

Solvers and Eigensolvers for Multicore Processors

Paolo Bientinesi

AICES, RWTH Aachen
pauldj@aices.rwth-aachen.de

Max-Planck-Institute für biologische Kybernetik
March 18th, 2011
Tübingen, Germany



- 1 Introduction
- 2 Part #1: Solvers
- 3 Part #2: Eigensolvers
- 4 Conclusions

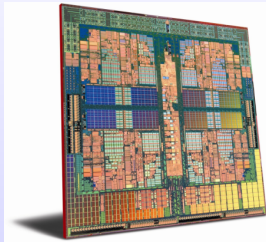
- Matrix-matrix multiplication (GEMM): $C \leftarrow C + AB$
- Factorizations: $LU = A, LL^T = A, QR = A$
- Linear system: $AX = B$
- Transformations: $QAQ^T = T, QA = H$
- Matrix Equations: $AX + XB = C$

- Eigenproblems: $AZ = Z\Lambda$
- Generalized eigenproblems: $AZ = BZ\Lambda$
- ...

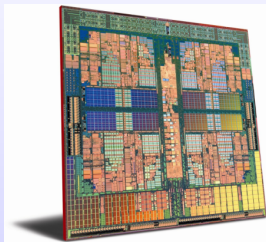
- Matrix-matrix multiplication (GEMM): $C \leftarrow C + AB$
- Factorizations: $LU = A, LL^T = A, QR = A$
- Linear system: $AX = B$
- Transformations: $QAAQ^T = T, QA = H$
- Matrix Equations: $AX + XB = C$
- Eigenproblems: $AZ = Z\Lambda$
- Generalized eigenproblems: $AZ = BZA$
- ...

Objective: High-performance

- Numerical stability
- Multiple algorithmic variants
- Multiple implementations
- Multicore → **Multiple types of parallelism!**



Multi-cores invasion: 499/500 entries of the Top 500



Multi-cores invasion: 499/500 entries of the Top 500

4 cores, 8 cores, ... 24 cores, ...

- More parallelism than we know what to do with?
- Is multi-threaded BLAS the solution for LA libs?
- Linear solvers \neq Eigensolvers

Algorithms expressed in terms of simpler linear algebra operations

Algorithms expressed in terms of simpler linear algebra operations

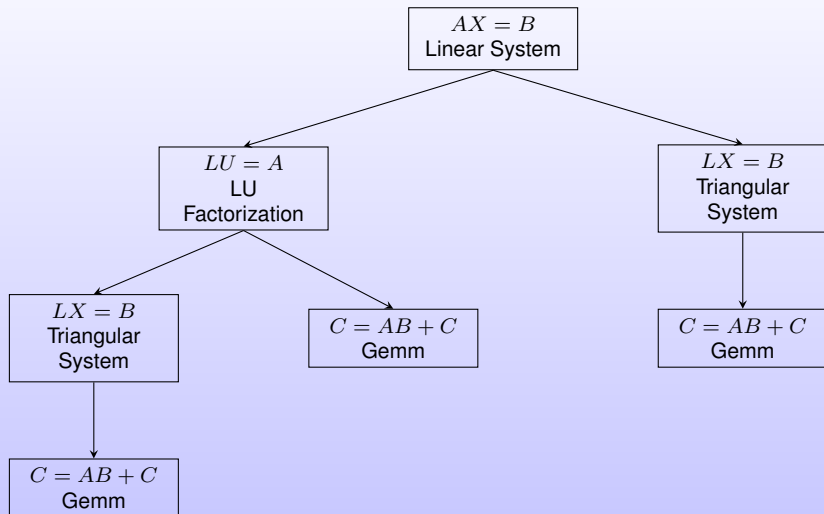
BLAS: Basic Linear Algebra Subroutines

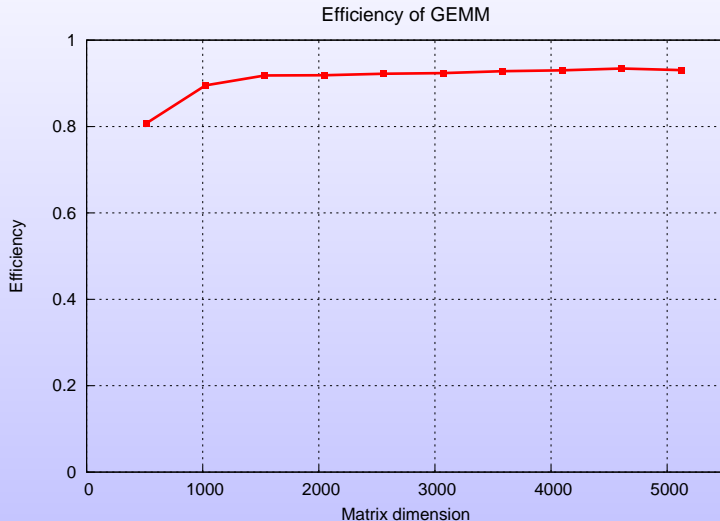
$$\text{BLAS-1: } \begin{array}{l} y := y + \alpha x \\ \beta := \alpha + x^T y \end{array} \quad x, y \in \mathbb{R}^n$$

