

Matrix computations: Going beyond libraries

Paolo Bientinesi

Umeå Universitet

pauldj@cs.umu.se

October 5, 2022

Swedish e-Science Academy



High Performance and
Automatic Computing



UMEÅ UNIVERSITY

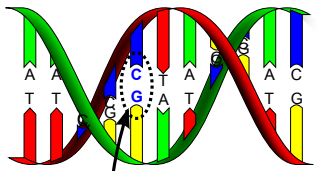


HPC2N

10-second lecture

- ▶ Matrix computations are ubiquitous in science and engineering
- ▶ Many excellent linear algebra libraries exist
- ▶ Are they used (properly)? Not so much

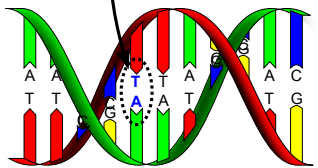
Matrix computations everywhere



1

SNP

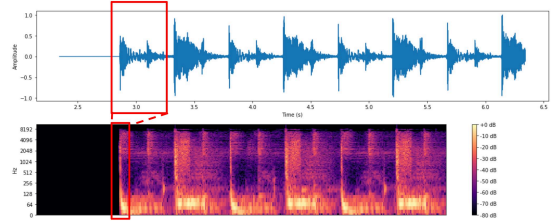
2



GENOME-WIDE ASSOCIATION STUDIES



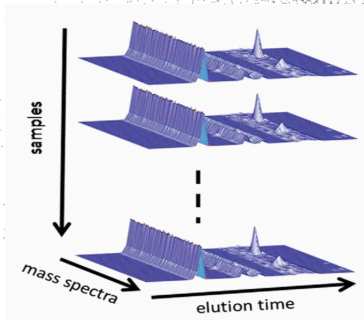
ESTIMATION OF EARTH'S GRAVITATIONAL FIELD



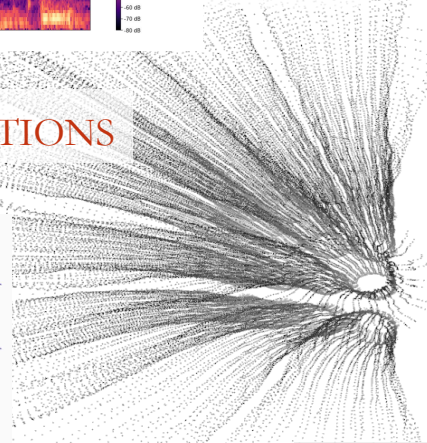
MUSIC
INFORMATION
RETRIEVAL

DIGITAL SIGNAL
PROCESSING

MATRIX COMPUTATIONS



CHEMOMETRICS



SIMULATION
OF CARBON
NANOTUBES

Libraries since the '70s!

- ▶ 1971: NAG Fortran Library – scientific computing
- ▶ 1972: EISPACK – eigenproblems
- ▶ 1973: BLAS 1 – basic vector ops
- ▶ 1975: LINPACK – linear systems
- ▶ ...
- ▶ ...

Hierarchy of Libraries: 1980s

BLAS-2: Mat-vec ops, ACM TOMS 1988

BLAS-3: mat-mat ops, ACM TOMS 1990

BLAS-1, BLAS-2, BLAS-3

Hierarchy of Libraries: 1990s

Solvers & eigensolvers, 1992

LAPACK

BLAS-1, BLAS-2, BLAS-3

Hierarchy of Libraries: 1990s

Distributed memory, 1995, 1997

ScaLAPACK, PLAPACK, ...

LAPACK

BLAS-1, BLAS-2, BLAS-3

Hierarchy of Libraries: 1990s

Dense & sparse, 1997

PETSc, ...

ScaLAPACK, PLAPACK, ...

LAPACK

BLAS-1, BLAS-2, BLAS-3

Many libraries! ... (and many more iterative ones)

PETSc, Trilinos, ...

ScaLAPACK, PLAPACK, Elemental, ...

LAPACK, Plasma, SuperMatrix, Magma, ...

BLAS-1, BLAS-2, BLAS-3, ATLAS, BTO-BLAS, BLIS, ...



