

High-Performance & Automatic Computing

Paolo Bientinesi
pauldj@cs.umu.se

December 5, 2018
Umeå Universitet



**High Performance and
Automatic Computing**

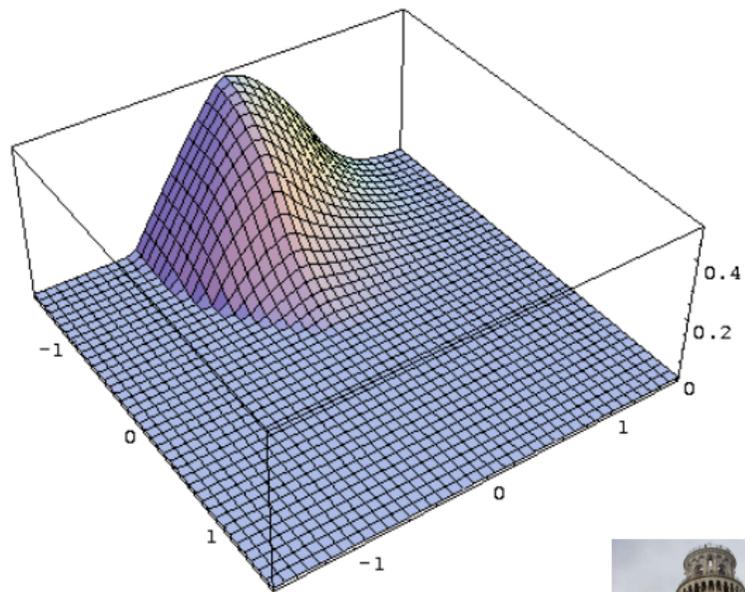
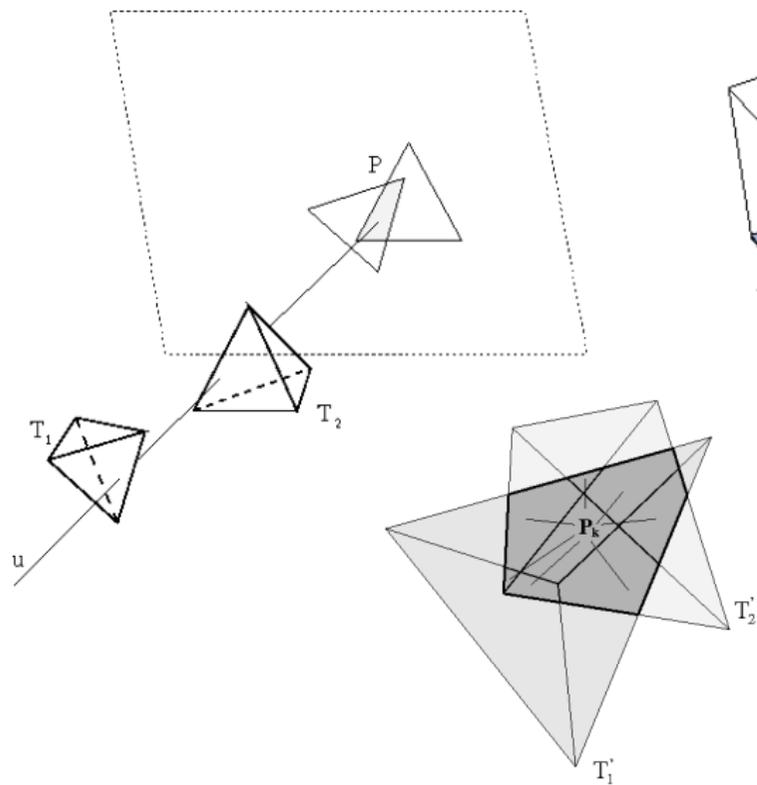
- ▶ **About me**
- ▶ **Molecular dynamics**
- ▶ **Mixed precision calculations**
- ▶ **Accuracy & performance modeling**
- ▶ **Linear Algebra, applications**
- ▶ **Tensor operations**

About me: Italy



Tuscany





Algorithm $LU = A$

Partition $A \rightarrow \left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right), L \rightarrow \left(\begin{array}{c|c} L_{TL} & 0 \\ \hline L_{BL} & L_{BR} \end{array} \right), U \rightarrow \left(\begin{array}{c|c} U_{TL} & U_{TR} \\ \hline 0 & U_{BR} \end{array} \right)$

where A_{TL}, L_{TL}, U_{TL} , are 0×0

While $m(A_{TL}) \leq m(A)$ **do**

Repartition

$$\left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right), \left(\begin{array}{c|c} L_{TL} & 0 \\ \hline L_{BL} & L_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c|c|c} L_{00} & 0 & 0 \\ \hline l_{10}^T & 1 & 0 \\ \hline L_{20} & l_{21} & L_{22} \end{array} \right),$$

$$\left(\begin{array}{c|c} U_{TL} & U_{TR} \\ \hline 0 & U_{BR} \end{array} \right) \rightarrow \left(\begin{array}{c|c|c} U_{00} & u_{01} & U_{02} \\ \hline 0 & v_{11} & u_{12}^T \\ \hline 0 & 0 & U_{22} \end{array} \right)$$

where $\alpha_{11}, 1, v_{11}$ are scalars

$$v_{11} := \alpha_{11} - l_{10}^T u_{01}$$

$$u_{12}^T := a_{12}^T - l_{10}^T U_{02}$$

$$l_{21} := (a_{21} - L_{20} u_{01}) / v_{11}$$

Continue with

$$\left(\begin{array}{c|c} A_{TL} & A_{TR} \\ \hline A_{BL} & A_{BR} \end{array} \right) \leftarrow \left(\begin{array}{c|c|c} A_{00} & a_{01} & A_{02} \\ \hline a_{10}^T & \alpha_{11} & a_{12}^T \\ \hline A_{20} & a_{21} & A_{22} \end{array} \right), \left(\begin{array}{c|c} L_{TL} & 0 \\ \hline L_{BL} & L_{BR} \end{array} \right) \leftarrow \left(\begin{array}{c|c|c} L_{00} & 0 & 0 \\ \hline l_{10}^T & 1 & 0 \\ \hline L_{20} & l_{21} & L_{22} \end{array} \right), \dots$$

endwhile

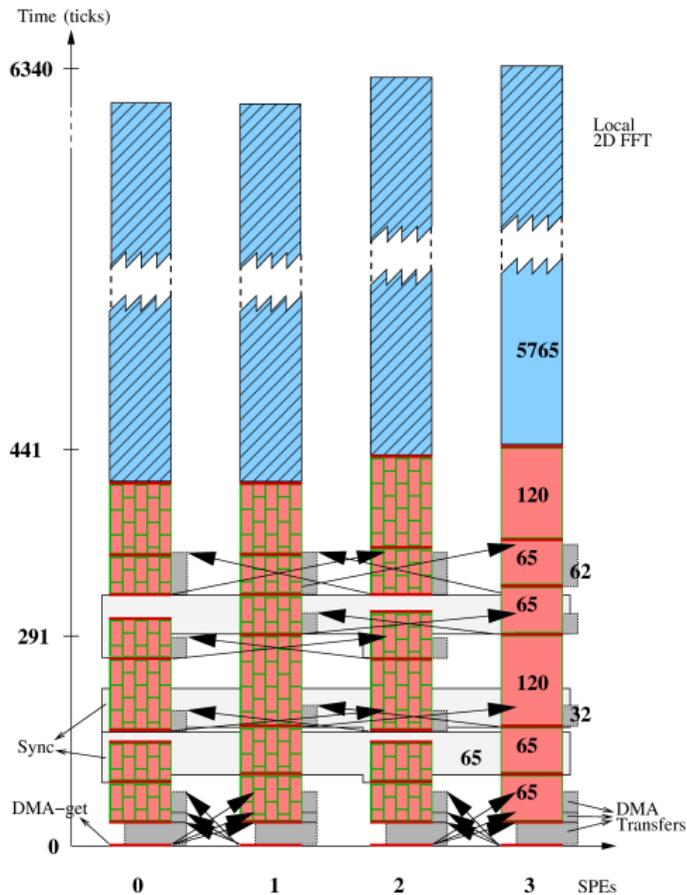
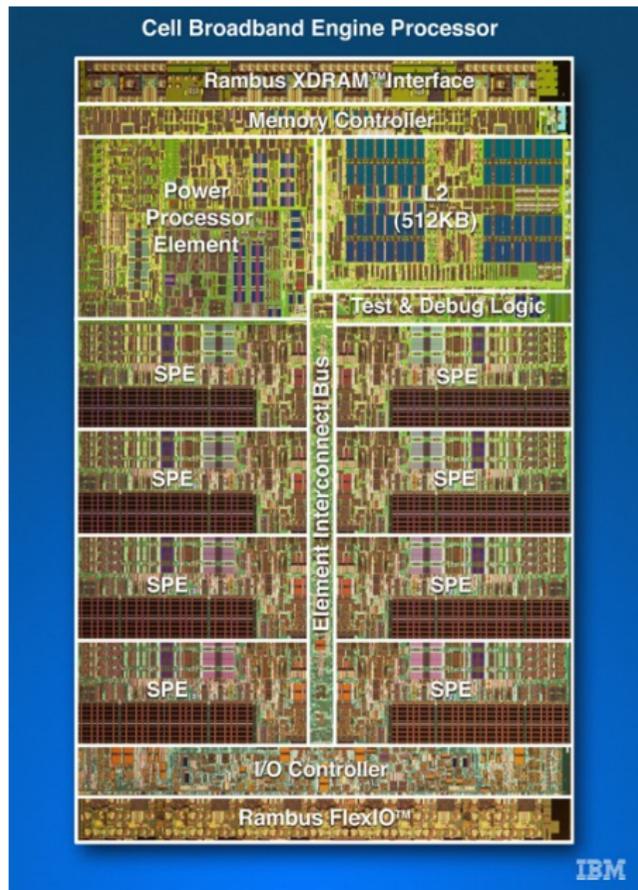


► Automatic backward analysis?

$v_{11} := \alpha_{11} - l_{10}^T u_{01}$	$\check{v}_{11} + \delta\alpha_{11} = \alpha_{11} - l_{10}^T u_{01}$ $\wedge \delta\alpha_{11} \leq \gamma_{k+1} (l_{10}^T u_{01} + \check{v}_{11}) \quad \text{Th. 3.9 R2-F}$
$u_{12}^T := a_{12}^T - l_{10}^T U_{02}$	$\check{u}_{12}^T + \delta\alpha_{12}^T = a_{12}^T - l_{10}^T U_{02}$ $\wedge \delta\alpha_{12}^T \leq \gamma_{k+1} (l_{10}^T U_{02} + \check{u}_{12}^T) \quad \text{Th. 5.1 R2-F}$
$l_{21} := (a_{21} - L_{20} u_{01}) / v_{11}$	$\check{l}_{21} v_{11} + \delta\alpha_{21} = a_{21} - L_{20} u_{01}$ $\wedge \delta\alpha_{21} \leq \gamma_{k+1} (\check{L}_{20} \check{u}_{01} + \check{l}_{21} v_{11}) \quad \text{Th. 5.1 R4-F}$

- Symmetric eigenproblem $AX = X\Lambda$
Multiple Relatively Robust Representations (MRRR)

Duke University – Cell Processor & FFTs



Germany – RWTH Aachen University



The world of scientific computing

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

LENNARD-JONES POTENTIAL

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

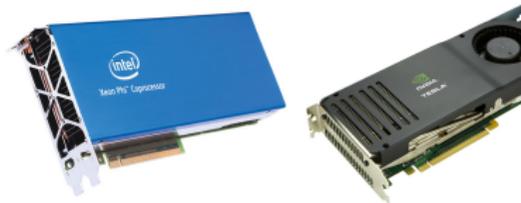
SCHRÖDINGER EQN.

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

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

LINEAR MIXED MODELS

⋮



The world of scientific computing

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

LENNARD-JONES POTENTIAL

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

SCHRÖDINGER EQN.

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

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

LINEAR MIXED MODELS

⋮



The world of scientific computing

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

LENNARD-JONES POTENTIAL

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

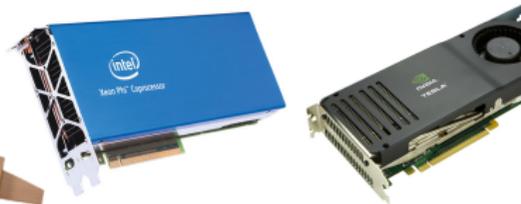
SCHRÖDINGER EQN.

