

Current state of programming languages for linear algebra computations

Paolo Bientinesi
pauldj@cs.umu.se

DCSE Summerschool
Numerical Linear Algebra on High Performance Computers
June 07, 2023
TU Delft



High Performance and
Automatic Computing



UMEÅ UNIVERSITY

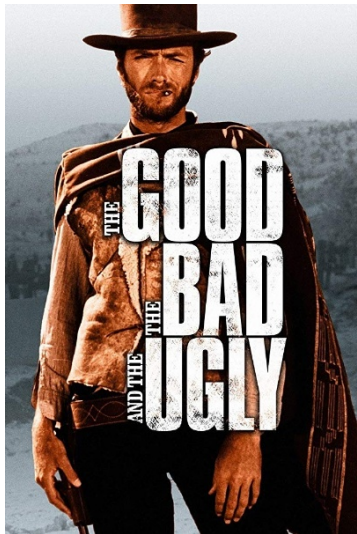


HPC2N

Preview

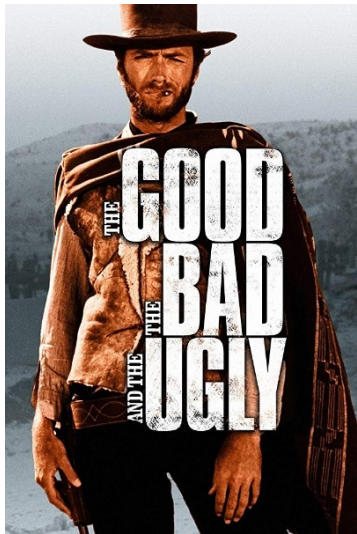


Preview



▶ GOOD: Plenty of excellent libraries for matrix computations

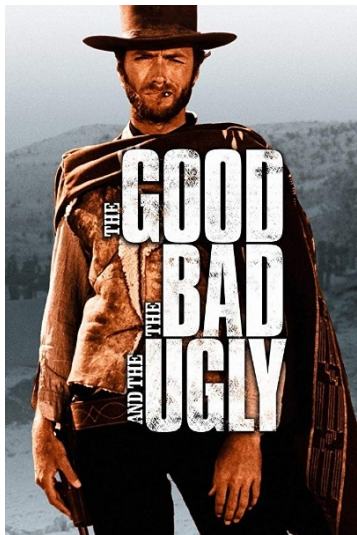
Preview



▶ GOOD: Plenty of excellent libraries for matrix computations

▶ BAD: Are they used (properly)? Not so much

Preview



▶ GOOD: Plenty of excellent libraries for matrix computations

▶ BAD: Are they used (properly)? Not so much

▶ UGLY: Current state of high-level programming languages

Hierarchy of Libraries: 1979–1990

“Basic Linear Algebra Subprograms”

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–

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 sparse and iterative ones

PETSc, Trilinos, ...

ScaLAPACK, PLAPACK, Elemental, ...

LAPACK, Plasma, SuperMatrix, Magma, ...

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

2021 ACM Turing Award



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

- Mathematical expression: $C := C + A * B^T + B * A^T$

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

- Mathematical expression: $C := C + A * B^T + B * A^T$
- Corresponding “naive” 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 )
```

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

- Mathematical expression: $C := C + A * B^T + B * A^T$

- Corresponding “naive” 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” call:

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

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

- Mathematical expression: $C := C + A * B^T + B * A^T$

- Corresponding “naive” 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” call:

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

- 12 args, explicit indexing, not user friendly. Certainly not users' preference (anymore).

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

- Mathematical expression: $C := C + A * B^T + B * A^T$

- Corresponding “naive” 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” call:

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

- 12 args, explicit indexing, not user friendly. Certainly not users' preference (anymore).
Did you catch the indexing mistake?

Users' preference(s)?

$$\mathbf{C} := \mathbf{C} + \mathbf{A} * \mathbf{B}^T + \mathbf{B} * \mathbf{A}^T$$

```
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

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

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

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

Image Restoration

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

Rand. Matrix Inversion

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

Generalized Least Squares

$$b := (X^T M^{-1} X)^{-1} X^T M^{-1} y \quad n > m; M \in \mathbb{R}^{n \times n}, \text{SPD}; X \in \mathbb{R}^{n \times m}; y \in \mathbb{R}^{n \times 1}$$