$$\text{BLAS-2: } \begin{array}{l} y := y + Ax \\ y := A^{-1}x \end{array} \quad A \in \mathbb{R}^{n \times n}, x, y \in \mathbb{R}^n$$

$$\text{BLAS-3: } \begin{array}{l} C := C + AB \\ C := A^{-1}B \end{array} \quad A, B, C \in \mathbb{R}^{n \times n}$$

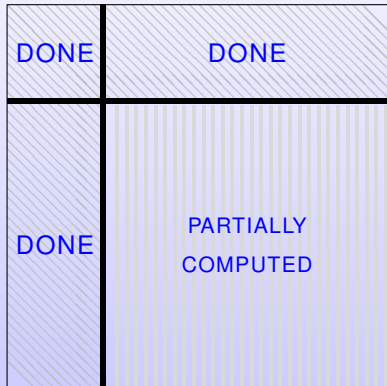
ESSL (IBM), MKL (Intel), ATLAS, GotoBLAS, ...



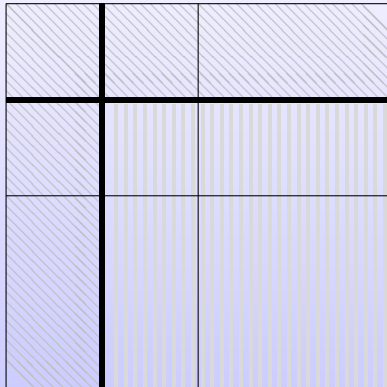


- 1 Introduction
- 2 Part #1: Solvers**
- 3 Part #2: Eigensolvers
- 4 Conclusions

Iteration i: completed

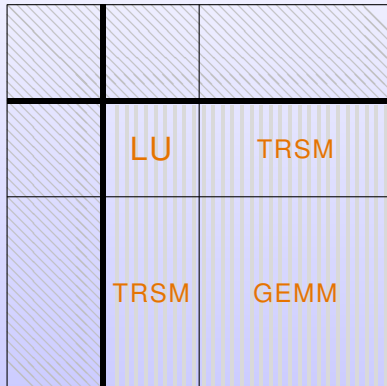


Iteration $i+1$: repartitioning



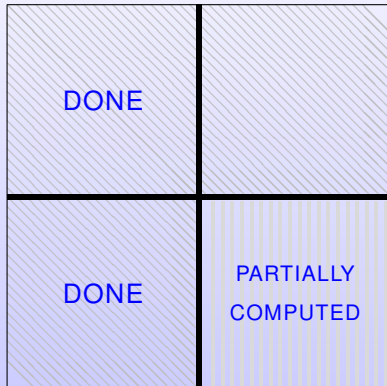
LU factorization: loop-based algorithm

Iteration $i+1$: computation

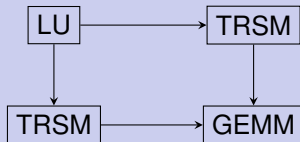


LU factorization: loop-based algorithm

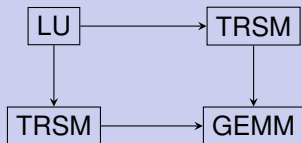
Iteration $i+1$: completed (boundary shift)



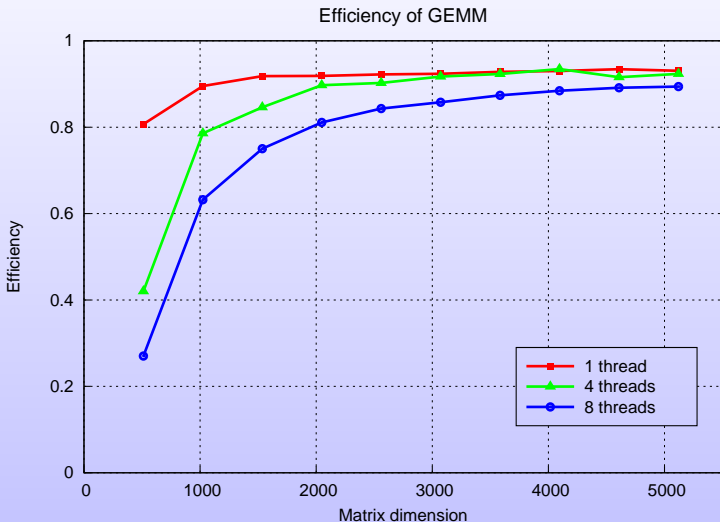
Solution #1: Multithreaded BLAS



Solution #1: Multithreaded BLAS



- Advantage: ease of use. Legacy code!
- Drawback: synchronization.