JACK DONGARRA

United States – 2021

CITATION

For his pioneering contributions to numerical algorithms and libraries that enabled high performance computational software to keep pace with exponential hardware improvements for over four decades



RESEARCH
SUBJECTS

But... Issue #1: expression vs. code

But... Issue #1: expression vs. code

$$C := C + A * B^T + B * A^T$$

Naive translation to code:

```
CALL DGEMM( 'N', 'Y', N-K-KB+1, N-K-KB+1, KB, ONE, A( K+KB, K+KB ), LDA,  
$      B( K, K+KB ), LDB, ONE, C( K+KB, K+KB ), LDC )  
CALL DGEMM( 'N', 'Y', N-K-KB+1, N-K-KB+1, KB, ONE, B( K+KB, K+KB ), LDA,  
$      A( K, K+KB ), LDB, ONE, C( K+KB, K+KB ), LDC )
```

The “right” way:

```
CALL DSYR2K( UPLO, 'Transpose', N-K-KB+1, KB, -ONE, A( K, K+KB ), LDA,  
$      B( K, K+KB ), LDB, ONE, A( K+KB, K+KB ), LDA )
```

But... Issue #1: expression vs. code

$$C := C + A * B^T + B * A^T$$

Naive translation to code:

```
CALL DGEMM( 'N', 'Y', N-K-KB+1, N-K-KB+1, KB, ONE, A( K+KB, K+KB ), LDA,  
$          B( K, K+KB ), LDB, ONE, C( K+KB, K+KB ), LDC )  
CALL DGEMM( 'N', 'Y', N-K-KB+1, N-K-KB+1, KB, ONE, B( K+KB, K+KB ), LDA,  
$          A( K, K+KB ), LDB, ONE, C( K+KB, K+KB ), LDC )
```

The “right” way:

```
CALL DSYR2K( UPLO, 'Transpose', N-K-KB+1, KB, -ONE, A( K, K+KB ), LDA,  
$          B( K, K+KB ), LDB, ONE, A( K+KB, K+KB ), LDA )
```

Not exactly user friendly. Certainly not users' preference (anymore)

Users preference(s)

$$C := C + A * B^T + B * A^T$$

```
C = A * B' + B * A' + C; // Matlab, Octave
```

```
C = A * transpose(B) + B * transpose(A) + C // Julia
```

```
C = A * trans(B) + B * trans(A) + C; // Armadillo
```

```
C = A * B.transpose() + B * A.transpose() + C; // Eigen
```

```
ct = at @ bt.T + bt @ at.T + ct // NumPy
```

```
ct <- at %*% t(bt) + bt %*% t(at) + ct // R
```

```
C = A@tf.transpose(B) + B@tf.transpose(A) + C // TensorFlow
```

```
C = A@torch.t(B) + B@torch.t(A) + C // PyTorch
```

Issue #2: building blocks vs. applications

- ▶ BLAS: **B**asic Linear Algebra Subprograms

$$\mathbf{w} := \mathbf{x}^T \mathbf{y} \text{ (DOT)}, \quad \mathbf{y} := \alpha \mathbf{x} + \mathbf{y} \text{ (AXPY)}, \quad \eta := (\mathbf{x}^T \mathbf{x})^{1/2} \text{ (NORM)}$$

$$\mathbf{y} := \mathbf{A} \mathbf{x} \text{ (GEMV)}, \quad \mathbf{A} \mathbf{x} = \mathbf{b} \text{ (TRSV)}$$

$$\mathbf{C} := \mathbf{A} \mathbf{B} \text{ (GEMM)}, \quad \mathbf{A} \mathbf{X} = \mathbf{B} \text{ (TRSM)}$$

- ▶ LAPACK: Linear Algebra Package

$$\mathbf{A} \mathbf{x} = \lambda \mathbf{x} \quad \text{eigenproblems}$$

$$\mathbf{A} \mathbf{x} = \mathbf{b} \quad \text{linear systems, least squares problems}$$

$$\mathbf{Q} \mathbf{R} = \mathbf{A} \quad \text{matrix factorizations, ...}$$

Applications

Signal Processing

$$x := (A^{-T}B^TBA^{-1} + R^T L R)^{-1} A^{-T}B^TBA^{-1}y \quad R \in \mathbb{R}^{n-1 \times n}, \text{ UT}; L \in \mathbb{R}^{n-1 \times n-1}, \text{ DI}$$

Kalman Filter

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; x_k^a := x_k^b + K_k (z_k - H x_k^b); P_k^a := (I - K_k H) P_k^b$$

Ensemble Kalman Filter

$$X^a := X^b + (B^{-1} + H^T R^{-1} H)^{-1} (Y - H X^b) \quad B \in \mathbb{R}^{N \times N} \text{ SSPD}; R \in \mathbb{R}^{m \times m}, \text{ SSPD}$$

Ensemble Kalman Filter

$$\delta X := (B^{-1} + H^T R^{-1} H)^{-1} H^T R^{-1} (Y - H X^b)$$

Ensemble Kalman Filter

$$\delta X := X V^T (R + H X (H X)^T)^{-1} (Y - H X^b)$$

Image Restoration

$$x_k := (H^T H + \lambda \sigma^2 I_n)^{-1} (H^T y + \lambda \sigma^2 (v_{k-1} - u_{k-1}))$$