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

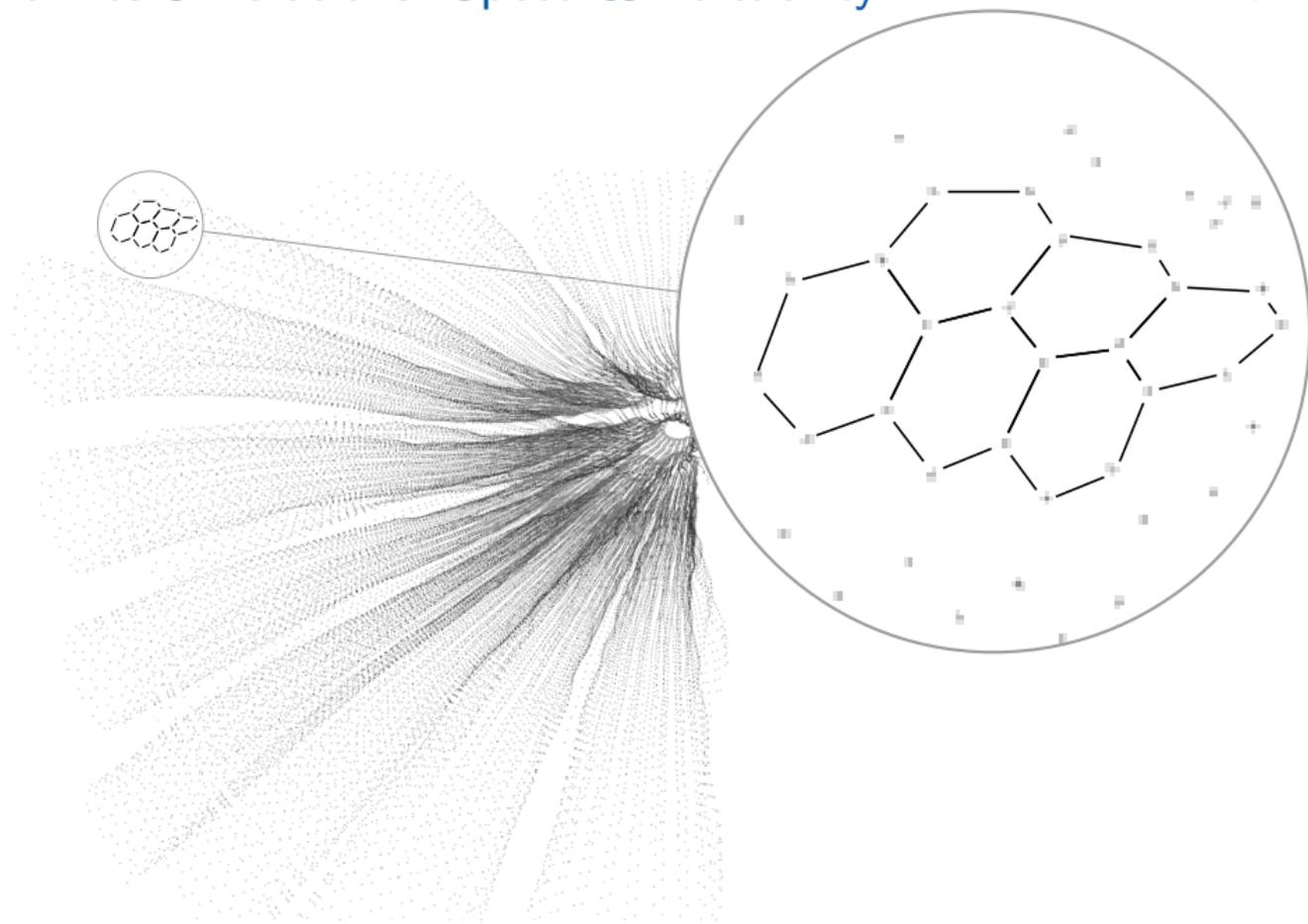
$$\min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|^2 + \|\boldsymbol{\Gamma}\mathbf{x}\|^2$$

LINEAR MIXED MODELS

⋮



Domain Specific Languages, Compilers, Libraries



Pair vs Many-body potentials

$$V = \sum_i \sum_j V(i, j)$$

$$F_i = -\nabla_{x_i} V$$

```
for i in atoms:
    for j in neighbors(i):
        V += V(i, j)
        F[i] -= dVdi(i, j)
        F[j] -= dVdj(i, j)
```

$$V = \sum_i \sum_j V(i, j, \sum_k f(i, j, k))$$

```
for i in atoms:
    for j in neighbors(i):
        tmp = 0
        for k in neighbors(i):
            tmp += f(i, j, k)
        V += V(i, j, tmp)
        F[i] -= dVdi(i, j, tmp)
        F[j] -= dVdj(i, j, tmp)
        tmp = dVdf(i, j, tmp)
        for k in neighbors(i):
            F[i] -= tmp * dfdi(i, j, k)
            F[j] -= tmp * dfdj(i, j, k)
            F[k] -= tmp * dfdk(i, j, k)
```

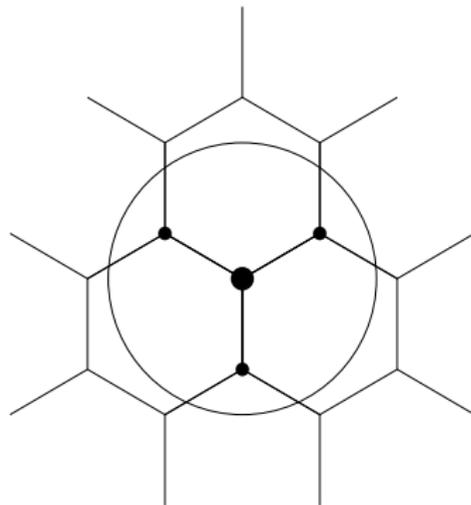
Vectorization

x86	SSE	128 bit
	AVX(2)	256 bit
	AVX-512 (IMCI)	512 bit
ARM	NEON	128 bit
	SVE	up to 2048 bit
POWER	AltiVec/VMX/VSX	128 bit
	QPX	256 bit
SPARC	HPC-ACE	128 bit
	HPC-ACE2	256 bit





$$V = \sum_i \sum_j V(i, j, \sum_k f(i, j, k))$$



The Vectorization of the Tersoff Multi-Body Potential: An Exercise in Performance Portability

Markus Höhnerbach
RWTH Aachen University

Ahmed E. Ismail
RWTH Aachen University,
West Virginia University

Paolo Bientinesi
RWTH Aachen University



Manual Optimization: AIREBO

$$\begin{aligned} V &= \sum_i \sum_j^i V(i, j, \sum_k^i f(i, k), \sum_l^j f(j, l), \sum_k^i \sum_l^j g(i, j, k, l)) \\ &+ \sum_i \sum_j^i \sum_k^i \sum_l^j V'(i, j, k, l) \\ &+ \sum_i \sum_j^{i'} V''(i, j, \max_k g(i, j, k, l) \forall k, l, \sum_k^i f(i, k), \sum_l^j f(j, l), \sum_k^i \sum_l^j g(i, j, k, l)) \end{aligned}$$

Optimizing AIREBO:

Navigating the Journey from Complex Legacy Code to High Performance

Markus Höhnerbach, Paolo Bientinesi

submitted

PotC: A domain-specific language + compiler

Goal: Good –not best– performance on different platforms and arbitrary potentials

```
energy 1 / 2 * sum(i : all_atoms)
  sum(j : neighbors(i, R(i, j, j) + D(i, j, j))) V(i, j);
function V(i : atom; j : atom) = f_C(i, j, r(i, j)) *
  (f_R(i, j, r(i, j)) + b(i, j) * f_A(i, j, r(i, j)));
function f_C(i : atom_type; j : atom_type; k : atom_type; r : distance) =
  implicit(i : i; j : j; k : k) piecewise(r <= R - D : 1;
  R - D < r < R + D: 1 / 2 - 1 / 2 * sin(pi / 2 * (r - R) / D); r >= R + D : 0);
# ...
function b(i : atom; j : atom) = (1 + beta(i, j) ^ n(i, j) *
  zeta(i, j) ^ n(i, j)) ^ (-1 / (2 * n(i, j)));
function zeta(i : atom; j : atom) =
  sum(k : neighbors(i, R(i, j, k) + D(i, j, k), j))
  implicit(i : i; j : j; k : k) f_C(r(i, k)) * g(theta(j, i, k)) *
    exp(lambda_3 ^ m * (r(i, j) - r(i, k)) ^ m);
parameter A(i : atom_type; j : atom_type) = file(1);
# ...
parameter gamma(i : atom_type; j : atom_type; k : atom_type) = file(5);
```

- ▶ About me
- ▶ Molecular dynamics
- ▶ **Mixed precision calculations**
- ▶ Accuracy & performance modeling
- ▶ Linear Algebra, applications
- ▶ Tensor operations

Tridiagonal Eigenproblem

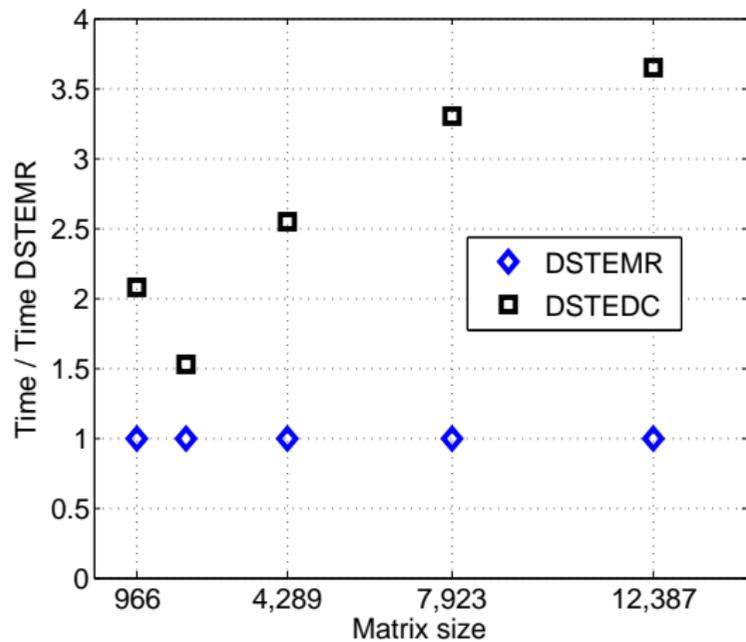
Algorithms			
1958	Bisection + Inverse Iteration (BI)	subsets	$O(kn^2)$
1961	QR	robust, accurate	$O(n^3)$
1981	Divide & Conquer (DC)	BLAS3, accurate	$O(n^3)$
1997	MR3 / MRRR	subsets, fast	$O(kn)$

Tridiagonal Eigenproblem

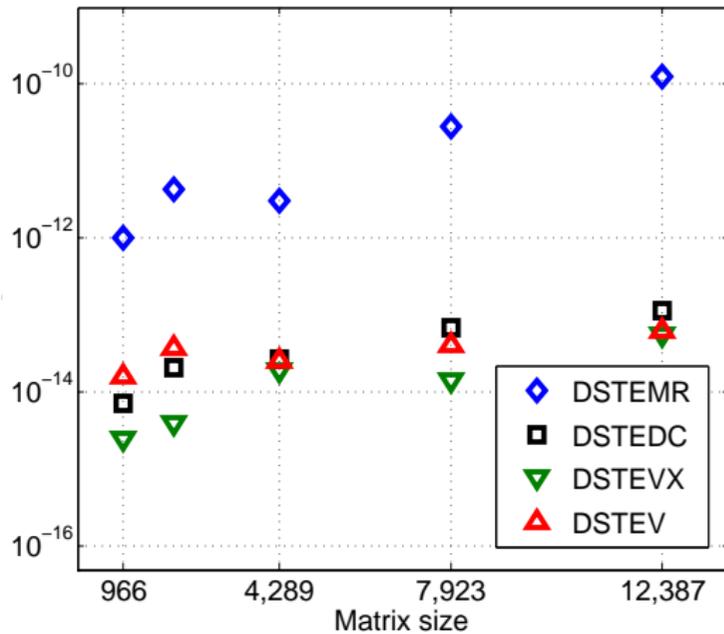
Algorithms			
1958	Bisection + Inverse Iteration (BI)	subsets	$O(kn^2)$
1961	QR	robust, accurate	$O(n^3)$
1981	Divide & Conquer (DC)	BLAS3, accurate	$O(n^3)$
1997	MR3 / MRRR	subsets, fast	$O(kn)$

