

Improving the Performance of the GMRES method using Mixed Precision Techniques

Neil Lindquist, Piotr Luszczek, Jack Dongarra

Smoky Mountains Conference

August 27th, 2020



GMRES

- General purpose, sparse linear solver
 - Iterative, Krylov solver
- Memory bound performance
 - Mix single and double precision

GMRES Algorithm

GMRES_{res}(A, x_0, b, M^{-1})

for $k = 0, 1, 2, \dots$

$$r_k \leftarrow b - Ax_k$$

$$z_k \leftarrow M^{-1}r_k$$

$$\beta \leftarrow \|z_k\|_2$$

$$V_{:,0} \leftarrow z_k/\beta$$

$$s \leftarrow [\beta, 0, 0, \dots, 0]^T$$

for $j = 0, 1, 2, \dots$

$$w \leftarrow M^{-1}AV_{:,j}$$

$$w, H_{:,j} \leftarrow \text{orthogonalize}(w, V_{:,j})$$

$$H_{j+1,j} \leftarrow \|w\|_2$$

$$V_{:,j+1} \leftarrow w/\|w\|_2$$

$$H_{:,j} \leftarrow G_0 G_1 \dots G_{j-1} H_{:,j}$$

$$G_j \leftarrow \text{rotation_matrix}(H_{:,j})$$

$$H_{:,j} \leftarrow G_j H_{:,j}$$

$$s \leftarrow G_j s$$

$$u_k \leftarrow VH^{-1}s$$

$$x_{k+1} \leftarrow x_k + u_k$$

Computing $Ax = b$. $A^{-1} \approx M^{-1}$
Restarts

Iteration count

GMRES Algorithm

GMRES_{res}(A, x_0, b, M^{-1})

for $k = 0, 1, 2, \dots$

Computing $Ax = b$. $A^{-1} \approx M^{-1}$

Restarts

Double:

$$r_k \leftarrow b - Ax_k$$

$$z_k \leftarrow M^{-1}r_k$$

$$\beta \leftarrow \|z_k\|_2$$

$$V_{:,0} \leftarrow z_k/\beta$$

$$s \leftarrow [\beta, 0, 0, \dots, 0]^T$$

for $j = 0, 1, 2, \dots$

$$w \leftarrow M^{-1}AV_{:,j}$$

$$w, H_{:,j} \leftarrow \text{orthogonalize}(w, V_{:,j})$$

$$H_{j+1,j} \leftarrow \|w\|_2$$

$$V_{:,j+1} \leftarrow w/\|w\|_2$$

$$H_{:,j} \leftarrow G_0 G_1 \dots G_{j-1} H_{:,j}$$

$$G_j \leftarrow \text{rotation_matrix}(H_{:,j})$$

$$H_{:,j} \leftarrow G_j H_{:,j}$$

$$s \leftarrow G_j s$$

$$u_k \leftarrow VH^{-1}s$$

Iteration count

Single:

Double:

$$x_{k+1} \leftarrow x_k + u_k$$

Effect on Memory Allocation

- Double: $8kn + 12n_{nz} + O(n + k^2)$ bytes
- Mixed: $4kn + 16n_{nz} + O(n + k^2)$ bytes
 - Including A, x, b
 - Excluding preconditioner
 - At most k inner iterations per restart
 - CSR matrix format

GMRES Algorithm

GMRES_{res}(A, x_0, b, M^{-1})

for $k = 0, 1, 2, \dots$

Computing $Ax = b$. $A^{-1} \approx M^{-1}$

Restarts

Double:

$$r_k \leftarrow b - Ax_k$$

$$z_k \leftarrow M^{-1}r_k$$

$$\beta \leftarrow \|z_k\|_2$$

$$V_{:,0} \leftarrow z_k/\beta$$

$$s \leftarrow [\beta, 0, 0, \dots, 0]^T$$

for $j = 0, 1, 2, \dots$

$$w \leftarrow M^{-1}AV_{:,j}$$

$$w, H_{:,j} \leftarrow \text{orthogonalize}(w, V_{:,j})$$

$$H_{j+1,j} \leftarrow \|w\|_2$$

$$V_{:,j+1} \leftarrow w/\|w\|_2$$

$$H_{:,j} \leftarrow G_0 G_1 \dots G_{j-1} H_{:,j}$$

$$G_j \leftarrow \text{rotation_matrix}(H_{:,j})$$

$$H_{:,j} \leftarrow G_j H_{:,j}$$

$$s \leftarrow G_j s$$

$$u_k \leftarrow VH^{-1}s$$

Iteration count

Single:

Double:

$$x_{k+1} \leftarrow x_k + u_k$$

GMRES Simplified Algorithm

$\text{GMRES}_{res}(A, x_0, b, M^{-1})$

for $k = 0, 1, 2, \dots$

Double: $r_k \leftarrow b - Ax_k$

Single: $u_k \leftarrow \text{GMRES}_{no\ res}(A, \vec{0}, r_k, M^{-1})$

Double: $x_{k+1} \leftarrow x_k + u_k$

GMRES Simplified Algorithm

$\text{GMRES}_{res}(A, x_0, b, M^{-1})$

for $k = 0, 1, 2, \dots$

Double: $r_k \leftarrow b - Ax_k$

Single: $u_k \leftarrow A^{-1} r_k$

Double: $x_{k+1} \leftarrow x_k + u_k$

Restart Strategies

1. Fixed # iterations

Restart Strategies

1. Fixed # iterations
2. Approximation of residual norm
 1. Threshold
 2. Improvement stops

Restart Strategies

1. Fixed # iterations
2. Approximation of residual norm
 1. Threshold
 2. Improvement stops
3. Detect basis losing independence

Restart Strategies

1. Fixed # iterations
2. Approximation of residual norm
 1. Threshold
 2. Improvement stops
3. Detect basis losing independence
4. # iterations for first restart