Image Restoration

$$H^\dagger := H^T (H H^T)^{-1}; y_k := H^\dagger y + (I_n - H^\dagger H) x_k$$

Rand. Matrix Inversion

$$X_{k+1} := S(S^T A S)^{-1} S^T + (I_n - S(S^T A S)^{-1} S^T A) X_k (I_n - A S(S^T A S)^{-1} S^T)$$

Rand. Matrix Inversion

$$X_{k+1} := X_k + W A^T S(S^T A W A^T S)^{-1} S^T (I_n - A X_k) \quad W \in \mathbb{R}^{n \times n}, \text{ SPD}$$

Rand. Matrix Inversion

$$X_{k+1} := X_k + (I_n - X_k A^T) S(S^T A^T W A S)^{-1} S^T A^T W$$

Rand. Matrix Inversion

$$\Lambda := S(S^T A W A S)^{-1} S^T; \Theta := \Lambda A W; M_k := X_k A - I \\ X_{k+1} := X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta$$

Generalized Least Squares	$b := (X^T M^{-1} X)^{-1} X^T M^{-1} y$	$n > m; M \in \mathbb{R}^{n \times n}, \text{SPD}; X \in \mathbb{R}^{n \times m}; y \in \mathbb{R}^{n \times 1}$
Stochastic Newton	$B_k := \frac{k}{k-1} B_{k-1} (I_n - A^T W_k ((k-1)I_l + W_k^T A B_{k-1} A^T W_k)^{-1} W_k^T A B_{k-1})$	
Optimization	$x_f := W A^T (A W A^T)^{-1} (b - A x); x_o := W (A^T (A W A^T)^{-1} A x - c)$	
Optimization	$x := W (A^T (A W A^T)^{-1} b - c)$	
Triangular Matrix Inv.	$X_{10} := L_{10} L_{00}^{-1}; X_{20} := L_{20} + L_{22}^{-1} L_{21} L_{11}^{-1} L_{10}; X_{11} := L_{11}^{-1}; X_{21} := -L_{22}^{-1} L_{21}$	
Tikhonov Regularization	$x := (A^T A + \Gamma^T \Gamma)^{-1} A^T b$	$A \in \mathbb{R}^{n \times m}; \Gamma \in \mathbb{R}^{m \times m}; b \in \mathbb{R}^{n \times 1}$
Tikhonov Regularization	$x := (A^T A + \alpha^2 I)^{-1} A^T b$	
Gen. Tikhonov Reg.	$x := (A^T P A + Q)^{-1} (A^T P b + Q x_0)$	$P \in \mathbb{R}^{n \times n}, \text{SSPD}; Q \in \mathbb{R}^{m \times m}, \text{SSPD}; x_0 \in \mathbb{R}^{m \times 1}$
Gen. Tikhonov reg.	$x := x_0 + (A^T P A + Q)^{-1} (A^T P (b - A x_0))$	
LMMSE estimator	$K_{t+1} := C_t A^T (A C_t A^T + C_z)^{-1}; x_{t+1} := x_t + K_{t+1} (y - A x_t); C_{t+1} := (I - K_{t+1} A) C_t$	
LMMSE estimator	$x_{\text{out}} = C_X A^T (A C_X A^T + C_Z)^{-1} (y - A x) + x$	
LMMSE estimator	$x_{\text{out}} := (A^T C_Z^{-1} A + C_X^{-1})^{-1} A^T C_Z^{-1} (y - A x) + x$	

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k (z_k - H x_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_{\dagger} := P C P^T + Q \\ K := C_{\dagger} H^T (H C_{\dagger} H^T)^{-1} \end{cases}$$

$$\Lambda := S(S^T A W A S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I \\ X_{k+1} := X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U (I + U^T Q^{-1} U)^{-1} U^T$$

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k (z_k - H x_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_{\dagger} := P C P^T + Q \\ K := C_{\dagger} H^T (H C_{\dagger} H^T)^{-1} \end{cases}$$

$$\Lambda := S(S^T A W A S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I \\ X_{k+1} := X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U (I + U^T Q^{-1} U)^{-1} U^T$$

$$\boxed{y := \alpha x + y} \quad \boxed{LU = A} \quad \dots \quad \boxed{C := \alpha AB + \beta C}$$

$$\boxed{X := A^{-1} B} \quad \boxed{C := AB^T + BA^T + C} \quad \boxed{X := L^{-1} M L^{-T}} \quad \boxed{QR = A}$$