- Inversion of a Symmetric Positive Definite matrix
- Covariance matrix
- Very large dense problems

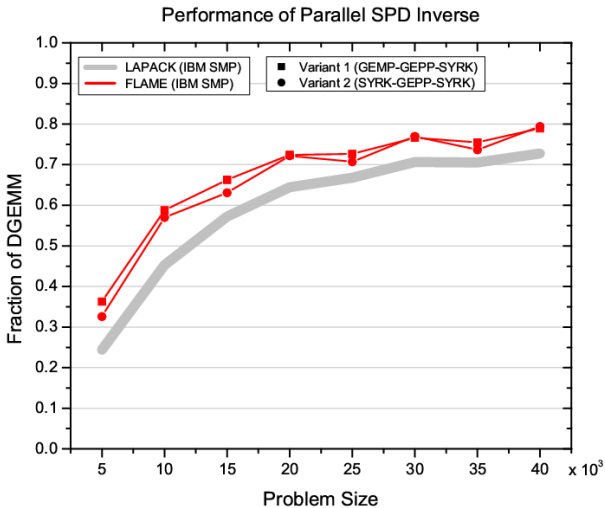
Cholesky factorization



Triangular inversion



Matrix-matrix multiplication

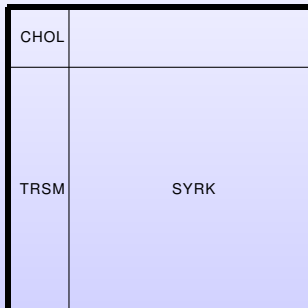


Solution #2: Algorithms by blocks

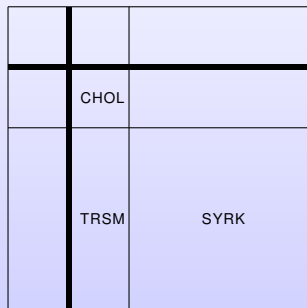
- Advantage: out of order execution.
- Advantage: parallelism limited only by the data dependencies between operations.
- Drawback: plateaux.

