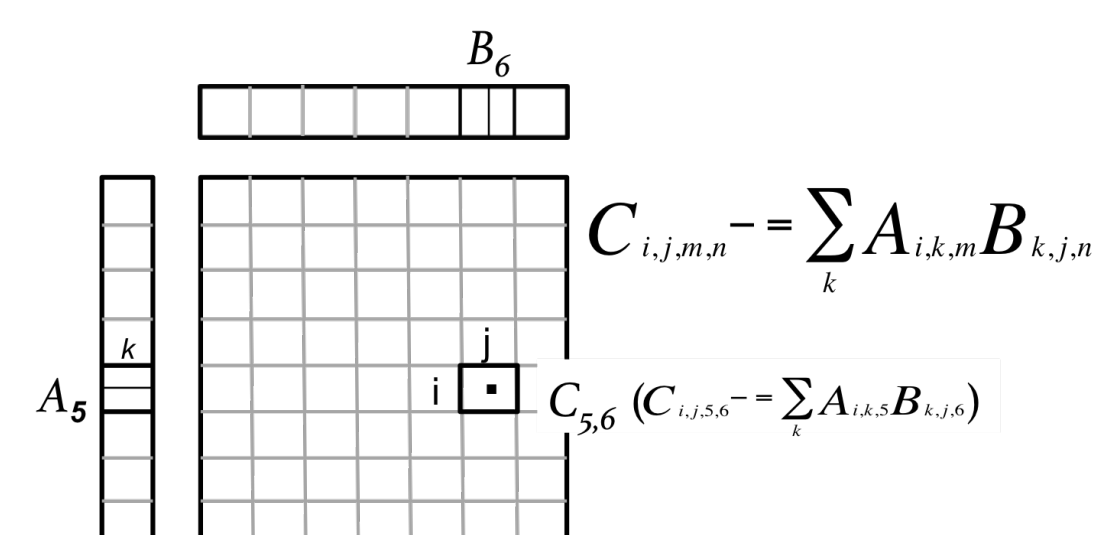
Numerous important applications can be expressed through tensors:

- High-order FEM simulations
- Data Mining
- Signal Processing
- Deep Learning
- Numerical Linear Algebra
- Graph Analysis
- Numerical Analysis
- Neuroscience and more

The goal is to design a:

- High-performance package for Tensor algebra
- Built-in architecture-awareness (GPU, Xeon Phi, multicore)
- User-friendly interface

Example cases



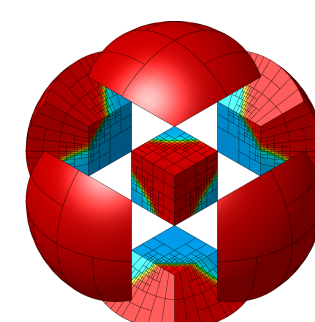
Numerical linear algebra:

- A 4-dimensional tensor contraction
- rank-k update on matrices in tile format (k can be small, e.g., sub-vector/warp size)
- Must determine (in software) if possible to do it through batched GEMM kernels

Lagrangian Hydrodynamics in the BLAST code^[1]

On semi-discrete level our method can be written as

$$\begin{aligned} \text{Momentum Conservation: } \frac{dv}{dt} &= -M_v^{-1} F \cdot \mathbf{1} \\ \text{Energy Conservation: } \frac{de}{dt} &= M_e^{-1} F^T \cdot \mathbf{v} \\ \text{Equation of Motion: } \frac{dx}{dt} &= \mathbf{v} \end{aligned}$$



where \mathbf{v} , e , and \mathbf{x} are the unknown velocity, specific internal energy, and grid position, respectively; M_v and M_e are independent of time velocity and energy mass matrices; and \mathbf{F} is the generalized corner force matrix depending on $(\mathbf{v}, e, \mathbf{x})$ that needs to be evaluated at every time step.

