

Fast and Scalable Eigensolvers for Multicore and Hybrid Architectures

Paolo Bientinesi

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



- 1 The Problem
- 2 Multicore Processors: MR³-SMP
- 3 Distributed Memory Architectures: PMRRR
- 4 GPUs

Symmetric Dense Eigenproblem

$$AX = X\Lambda$$

$$AX = XB\Lambda$$

STDEIG

GENEIG

- Input:

$$\begin{array}{ll} A \in \mathcal{C}^{n \times n}, & A^H = A \\ B \in \mathcal{C}^{n \times n}, & \text{SPD} \\ k, 1 \leq k \leq n & \text{\#eigenpairs} \end{array}$$

- Output:

$$\begin{array}{ll} X \in \mathcal{C}^{n \times k}, & \text{eigenvectors} \\ \Lambda \in \mathcal{R}^{k \times k}, & \text{eigenvalues} \end{array}$$

- Accuracy:

$$\begin{array}{ll} \|AX - X\Lambda\|, & \text{residual} \\ \|X^H X - I\|, & \text{orthogonality} \end{array}$$

GENEIG $AX = XBA\Lambda$

1	$LL^H = B$	Cholesky factorization	$O(n^3)$
2	$M \leftarrow L^{-1}AL^{-H}$	Reduction to standard form	$O(n^3)$
3	$T = Q^HMQ$	Reduction to tridiagonal form	$O(n^3)$
4	$TZ = Z\Lambda$	Tridiagonal eigenproblem	$O(kn) - O(n^3)$
5	$Y = QZ$	Backtransformation #1	$O(kn^2)$
6	$X = L^{-H}Y$	Backtransformation #2	$O(kn^2)$

Nested Eigensolvers

GENEIG \rightarrow STDEIG \rightarrow TRDEIG

- | | | | |
|---|------------------------------|-------------------------------|------------------|
| 1 | $LL^H = B$ | Cholesky factorization | $O(n^3)$ |
| 2 | $M \leftarrow L^{-1}AL^{-H}$ | Reduction to standard form | $O(n^3)$ |
| 3 | $T = Q^HMQ$ | Reduction to tridiagonal form | $O(n^3)$ |
| 4 | $TZ = Z\Lambda$ | Tridiagonal eigenproblem | $O(kn) - O(n^3)$ |
| 5 | $Y = QZ$ | Backtransformation #1 | $O(kn^2)$ |
| 6 | $X = L^{-H}Y$ | Backtransformation #2 | $O(kn^2)$ |

Stage 4: TRDEIG

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

Stage 4: TRDEIG

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

Stage 3: Reduction to TRDEIG

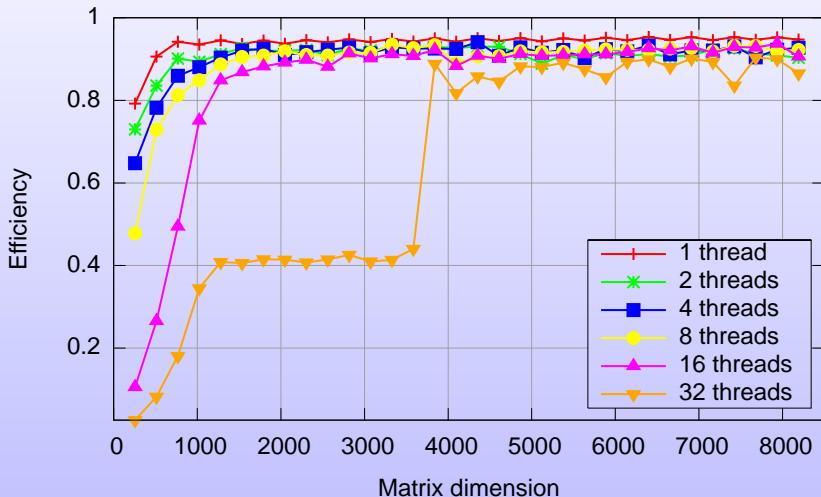
- 1-stage Householder
- Successive Banded Reduction

