

# **Multiprecision Approach in GMRES and its Effects on Performance**

Neil Lindquist, Piotr Luszczek, Jack Dongarra

SIAM LA21

May 18<sup>th</sup>, 2021

# GMRES

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

# GMRES Algorithm

GMRES<sub>res</sub>( $A, x_0, b, M^{-1}$ )

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<sub>res</sub>( $A, x_0, b, M^{-1}$ )

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



# Performance – Precision Choices

- FP64 GMRES
- Compressed basis GMRES
- Our mixed precision GMRES
- FP32 GMRES

# Performance – Test Setup

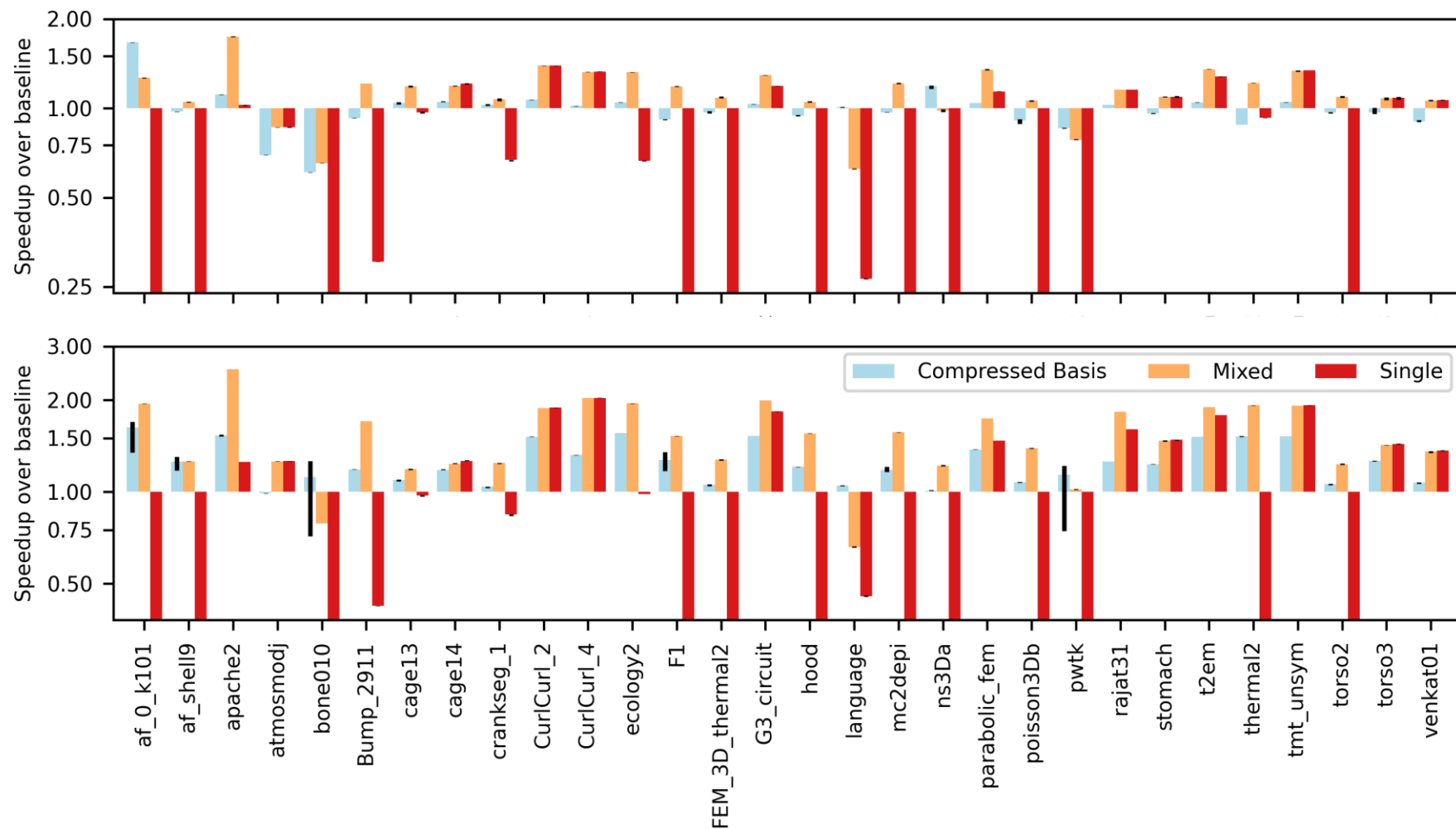
- Target accuracy  $10^{-10} = \frac{\|b - Ax\|_2}{\|A\|_F \|x\|_2 + \|b\|_2}$
- Restart strategies:
  - I. 100 inner iterations
  - II. 100 inner iterations or residual estimate of  $10^{-10}$
  - III. First: 100 inner iterations or residual estimate of  $10^{-6}$   
Then: same number of inner iterations
- 20-core Haswell node with NVIDIA V100 GPU
  - cuSparse, cuBLAS, Kokkos
- CSR matrix format



