# A Compiler for Linear Algebra Operations

## Henrik Barthels

| Introduction   | Encoded Linear Algebra Knowledge  |                        |
|--|---|------------------------|
| <ul> <li>Translating the mathematical description of a linear algebra problem to efficient<br/>code is a difficult task. Examples of such problems are:</li> </ul>   | Knowledge about linear algebra is used to simplify and rewrite expressions, as well as to infer properties. Matrix properties are crucial to select the most suitable kornels, as well as for simplifications.  |                        |
| b := $(X^{T}X)^{-1}X^{T}y$<br>$x := (A^{-T}B^{T}BA^{-1} + R^{T}[\Lambda(Rz)]R)^{-1}A^{-T}B^{T}BA^{-1}y$  | Inference of Properties Simplifications   |                        |
| <ul> <li>Languages such as Matlab and Julia are easy to use, but performance is usually suboptimal.</li> <li>We are developing a compiler that offers the ease-of-use, and thus, productivity of Matlab paired with performance that comes close to what a human expert achieves.</li> </ul> | $\begin{array}{cccc} A & \rightarrow & A^{T} \\ \hline A & B & \rightarrow & AB \\ \hline A & B & \rightarrow & AB \\ \hline SPD(S) & \rightarrow & SPD(S^{-1}) \\ expr & = expr^{T} & \rightarrow & Symmetric(expr) \end{array} \begin{array}{ccc} (AB)^{T} & \rightarrow & B^{T}A^{T} \\ \hline A & B & \rightarrow & AB \\ \hline Q^{T}Q & \rightarrow & I \\ Q^{-1} & \rightarrow & Q^{T} \\ \hline \alpha & A + & \beta & A \\ \hline \alpha & A + & \beta & A \\ \hline A & B & \rightarrow & AB \\ \hline Q^{T}Q & \rightarrow & I \\ \hline Q^{-1} & \rightarrow & Q^{T} \\ \hline \alpha & A + & \beta & A \\ \hline A & B & \rightarrow & AB \\ \hline Q^{T}Q & \rightarrow & I \\ \hline A & B & \rightarrow & AB \\ \hline Q^{T}Q & \rightarrow & I \\ \hline A & B & \rightarrow & AB \\ \hline Q^{T}Q & \rightarrow & I \\ \hline A & B & \rightarrow & AB \\ \hline Q^{T}Q & \rightarrow & I \\ \hline A & B & \rightarrow & AB \\ \hline Q^{T}Q & \rightarrow & I \\ \hline A & B & \rightarrow & AB \\ \hline A & B & AB \\ \hline A & $ | c(A)<br>al(Q)<br>al(Q) |
| Compute $X^T X$ directly<br>or apply factorization   | Derivation Graph  |                        |

Algorithms are derived by repeatedly applying different derivation steps.



Input  $\rightarrow$  Simplifications  $\rightarrow$  Common Subexpr.  $\longrightarrow$  Compute  $\rightarrow$  Algorithms

Each of these steps yields one or more results. We use a graph to represent all algorithms. The root of the graph is the input expression, and each path in the graph corresponds to one algorithm. Below, a part of such a graph is shown. The full graph has 83 nodes.



size  $\rightarrow$  (*rows*, *columns*) property  $\rightarrow$  LowerTriangular | Orthogonal | ...  $expr^{-1} | expr^{T} | expr^{-T}$ operand  $\rightarrow$  *symbol* | *symbol*<sub>indices</sub> indices  $\rightarrow$  *index*<sup>+</sup>

### **Instruction Set**

As the instruction set, we use optimized kernels as offered by linear algebra libraries.

**BLAS** [3]

•

LAPACK [1]

•  $\mathbf{y} \leftarrow \alpha \mathbf{A}\mathbf{x} + \beta \mathbf{y}$ 

- $C \leftarrow \alpha AB + \beta C$
- $\mathbf{B} \leftarrow \alpha \mathbf{A}^{-1}\mathbf{B}$  •
- Matrix factorizations.
- Eigensolvers.

Solvers for linear systems with specific properties.

## Optimizations

#### **Common Subexpression Elimination**

- FLOPS can be saved by reusing subexpressions that appear more than once, for example AB in AB + ABC.
- Operators such as transposition and inversion make the detection of common subexpression more complicated:

 $\underbrace{AB^{-T}}_{}\underbrace{B^{-1}A^{T}}_{} \rightarrow C_{1}C_{1}^{T} \qquad C_{1} = AB^{-T} = \left(B^{-1}A^{T}\right)^{T} = C_{2}^{T}$ 

## Results

• Example:  $z := (X^T M^{-1} X)^{-1} X^T M^{-1} y$ ,  $M \in \mathbb{R}^{n \times n}$ ,  $X \in \mathbb{R}^{n \times m}$ ,  $y \in \mathbb{R}^n$ ,  $n \ge m$ . *M* is symmetric positive definite.

- The compiler (written in Python) finds 23 algorithms in less than a second.
- The algorithm shown below is identified as the cheapest one in terms of FLOPS.
- It is compared to two other implementations that rely on Matlab's strategy to evaluate expressions.

Naive implementation
z = inv(X'\*inv(M)\*X)\*X'\*inv(M)\*y;

Recommended implementation  $z = (X'*(M\backslash X))\backslash X'*(M\backslash y);$ 

Compiler implementation
L1 = chol(M, 'lower');
t1 = linsolve(L1, X, opts1);





• We developed algorithms to detect such common subexpressions.

#### **Generalized Matrix Chain Problem**

- The cost of a product  $M_0 \cdots M_n$  highly depends on the parenthesization.
- We developed a generalized version of the  $O(n^3)$  matrix chain algorithm [2].
- It finds the optimal parenthesization for matrix chains containing transposition and inversion, for example  $X := AB^{-T}C^{-1}D$ .
- This algorithm also takes properties into account.

t2 = t1'\*t1; L2 = chol(t2, 'lower'); t3 = linsolve(L1, y, opts1); t4 = t1'\*t3; t5 = linsolve(L2, t4, opts1); z = linsolve(L2, t5, opts2);

#### References

[1] E. Anderson, Z. Bai, et al. LAPACK Users' guide, volume 9. SIAM, 1999.
[2] T. H. Cormen, R. L. Rivest, and C. E. Leiserson. Introduction to Algorithms. McGraw-Hill, Inc., 1990.
[3] J. J. Dongarra, J. Du Croz, et al. A set of Level 3 Basic Linear Algebra Subprograms. ACM TOMS, 16(1):1–17, 1990.

#### Henrik Barthels

Aachen Institute for Advanced Study in Computational Engineering Science (AICES)

RWTH Aachen, Germany

http://hpac.rwth-aachen.de

barthels@aices.rwth-aachen.de