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k (z_k - H x_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_{\dagger} := P C P^T + Q \\ K := C_{\dagger} H^T (H C_{\dagger} H^T)^{-1} \end{cases}$$

$$\Lambda := S(S^T A W A S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I \\ X_{k+1} := X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U (I + U^T Q^{-1} U)^{-1} U^T$$



$$\begin{array}{cccc} \boxed{y := \alpha x + y} & \boxed{LU = A} & \dots & \boxed{C := \alpha AB + \beta C} \\ \boxed{X := A^{-1} B} & \boxed{C := AB^T + BA^T + C} & \boxed{X := L^{-1} M L^{-T}} & \boxed{QR = A} \end{array}$$

$$K_k := P_k^b H^T (H P_k^b H^T + R)^{-1}; \quad x_k^a := x_k^b + K_k (z_k - H x_k^b); \quad P_k^a := (I - K_k H) P_k^b$$

$$\begin{cases} C_{\dagger} := P C P^T + Q \\ K := C_{\dagger} H^T (H C_{\dagger} H^T)^{-1} \end{cases}$$

$$\Lambda := S(S^T A W A S)^{-1} S^T; \quad \Theta := \Lambda A W; \quad M_k := X_k A - I \\ X_{k+1} := X_k - M_k \Theta - (M_k \Theta)^T + \Theta^T (A X_k A - A) \Theta$$

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y \quad \dots \quad E := Q^{-1} U (I + U^T Q^{-1} U)^{-1} U^T$$

LINEAR ALGEBRA MAPPING PROBLEM

C. Psarras, H. Barthels, P. Bientinesi, [arXiv:1911.09421]

"The Linear Algebra Mapping Problem. Current state of linear algebra languages and libraries", ACM TOMS, 2022

A. Sankaran, N.A. Alashti, C. Psarras, P. Bientinesi, [arXiv:2202.09888]

"Benchmarking the Linear Algebra Awareness of TensorFlow and PyTorch", iWAPT-22

H. Barthels, C. Psarras, P. Bientinesi, [arXiv:1912.12924]

"Linnea: Automatic Generation of Efficient Linear Algebra Programs", ACM TOMS, 2021

Linear Algebra Mapping Problem (“LAMP”)

Linear Algebra Mapping Problem (“LAMP”)

- ▶ \mathcal{E} : a sequence of explicit assignments $var_i := EXP_i$
- ▶ \mathcal{K} : a set of available computational building blocks e.g., BLAS, LAPACK, ...
- ▶ \mathcal{M} : a cost function defined over \mathcal{K}^+ #FLOPs, exec. time, #mem.ops, stability

Linear Algebra Mapping Problem (“LAMP”)

- ▶ \mathcal{E} : a sequence of explicit assignments $var_i := EXP_i$
- ▶ \mathcal{K} : a set of available computational building blocks e.g., BLAS, LAPACK, ...
- ▶ \mathcal{M} : a cost function defined over \mathcal{K}^+ #FLOPs, exec. time, #mem.ops, stability

LAMP:

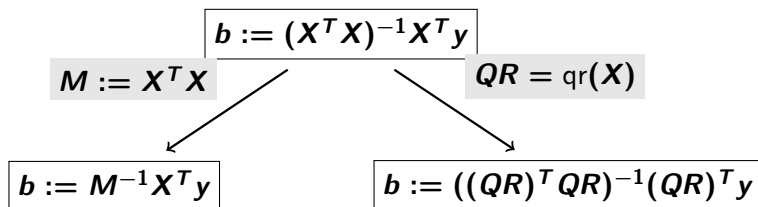
Find a sequence of calls to building blocks in \mathcal{K} , optimal according to \mathcal{M} , that computes all the assignments in \mathcal{E} .

- ▶ Suboptimal solution \rightarrow easy
- ▶ Optimality \rightarrow NP complete \leftarrow reduction from Ensemble Computation

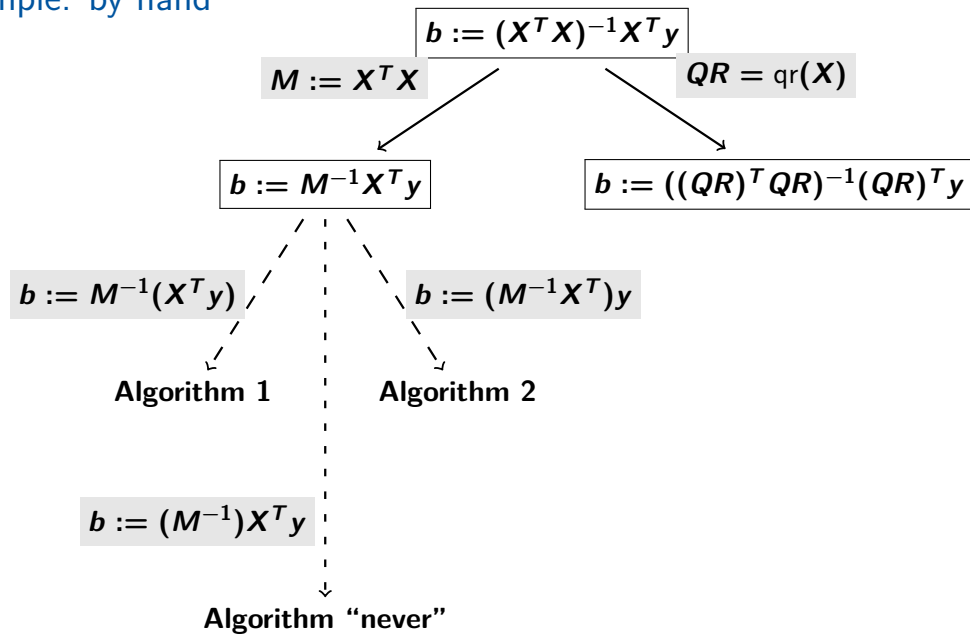
Example:

$$\mathbf{b} := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

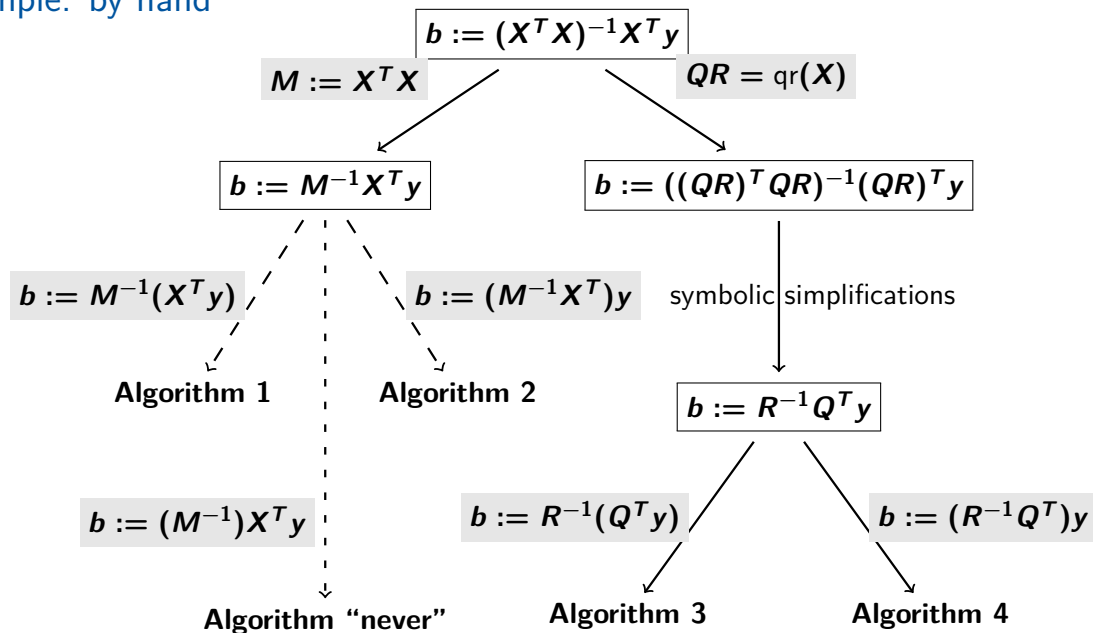
Example: by hand



Example: by hand



Example: by hand



Example: high-level language

$$b := (X^T X)^{-1} X^T y$$

```
b = inv(X'*X)*X'*y           // Matlab (naive)
```

```
b = (X'*X)\X'*y             // Matlab (recommended)
```

Example: $\mathbf{b} := (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

1) **By hand:** C/Fortran Required: expertise, time, patience

computer efficiency! 😊 ... but human productivity 😞

2) **High-level languages:** Matlab, Julia, R, Armadillo (C++), NumPy, TensorFlow, ...

human productivity! 😊 ... computer efficiency?

Investigation: How well do high-level languages solve LAMPs?

High-level notation

- ▶ Matlab
- ▶ Octave
- ▶ Julia
- ▶ C++ with Armadillo
- ▶ C++ with Eigen
- ▶ NumPy
- ▶ R
- ▶ (TensorFlow)
- ▶ (PyTorch)

Q1: Do they map? matrix products

input

Matlab

Octave

Julia

R

Eigen

Armad.