- 1 The Problem
- 2 Multicore Processors: MR³-SMP**
- 3 Distributed Memory Architectures: PMRRR
- 4 GPUs

Multi-threaded BLAS

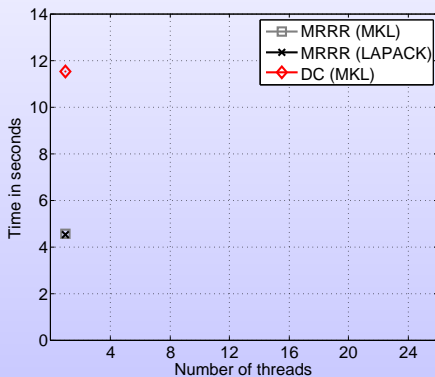
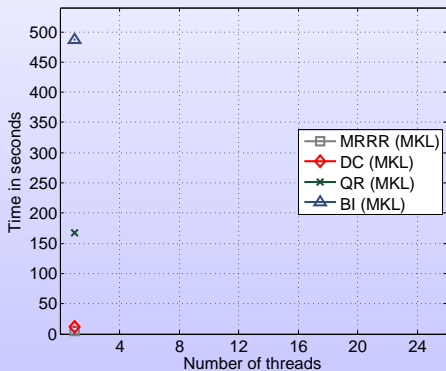
Xeon, 32 physical cores

Efficiency of GEMM



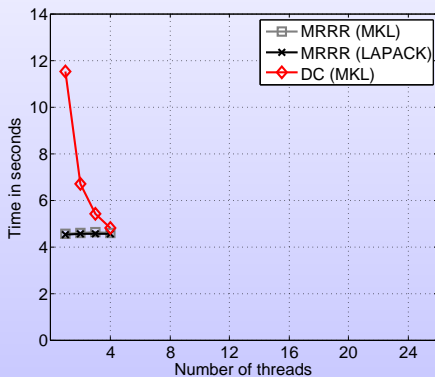
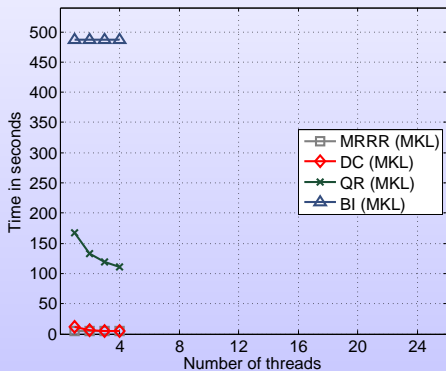
Multi-threaded BLAS for TRDEIG?

Tridiagonal eigensolvers. Matrix size=4289, from DFT.



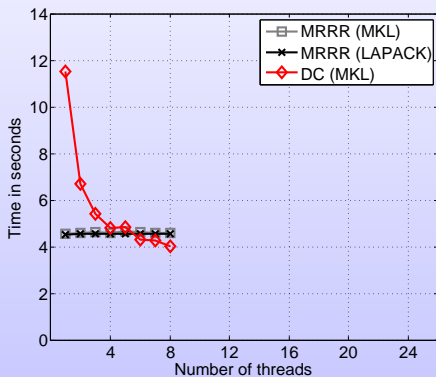
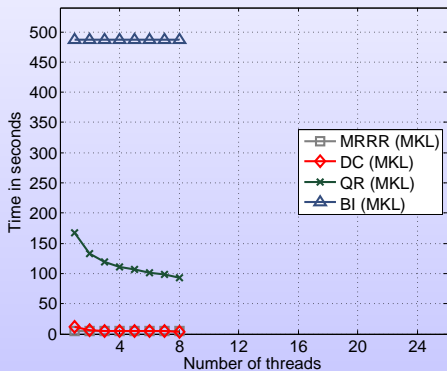
Multi-threaded BLAS for TRDEIG?

Tridiagonal eigensolvers. Matrix size=4289, from DFT.