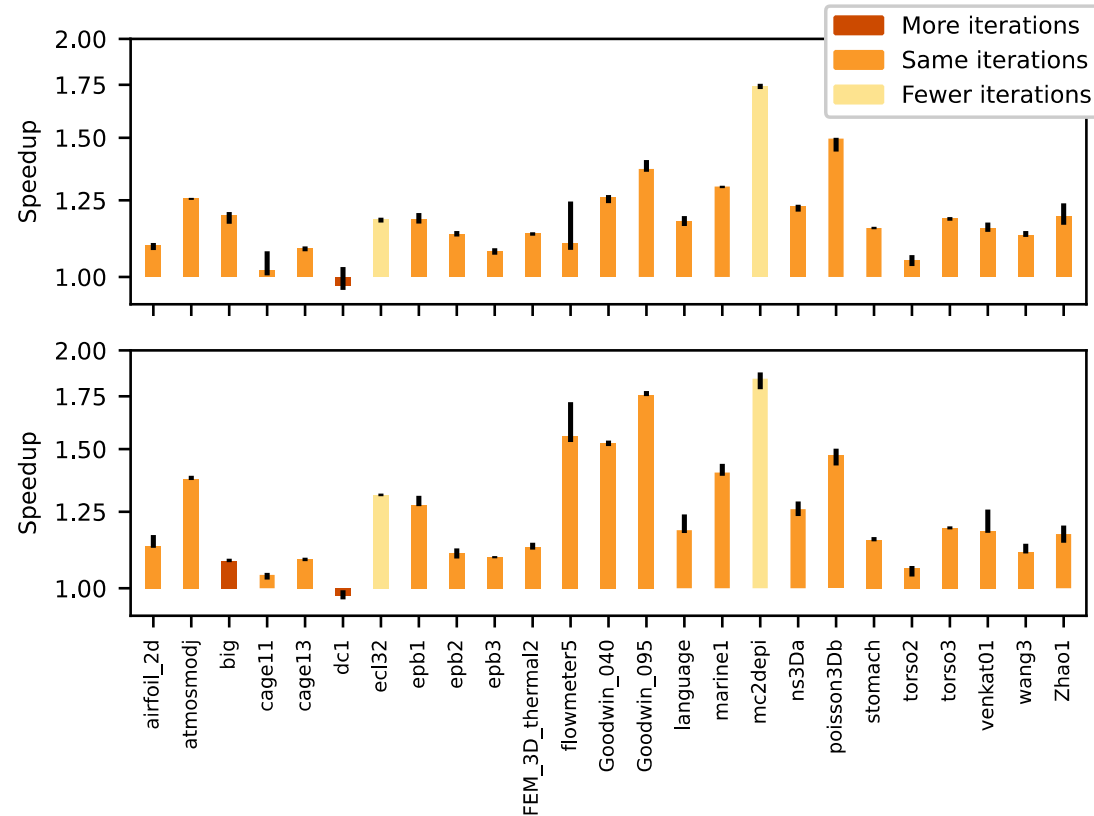
Effect on Performance: Configuration

- ILU(0) preconditioner (M^{-1})
- CSR matrix format
- KokkosKernels with Intel MKL
- 20-core Haswell node
 - 2x Intel® Xeon® E5-2650 v3 processors

Performance - Restarted

- Target accuracy $10^{-10} = \frac{\|b - Ax\|_2}{\|A\|_F \|x\|_2 + \|b\|_2}$
- Restart: half the iterations needed for double-precision implementation

Performance - Restarted

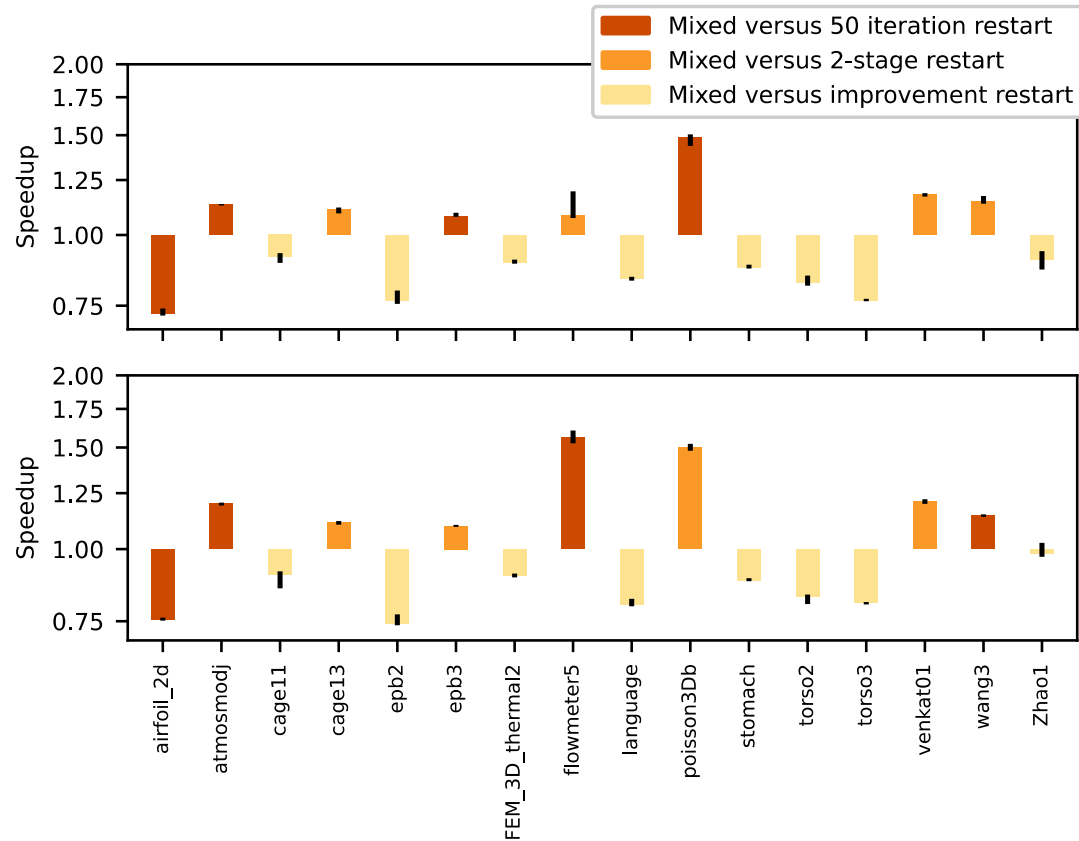


- Speedups of median times
 - 5 runs
 - Error bars: mins and maxes
- Geometric mean of speedup
 - MGS: 19%
 - CGSR: 24%