NumPy

C

$C = A*B$

Q1: Do they map? matrix products

input	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$C = A*B$	0.29	0.28	0.30	0.31	0.29	0.29	0.29	0.27

Q1: Do they map? matrix products

input	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$C = A*B$	0.29	0.28	0.30	0.31	0.29	0.29	0.29	0.27
GEMM	✓	✓	✓	✓	✓	✓	✓	

Q1: Do they map? matrix products

input	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$C = A*B$	0.29	0.28	0.30	0.31	0.29	0.29	0.29	0.27
GEMM	✓	✓	✓	✓	✓	✓	✓	

$C = C+A*A'$

Q1: Do they map? matrix products

input	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$C = A*B$	0.29	0.28	0.30	0.31	0.29	0.29	0.29	0.27
GEMM	✓	✓	✓	✓	✓	✓	✓	
$C = C+A*A'$	0.18	0.17	0.21	0.32	0.29	0.17	0.18	0.14

Q1: Do they map? matrix products

input	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$C = A*B$	0.29	0.28	0.30	0.31	0.29	0.29	0.29	0.27
GEMM	✓	✓	✓	✓	✓	✓	✓	
$C = C+A*A'$	0.18	0.17	0.21	0.32	0.29	0.17	0.18	0.14
SYRK	✓	✓	✓	×	×	✓	✓	

Q1: Do they map? matrix products

input	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$C = A*B$	0.29	0.28	0.30	0.31	0.29	0.29	0.29	0.27
GEMM	✓	✓	✓	✓	✓	✓	✓	
$C = C+A*A'$	0.18	0.17	0.21	0.32	0.29	0.17	0.18	0.14
SYRK	✓	✓	✓	×	×	✓	✓	
$C = C+AB'+BA'$								

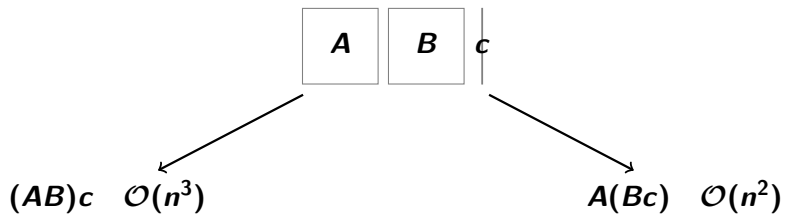
Q1: Do they map? matrix products

input	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$C = A*B$	0.29	0.28	0.30	0.31	0.29	0.29	0.29	0.27
GEMM	✓	✓	✓	✓	✓	✓	✓	
$C = C+A*A'$	0.18	0.17	0.21	0.32	0.29	0.17	0.18	0.14
SYRK	✓	✓	✓	×	×	✓	✓	
$C = C+AB'+BA'$	0.57	0.59	0.69	0.59	0.58	0.57	0.58	0.28

Q1: Do they map? matrix products

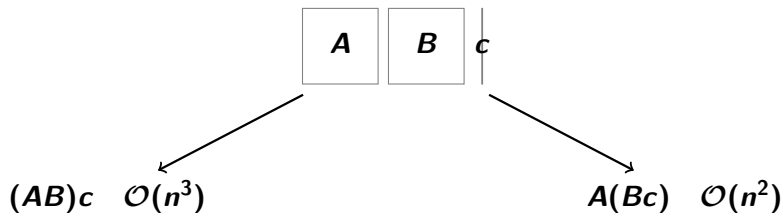
input	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$C = A*B$ GEMM	0.29 ✓	0.28 ✓	0.30 ✓	0.31 ✓	0.29 ✓	0.29 ✓	0.29 ✓	0.27
$C = C+A*A'$ SYRK	0.18 ✓	0.17 ✓	0.21 ✓	0.32 ×	0.29 ×	0.17 ✓	0.18 ✓	0.14
$C = C+AB'+BA'$ SYR2K	0.57 ×	0.59 ×	0.69 ×	0.59 ×	0.58 ×	0.57 ×	0.58 ×	0.28

Optimal parenthesisation



Matrix product is associative, but its cost is not

Optimal parenthesisation



Matrix product is associative, but its cost is not

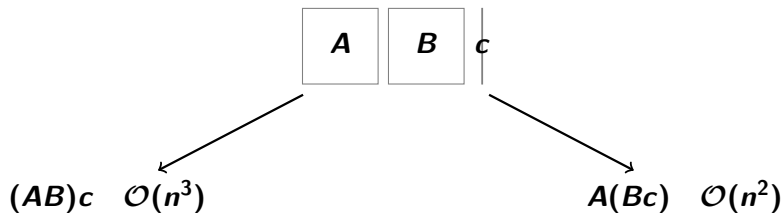
Matrix Chain(k)

Catalan($k - 1$) parenthesisizations

if parallelism, then $k!$ parenthesisizations

Which is the best one?

Optimal parenthesisation



Matrix product is associative, but its cost is not

Matrix Chain(k)

Catalan($k - 1$) parenthesizations
if parallelism, then $k!$ parenthesizations

Which is the best one?

Matrix Chain Algorithm(k)

$\mathcal{O}(k \log k)$ Hu & Shing 1982
 $\mathcal{O}(k^3)$ dynamic programming

Q2: Matrix Chain?

Chain	Optimal Evaluation
1) "left-to-right"	$((A B) C)$
2) "right-to-left"	$(A (B C))$
3) "mixed"	$((A B) (C D))$

Q2: Matrix Chain?