- ▶ I.Dhillon & B.Parlett (U.C. Berkeley)
- ▶ no reorthogonalization $\Rightarrow k$ eigenpairs in $O(nk)$ operations
- ▶ 1) eigenvalues (*dqds*, *bisection*)
2) eigenvectors + eigenvalues

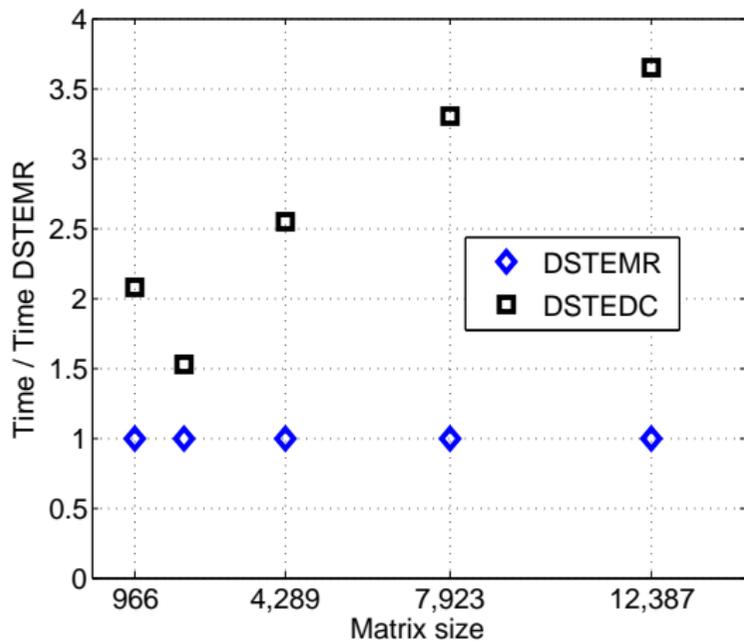
MR3 speedups



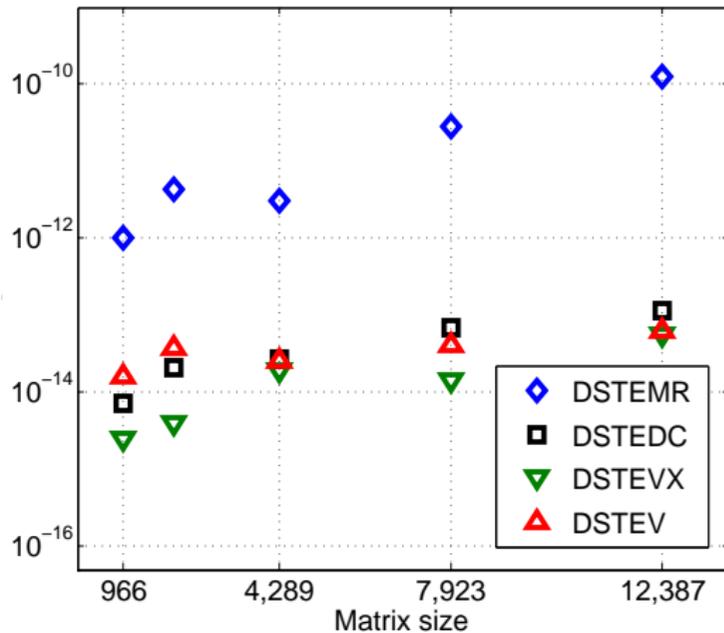
$$\max_{i \neq j} |z_i^H z_j|$$



MR3 speedups



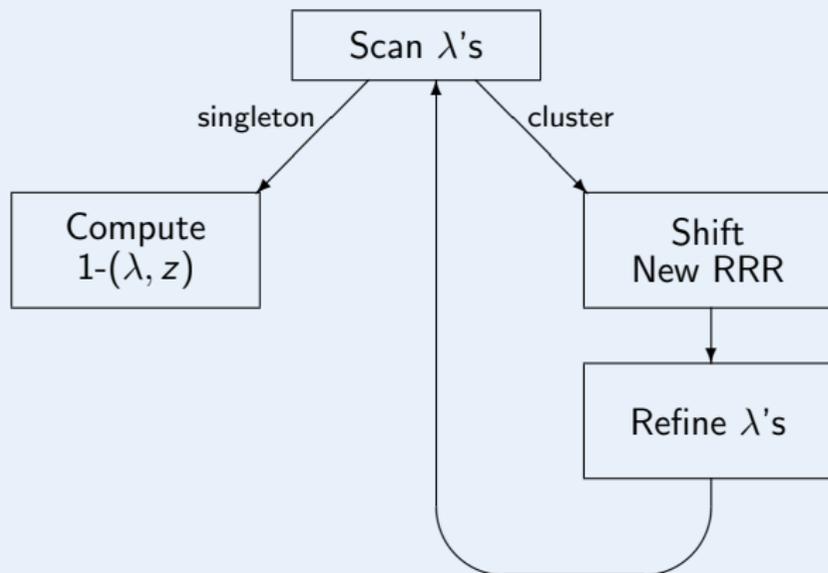
$$\max_{i \neq j} |z_i^H z_j|$$



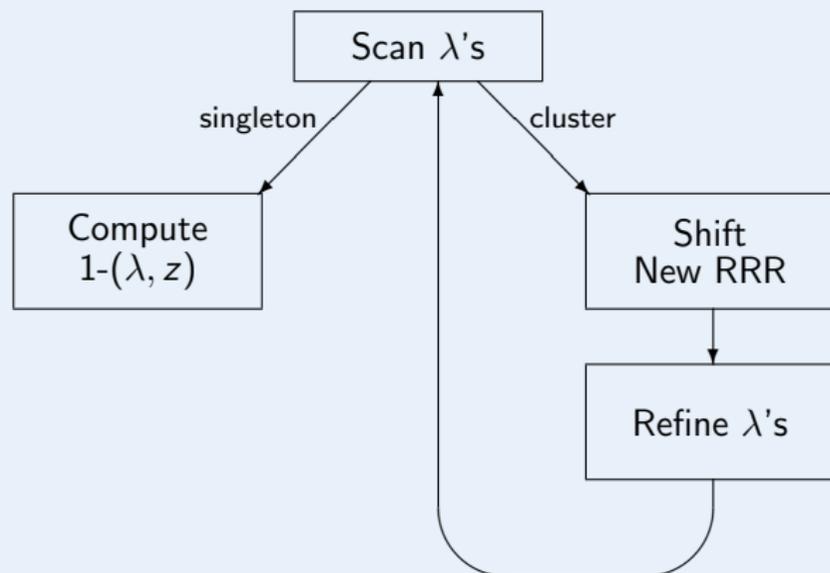
Idea: Trade speed for accuracy

...diametrically opposite approach wrt to current trends

Workflow



Workflow



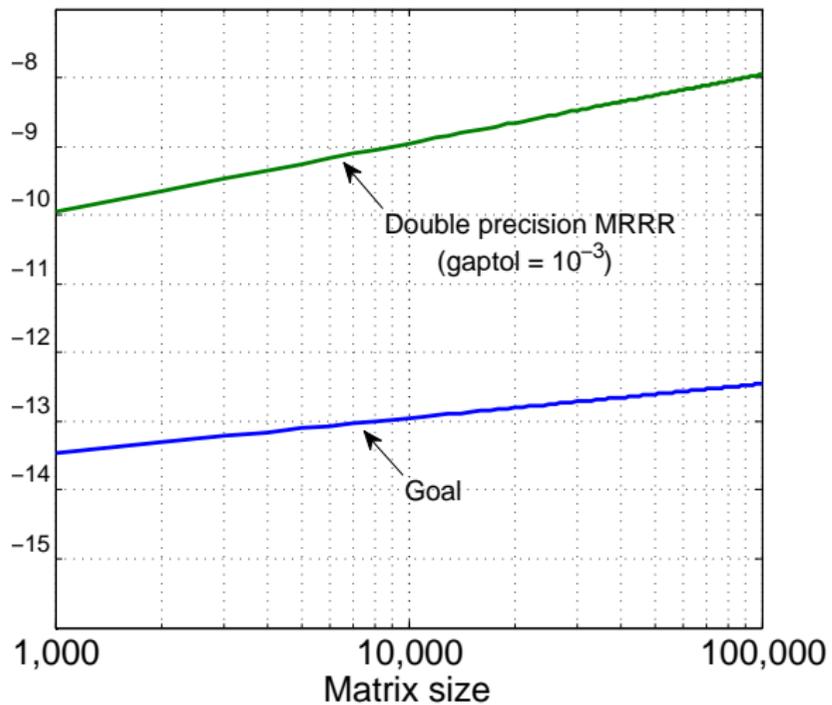
if $\min_{j \neq i} \frac{|\lambda_i - \lambda_j|}{|\lambda_i|} \geq \text{gaptol}$
 then λ_i is a singleton

- ▶ gaptol : “small” (loose) \Rightarrow less clustering, better robustness, more parallelism, more BX
- ▶ gaptol : “large” (strict) \Rightarrow more work, deeper trees, better orthogonality (?), more failures

$$\text{orthogonality} \lesssim \frac{n \varepsilon}{\text{gaptol}}$$

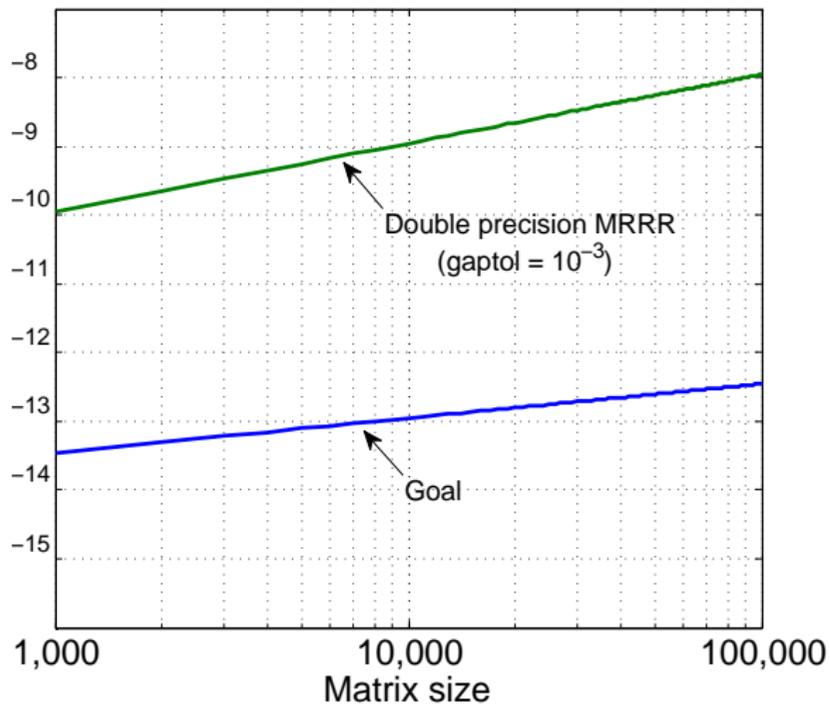
$$\text{orthogonality} \lesssim \frac{n \varepsilon}{\text{gaptol}}$$