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

$$x_f := W A^T (A W A^T)^{-1} (b - A x); x_o := W (A^T (A W A^T)^{-1} A x - c)$$

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

$$x := (A^T A + \Gamma^T \Gamma)^{-1} A^T b \quad A \in \mathbb{R}^{n \times m}; \Gamma \in \mathbb{R}^{m \times m}; b \in \mathbb{R}^{n \times 1}$$

Gen. Tikhonov reg.

$$x := x_0 + (A^T P A + Q)^{-1} (A^T P (b - A x_0))$$

LMMSE estimator

$$K_{t+1} := C_t A^T (A C_t A^T + C_z)^{-1}; x_{t+1} := x_t + K_{t+1} (y - A x_t); C_{t+1} := (I - K_{t+1} A) C_t$$

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

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

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

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



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

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

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

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

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

**LINEAR ALGEBRA
MAPPING PROBLEM
“LAMP”**

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

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

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

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

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

**LINEAR ALGEBRA
MAPPING PROBLEM
“LAMP”**

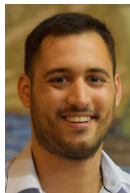
$$\begin{array}{ccccccc} \boxed{y := \alpha x + y} & \boxed{LU = A} & \dots & \boxed{C := \alpha AB + \beta C} \\ \boxed{X := A^{-1} B} & \boxed{C := AB^T + BA^T + C} & \boxed{X := L^{-1} M L^{-T}} & \boxed{QR = A} \end{array}$$

All the aforementioned high-level languages solve LAMPs

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

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

1. Matlab
2. Octave
3. Julia
4. C++ with Armadillo
5. C++ with Eigen
6. NumPy
7. R
8. (TensorFlow)
9. (PyTorch)



Christos Psarras



Aravind Sankaran

- ▶ 12 experiments, each exposing one specific optimization
- ▶ Not a ranking of languages!
- ▶ **“The Linear Algebra Mapping Problem. Current state of linear algebra languages and libraries”**, ACM Transactions on Mathematical Software, Vol. 48(3), pp.1–30, September 2022. [arXiv:1911.09421]
- ▶ *“Benchmarking the Linear Algebra Awareness of TensorFlow and PyTorch”*, Proceedings of iWAPT-22. [arXiv:2202.09888]

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	✓	✓	✓	✓	✓	✓	✓	
<hr/>								
$C = C+A*A'$								

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
$C = A*B$ 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

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
$C = C+AB'+BA'$								

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
$C = C+AB'+BA'$	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

Do they map? $\mathbf{Ax} = \mathbf{b} \quad \equiv \quad \mathbf{x} := \mathbf{A} \backslash \mathbf{b} \quad \not\equiv \quad \text{inv}(\mathbf{A}) * \mathbf{b}$

input

Matlab

Octave

Julia

R

Eigen

Armad.

NumPy

C

$\mathbf{x} := \mathbf{A} \backslash \mathbf{b}$

Do they map? $\mathbf{Ax} = \mathbf{b} \equiv \mathbf{x} := \mathbf{A} \backslash \mathbf{b} \not\equiv \text{inv}(\mathbf{A}) * \mathbf{b}$

input	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$\mathbf{x} := \mathbf{A} \backslash \mathbf{b}$	0.71	0.72	0.63	0.68	0.64	0.63	0.72	0.61

Do they map? $\mathbf{Ax} = \mathbf{b} \equiv \mathbf{x} := \mathbf{A} \backslash \mathbf{b} \not\equiv \text{inv}(\mathbf{A}) * \mathbf{b}$

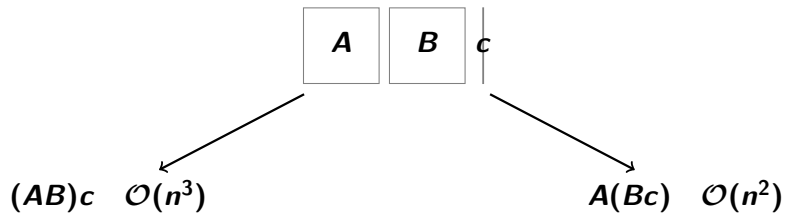
input	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$\mathbf{x} := \mathbf{A} \backslash \mathbf{b}$	0.71	0.72	0.63	0.68	0.64	0.63	0.72	0.61
$\text{inv}(\mathbf{A}) * \mathbf{b}$	1.76	1.82	1.69	2.20	2.21	0.63	2.49	1.71