Chain	Optimal Evaluation
1) "left-to-right"	$((A B) C)$
2) "right-to-left"	$(A (B C))$
3) "mixed"	$((A B) (C D))$

	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy
1) $A*B*C$	0.056	0.056	0.055	0.061	0.058	0.056	0.055
$(A*B)*C$	0.056	0.056	0.055	0.061	0.058	0.056	0.055

Q2: Matrix Chain?

Chain	Optimal Evaluation
1) "left-to-right"	$((A B) C)$
2) "right-to-left"	$(A (B C))$
3) "mixed"	$((A B) (C D))$

	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy
1) $A*B*C$	0.056	0.056	0.055	0.061	0.058	0.056	0.055
1) $(A*B)*C$	0.056	0.056	0.055	0.061	0.058	0.056	0.055
2) $A*B*C$	0.42	0.43	0.42	0.44	0.42	0.055	0.42
2) $A*(B*C)$	0.055	0.056	0.054	0.059	0.056	0.055	0.056

Q2: Matrix Chain?

Chain	Optimal Evaluation
1) "left-to-right"	$((A B) C)$
2) "right-to-left"	$(A (B C))$
3) "mixed"	$((A B) (C D))$

	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy
1) $A*B*C$	0.056	0.056	0.055	0.061	0.058	0.056	0.055
$(A*B)*C$	0.056	0.056	0.055	0.061	0.058	0.056	0.055
2) $A*B*C$	0.42	0.43	0.42	0.44	0.42	0.055	0.42
$A*(B*C)$	0.055	0.056	0.054	0.059	0.056	0.055	0.056
3) $A*B*C*D$	0.32	0.33	0.33	0.33	0.35	0.31	0.33
$(A*B)*(C*D)$	0.21	0.22	0.22	0.22	0.23	0.20	0.22

Q2: Matrix Chain?

Chain	Optimal Evaluation
1) "left-to-right"	$((A B) C)$
2) "right-to-left"	$(A (B C))$
3) "mixed"	$((A B) (C D))$

	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy
1)	$A*B*C$	0.056	0.056	0.055	0.061	0.058	0.055
	$(A*B)*C$	0.056	0.056	0.055	0.061	0.058	0.055
2)	$A*B*C$	0.42	0.43	0.42	0.44	0.42	0.055
	$A*(B*C)$	0.055	0.056	0.054	0.059	0.056	0.055
3)	$A*B*C*D$	0.32	0.33	0.33	0.33	0.35	0.31
	$(A*B)*(C*D)$	0.21	0.22	0.22	0.22	0.23	0.20
Matrix chains	×	×	×	×	×	≈	×

Wait! They read our paper!

Chain	Optimal Evaluation
1) "left-to-right"	$((A B) C)$
2) "right-to-left"	$(A (B C))$
3) "mixed"	$((A B) (C D))$

	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	
1)	$A*B*C$	0.056	0.056	0.055	0.061	0.058	0.056	0.055
	$(A*B)*C$	0.056	0.056	0.055	0.061	0.058	0.056	0.055
2)	$A*B*C$	0.42	0.43	0.054	0.44	0.42	0.055	0.42
	$A*(B*C)$	0.055	0.056	0.054	0.059	0.056	0.055	0.056
3)	$A*B*C*D$	0.32	0.33	0.22	0.33	0.35	0.31	0.33
	$(A*B)*(C*D)$	0.21	0.22	0.22	0.22	0.23	0.20	0.22
	Matrix chains	×	×	✓	×	×	≈	×

Matrix chains in practice

- ▶ Unary operators: transposition, inversion
- ▶ Overlapping kernels
- ▶ Decompositions
- ▶ Properties & specialized kernels

$$(\mathbf{X} := \mathbf{A}\mathbf{B}^T\mathbf{C}^{-T}\mathbf{D} + \dots)$$

$$(\text{e.g., } \mathbf{L} \leftarrow \mathbf{L}^{-1}, \mathbf{X} = \mathbf{A}^{-1}\mathbf{B})$$

$$(\text{e.g., } \mathbf{A} \rightarrow \mathbf{Q}^T\mathbf{D}\mathbf{Q}, \mathbf{A} \rightarrow \mathbf{L}\mathbf{U})$$

$$(\text{GEMM, TRMM, SYMM, } \dots)$$

#FLOPs vs. execution time (vs. numerical stability)

$$\underset{\mathcal{A}}{\operatorname{argmin}} (\operatorname{FLOPs}(\mathcal{A})) \neq \underset{\mathcal{A}}{\operatorname{argmin}} (\operatorname{time}(\mathcal{A}))$$

⇒ **Performance prediction:** efficiency

Q3: Properties?

Operation	Property	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$A*X = B$	-	0.71	0.74	0.62	0.67	0.63	0.62	0.65	0.61

Q3: Properties?