Performance – Non-restarted Option

- Target accuracy: $10^{-10} = \frac{\|b - Ax\|_2}{\|A\|_F \|x\|_2 + \|b\|_2}$
- Matrices restarted ≤ 50 iterations in last test
- Restart: 50 iterations plus one of:
 1. Residual improved by 10^{-6} for 1st restart, then same # iterations
 2. Residual improved by 10^{-8}
 - Non-restarted
 - only baseline
 3. Nothing
 - only baseline

Performance – Non-restarted Option



- Speedups of median times
 - 5 runs
 - Error bars: mins and maxs
- Geometric mean of speedup
 - MGS: -4%
 - CGSR: 0%
- Speedup \Leftrightarrow baseline restarted

Future Directions

- Choice of low-precision
 - Half, Bfloat16, integers
 - Compression
- Hardware
 - GPUs
 - Distributed
- Algorithms
 - Communication hiding/avoiding

Conclusions

- With restarts, mixed-precision GMRES
 - Converges to double-precision accuracy
 - Outperforms restarted, double-precision GMRES

Extra Slides

Effect on Convergence: Configuration

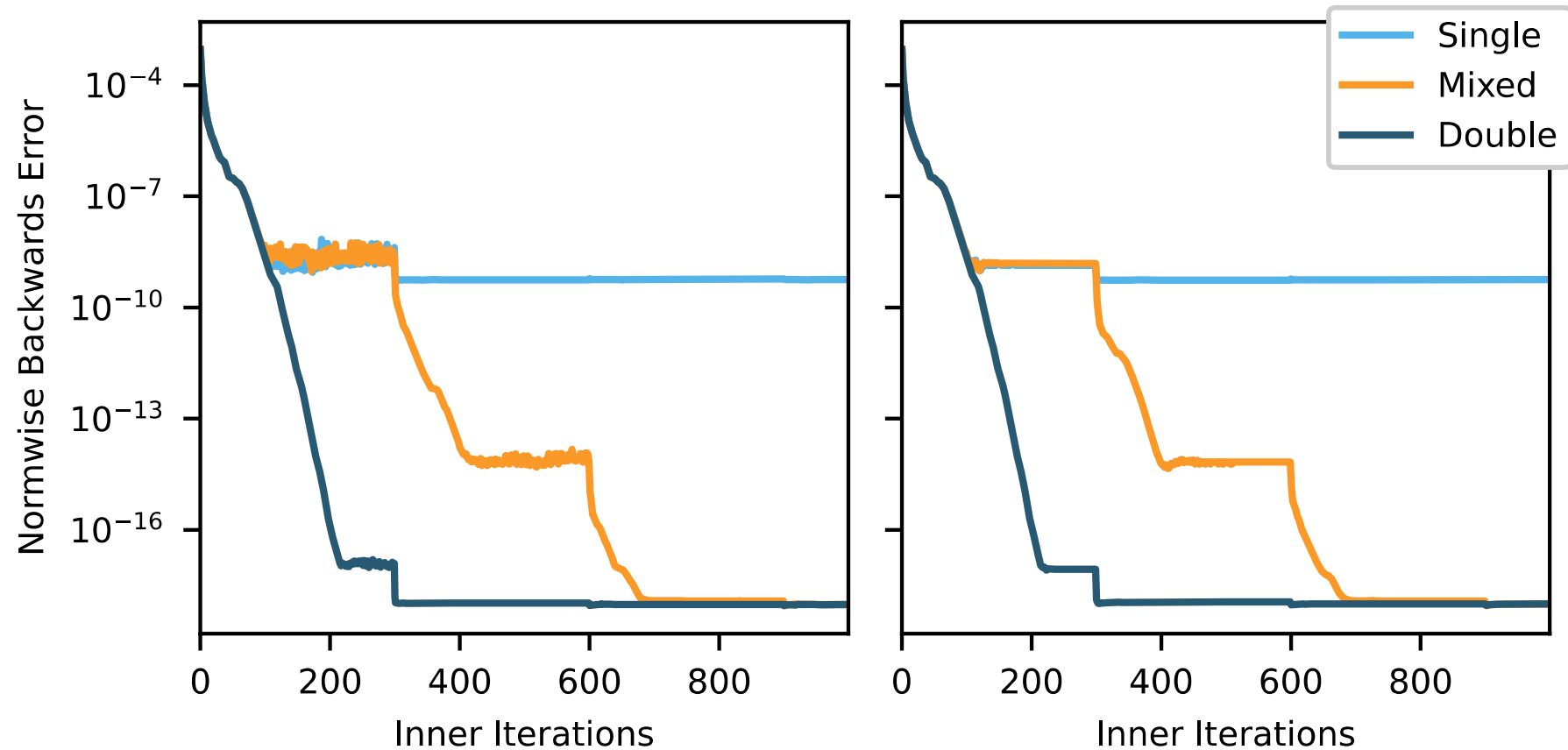
- ILU(0) preconditioner (M^{-1})
- CSR matrix format
- Custom, mixed precision kernels w/ Kokkos
- 20-core Haswell node
 - 2x Intel® Xeon® E5-2650 v3 processors
- 2 orthogonalization schemes
 - Modified Gram-Schmidt (MGS)
 - Classical Gram-Schmidt with Reorthogonalization (CGSR)