# Performance – Plots

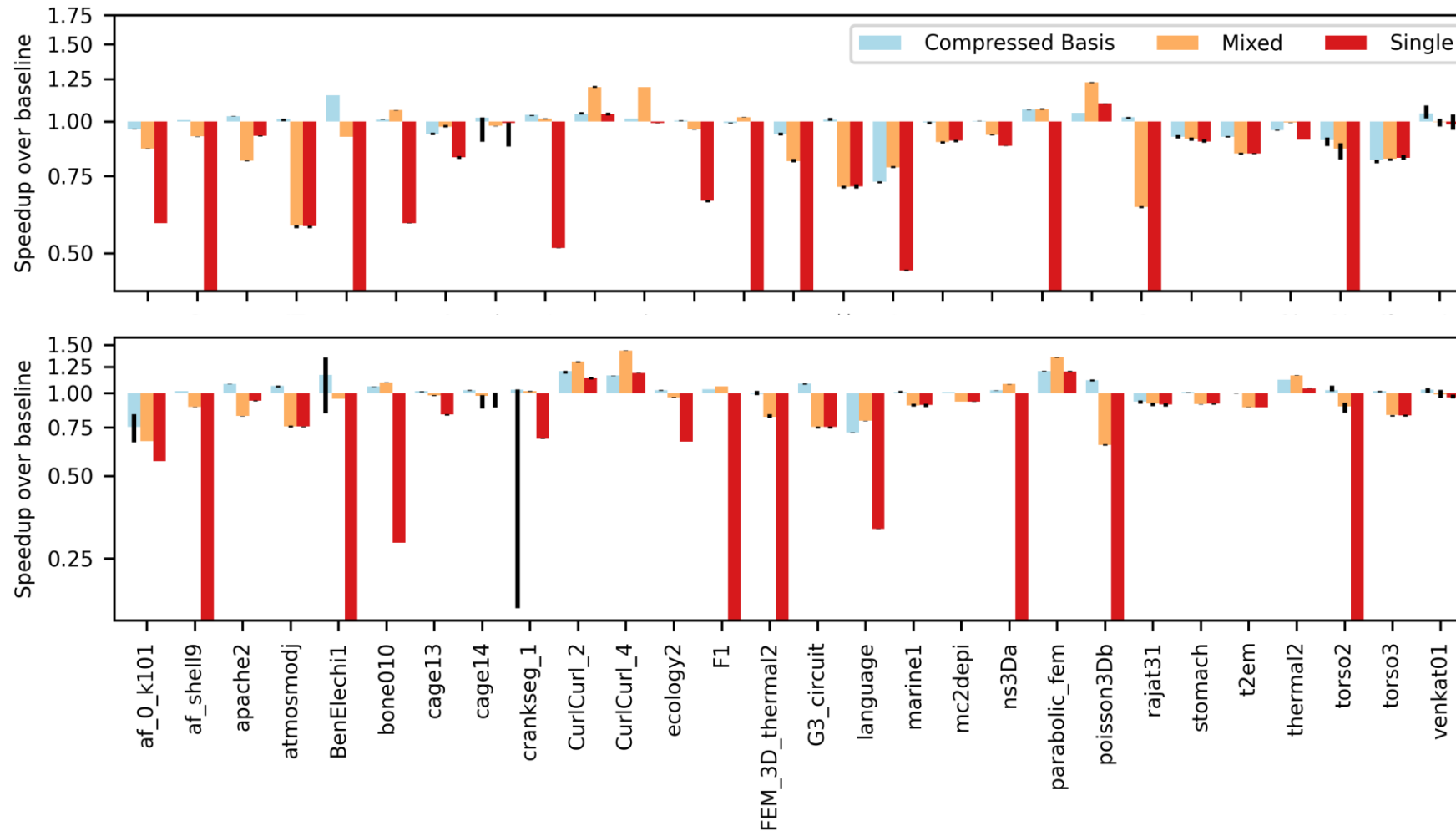
- 3 runs each – take median runtime
- Plotted speedups over FP64
  - Error bars – min and max speedup
- Performance summarized w/ geometric mean