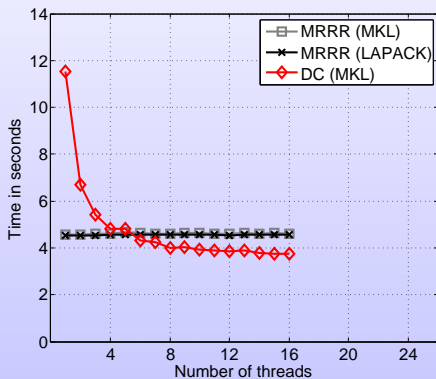
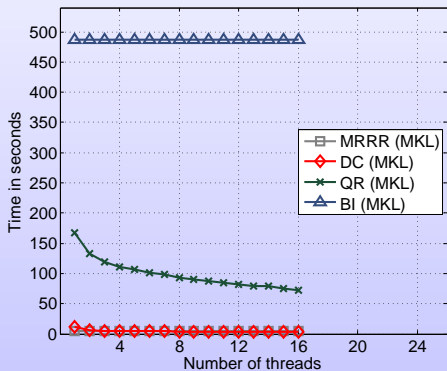
Multi-threaded BLAS for TRDEIG?

Tridiagonal eigensolvers. Matrix size=4289, from DFT.



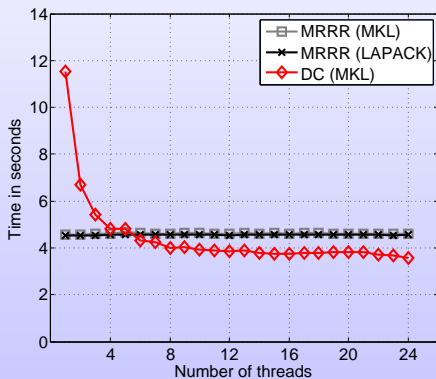
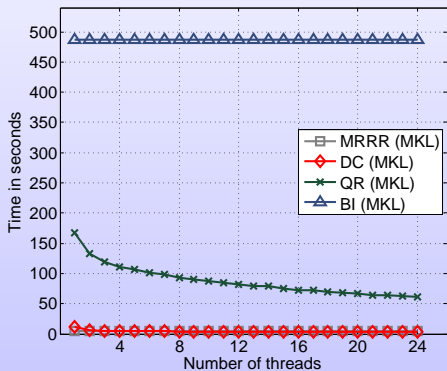
Multi-threaded BLAS for TRDEIG?

Tridiagonal eigensolvers. Matrix size=4289, from DFT.



Multi-threaded BLAS for TRDEIG?

Tridiagonal eigensolvers. Matrix size=4289, from DFT.

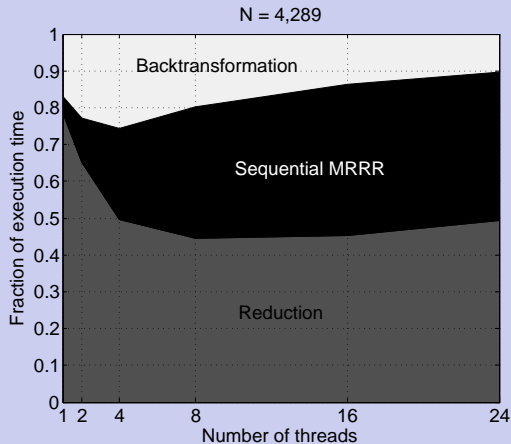


More motivation?

“MR3 is $O(n^2)$ anyway. . .”

More motivation?

“MR3 is $O(n^2)$ anyway...”