[1] V. Dobrev, T. Kolev, R. Rieben. *High order curvilinear finite element methods for Lagrangian hydrodynamics*. SIAM J.Sci.Comp.34(5), B606–B641. (36 pages)

Tensor operations in high-order FEM

Consider the FE mass matrix \mathbf{M}_E for an element/zone E with weight ρ , as a 2-dimensional tensor:

$$(M_E)_{ij} = \sum_{k=1}^{nq} \alpha_k \rho(q_k) \varphi_i(q_k) \varphi_j(q_k) |J_E(q_k)|,$$

$i, j = 1, \dots, nd$, where

- nd is the number of FE degrees of freedom (dofs)
- nq is the number of quadrature points
- $\{\varphi_i\}_{i=1}^{nd}$ are the FE basis functions on the reference element
- $|J_E|$ is the determinant of the element transformation
- $\{q_k\}_{k=1}^{nq}$ and $\{\alpha_k\}_{k=1}^{nq}$ are the points and weights of the quadrature rule

Take the $nq \times nd$ matrix $B_{ki} = \varphi_i(q_k)$, and $(D_E)_{kk} = \alpha_k \rho(q_k) |J_E(q_k)|$. Then, $(M_E)_{ij} = \sum_{k=1}^{nq} B_{ki} (D_E)_{kk} B_{kj}$, or omitting the E subscript $M = B^T D B$.

Using FE of order p , we have $nd = O(p^d)$ and $nq = O(p^d)$, so B is dense $O(p^d) \times O(p^d)$ matrix.

If the FE basis and the quadrature rule have tensor product structure, we can decompose dofs and quadrature point indices in logical coordinate axes

$$\mathbf{i} = (i_1, \dots, i_d), \quad \mathbf{j} = (j_1, \dots, j_d), \quad \mathbf{k} = (k_1, \dots, k_d)$$

so M_{ij} can be viewed as 2d-dimensional tensor $M_{i_1, \dots, i_d, j_1, \dots, j_d}$

Summary of kernels needed:

- Assembly of M , referred as equations (1) & (2) below
- Evaluations of M times V , referred as equations (3) & (4) below

stored components	FLOPs for assembly	amount of storage	FLOPs for matvec	numerical kernels
full assembly				
M	$O(p^{2d})$	$O(p^{2d})$	$O(p^{2d})$	$B, D \mapsto B^T D B, x \mapsto Mx$
decomposed evaluation				
B, D	$O(p^{2d})$	$O(p^{2d})$	$O(p^{2d})$	$x \mapsto Bx, x \mapsto B^T x, x \mapsto Dx$
near-optimal assembly – equations (1) and (2)				
M_{i_1, \dots, i_d}	$O(p^{2d+1})$	$O(p^{2d})$	$O(p^{2d})$	$A_{i_1, k_1, j_1} = \sum_{k_2} B_{k_1, k_2}^{i_1} B_{k_2, j_1}^{i_1} D_{k_1, k_2}$ (1a)
				$A_{i_1, k_1, j_1, j_2} = \sum_{k_2} B_{k_1, k_2}^{i_1} B_{k_2, j_1}^{i_1} C_{i_1, k_1, j_2}$ (1b)
				$A_{i_1, k_1, k_2, j_1} = \sum_{k_3} B_{k_1, k_3}^{i_1} B_{k_3, j_1}^{i_1} D_{k_1, k_3, k_2}$ (2a)
				$A_{i_1, k_1, k_2, j_1, j_2} = \sum_{k_3} B_{k_1, k_3}^{i_1} B_{k_3, j_1}^{i_1} C_{i_1, k_1, k_2, j_2}$ (2b)
				$A_{i_1, k_1, k_2, k_3, j_1, j_2} = \sum_{k_4} B_{k_1, k_4}^{i_1} B_{k_4, j_1}^{i_1} C_{i_1, k_1, k_2, k_3, j_2}$ (2c)
near-optimal evaluation (partial assembly) – equations (3) and (4)				
B^{i_1}, D	$O(p^d)$	$O(p^d)$	$O(p^{d+1})$	$A_{j_1, k_1} = \sum_{j_2} B_{k_1, j_2}^{i_1} V_{j_1, j_2}$ (3a)
				$A_{k_1, j_2} = \sum_{j_1} B_{k_1, j_1}^{i_1} C_{j_1, j_2}$ (3b)
				$A_{k_1, j_2} = \sum_{k_2} B_{k_1, k_2}^{i_1} C_{k_2, j_2}$ (3c)
				$A_{i_1, j_2} = \sum_{k_1} B_{k_1, j_2}^{i_1} C_{k_1, j_2}$ (3d)
				$A_{j_1, j_2, k_1} = \sum_{j_3} B_{k_1, j_3}^{i_1} V_{j_1, j_3}$ (4a)
				$A_{j_1, k_1, k_2} = \sum_{j_3} B_{k_1, j_3}^{i_1} C_{j_3, k_2}$ (4b)
				$A_{k_1, k_2, k_3} = \sum_{j_1} B_{k_1, j_1}^{i_1} C_{j_1, k_3}$ (4c)
				$A_{k_1, k_2, k_3} = \sum_{k_4} B_{k_1, k_4}^{i_1} C_{k_4, k_3}$ (4d)
				$A_{k_1, k_2, k_3} = \sum_{k_4} B_{k_1, k_4}^{i_1} C_{k_4, k_2, k_3}$ (4e)
				$A_{i_1, k_1, k_2} = \sum_{k_3} B_{k_1, k_3}^{i_1} C_{k_3, k_2}$ (4f)
matrix-free evaluation				
none	none	none	$O(p^{d+1})$	evaluating entries of $B^{i_1}, D, (3a)-(4f)$ sums