# Performance – Scalar Jacobi



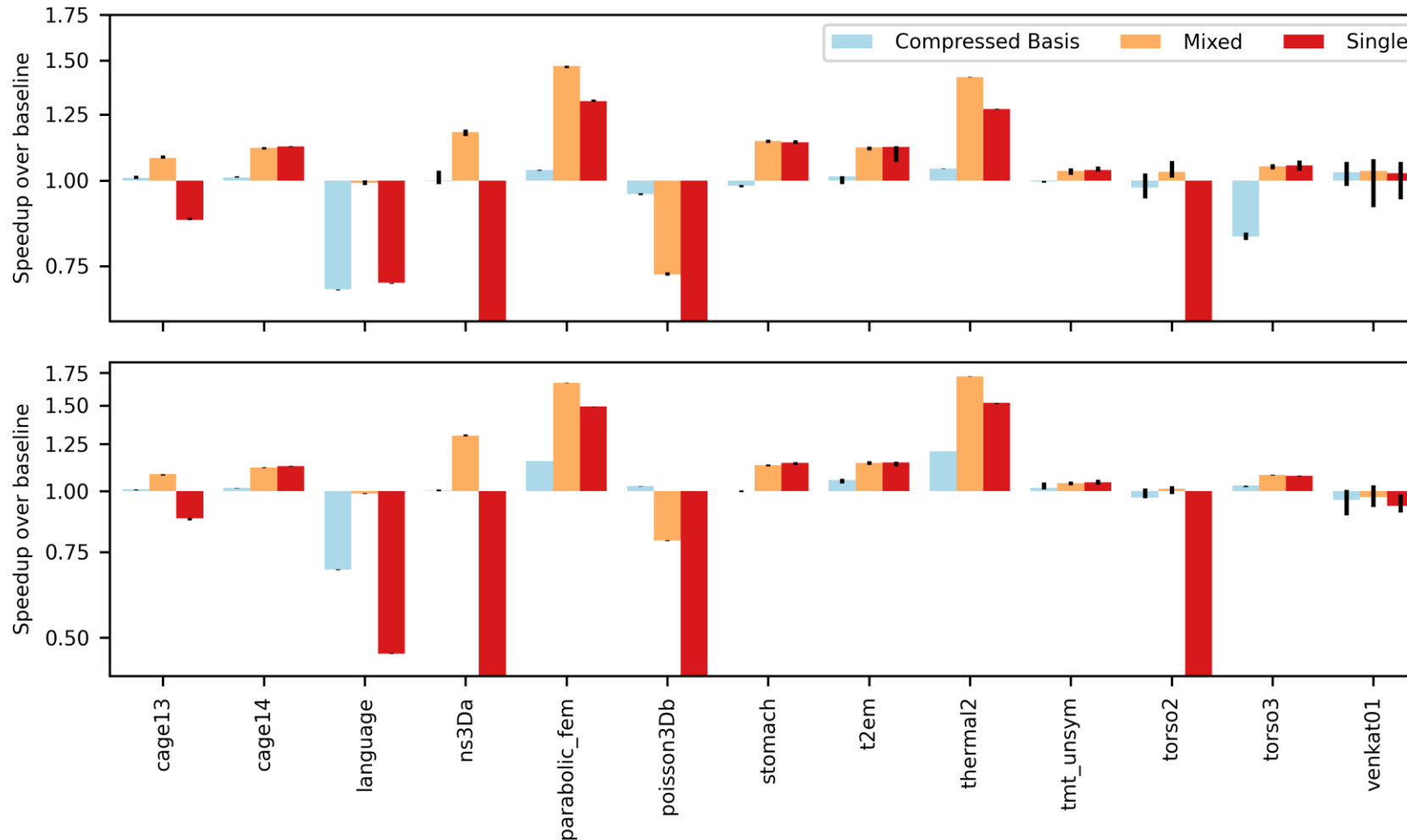
- C. Basis
  - MGS: -2%
  - CGSR: 24%
- Mixed
  - MGS: 12%
  - CGSR: 48%
- Single
  - MGS: -8%
  - CGSR: 24%

# Performance – ILU(0)



- C. Basis
  - MGS: -2%
  - CGSR: 3%
- Mixed
  - MGS: -9%
  - CGSR: -6%
- Single
  - MGS: -21%
  - CGSR: -20%

# Performance – ILU(0) w/ Jacobi Solves



- C. Basis
  - MGS: -4%
  - CGSR: 0%
- Mixed
  - MGS: 9%
  - CGSR: 13%
- Single
  - MGS: 5%
  - CGSR: 4%

# Future Directions

- Choice of low-precision
  - Half, Bfloat16
  - Compression
- Distributed systems
- Other Krylov methods
- Applications

# Conclusions

- When restarted, mixed-precision GMRES can provide speedups
  - Depending on the preconditioner



# Extra Slides

# Test Configuration Details

- CUDA 10.2.199, Kokkos 3.1.01, GCC 7.3.0
- <https://bitbucket.org/icl/mixed-precision-gmres>
  - tag [TPDS](#)

# Publications

- N. Lindquist, P. Luszczyk, and J. Dongarra, “Improving the performance of the GMRES method using mixed-precision techniques,” in Driving Scientific and Engineering Discoveries through the Convergence of HPC, Big Data and AI. DOI: [10.1007/978-3-030-63393-6\\_4](https://doi.org/10.1007/978-3-030-63393-6_4)
- [Submitted] N. Lindquist, P. Luszczyk, and J. Dongarra, “Accelerating restarted GMRES with mixed precision arithmetic.”