Cholesky factorization

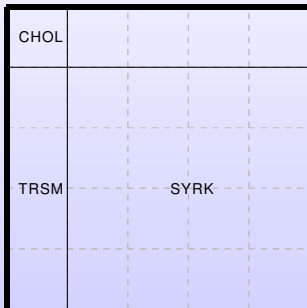
$$LL^T = A$$



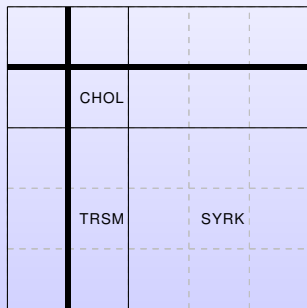
Iteration 1



Iteration 2



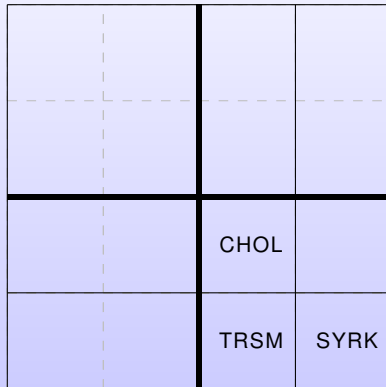
Iteration 1

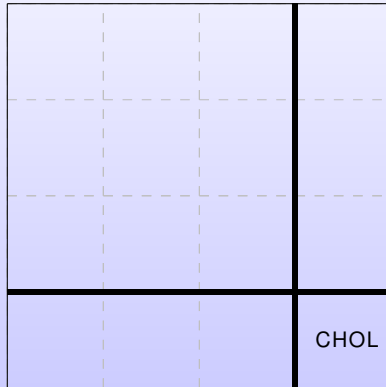


Iteration 2

CHOL			
TRSM	SYRK		
TRSM	GEMM	SYRK	
TRSM	GEMM	GEMM	SYRK

	CHOL		
	TRSM	SYRK	
	TRSM	GEMM	SYRK





CHOL			
TRSM	SYRK		
TRSM	GEMM	SYRK	
TRSM	GEMM	GEMM	SYRK

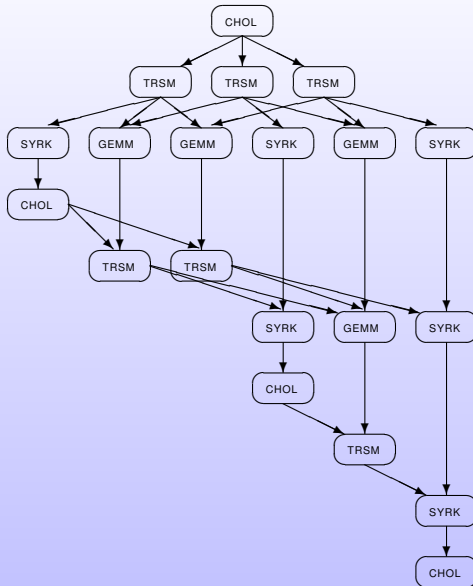
CHOL			
TRSM	SYRK		
TRSM	GEMM	SYRK	
TRSM	GEMM	GEMM	SYRK

CHOL			
TRSM	SYRK		
TRSM	GEMM	SYRK	
TRSM	GEMM	GEMM	SYRK

CHOL			
TRSM	SYRK		
TRSM	GEMM	SYRK	
TRSM	GEMM	GEMM	SYRK

CHOL			
TRSM	SYRK		
TRSM	GEMM	SYRK	
TRSM	GEMM	GEMM	SYRK

4×4 -tile matrix



4 × 4-tile matrix

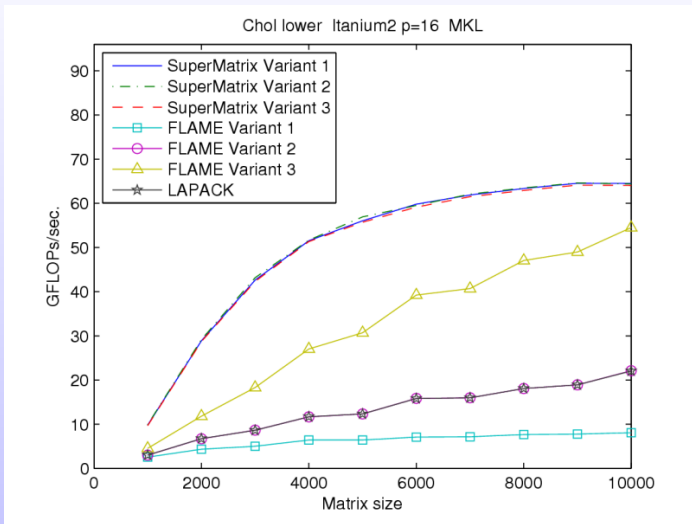
Stage	Scheduled Tasks			
1	CHOL			
2	TRSM	TRSM	TRSM	TRSM
3	SYRK	GEMM	SYRK	GEMM
4	GEMM	SYRK	GEMM	GEMM
5	GEMM	SYRK	CHOL	
6	TRSM	TRSM	TRSM	
7	SYRK	GEMM	SYRK	GEMM
8	GEMM	SYRK	CHOL	
9	TRSM	TRSM		
10	SYRK	GEMM	SYRK	
11	CHOL			
12	TRSM			
13	SYRK			
14	CHOL			

