

Programming languages for matrix computations

Paolo Bientinesi, Henrik Barthels, Christos Psarras

Umeå Universitet

May 10, 2019

University of Manchester



**High Performance and
Automatic Computing**

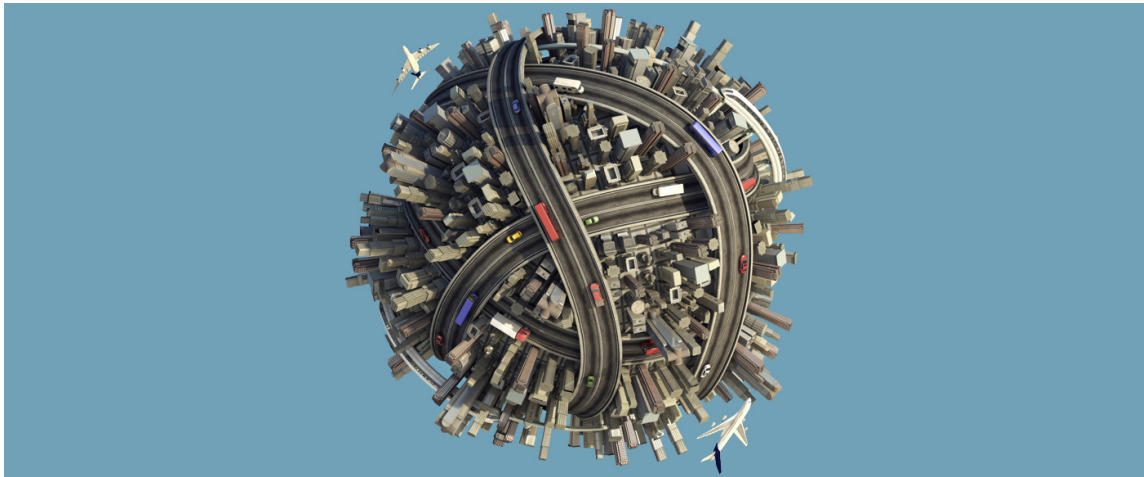


1950s & 1960s

- ▶ Computers difficult to program

BUT

- ▶ “simple” architectures: no memory hierarchy
- ▶ $\text{Cost}(\mathcal{A}lg) \equiv \#\text{operations}(\mathcal{A}lg)$
- ▶ Programs “easy” to optimize — smart code



Cheong Suk-Wai, *A borderless world*

Since the 1970s

- ▶ Libraries
- ▶ Identification, analysis, optimization of **building blocks**
- ▶ EISPACK, LINPACK, BLAS, LAPACK, ...
FFT, numerical integration, ...
- ▶ Convenience, portability, **separation** of concerns

Since the 1980s

- ▶ Memory hierarchies
- ▶ Increasingly complex architectures — caching, prefetching, locality, ...
- ▶ $\text{Cost}(\mathcal{A}lg) \not\equiv \#operations(\mathcal{A}lg)$
- ▶ Programs difficult to optimize — complex, non-portable code
- ▶ Libraries: necessity (wrt performance)

Present: the world of scientific computing

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$$

$$\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 + \|\boldsymbol{\Gamma}\mathbf{x}\|^2$$

LINEAR MIXED MODELS

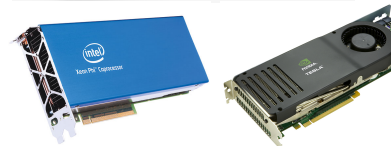
$$V_{LJ} = 4\epsilon \left[\left(\frac{\sigma}{r} \right)^{12} - \left(\frac{\sigma}{r} \right)^6 \right]$$

LENNARD-JONES POTENTIAL

$$i\hbar \frac{\partial}{\partial t} \Psi(\mathbf{r}, t) = \left[\frac{-2\hbar}{2\mu} \nabla^2 + V(\mathbf{r}, t) \right] \Psi(\mathbf{r}, t)$$

SCHRÖDINGER EQN.