Effect on Convergence: Configuration

- airfoil_2d from SuiteSparse collection
 - $n = 14,214$
 - $nnz = 259,688$
 - $\kappa_2 = 1.8 \cdot 10^6$
- Error if GMRES stopped

$$\frac{\|b - Ax\|_2}{\|A\|_F \|x\|_2 + \|b\|_2}$$

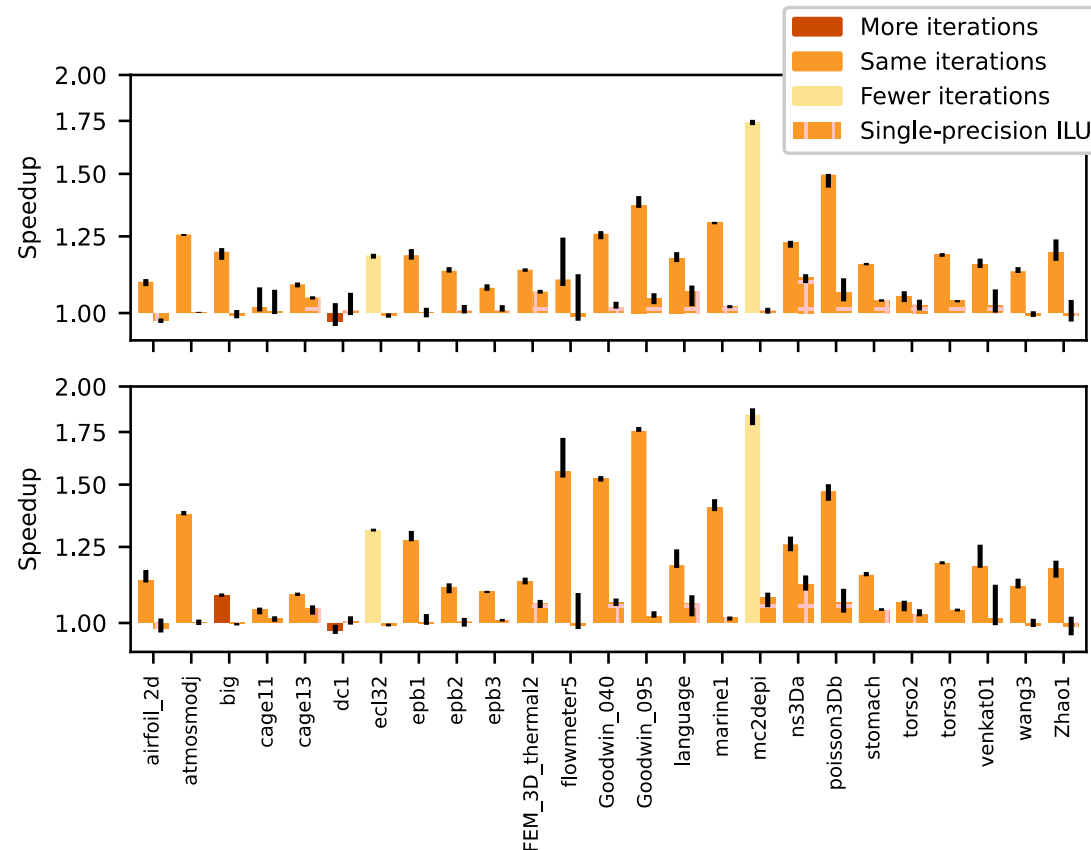
Accuracy results



Modified Gram-Schmidt
Orthogonalization (MGS)