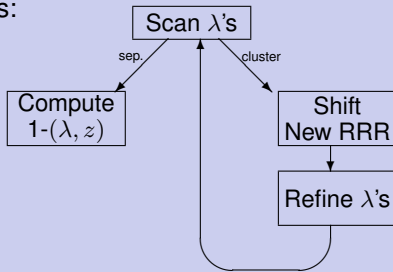


Multiple Relatively Robust Representations

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

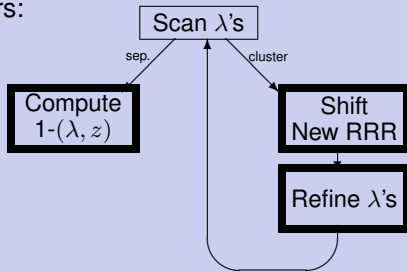
Multiple Relatively Robust Representations

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

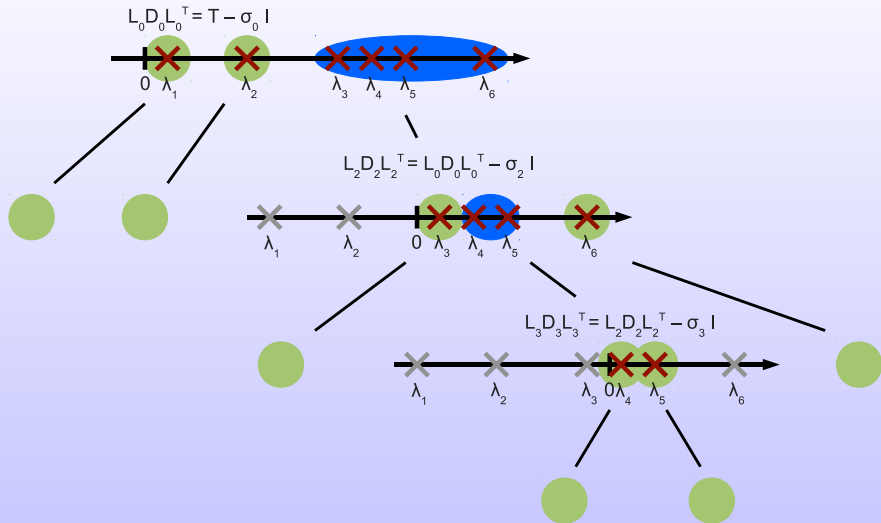


Multiple Relatively Robust Representations

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

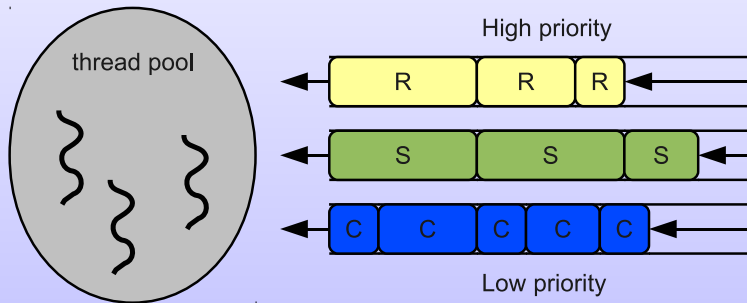


Representation Tree



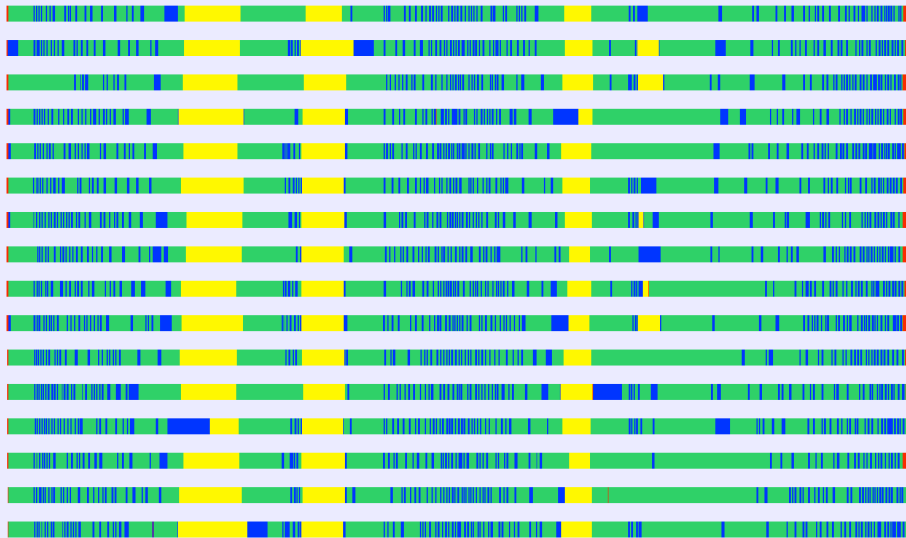
MR³-SMP: the work queue

- Tasks:
- a) Singleton ⇒ **S**: Eigenvector computation
 - b) Cluster ⇒ **C**: Shift + new representation (RRR)
 - c) New RRR ⇒ **R**: Eigenvalues refinement