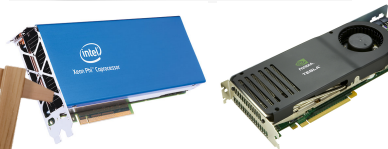
⋮



This talk

Matrix
computations

Matlab, Julia, R,
Eigen, Armadillo, ...



Matrix computations

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

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

Matrix computations (2)

Generalized Least Squares $b := (X^T M^{-1} X)^{-1} X^T M^{-1} y$ $n > m; M \in \mathbb{R}^{n \times n}$, 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. $x := (A^T P A + Q)^{-1} (A^T P b + Q x_0)$ $P \in \mathbb{R}^{n \times n}$, SSPD; $Q \in \mathbb{R}^{m \times m}$, 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 $K_{t+1} := C_t A^T (A C_t A^T + C_z)^{-1}$; $x_{t+1} := x_t + K_{t+1} (y - A x_t)$; $C_{t+1} := (I - K_{t+1} A) C_t$

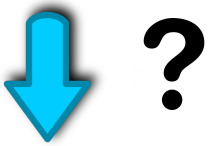
LMMSE estimator $x_{\text{out}} = C_X A^T (A C_X A^T + C_Z)^{-1} (y - A x) + x$

LMMSE estimator $x_{\text{out}} := (A^T C_Z^{-1} A + C_X^{-1})^{-1} A^T C_Z^{-1} (y - A x) + x$

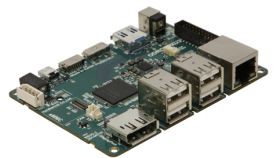
$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y$$

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

$$E := Q^{-1} U (I + U^T Q^{-1} U)^{-1} U^T \quad \dots$$



2-step solution



- MUL
- ADD
- MOV
- MOVAPD
- VFMADDPD ...

$$x := A(B^T B + A^T R^T \Lambda R A)^{-1} B^T B A^{-1} y$$

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

$$E := Q^{-1} U (I + U^T Q^{-1} U)^{-1} U^T \quad \dots$$

**LINEAR
ALGEBRA
MAPPING
PROBLEM
(LAMP)**

$$y := \alpha x + y$$

$$:= \alpha A B + \beta C$$

$$X := A^{-1} B$$

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

$$X := L^{-1} M L^{-T}$$

$$Q R = A$$

BLAS



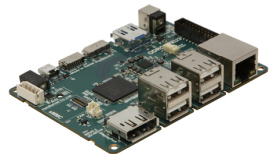
LAPACK



MKL



...



MUL ADD MOV

MOVAPD

VFMADDPD ...

Linear Algebra Mapping Problem

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

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 easy
- ▶ Optimality NP complete reduction from Ensemble Computation

LAMPs are everywhere

- ▶ Kernels

$\mathcal{E} : \{C := A * B + C\}$ $\mathcal{K} : \text{processor's ISA}$ $\mathcal{M} : \text{execution time}$

- ▶ Automatic code generation

ATLAS [Whaley 2001], FLAME [Gunnels 2001], Build-To-Order BLAS [Siek 2008], LGEN [Spampinato 2014], ...

- ▶ “Matrix Chain Problem”

$\mathcal{E} : \{X := M_1 M_2 \cdots M_k\}$ $\mathcal{K} : \{C := A * B\}$ $\mathcal{M} : \text{\#flops}$

- ▶ Applications

$\mathcal{E} : \text{explicit assignments}$ $\mathcal{K} : \text{libraries}$ $\mathcal{M} : \text{execution time} + \dots$

Problem acknowledged, yet overlooked

Libraries exist (a myriad of them) — How do I express my problem to use them?

- ▶ By hand. C/Fortran. Patience. Expertize.

computer efficiency! ... but human productivity?