APPROACH AND RESULTS

User-friendly interface

To provide various interfaces, including one using C++11.

Top level design to provide features similar to the

mshadow library. <https://github.com/dmlc/mshadow>

```
// Our current interface :
// create a 2 x 5 x 2 float tensor, default locality is cpu using std::vector as default backend for data
Tensor<2,5,2> ts;
// create a 2 x 5 x 2 tensor on the gpu using thrust as the default backend for data
Tensor<2,5,2> gpu_ts;
// Call a thrust function to set values to 9
thrust::fill(ts.begin(), ts.end(), 9);
// Send back values to the cpu tensor
ts = gpu_ts;
// Reshape the 2 x 5 x 2 tensor to a matrix 2 x 10 using views
view<2,10> mat = ts;
```

Index reordering/reshape

If we store tensors as column-wise 1D arrays,

$$M_{i_1, i_2, j_1, j_2}^{nd_1 \times nd_2 \times nd_1 \times nd_2} = M_{i, j}^{nd \times nd} = M_{i_1 + nd_1 i_2 + nd(j_1 + nd_1 j_2)}^{nd^2}$$

, i.e., M can be interpreted as a 4th order tensor, a $nd \times nd$ matrix, or a vector of size nd^2 , without changing the storage. We can define

$$Reshape(T)_{j_1, \dots, j_q}^{m_1 \times \dots \times m_q} = T_{i_1, \dots, i_r}^{n_1 \times \dots \times n_r}$$