Double Precision

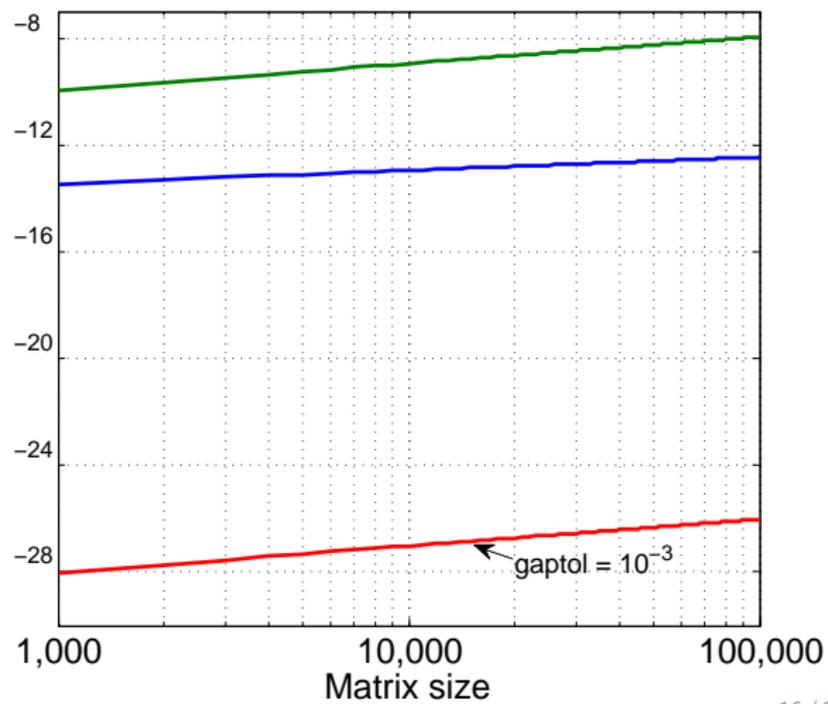


$$\text{orthogonality} \lesssim \frac{n \varepsilon}{\text{gaptol}}$$

Double Precision

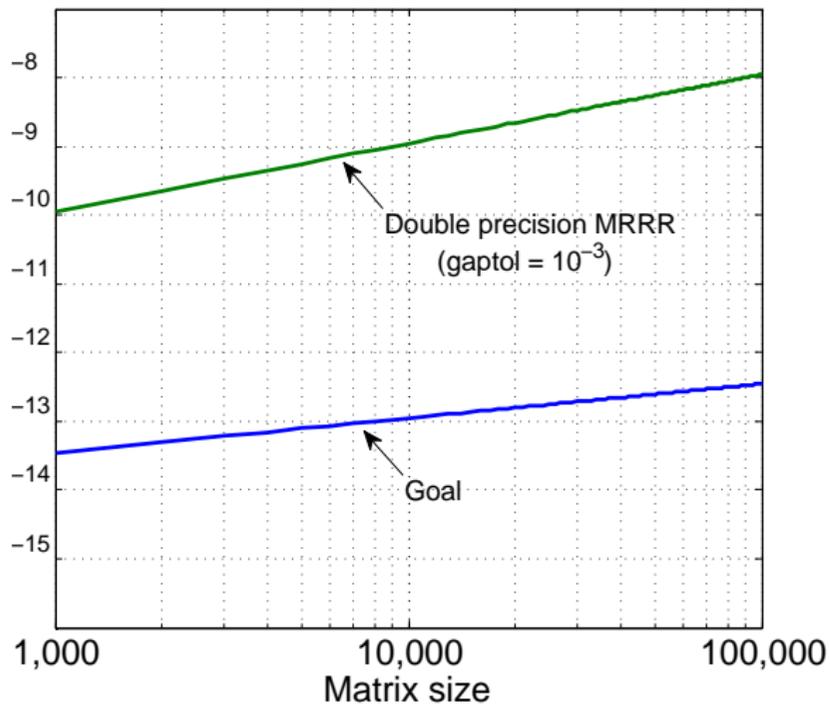


Quad Precision

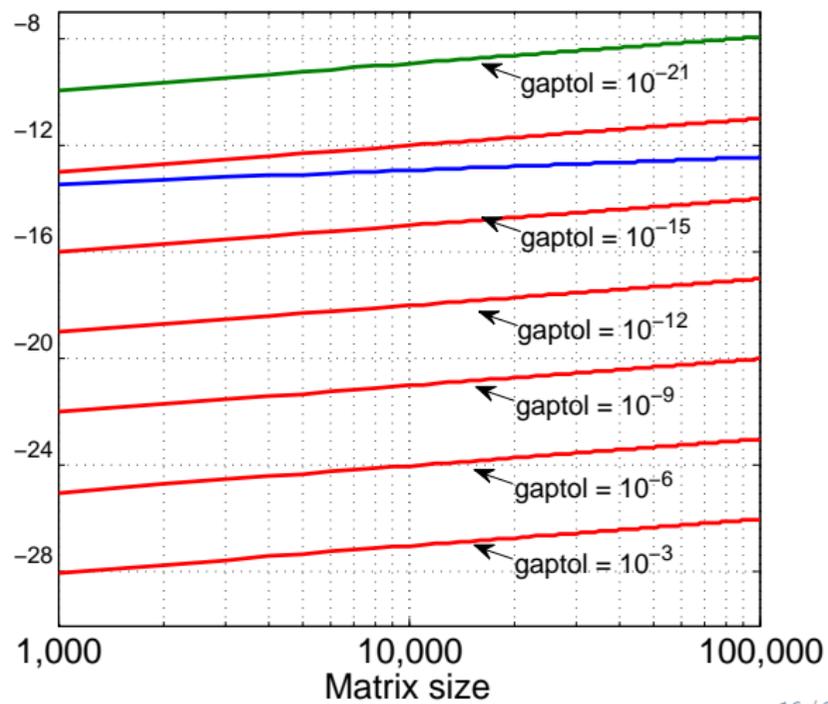


$$\text{orthogonality} \lesssim \frac{n \varepsilon}{\text{gaptol}}$$

Double Precision

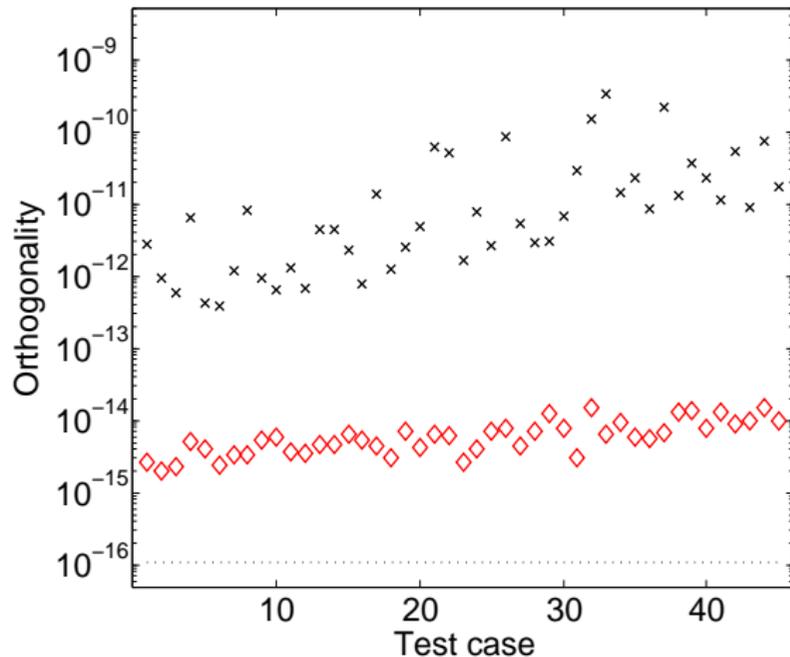
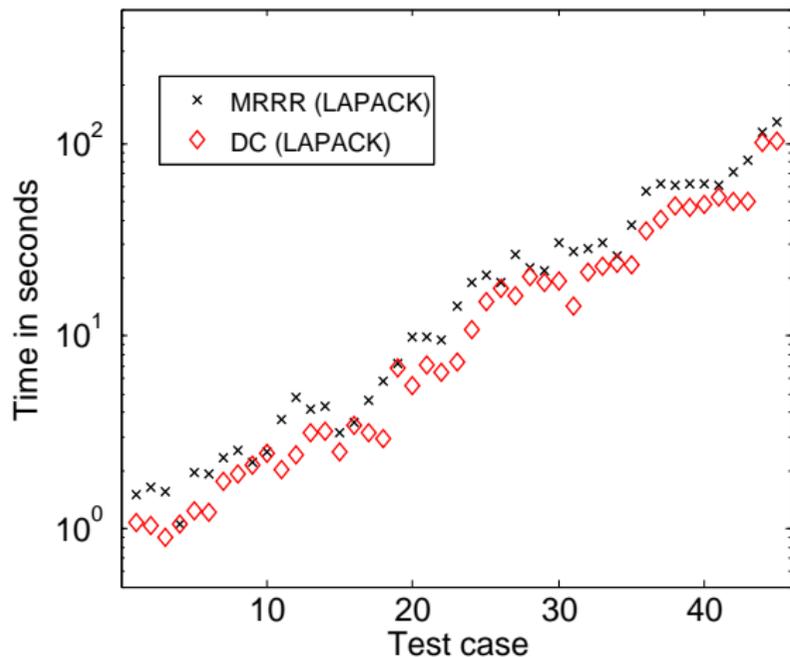


Quad Precision



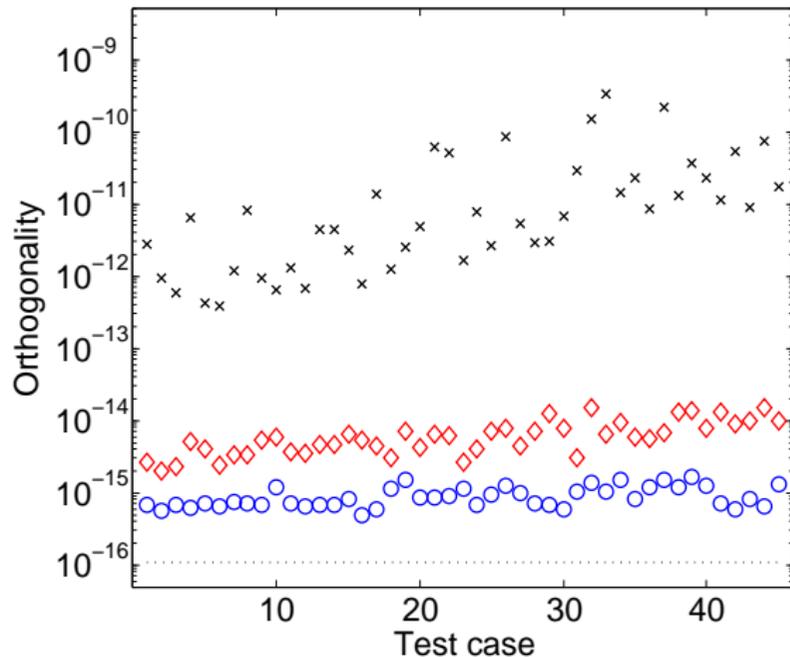
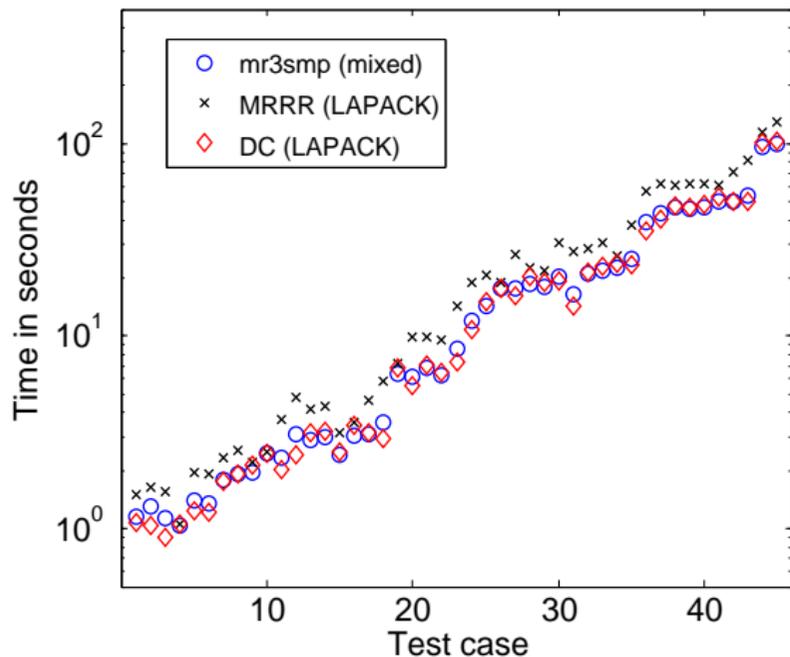
Results: Dense Real Symmetric Eigenproblems

32 cores — 45 application matrices $n \in [1000..8000]$



Results: Dense Real Symmetric Eigenproblems

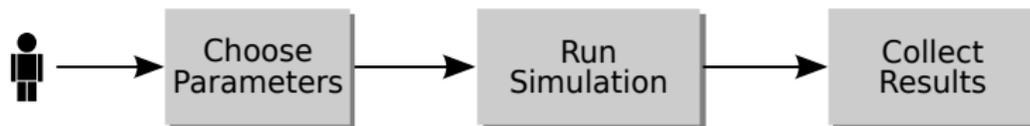
32 cores — 45 application matrices $n \in [1000..8000]$



- ▶ About me
- ▶ Molecular dynamics
- ▶ Mixed precision calculations
- ▶ **Accuracy & performance modeling**
- ▶ Linear Algebra, applications
- ▶ Tensor operations