- ▶ High-level language. Matlab, Julia, R, Eigen, Armadillo, NumPy, ... Quick prototyping.

human productivity! ... but computer efficiency?

Slow solutions → “2-language problem”

2-language problem

- ▶ [arXiv:1904.12380]: Softmax function.
Original implementation: Eigen.
Reimplemented with explicit calls to BLAS and MKL.

- ▶ [F1000Res.2016]: The GenABEL Project for statistical genomics.
Original implementation: R.
Kernels reimplemented in C.

- ▶ Convex, non-smooth optimization in image processing.
Original implementation: Python.
Looking for alternatives.

State of the art

Many languages & environments that allow a high level of abstraction

Matlab, Julia, R, Eigen, Armadillo, NumPy

Why? Popularity, expressiveness, (performance...)

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

```
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
```

Do they map?

matrix products

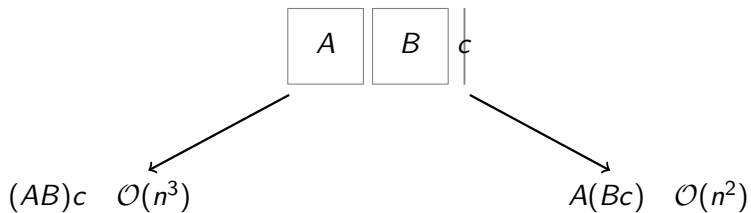
	Matlab	Julia	R	Eigen	Armad.	NumPy	C
$C = AB+C$	0.28	0.31	0.30	0.29	0.28	0.29	0.26
$C = AB$...						
$C = \alpha AB$							
$C = \alpha AB + \beta C$							
GEMM	✓	✓	✓	✓	✓	✓	
$C = C+AA'$	0.17	0.22	0.31	0.29	0.17	0.18	0.13
SYRK	✓	✓	×	×	✓	✓	
$C = C+AB'+BA'$	0.56	0.70	0.61	0.57	0.56	0.58	0.27
SYR2K	×	×	×	×	×	×	

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

	Matlab	Julia	R	Eigen	Armad.	NumPy	C
$x := A \setminus b$	0.70	0.62	0.67	0.63	0.62	0.65	0.61
$\text{inv}(A) * b$	1.74	1.45	2.20	2.20	0.62	2.22	1.71
LinSolve	-	-	-	-	✓	-	-

... should they map?

Parenthesisation



Product is associative, but cost is not

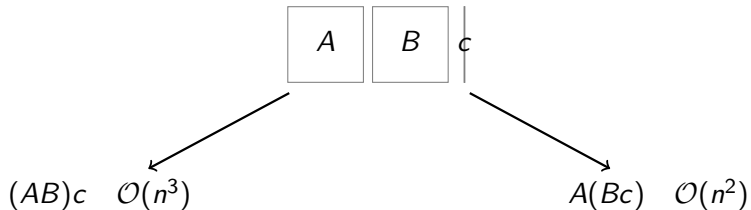
Matrix Chain Algorithm

$O(k \log k)$ Hu & Shing 1982; $O(k^3)$ dynamic programming

Matrix Chain?

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

	Matlab	Julia	R	Eigen	Armad.	NumPy
LtR no par.	0.056	0.055	0.061	0.058	0.056	0.055
LtR guided	0.056	0.055	0.061	0.058	0.056	0.055
RtL no par.	0.42	0.42	0.44	0.42	0.055	0.42
RtL guided	0.055	0.054	0.059	0.056	0.055	0.056
Mix no par.	0.32	0.33	0.33	0.35	0.31	0.33
Mix guided	0.21	0.22	0.22	0.23	0.20	0.22
Matrix chains	×	×	×	×	≈	×



In practice

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

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

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

(e.g., $A \rightarrow Q^T D Q$, $A \rightarrow LU$)

(GEMM, TRMM, SYMM, ...)