Operation	Property	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$A * X = B$	-	0.71	0.74	0.62	0.67	0.63	0.62	0.65	0.61
	Symmetric	0.71	0.73	0.62	0.69	N/A	0.62	0.65	0.46

Q3: Properties?

Operation	Property	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$A*X = B$	-	0.71	0.74	0.62	0.67	0.63	0.62	0.65	0.61
	Symmetric	0.71	0.73	0.62	0.69	N/A	0.62	0.65	0.46
	SPD	0.41	0.40	0.60	0.63	N/A	0.34	0.62	0.31

Q3: Properties?

Operation	Property	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$A * X = B$	-	0.71	0.74	0.62	0.67	0.63	0.62	0.65	0.61
	Symmetric	0.71	0.73	0.62	0.69	N/A	0.62	0.65	0.46
	SPD	0.41	0.40	0.60	0.63	N/A	0.34	0.62	0.31
	Triangular	0.03	0.04	0.03	0.63	N/A	0.62	0.65	0.03

Q3: Properties?

Operation	Property	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
A*X = B	-	0.71	0.74	0.62	0.67	0.63	0.62	0.65	0.61
	Symmetric	0.71	0.73	0.62	0.69	N/A	0.62	0.65	0.46
	SPD	0.41	0.40	0.60	0.63	N/A	0.34	0.62	0.31
	Triangular	0.03	0.04	0.03	0.63	N/A	0.62	0.65	0.03
	Diagonal	0.03	0.05	0.01	0.63	N/A	0.03	0.62	0.001
		≈	≈	≈	×	×	≈	×	

Q3: Properties?

Operation	Property	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
A*X = B	-	0.71	0.74	0.62	0.67	0.63	0.62	0.65	0.61
	Symmetric	0.71	0.73	0.62	0.69	N/A	0.62	0.65	0.46
	SPD	0.41	0.40	0.60	0.63	N/A	0.34	0.62	0.31
	Triangular	0.03	0.04	0.03	0.63	N/A	0.62	0.65	0.03
	Diagonal	0.03	0.05	0.01	0.63	N/A	0.03	0.62	0.001
		≈	≈	≈	×	×	≈	×	
C = A*B	-	1.44	1.48	1.47	1.47	1.45	1.44	1.44	1.46
	Triangular	1.44	1.48	0.75	1.47	1.45	1.44	1.44	0.74
	Diagonal	1.44	1.48	0.03	1.47	1.45	1.42	1.44	0.06
		×	×	✓	×	×	×	×	

Q4: Common Subexpressions?

$$\begin{cases} X := AB \\ Y := AB \end{cases} \rightarrow \begin{cases} X := AB \\ Y := X \end{cases}$$

	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy
copy	0.27	0.31	0.36	0.30	0.30	0.26	0.30
direct	0.54	0.6	0.61	0.56	0.58	0.52	0.55
	×	×	×	×	×	×	×

$$\begin{cases} X := AB^{-T}C \\ Y := B^{-1}A^T D \end{cases} \rightarrow \begin{cases} Z := AB^{-T} \\ X := ZC \\ Y := Z^T D \end{cases}$$

BUT

$$X := ABABv \not\rightarrow \begin{cases} Z := AB \\ X := ZZv \end{cases}$$

Q6: Other features?

▶ Code motion

```
for i = 1:n,  
    C = A*B;  
    d[i] = C[i,i];  
end
```

→ ?

```
C = A*B;  
for i = 1:n,  
    d[i] = C[i,i];  
end
```

×

▶ Blocked operands

$$M = \begin{bmatrix} M_1 & 0 \\ 0 & M_2 \end{bmatrix}$$

→ ?

$$\begin{cases} M_1 x_T = y_T \\ M_2 x_B = y_B \end{cases}, \quad \begin{cases} y_T := M_1 x_T \\ y_B := M_2 x_B \end{cases}$$

×, ×

▶ $\text{diag}(\mathbf{A} + \mathbf{B})$ vs. $\text{diag}(\mathbf{A}) + \text{diag}(\mathbf{B})$

→

Armadillo

$\text{diag}(\mathbf{AB})$ vs. ...

→

×

Summary

- ▶ Note: This evaluation was **NOT** a ranking of languages [arXiv:1911.09421]
- ▶ LAMPs are challenging — the optimal solution requires deep & interdisciplinary expertise
- ▶ To users:
Beware! Compilers & languages are great with scalars, not so much with matrices (yet?)
- ▶ To language developers:
Please pay attention to the optimizations we exposed. They arise frequently in the solution of LAMPs.
- ▶ Our approach: (a different talk) “**Linnea**: A linear algebra compiler”
<https://linnea.cs.umu.se/>
[arXiv:1912.12924]