# Effect on Convergence: Configuration

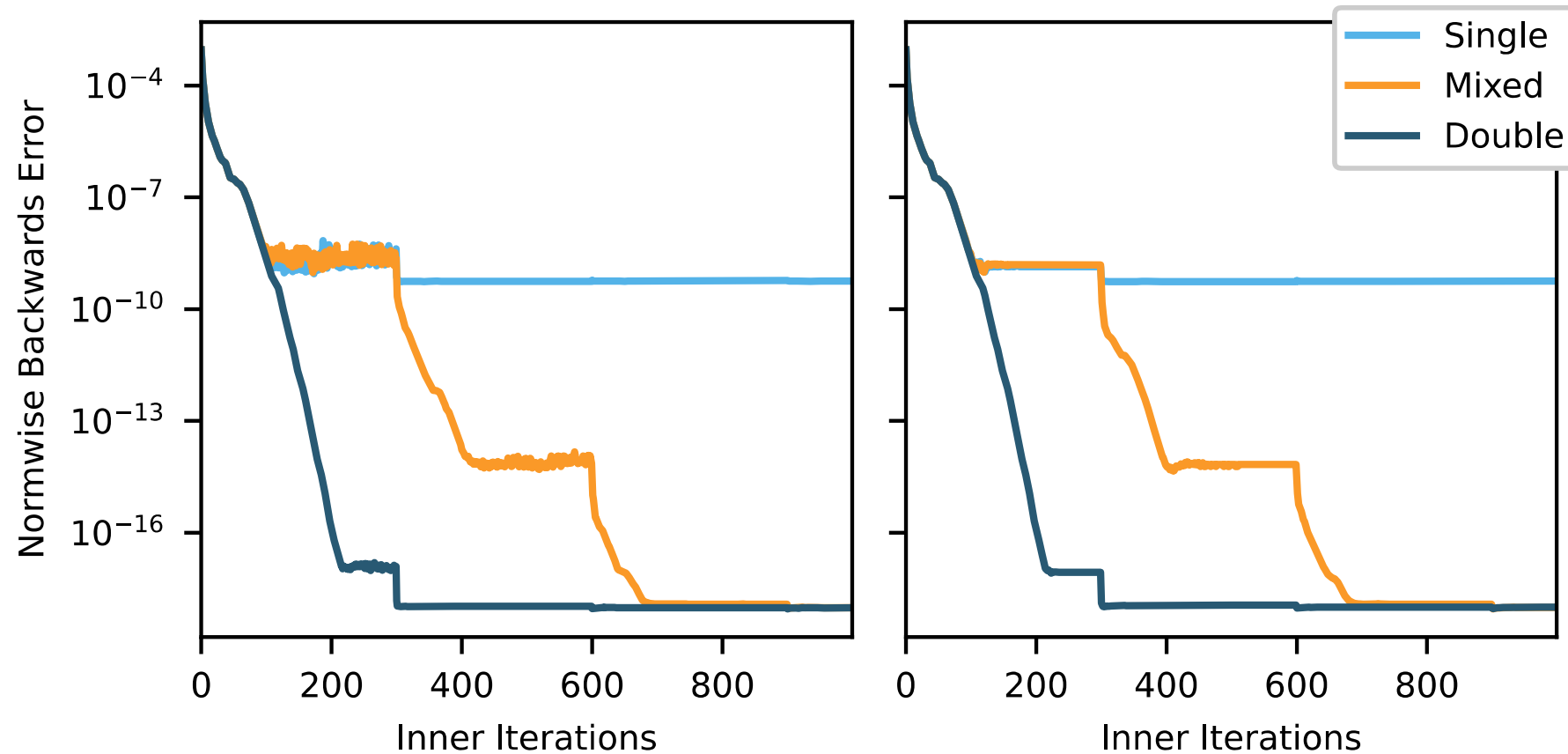
- ILU(0) preconditioner ( $M^{-1}$ )
- CSR matrix format
- Custom, mixed precision kernels w/ Kokkos
- 20-core Haswell node
  - 2x Intel® Xeon® E5-2650 v3 processors

# 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)