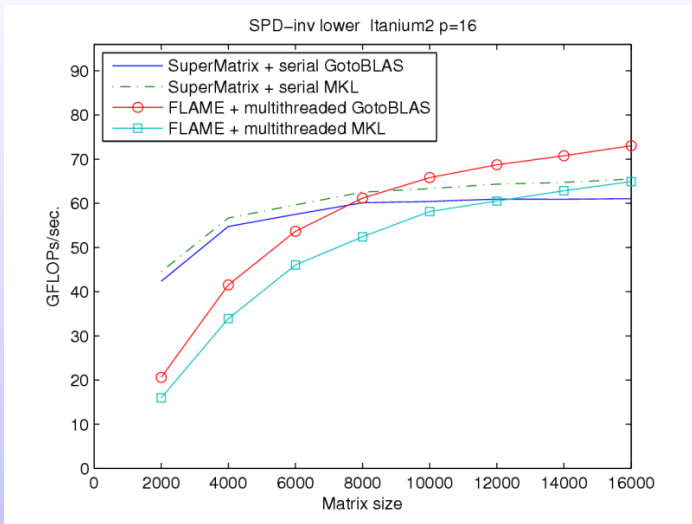
SPD Inverse again: Chol+Inv+GEMM

5×5 -tile matrix

5×5 -tile matrix

Stage	Scheduled Tasks			
1	CHOL			
2	TRSM	TRSM	TRSM	TRSM
3	SYRK	GEMM	SYRK	GEMM
4	GEMM	SYRK	GEMM	GEMM
5	GEMM	SYRK	CHOL	TRSM
6	TRSM	TRSM	TRSM	TRSM
7	TRSM	TRSM	TRINV	SYRK
8	GEMM	SYRK	GEMM	GEMM
9	SYRK	TTMM	CHOL	TRSM
10	TRSM	TRSM	TRSM	TRSM
11	GEMM	GEMM	GEMM	SYRK
12	GEMM	SYRK	TRSM	CHOL
13	TRSM	TRSM	TRINV	SYRK
14	TRSM	GEMM	GEMM	GEMM
15	GEMM	TRMM	SYRK	TRSM
16	TRSM	TTMM	CHOL	TRSM
17	SYRK	TRINV	GEMM	SYRK
18	GEMM	GEMM	GEMM	TRMM
19	TRMM	TRSM	TRSM	TRSM
20	TRSM	TRSM	TRSM	TRSM
21	TTMM	SYRK	GEMM	SYRK
22	TRINV	GEMM	GEMM	TRINV
23	SYRK	SYRK	GEMM	SYRK
24	TRMM	GEMM	TRMM	GEMM
25	TRMM	SYRK	GEMM	GEMM
26	TTMM	GEMM	TRMM	TRMM
27	SYRK	TRMM		
28	TRMM			
29	TTMM			





Multithreaded BLAS vs. Algorithms by blocks

No absolute winner: crossover!

- | | |
|-------------------|---|
| ✓ Ease of use | ✓ Out of order execution |
| ✗ Synchronization | ✓ Parallelism dictated by data dependencies |
| | ✗ Plateaux |

- 1 Introduction
- 2 Part #1: Solvers
- 3 Part #2: Eigensolvers**
- 4 Conclusions

$$AX = X\Lambda$$

- Input: $A \in \mathcal{C}^{n \times n}$, $A^H = A$; #eigenpairs: $1 \leq k \leq n$
- Output: $X \in \mathcal{C}^{n \times k}$ eigenvectors
 $\Lambda \in \mathcal{R}^{k \times k}$ eigenvalues

$$AX = X\Lambda$$

- Input: $A \in \mathcal{C}^{n \times n}$, $A^H = A$; #eigenpairs: $1 \leq k \leq n$
- Output: $X \in \mathcal{C}^{n \times k}$ eigenvectors
 $\Lambda \in \mathcal{R}^{k \times k}$ eigenvalues

Approach

- $T = Q^H A Q$ Reduction to tridiagonal form $O(n^3)$
- $T Z = Z \Lambda$ Tridiagonal eigenproblem $O(kn) - O(n^3)$
- $X = Q Z$ Backtransformation $O(kn^2)$

$$AX = X\Lambda$$

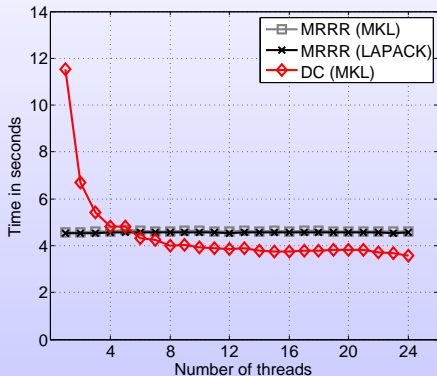
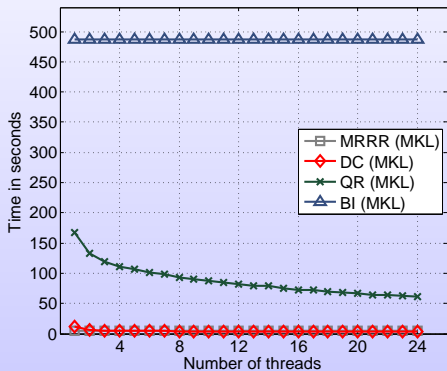
- Input: $A \in \mathcal{C}^{n \times n}$, $A^H = A$; #eigenpairs: $1 \leq k \leq n$
- Output: $X \in \mathcal{C}^{n \times k}$ eigenvectors
 $\Lambda \in \mathcal{R}^{k \times k}$ eigenvalues

Approach

- $T = Q^H A Q$ Reduction to tridiagonal form $O(n^3)$
- $T Z = Z \Lambda$ Tridiagonal eigenproblem $O(kn) - O(n^3)$
- $X = Q Z$ Backtransformation $O(kn^2)$