Do they map? $\mathbf{Ax} = \mathbf{b} \quad \equiv \quad \mathbf{x} := \mathbf{A} \backslash \mathbf{b} \quad \not\equiv \quad \text{inv}(\mathbf{A}) * \mathbf{b}$

input	Matlab	Octave	Julia	R	Eigen	Armad.	NumPy	C
$\mathbf{x} := \mathbf{A} \backslash \mathbf{b}$	0.71	0.72	0.63	0.68	0.64	0.63	0.72	0.61
$\text{inv}(\mathbf{A}) * \mathbf{b}$	1.76	1.82	1.69	2.20	2.21	0.63	2.49	1.71
LinSolve	-	-	-	-	-	✓	-	-

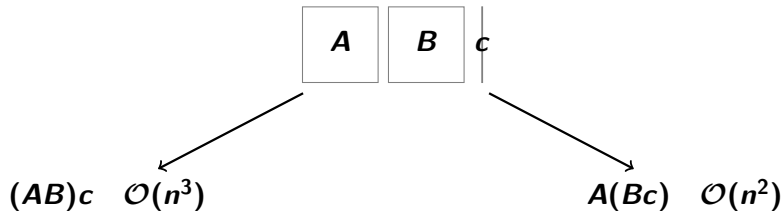
but ...should they map?

Optimal parenthesisation



Matrix product is associative, but its cost is not

Optimal parenthesisation



Matrix product is associative, but its cost is not

$$ABCD = (A(B(CD))) = (A((BC)D)) = (AB)(CD) = (((AB)C)D) = ((A(BC))D)$$

BUT the best parenthesisation depends on the sizes of the matrices A , B , C , and D

Q2: Matrix Chain?

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

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
	$(A*B)*C$	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

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	0.056	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
Matrix chains		×	×	×	×	×	≈	×

Wait!

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		×	×	✓	×	×	≈	×

In practice

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

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

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

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

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

Q3–4: 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

Q3–4: 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

Q3-4: 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

Q3–4: 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

Q3–4: 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
		≈	≈	≈	×	×	≈	×	

Q3-4: 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
		\approx	\approx	\approx	\times	\times	\approx	\times	
<hr/>									
$C = A * B$	-	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
		\times	\times	\checkmark	\times	\times	\times	\times	

Q5: Common Subexpressions?

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

Q5: 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
	×	×	×	×	×	×	×

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

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

BUT

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

Q6–8: 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
```

Q6–8: 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
```

×

Q6–8: 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

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

→ ?

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

Q6–8: 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

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

→ ?

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

×, ×

Q6–8: 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. $\text{diag}(\mathbf{A}) + \text{diag}(\mathbf{B})$

Q6–8: 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 ✓

Q6–8: 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

- ▶ LAMPs are challenging — the optimal solution requires expertise in LA and HPC
[arXiv:1911.09421]

Summary

- ▶ LAMPs are challenging — the optimal solution requires expertise in LA and HPC
[arXiv:1911.09421]
- ▶ To users: Beware!
Compilers & languages are great with scalars, not so much with matrices (yet?)

Summary

- ▶ LAMPs are challenging — the optimal solution requires expertise in LA and HPC
[arXiv:1911.09421]
- ▶ 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.

Summary

- ▶ LAMPs are challenging — the optimal solution requires expertise in LA and HPC
[arXiv:1911.09421]
- ▶ 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:
“*Linnea: Automatic Generation of Efficient Linear Algebra Programs*”, ACM TOMS, 2021
<https://linnea.cs.umu.se/> [arXiv:1912.12924]

Summary

- ▶ LAMPs are challenging — the optimal solution requires expertise in LA and HPC
[arXiv:1911.09421]
- ▶ 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:
“*Linnea: Automatic Generation of Efficient Linear Algebra Programs*”, ACM TOMS, 2021
<https://linnea.cs.umu.se/> [arXiv:1912.12924]

Thank you for your attention!