Numerical simulations

Typical workflow

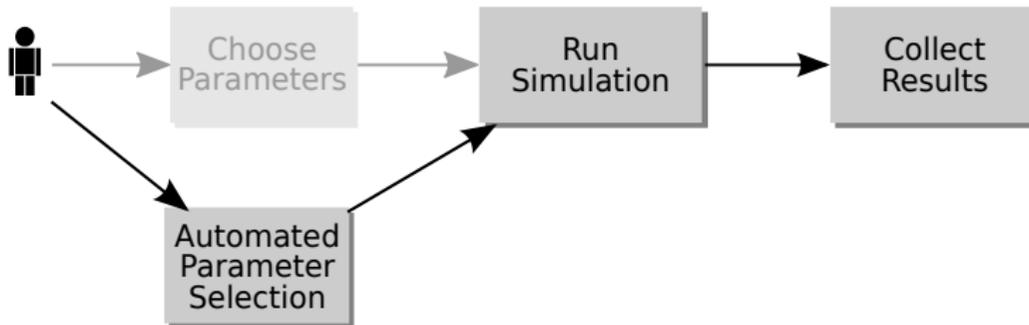


Numerical simulations

Typical workflow



Proposed workflow



Problem: Parameters selection

- ▶ Optimization problem — p : input parameters

$$\min_p \text{time}(\text{resources}, p) \quad \text{subject to} \quad \text{accurate}(p)$$

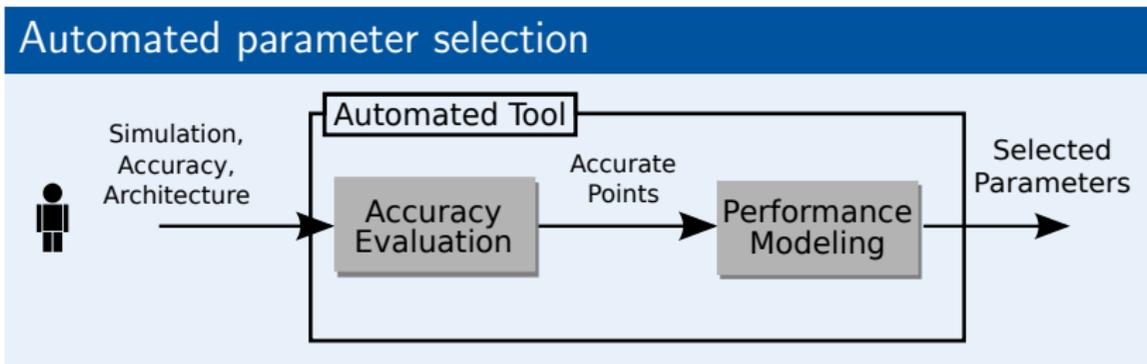
- ▶ Easily millions of combinations of parameters
- ▶ What is the “fastest” combination that leads to sufficiently accurate results?

Problem: Parameters selection

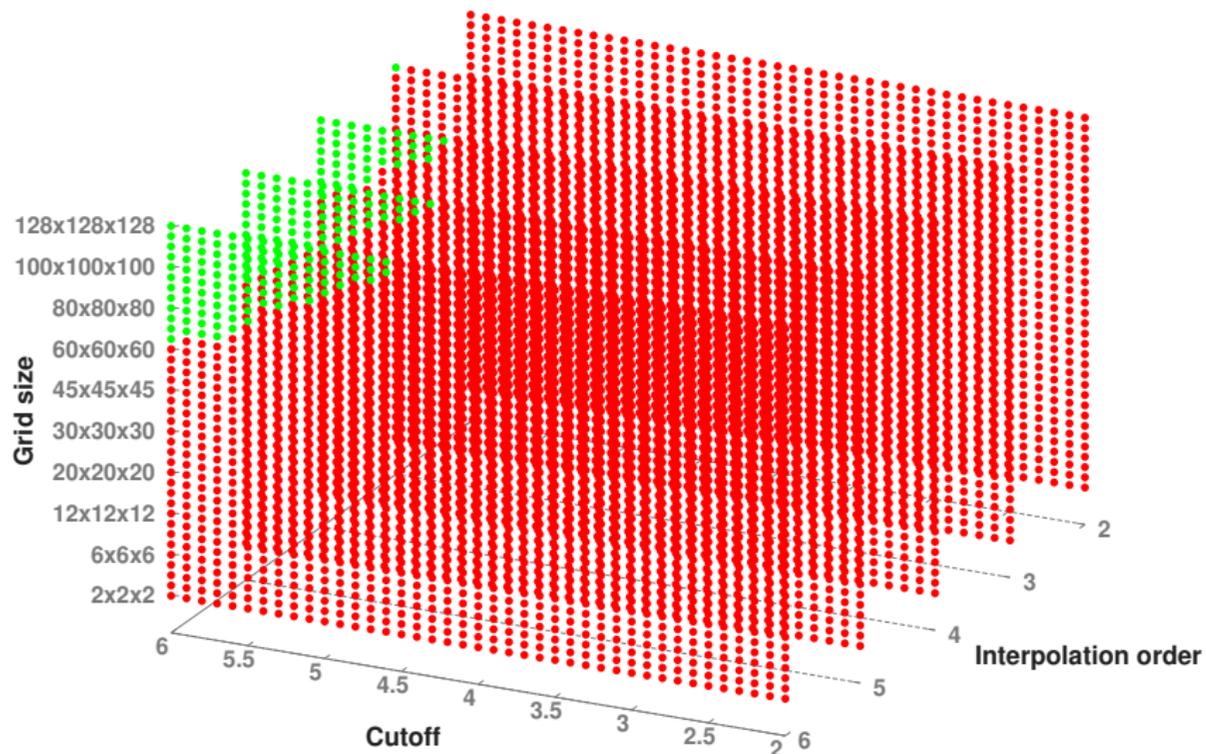
- ▶ Optimization problem — p : input parameters

$$\min_p \text{time}(\text{resources}, p) \quad \text{subject to} \quad \text{accurate}(p)$$

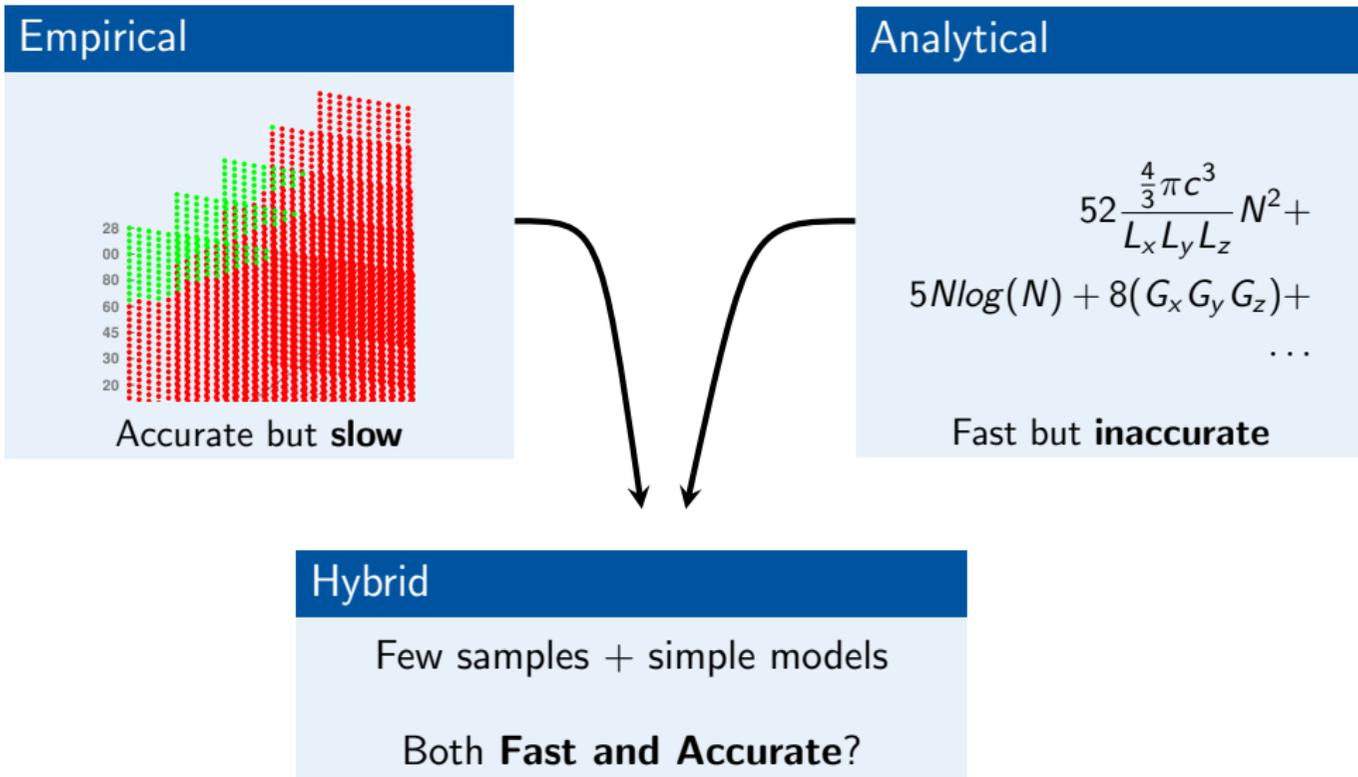
- ▶ Easily millions of combinations of parameters
- ▶ What is the “fastest” combination that leads to sufficiently accurate results?