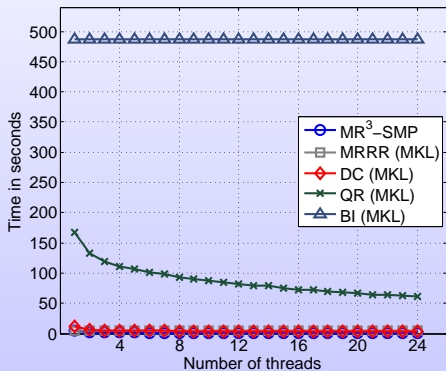
Example trace: 16 cores—eigenvectors

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



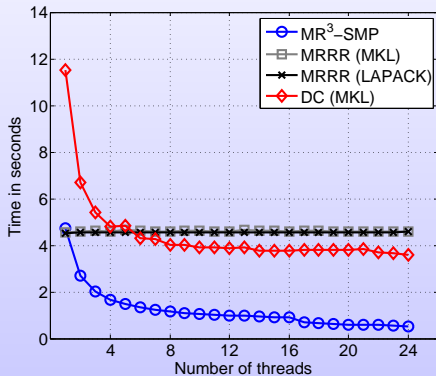
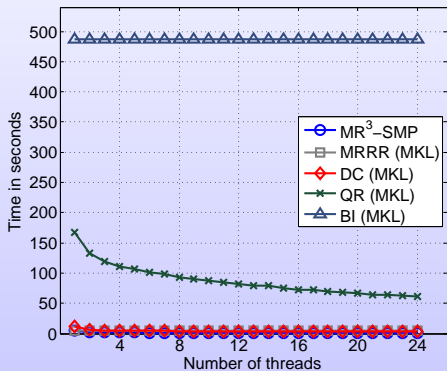
MR³-SMP: Timings

Matrix size=4289, from DFT.



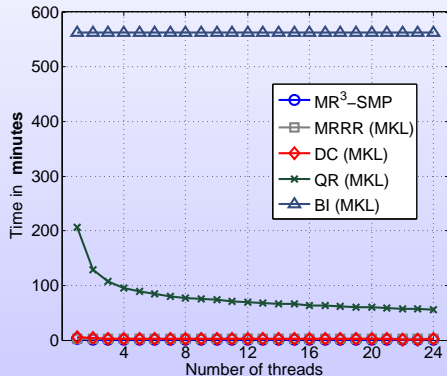
MR³-SMP: Timings

Matrix size=4289, from DFT.



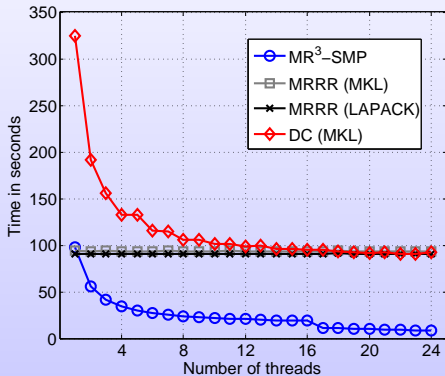
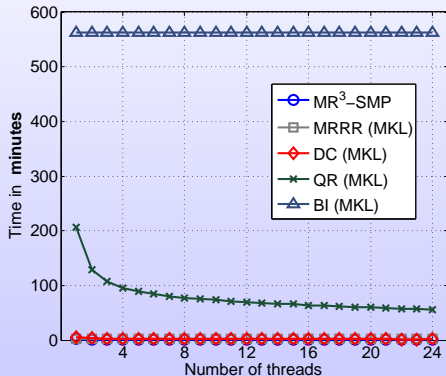
A larger example

Matrix size=16023; frequency response analysis of automobiles.