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

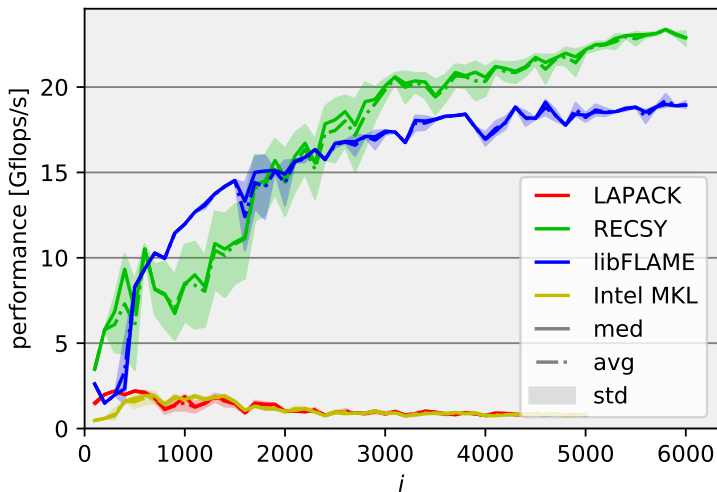
$$\operatorname{argmin}_{\mathcal{A}} (\text{FLOPs}(\mathcal{A})) \neq \operatorname{argmin}_{\mathcal{A}} (\text{time}(\mathcal{A}))$$

⇒ **Performance prediction:** efficiency

Challenge: Not all flops were created equal

[arXiv:1504.08035]

Triangular Sylvester equation – all libs performs the same # of flops



Challenge: Parallelism

multi-threaded libs, explicit multi-threading, runtime, hybrid, offloading

$$X := A((B^T C^{-T})D) \quad \text{vs.} \quad X := (AB^T)(C^{-T}D) \quad \text{vs.} \quad \dots$$

⇒ **Performance prediction:** efficiency, scalability

Properties?

Operation	Property	Matlab	Julia	R	Eigen	Armad.	NumPy	C
Linear System	-	0.70	0.62	0.67	0.63	0.62	0.65	0.61
	Symmetric	0.71	0.62	0.69	N/A	0.62	0.65	0.46
	SPD	0.41	0.60	0.63	N/A	0.34	0.62	0.31
	Triangular	0.03	0.03	0.63	N/A	0.62	0.65	0.03
	Diagonal	0.03	0.01	0.63	N/A	0.03	0.62	0.001
		≈	≈	×	×	≈	×	
Multiplication	Triangular	1.44	0.75	1.47	1.45	1.44	1.44	0.74
	Diagonal	1.44	0.03	1.47	1.45	1.42	1.44	0.06
		×	✓	×	×	×	×	

Challenge: Inference of properties

- ▶ **easy** $E := L_1 * U^T * L_2$ triangular(E) ?
- ▶ **hard** $\lambda(L^{-T}AL^{-1})$ symmetric($L^{-T}AL^{-1}$) ?
- ▶ **impossible?** $E := Q^{-1}U(I + U^TQ^{-1}U)^{-1}U^T$ properties($I + U^TQ^{-1}U$) ?

⇒ **Symbolic analysis:** pattern matching

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

⇒ **Pattern matching**

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

Common Subexpressions?

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

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

Other features

▶ Code motion

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

→ ?

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

×

▶ Blocked operands

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

→ ?

$$Mx = y, \quad y = Mx$$

×

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

→

Armadillo

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

→

×

Summary

- ▶ LAMP is challenging — lots of expertise needed; interdisciplinary
- ▶ Compilers are great with scalars, not so much with matrices
- ▶ **Linnea**: A linear algebra compiler

Linear algebra knowledge: operators, identities, theorems

- Distributivity, commutativity, partitionings, ...
- $((QR)^T QR)^{-1}(QR)^T y \rightarrow (R^T Q^T QR)^{-1}R^T Q^T y \rightarrow R^{-1}R^{-T}R^T Q^T y \rightarrow R^{-1}Q^T y$
- $\text{SPD}(A) \rightarrow \text{SPD}(A_{BR} - A_{BL}A_{TL}^{-1}A_{BL}^T)$ Schur complement
- ...