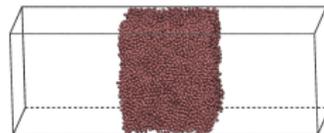
1) Accuracy boundaries



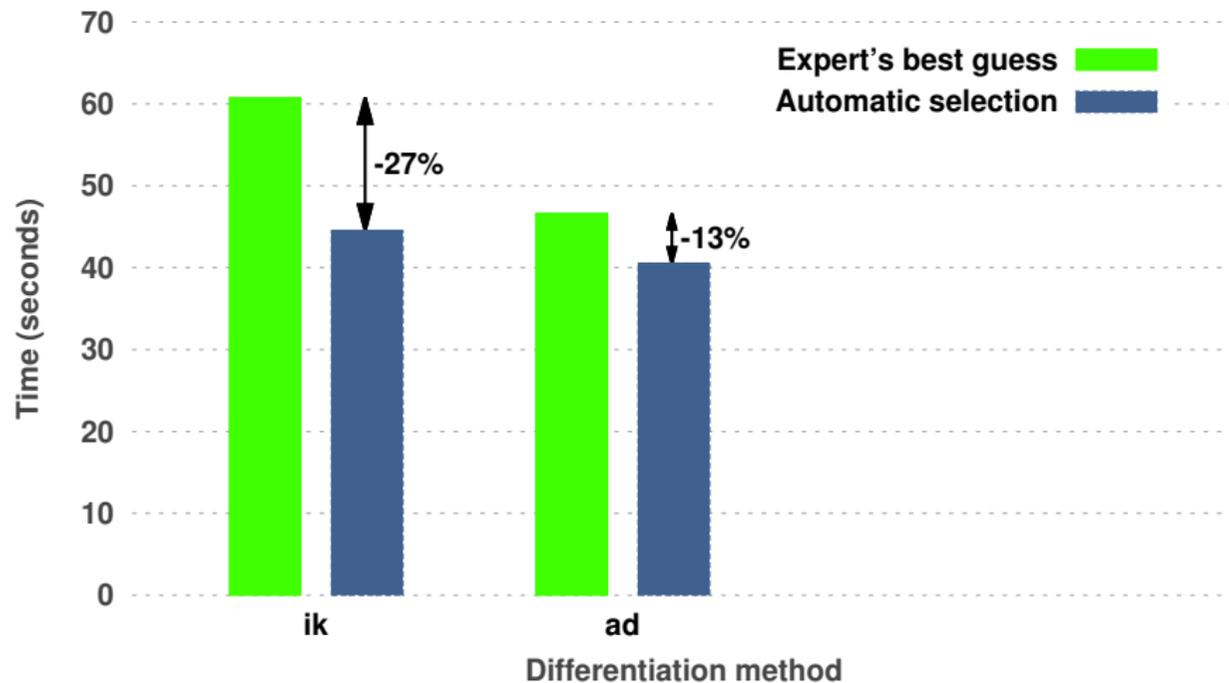
2) Performance models



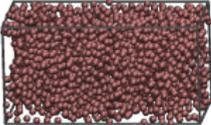
Results – Interfacial System



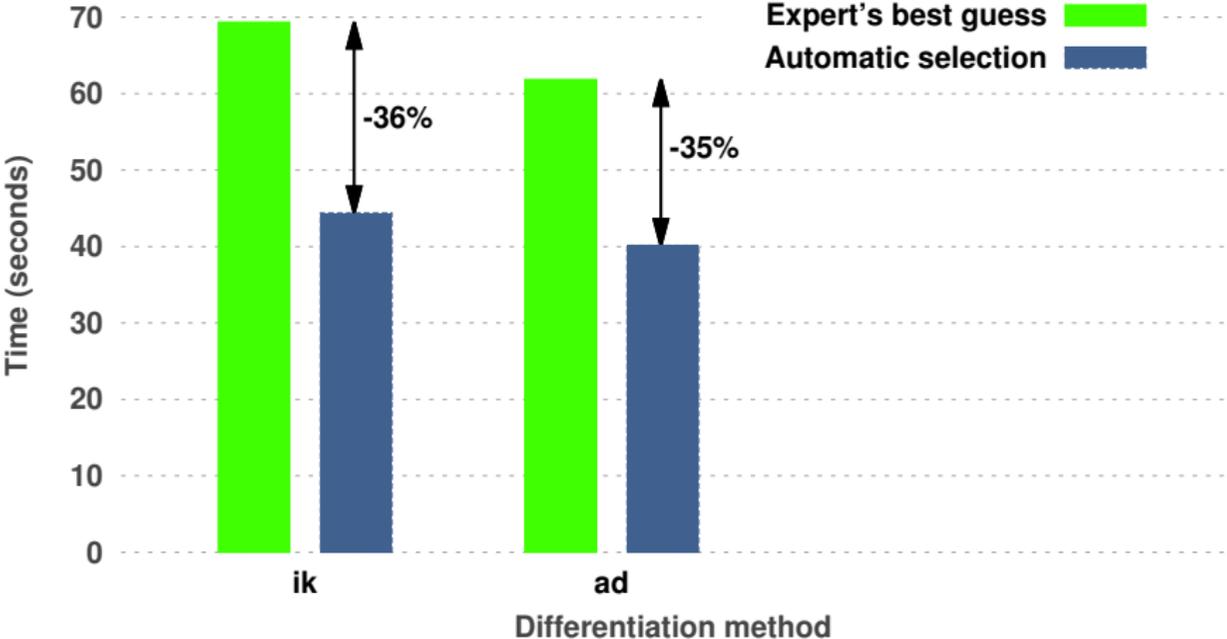
Reduction of time-to-solution



Results – Bulk System



Reduction of time-to-solution



- ▶ About me
- ▶ Molecular dynamics
- ▶ Mixed precision calculations
- ▶ Accuracy & performance modeling
- ▶ **Linear Algebra, applications**
- ▶ Tensor operations

Linear Algebra Applications

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

Signal Processing $x := (A^{-T} B^T B A^{-1} + R^T L R)^{-1} A^{-T} B^T B A^{-1} y$

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

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

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

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

Gen. Tikhonov Reg. $x := (A^T P A + Q)^{-1} (A^T P b + Q x_0)$ $P \in \mathbb{R}^{n \times n}, \text{SSPD}; Q \in \mathbb{R}^{m \times m}, \text{SSPD}; x_0 \in \mathbb{R}^{m \times 1}$

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

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

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

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



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

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

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

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

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

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

BLAS



...



LAPACK



...



MUL ADD MOV

MOVAPD

VFMADDPD ...

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

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

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



$y := \alpha x + y$	$LU = A$...	$C := \alpha AB + \beta C$
$X := A^{-1} B$	$C := AB^T + BA^T + C$	$X := L^{-1} M L^{-T}$	$QR = A$

BLAS

...

LAPACK

...

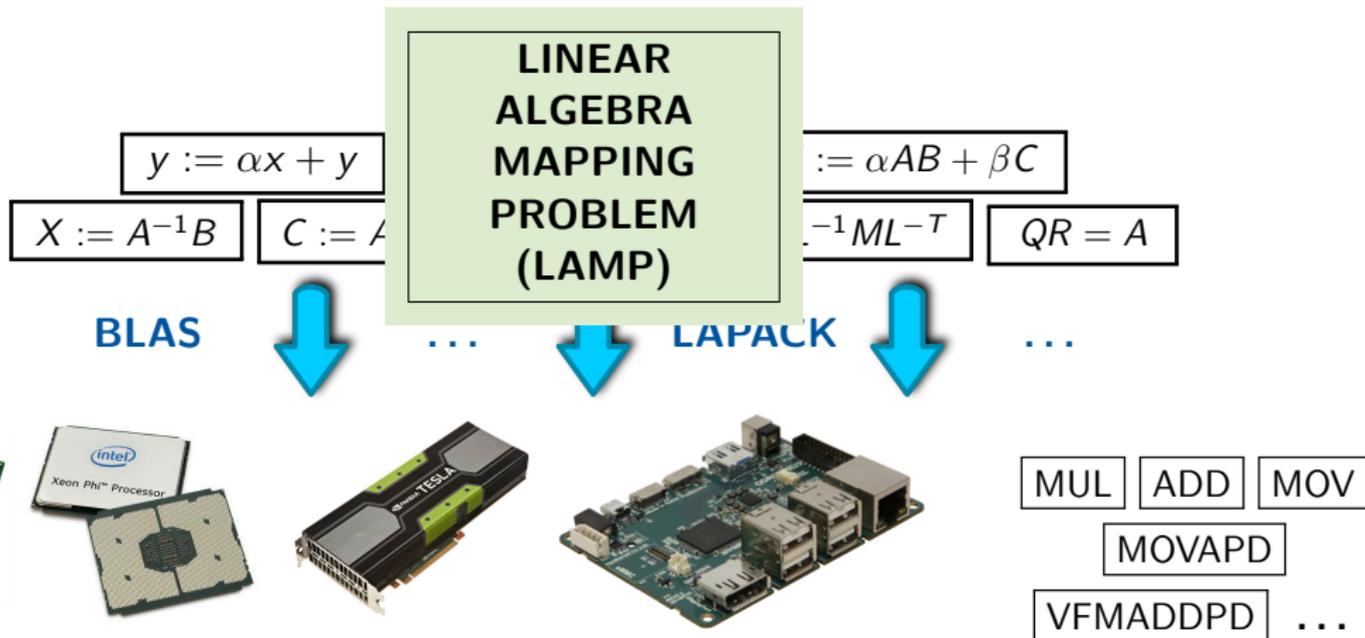


MUL	ADD	MOV
MOVAPD		
VFMADDPD	...	

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

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

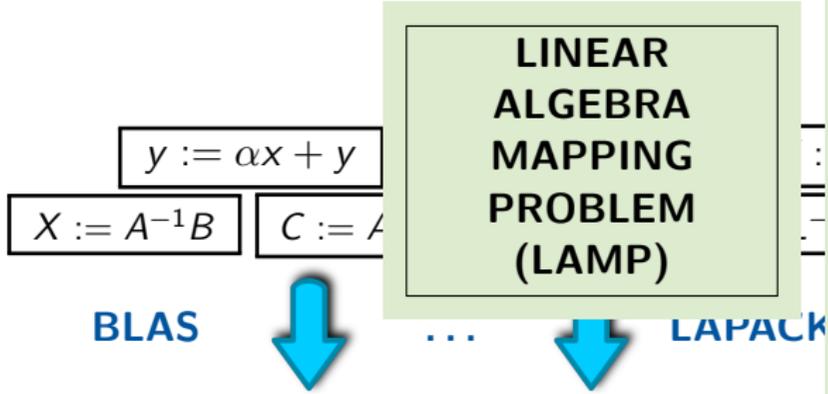
$$E := Q^{-1} U (I + U^T Q^{-1} U)^{-1} U^T \quad \dots$$



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

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

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

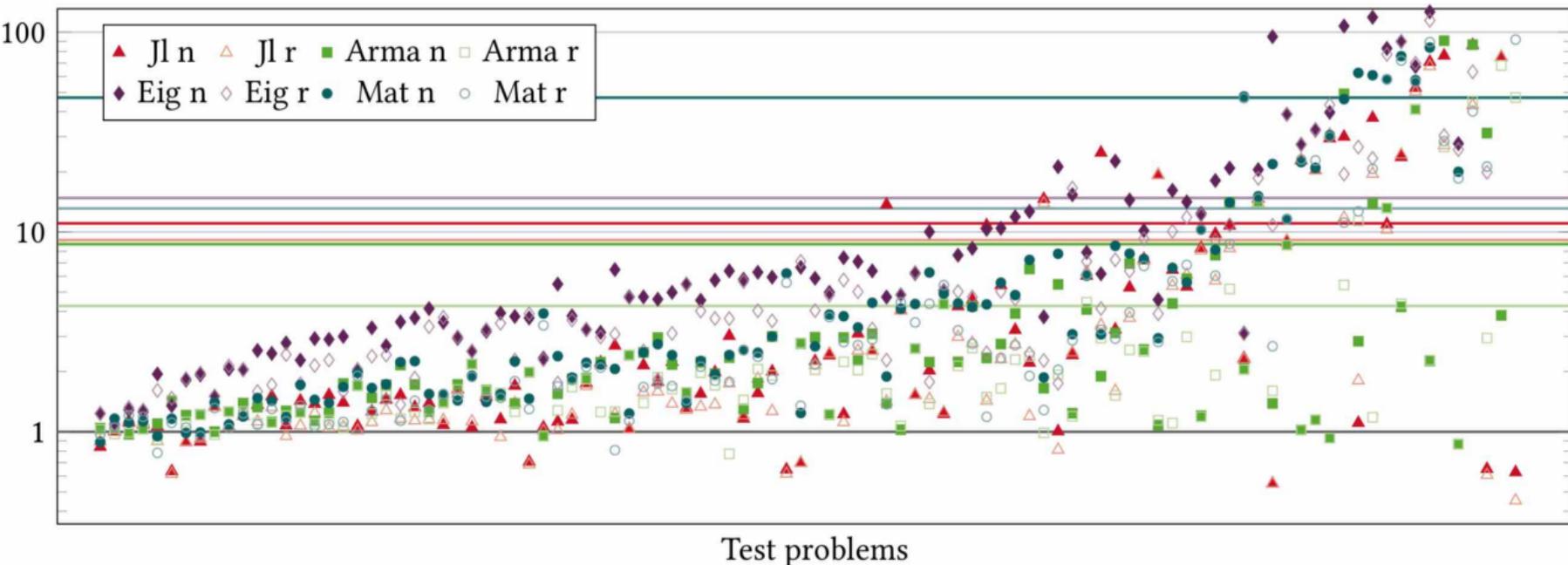


1. Associativity
Generalized matrix chain
2. Specialized kernels
Property propagation
3. Common subexpressions
4. Algebraic identities



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

Linnea's speedups



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

n/r: naive/recommended implementation.