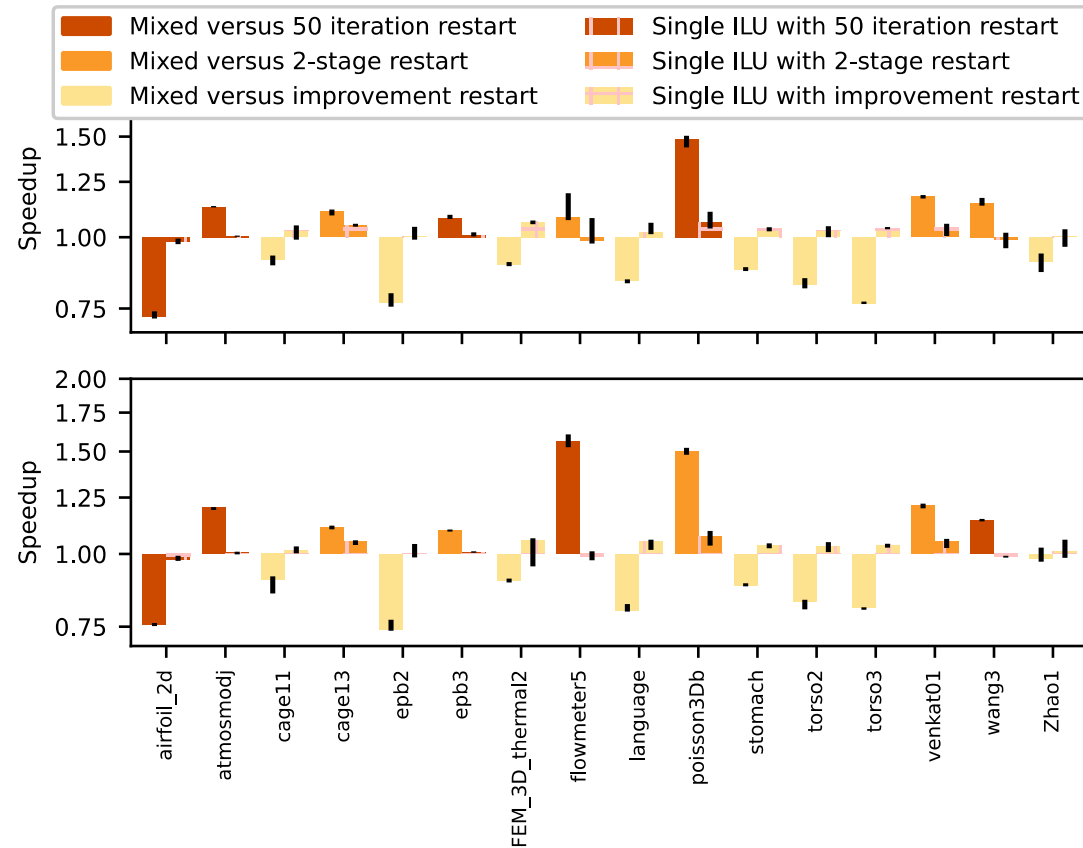
Classical Gram-Schmidt with
Reorthogonalization (CGSR)

Performance - Restarted



- Speedups of median times for 5 runs
 - Error bars for minimums and maximums
- Geometric mean of speedup
 - MGS: 19%
 - CGSR: 24%
 - Single-precision preconditioner: 2%

Performance - Restarted



- Speedups of median times for 5 runs
 - Error bars for minimums and maximums
- Geometric mean of speedup
 - MGS: -4%
 - CGSR: 0%
 - Single-precision preconditioner:
 - MGS: 2%
 - CGSR: 1%