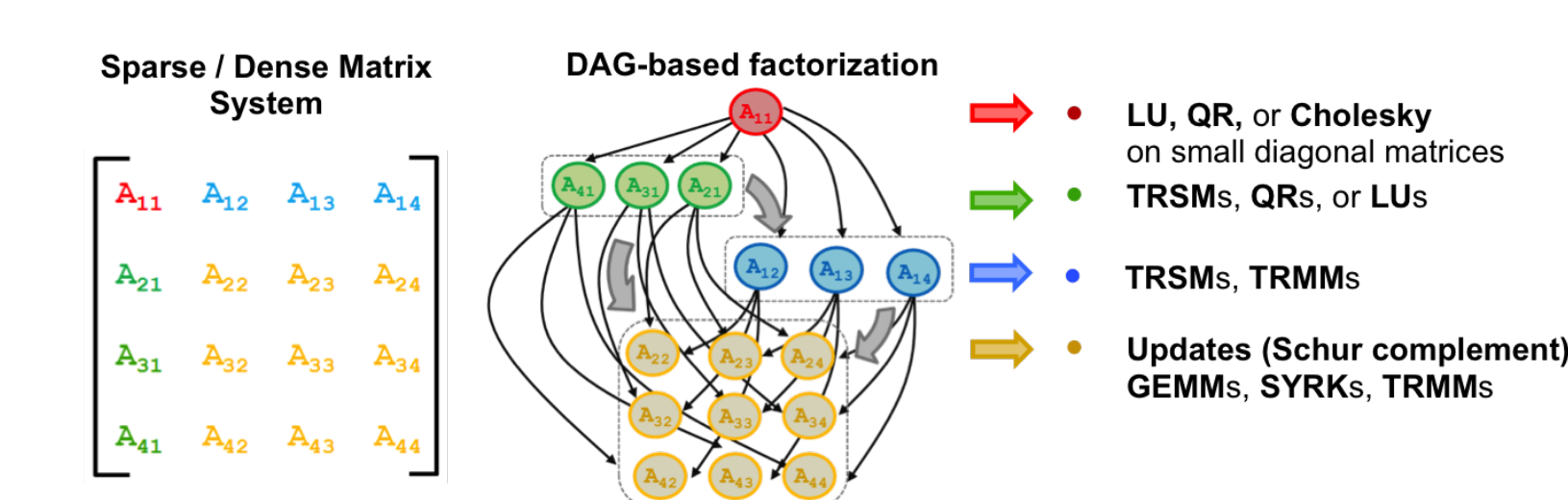
as long as $n_1 \dots n_r = m_1 \dots m_q$ and for every

$$i_1 \dots j_1 \dots i_1 + n_1 i_2 + \dots + n_1 n_2 \dots i_{r-1} i_r = j_1 + m_1 j_2 + \dots + m_1 m_2 \dots m_{q-1} j_q$$

Contractions can be implemented as a sequence of pairwise contractions. There is enough complexity here to search for something better: code generation, index reordering, and autotuning will be used, e.g., contractions (3a) - (4f) can be implemented as tensor index-reordering plus gemm $A, B \rightarrow A^T B$.

Batched LA

Tensor contractions are transformed through reshapes to batched LA operations, many of which available in MAGMA^[2] <http://icl.cs.utk.edu/magma/> (including LU, QR, Cholesky, GEMM, GEMV, TRSM, SYRK).



[2] A.Haidar, T.Dong, S.Tomov, P.Luszczek, and J.Dongarra. *A framework for batched and GPU-resident factorization algorithms applied to block Householder transformations*. ISC High Performance 2015, Frankfurt, Germany, July 12-16, 2015.

Conclusions and Future directions

- High-performance package on Tensor Algebra has the potential for high-impact on a number of important applications
- Multidisciplinary effort
- Current results show promising performance, where various components will be leveraged from autotuning MAGMA Batched linear algebra kernels, and BLAST from LLNL
- This is an ongoing work

Code Generation

C++11 features will be used as much as possible. Additional needs will be handled by defining a domain specific embedded language (DSEL). This technique is used in C++ to take advantage of DSL features while using the optimizations provided by a standard compiler. It will handle the generation of versions (index reordering, next) to be empirically evaluated and be part of the autotuning framework.

Autotuning

We are developing fixed-size gemm kernels for GPUs, Xeon Phi, and multicore (see Figure on Right for a single core intel Xeon E5-2620 and K40) through an autotuning framework. A number of generic versions are developed and parametrized for performance. The parameters are autotuned (empirically) to find "best" kernels for specific size.