- ▶ About me
- ▶ Molecular dynamics
- ▶ Mixed precision calculations
- ▶ Accuracy & performance modeling
- ▶ Linear Algebra, applications
- ▶ **Tensor operations**

Tensor Operations

Coupled-Cluster methods

$$\tau_{ij}^{ab} = t_{ij}^{ab} + \frac{1}{2} P_b^a P_j^i t_i^a t_j^b,$$

$$\tilde{F}_e^m = f_e^m + \sum_{fn} v_{ef}^{mn} t_n^f,$$

$$\tilde{F}_e^a = (1 - \delta_{ae}) f_e^a - \sum_m \tilde{F}_e^m t_m^a - \frac{1}{2} \sum_{mnf} v_{ef}^{mn} t_{mn}^{af} + \sum_{fn} v_{ef}^{an} t_n^f,$$

$$\tilde{F}_i^m = (1 - \delta_{mi}) f_i^m + \sum_e \tilde{F}_e^m t_i^e + \frac{1}{2} \sum_{nef} v_{ef}^{mn} t_{in}^{ef} + \sum_{fn} v_{if}^{mn} t_n^f,$$

$$\tilde{W}_{ei}^{mn} = v_{ei}^{mn} + \sum_f v_{ef}^{mn} t_i^f,$$

$$\tilde{W}_{ij}^{mn} = v_{ij}^{mn} + P_j^i \sum_e v_{ie}^{mn} t_j^e + \frac{1}{2} \sum_{ef} v_{ef}^{mn} \tau_{ij}^{ef},$$

$$\tilde{W}_{ie}^{am} = v_{ie}^{am} - \sum_n \tilde{W}_{ei}^{mn} t_n^a + \sum_f v_{ef}^{ma} t_i^f + \frac{1}{2} \sum_{nf} v_{ef}^{mn} t_{in}^{af},$$

$$\tilde{W}_{ij}^{am} = v_{ij}^{am} + P_j^i \sum_e v_{ie}^{am} t_j^e + \frac{1}{2} \sum_{ef} v_{ef}^{am} \tau_{ij}^{ef},$$

$$z_i^a = f_i^a - \sum_m \tilde{F}_i^m t_m^a + \sum_e f_e^a t_i^e + \sum_{em} v_{ei}^{ma} t_m^e + \sum_{em} v_{im}^{ae} \tilde{F}_e^m + \frac{1}{2} \sum_{efm} v_{ef}^{am} \tau_{ij}^{ef},$$

$$z_{ij}^{ab} = v_{ij}^{ab} + P_j^i \sum_e v_{ie}^{ab} t_j^e + P_b^a P_j^i \sum_{me} \tilde{W}_{ie}^{am} t_{mj}^{eb} - P_b^a \sum_m \tilde{W}_{ij}^{am} t_m^b + P_b^a \sum_m \tilde{F}_i^m t_m^b,$$

credits to D. Matthews, E. Solomonik, J. Stanton, and J. Gauss

Tensor Operations

Coupled-Cluster methods

$$\tau_{ij}^{ab} = t_{ij}^{ab} + \frac{1}{2} P_b^a P_j^i t_i^a t_j^b,$$

$$\tilde{F}_e^m = f_e^m + \sum_{fn} v_{ef}^{mn} t_n^f,$$

$$\tilde{F}_e^a = (1 - \delta_{ae}) f_e^a - \sum_m \tilde{F}_e^m t_m^a - \frac{1}{2} \sum_{mnf} v_{ef}^{mn} t_{mn}^{af} + \sum_{fn} v_{ef}^{an} t_n^f,$$

$$\tilde{F}_i^m = (1 - \delta_{mi}) f_i^m + \sum_e \tilde{F}_e^m t_i^e + \frac{1}{2} \sum_{nef} v_{ef}^{mn} t_{in}^{ef} + \sum_{fn} v_{if}^{mn} t_n^f,$$

$$\tilde{W}_{ei}^{mn} = v_{ei}^{mn} + \sum_f v_{ef}^{mn} t_i^f,$$

$$\tilde{W}_{ij}^{mn} = v_{ij}^{mn} + P_j^i \sum_e v_{ie}^{mn} t_j^e + \frac{1}{2} \sum_{ef} v_{ef}^{mn} \tau_{ij}^{ef},$$

$$\tilde{W}_{ie}^{am} = v_{ie}^{am} - \sum_n \tilde{W}_{ei}^{mn} t_n^a + \sum_f v_{ef}^{ma} t_i^f + \frac{1}{2} \sum_{nf} v_{ef}^{mn} t_{in}^{af},$$

$$\tilde{W}_{ij}^{am} = v_{ij}^{am} + P_j^i \sum_e v_{ie}^{am} t_j^e + \frac{1}{2} \sum_{ef} v_{ef}^{am} \tau_{ij}^{ef},$$

$$z_i^a = f_i^a - \sum_m \tilde{F}_i^m t_m^a + \sum_e f_e^a t_i^e + \sum_{em} v_{ei}^{ma} t_m^e + \sum_{em} v_{im}^{ae} \tilde{F}_e^m + \frac{1}{2} \sum_{efm} v_{ef}^{am} t_{im}^{ef},$$

$$z_{ij}^{ab} = v_{ij}^{ab} + P_j^i \sum_e v_{ie}^{ab} t_j^e + P_b^a P_j^i \sum_{me} \tilde{W}_{ie}^{am} t_{mj}^{eb} - P_b^a \sum_m \tilde{W}_{ij}^{am} t_m^b + P_b^a \sum_m \tilde{F}_i^m t_m^b,$$

credits to D. Matthews, E. Solomonik, J. Stanton, and J. Gauss

- ▶ **1) Kernels**
Identification, standardization
- ▶ **2) Decomposition / Mapping**
- ▶ Tensor Transpositions Compiler
TTC – github.com/HPAC/TTC
- ▶ High-Perf. Tensor Transp. Library
HPTT – github.com/HPAC/HPTT
- ▶ Tensor Contraction Code Generator
TCCG – github.com/HPAC/TCCG
- ▶ Tensor Contraction Library
TCL – github.com/springer13/tcl

▶ **Tensor Transpositions**

$$\mathcal{B}_{i_1, i_2, \dots, i_N} \leftarrow \alpha \cdot \mathcal{A}_{\pi(i_1, i_2, \dots, i_N)} + \beta \cdot \mathcal{B}_{i_1, i_2, \dots, i_N}$$

▶ **Summations** — linear summation over tensor transpositions

$$\mathcal{B}_{i_0 i_1 i_2} \leftarrow 2\mathcal{A}_{i_0 i_1 i_2} - \mathcal{A}_{i_2 i_1 i_0} - \mathcal{A}_{i_0 i_2 i_1}$$

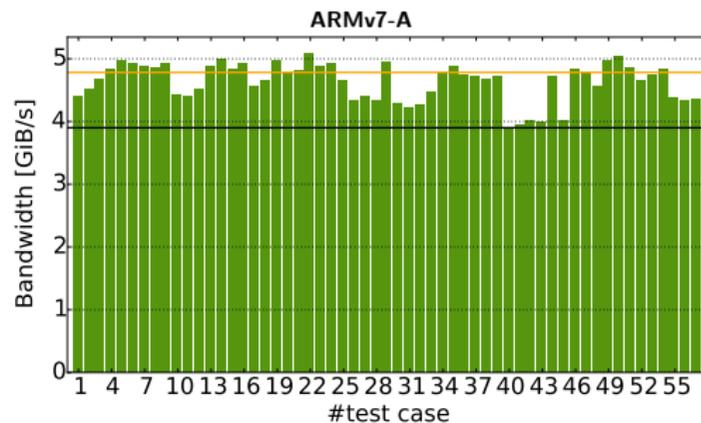
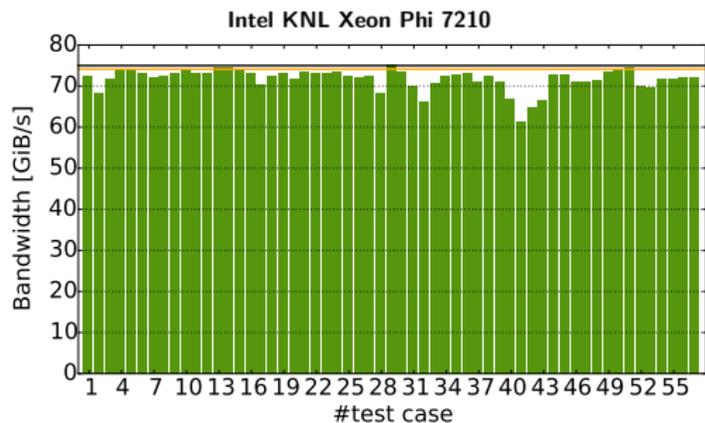
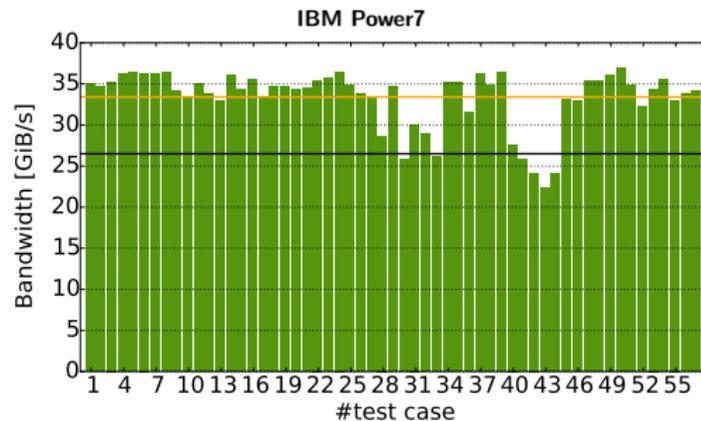
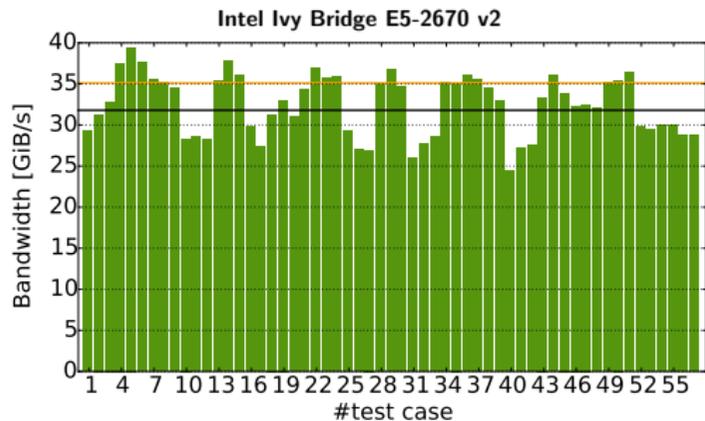
$$\mathcal{B}_{i_0 i_1 i_2} \leftarrow 4\mathcal{A}_{i_0 i_1 i_2} - 2\mathcal{A}_{i_1 i_0 i_2} - 2\mathcal{A}_{i_2 i_1 i_0} + \mathcal{A}_{i_1 i_2 i_0} - 2\mathcal{A}_{i_0 i_2 i_1} + \mathcal{A}_{i_2 i_0 i_1}$$

$$\mathcal{B}_{i_0 i_1 i_2 i_3} \leftarrow 2\mathcal{A}_{i_0 i_1 i_2 i_3} - \mathcal{A}_{i_2 i_1 i_0 i_3} - \mathcal{A}_{i_0 i_2 i_1 i_3} - \mathcal{A}_{i_0 i_1 i_3 i_2}$$

▶ **Tensor Contractions**

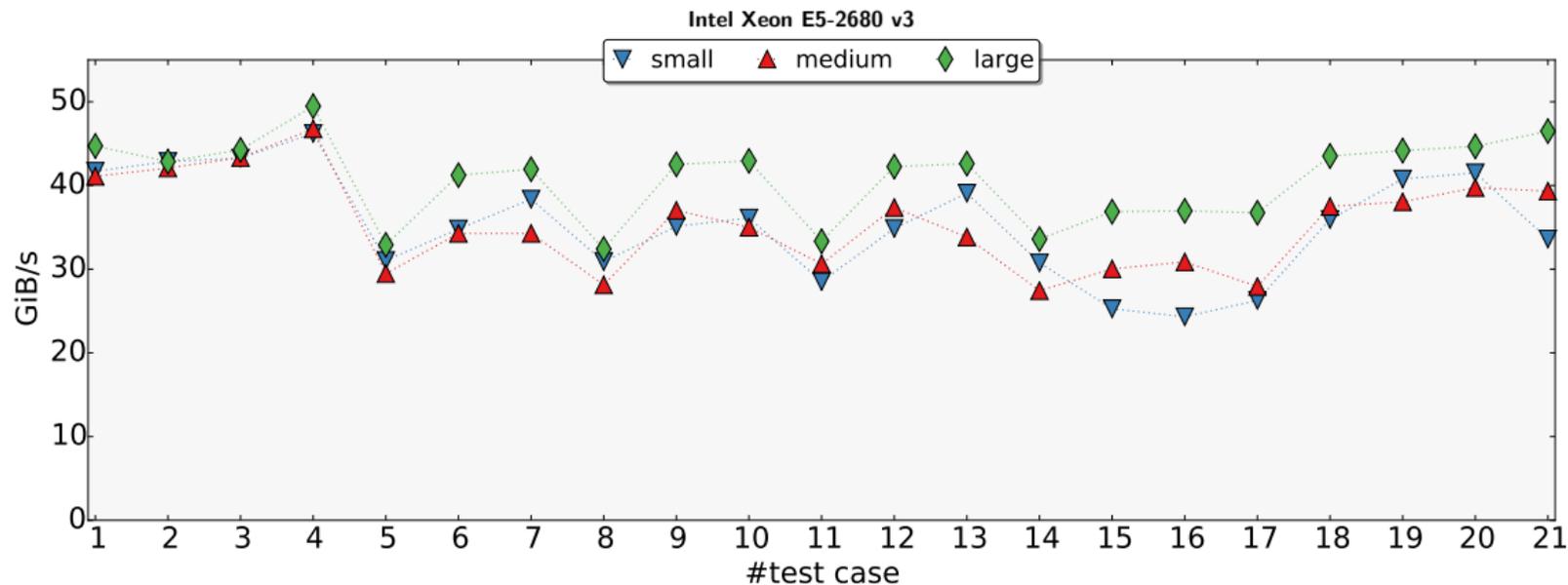
$$\mathcal{C}_{\pi^c(I_m \cup I_n)} \leftarrow \alpha \cdot \mathcal{A}_{\pi^a(I_m \cup I_k)} \times \mathcal{B}_{\pi^b(I_n \cup I_k)} + \beta \cdot \mathcal{C}_{\pi^c(I_m \cup I_n)}$$

HPTT – Transposition



Black and orange lines denote STREAM and AXPY bandwidth, respectively.