Algorithms

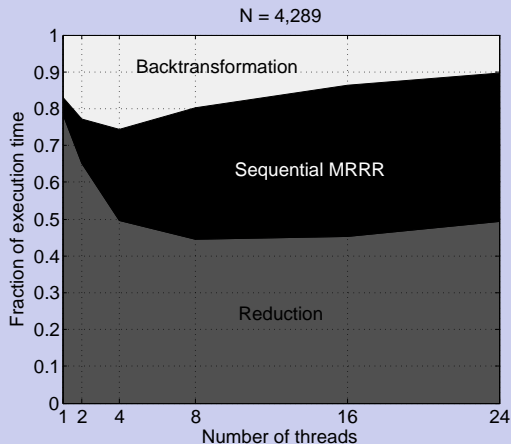
- Inverse Iteration (1958): subsets $O(kn^2)$
- QR (1961): high-accuracy $O(n^3)$
- Divide & Conquer (1981): parallel, BLAS3 $O(n^3)$
- MRRR (1997): subsets, no re-orth. $O(kn)$



Tridiagonal eigensolver, matrix size=4289, from DFT.

... it's $O(n^2)$ anyway

... it's $O(n^2)$ anyway



If not properly parallelized, even $O(n^2)$ dominates!

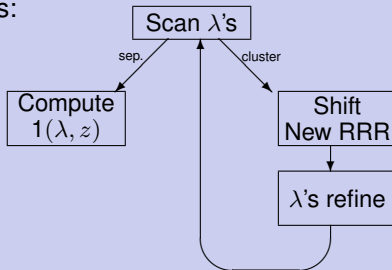
- 1 Introduction
- 2 Part #1: Solvers
- 3 Part #2: Eigensolvers**
- 4 Conclusions

Multiple Relatively Robust Representations

- first stable algorithm to compute k eigenpairs in $O(nk)$ ops
- no reorthogonalization

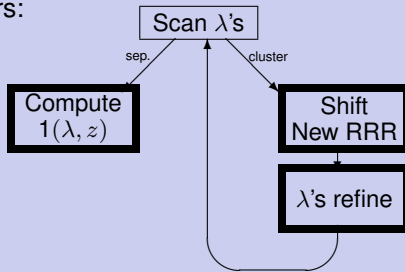
Multiple Relatively Robust Representations

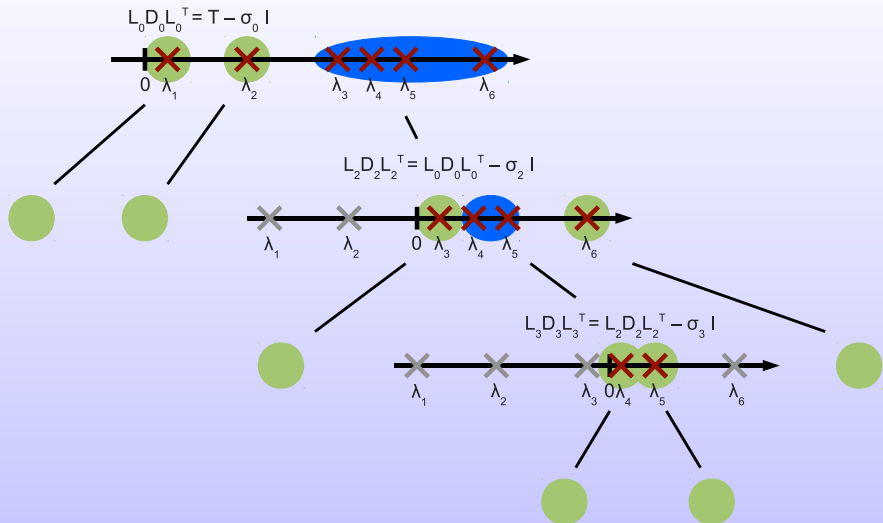
- first stable algorithm to compute k eigenpairs in $O(nk)$ ops
- no reorthogonalization
- 1) eigenvalues \rightarrow 2) eigenvectors + eigenvalues
- eigenvalues: *Bisection* or *dqds*
- eigenvectors:

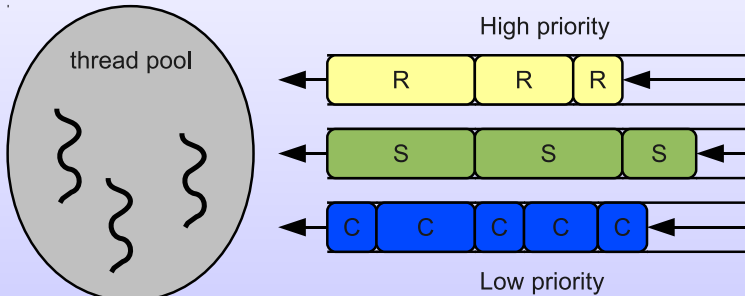


Multiple Relatively Robust Representations

- first stable algorithm to compute k eigenpairs in $O(nk)$ ops
- no reorthogonalization
- 1) eigenvalues \rightarrow 2) eigenvectors + eigenvalues
- eigenvalues: *Bisection* or *dqds*
- eigenvectors:

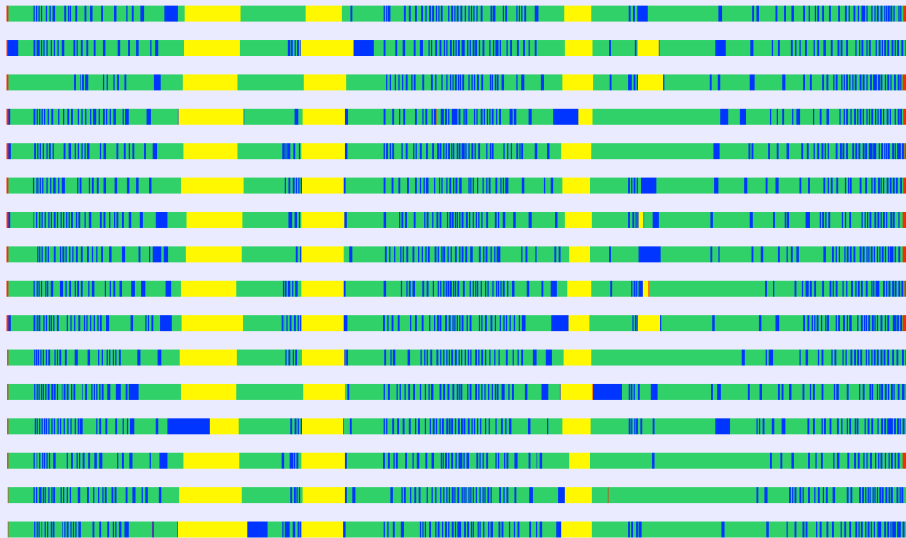


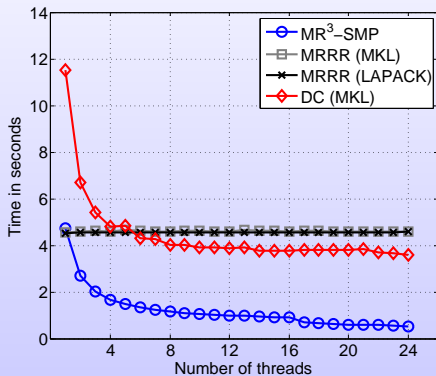
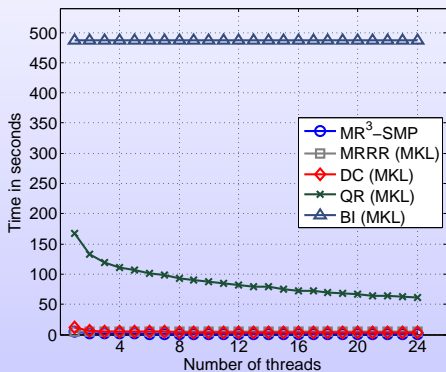




Example trace: 16 cores—eigenvectors

Matrix size: 12387 Execution time: 3.3s Sequential: 49.3s (LAPACK)

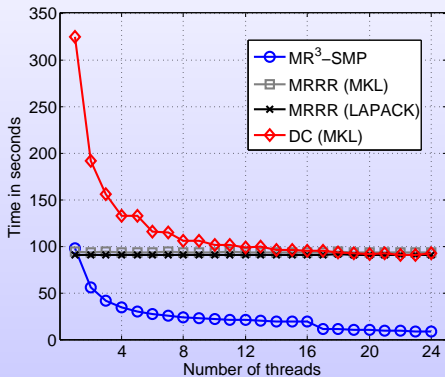
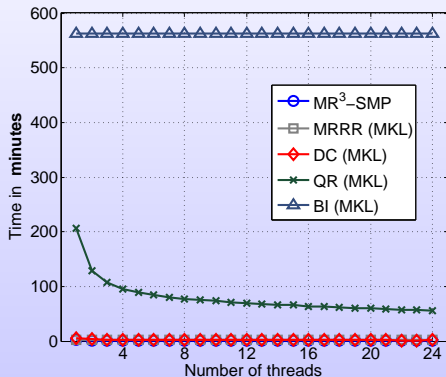