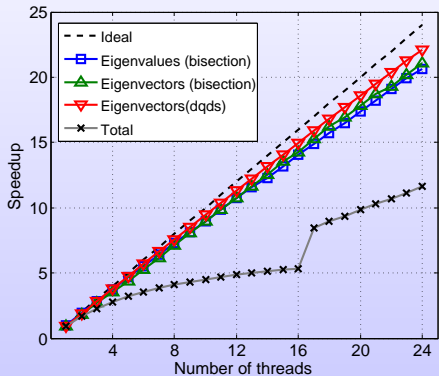
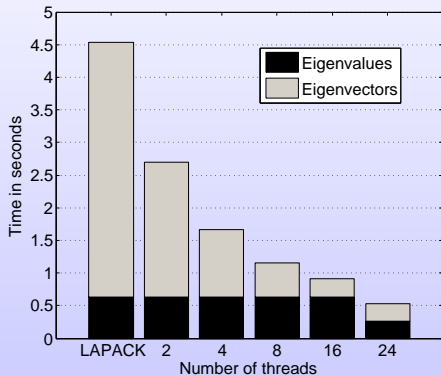
A larger example

Matrix size=16023; frequency response analysis of automobiles.

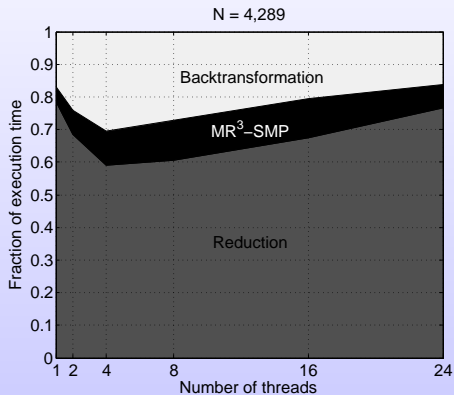
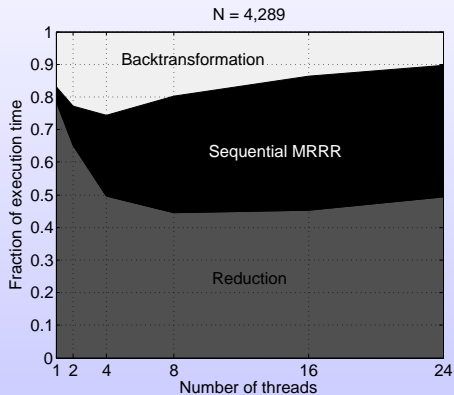


From 9+ hours to 8.3 seconds.

Speedups



3 stages: before and after



- 1 The Problem
- 2 Multicore Processors: MR³-SMP
- 3 Distributed Memory Architectures: PMRRR**
- 4 GPUs

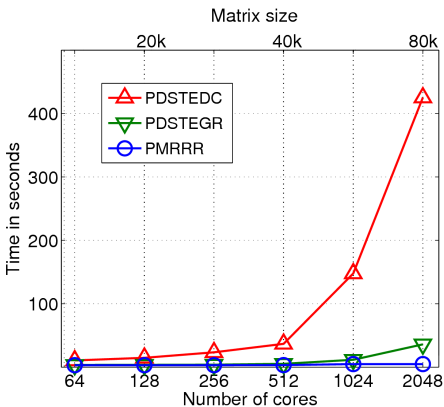
PMRRR, EleMRRR

- Static assignment of eigenpairs to nodes
- Multithreading
- Node-node communication: only eigenvalues
- PMRRR + Elemental \Rightarrow EleMRRR