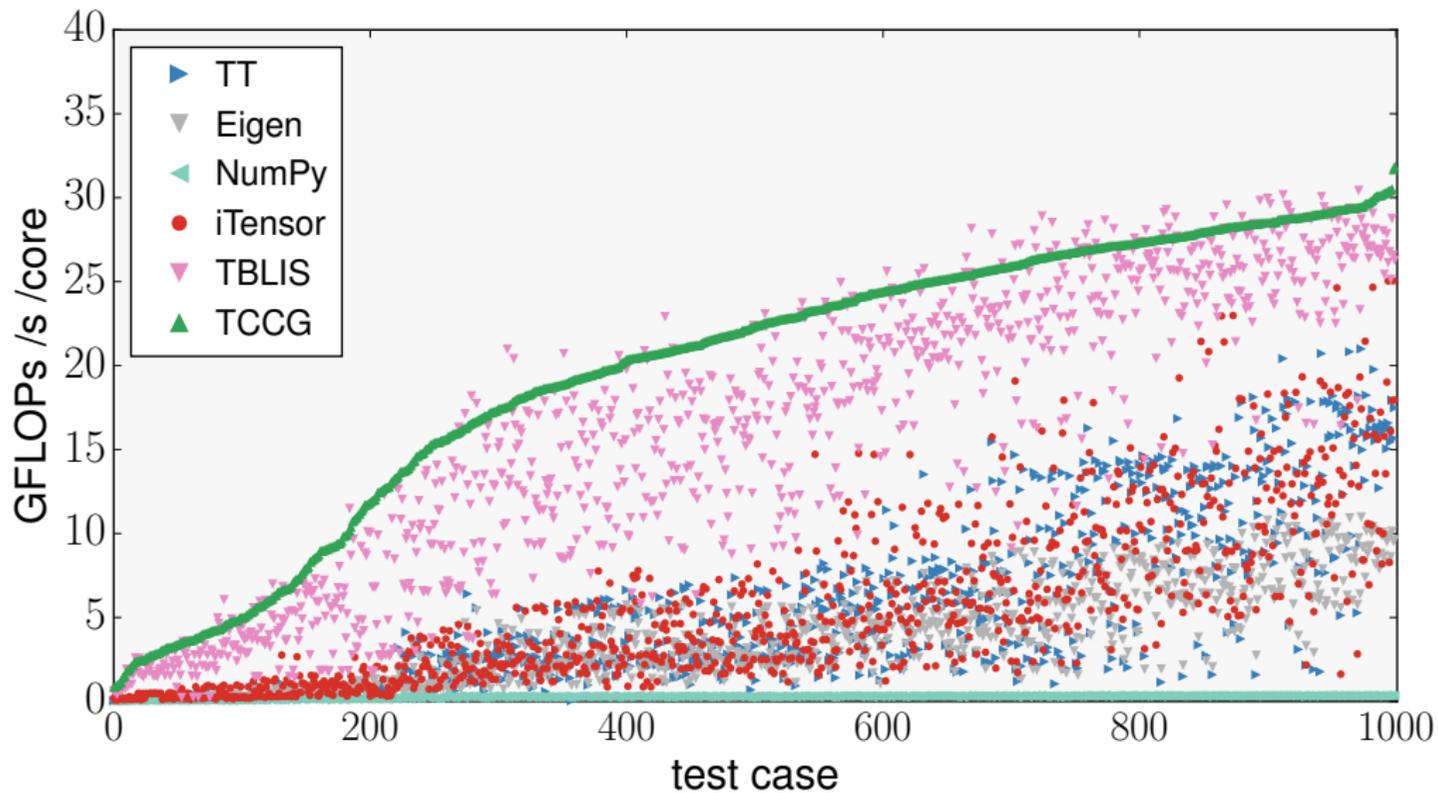
Summation



$$BW = \frac{2 \times N^d \times \text{sizeof}(\text{double})}{2^{30} \times \text{Time}} \text{ GiB/s}, \quad \text{TRIAD : 107.7 GiB/s}$$

Problem size: Small (70MBs), medium (320MBs), large (1200MBs).

TCCG – Contraction



DP tensor contractions. Host: 2×Haswell-EP E5-2680 v3 @ 24 threads

Mapping – LAMP (again)

Generation of Hamiltonian and Overlap Matrices

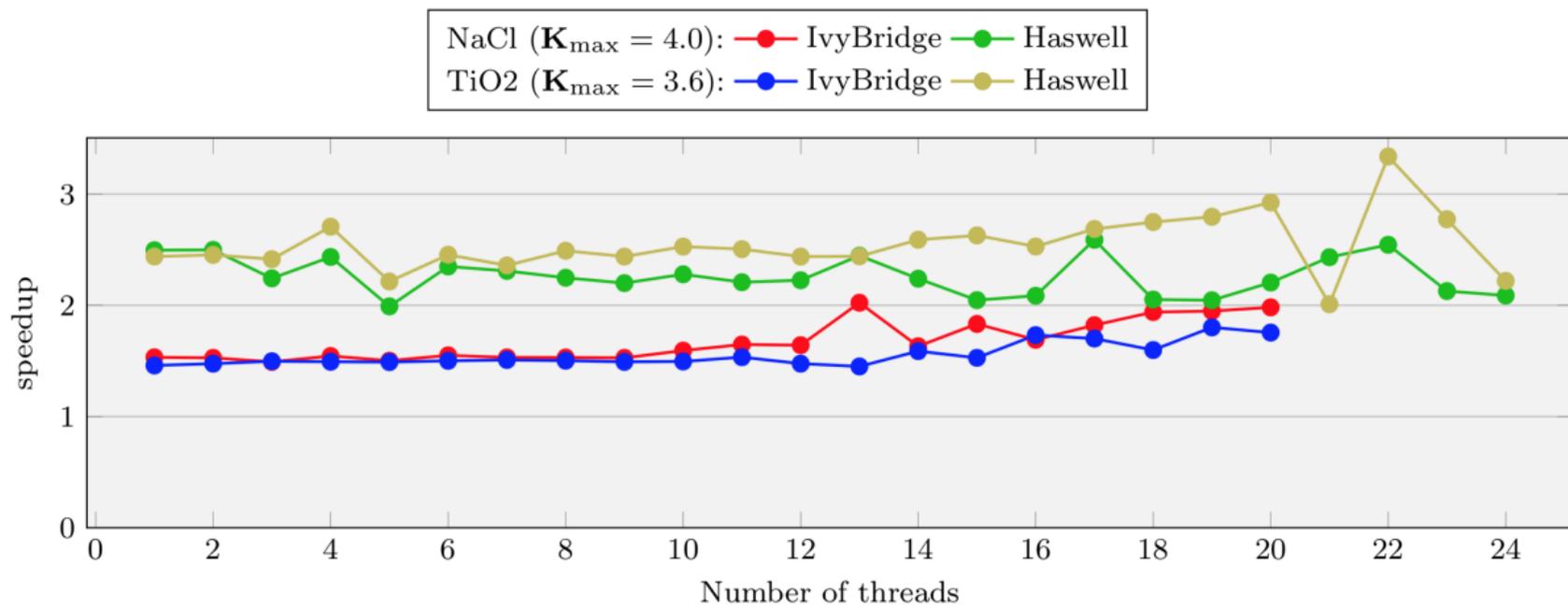
$$(S)_{G',G} = \sum_a \sum_{L=(l,m)} \left(A_L^{a,G'} \right)^* A_L^{a,G} + \left(B_L^{a,G'} \right)^* B_L^{a,G} \|\dot{u}_{l,a}\|^2$$

$$(H)_{G',G} = \sum_a \sum_{L',L} \left(A_{L',a,t'}^* T_{L',L;a}^{[AA]} A_{L,a,t} \right) + \left(A_{L',a,t'}^* T_{L',L;a}^{[AB]} B_{L,a,t} \right) \\ + \left(B_{L',a,t'}^* T_{L',L;a}^{[BA]} A_{L,a,t} \right) + \left(B_{L',a,t'}^* T_{L',L;a}^{[BB]} B_{L,a,t} \right)$$

“mappable” to BLAS and LAPACK kernels

Mapping – LAMP (again)

Generation of Hamiltonian and Overlap Matrices



More on Tensors?

- ▶ The “other” operations?

Data science — Decompositions, rank, compression, completion, ...

- ▶ Kernels?

- ▶ Common language? Interface?

- ▶ Mapping onto kernels

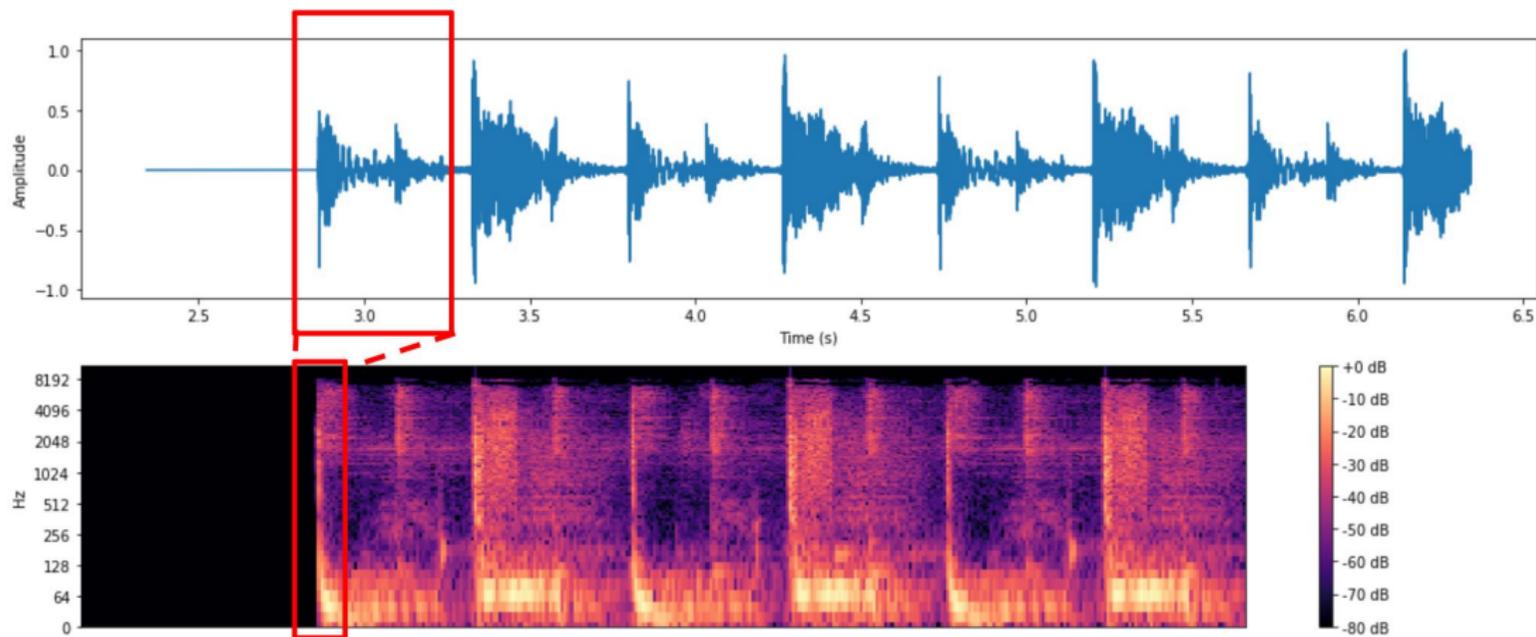
Two passions: mountains & music



Two passions: mountains & music



Two passions: mountains & music



Talk: “Automatic Dj-ing”, next Wednesday, 9am