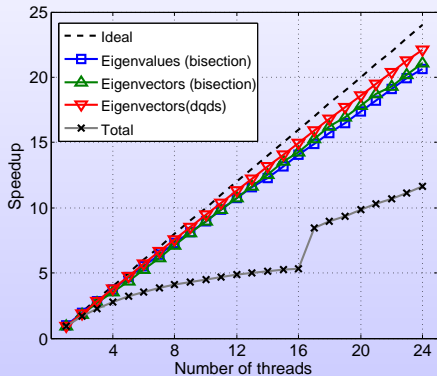
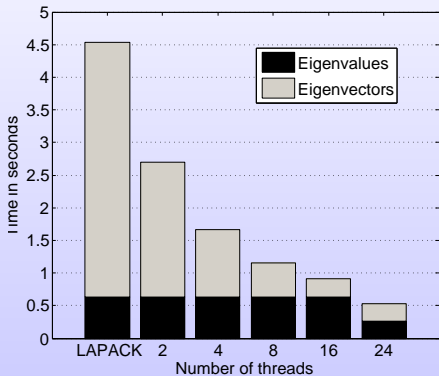
Matrix size: 4289.

A larger example: look at the scale!

Matrix size: 16023. Frequency response analysis of automobiles.

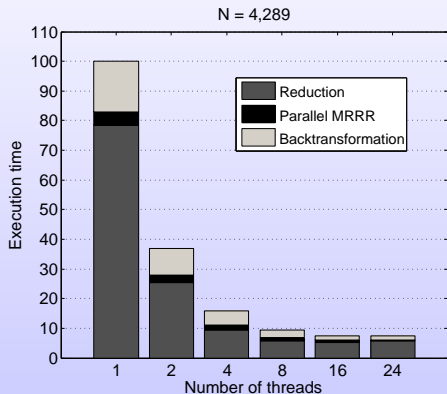
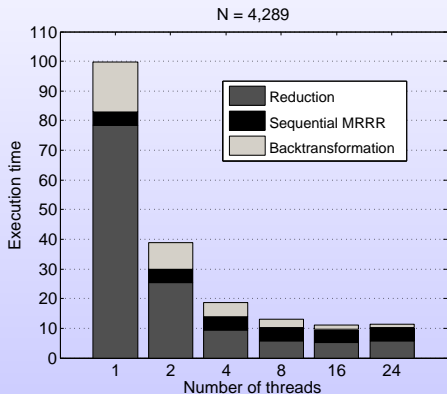


From almost 10 hours to 8.3 seconds.

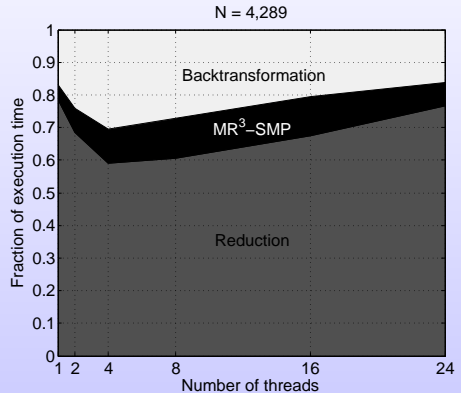
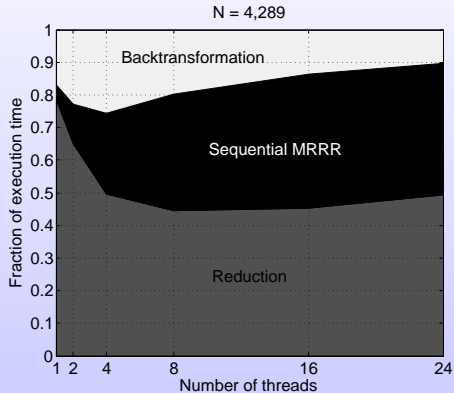


- 1 Introduction
- 2 Part #1: Solvers
- 3 Part #2: Eigensolvers
- 4 Conclusions

3 stages: before and after



3 stages: before and after



MRRR-SMP

- Matthias Petschow (AICES)
 - Eigensolver tailored for multi-cores
 - Almost perfect speedups
 - Routines are available
-
- Multi-threaded BLAS for solvers: nice and easy.
 - Multi-threaded BLAS for eigensolvers: not THAT good.

MRRR-SMP

- Matthias Petschow (AICES)
 - Eigensolver tailored for multi-cores
 - Almost perfect speedups
 - Routines are available
-
- Multi-threaded BLAS for solvers: nice and easy.
 - Multi-threaded BLAS for eigensolvers: not THAT good.

Thank you for the attention.

Deutsche
Forschungsgemeinschaft

DFG

Financial support from the Deutsche Forschungsgemeinschaft (German Research Association) through grant GSC 111 is gratefully acknowledged.