Example

$$w := AB^{-1}c, \quad \text{SPD}(B)$$

Naive

```
w = A*inv(B)*c
```

Recommended

```
w = A*(B\c)
```

Expert

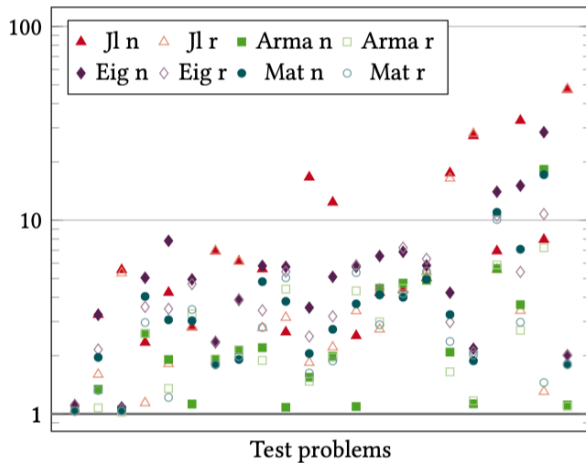
```
L = Chol(B)
w = A*(L'\(L\c))
```

Generated — “Linnea”

```
m10 = A; m11 = B; m12 = c;
potrf!('L', m11)
trsv!('L', 'N', 'N', m11, m12)
trsv!('L', 'T', 'N', m11, m12)
m13 = Array{Float64}(10)
gemv!('N', 1.0, m10, m12, 0.0, m13)
w = m13
```


Results — applications

Linnea's speedups

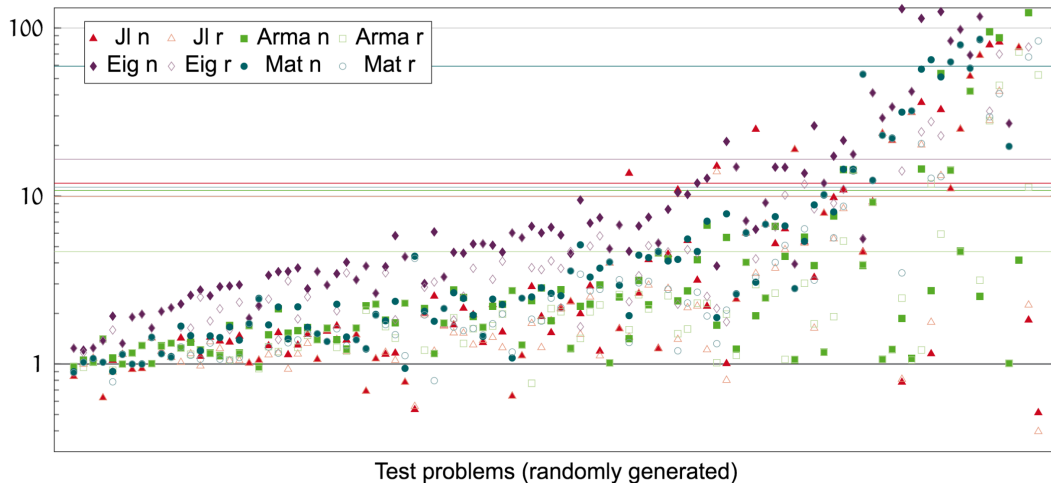


JI: Julia, **Arma:** Armadillo, **Eig:** Eigen, **Mat:** Matlab.

n/r: naive/recommended implementation

Results — random expressions

Linnea's speedups



JI: Julia, **Arma:** Armadillo, **Eig:** Eigen, **Mat:** Matlab.

n/r: naive/recommended implementation