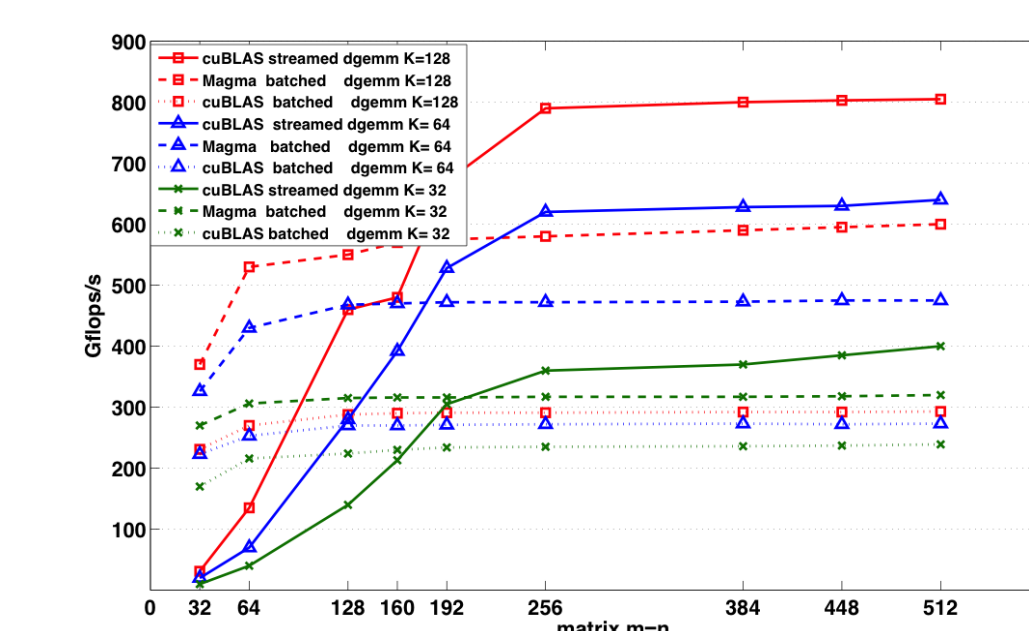
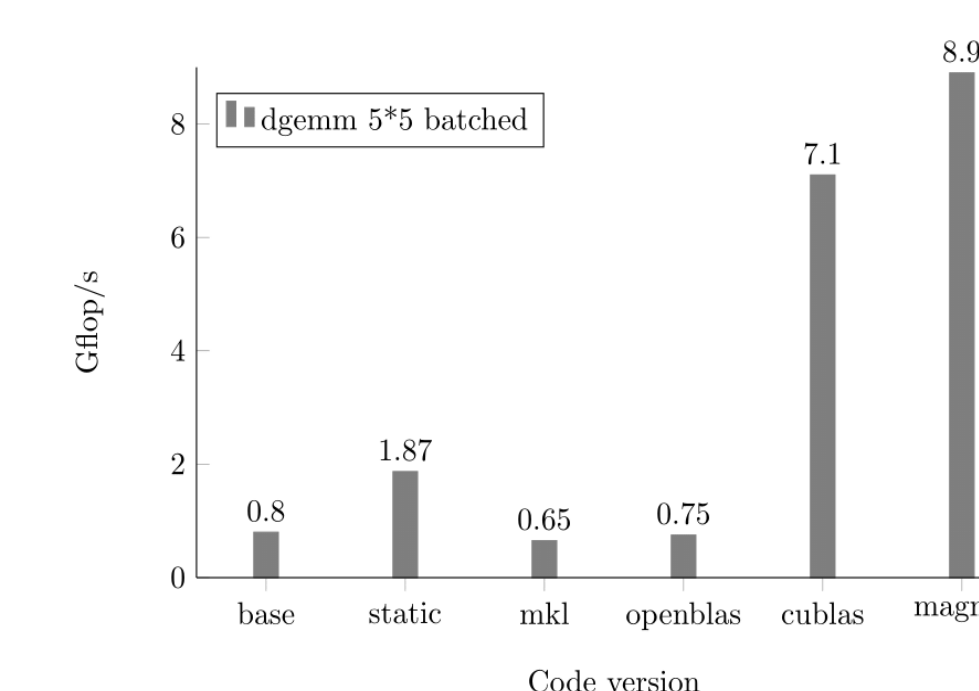


Figure: Batched dgemms on K40 GPU. Batch count is 2,000.

MAGMA exceeds in performance CUBLAS for "small" sizes, currently tuned for above 32. Current work is concentrated on kernels for fixed smaller (sub-warp) sizes.