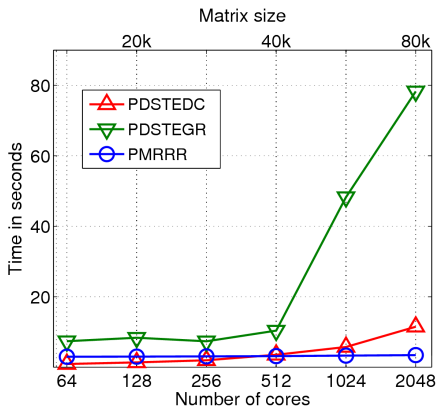
Generalized, standard and tridiagonal
hybrid eigensolvers

TRDEIG: PMRRR

1-2-1 matrix

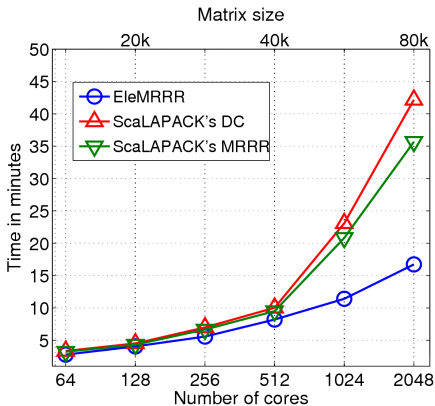


Wilkinson matrix

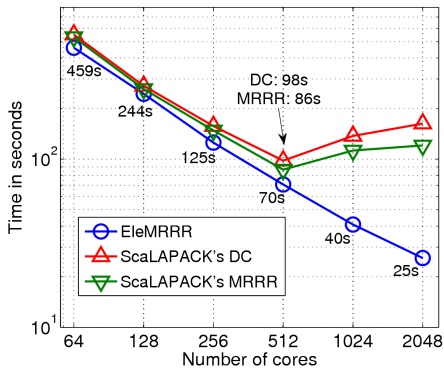


GENEIG: Weak & strong scaling

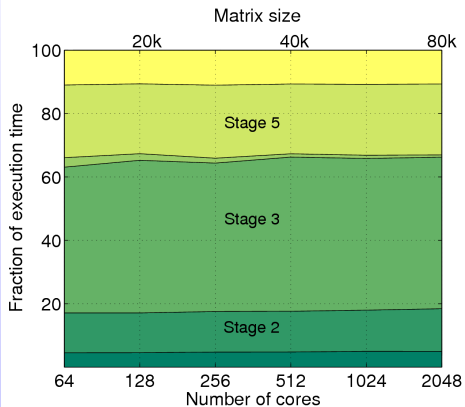
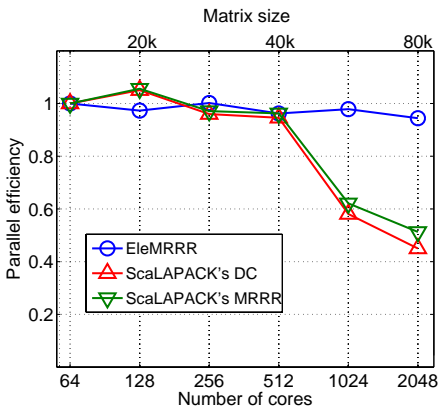
Weak scalability



Strong scalability, n=20000



GENEIG: Efficiency



- 1 The Problem
- 2 Multicore Processors: MR³-SMP
- 3 Distributed Memory Architectures: PMRRR
- 4 GPUs

mrrr_dp = data-parallel MRRR

n	rand(0,1)		rand(-1,1)	
	LAPACK	mrrr_dp	LAPACK	mrrr_dp
128	6.98	6.26	6.79	3.84
256	32.1	13.0	31.86	8.34
512	154.9	28.7	152.7	19.2
1024	656.1	60.2	647.6	54.0

Reduction to tridiagonal form

n	LAPACK	SBR	SBR + GPU
2048	0.23	0.6	0.58
6144	8.4	8.58	6.26
10240	40.5	30.4	20.32
24576	582.4	308.4	166.8

Reduction + backtransformation

n	LAPACK	SBR	SBR + GPU
2048	0.50	1.77	1.12
6144	13.5	29.0	12.7
10240	61.6	116.8	43.8
24576	845.1	1416.7	403.3