# Accelerating GMRES via Mixed Precision

**Neil Lindquist**, Piotr Luszczek, Jack Dongarra

12th JLESC Workshop

February 25th, 2020

# GMRES

- General purpose, sparse linear solver

  - Iterative, Krylov solver

- Memory bound performance

  - Mix single and double precision

# GMRES Algorithm

$\text{GMRES}_{res}(\boldsymbol{A}, \boldsymbol{x_0}, \boldsymbol{b}, \boldsymbol{M^{-1}})$    Computing $\boldsymbol{Ax = b}$. $\boldsymbol{A^{-1} \approx M^{-1}}$

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

$$\boldsymbol{r_k} \leftarrow \boldsymbol{b} - \boldsymbol{Ax_k}$$

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

$$\beta \leftarrow \|\boldsymbol{z_k}\|_2$$

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

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

$\quad\quad$ for j $= 0, 1, 2, \dots$    Iteration count

$$\boldsymbol{w} \leftarrow \boldsymbol{M^{-1} A V_{:,j}}$$

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

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

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

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

$$\boldsymbol{G_j} \leftarrow rotation\_matrix(\boldsymbol{H_{:,j}})$$

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

$$\boldsymbol{s} \leftarrow \boldsymbol{G_j s}$$

$$\boldsymbol{u_k} \leftarrow \boldsymbol{V H^{-1} s}$$

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

# GMRES Algorithm

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

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

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

Double:

$$r_k \leftarrow b - Ax_k$$

$$z_k \leftarrow M^{-1}r_k$$
$$\beta \leftarrow \|z_k\|_2$$
$$V_{:,0} \leftarrow z_k/\beta$$
$$s \leftarrow [\beta, 0, 0, \ldots, 0]^T$$

for j = 0, 1, 2, \ldots

Iteration count

$$w \leftarrow M^{-1}AV_{:,j}$$
$$w, H_{:,j} \leftarrow orthogonalize(w, V_{:,j})$$
$$H_{j+1,j} \leftarrow \|w\|_2$$
$$V_{:,j+1} \leftarrow w/\|w\|_2$$

Single:

$$H_{:,j} \leftarrow G_0 G_1 \ldots G_{j-1} H_{:,j}$$
$$G_j \leftarrow rotation\_matrix(H_{:,j})$$
$$H_{:,j} \leftarrow G_j H_{:,j}$$
$$s \leftarrow G_j s$$

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

Double:

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

# GMRES Simplified Algorithm

$$\text{GMRES}_{res}(A, x_0, b, M^{-1})$$
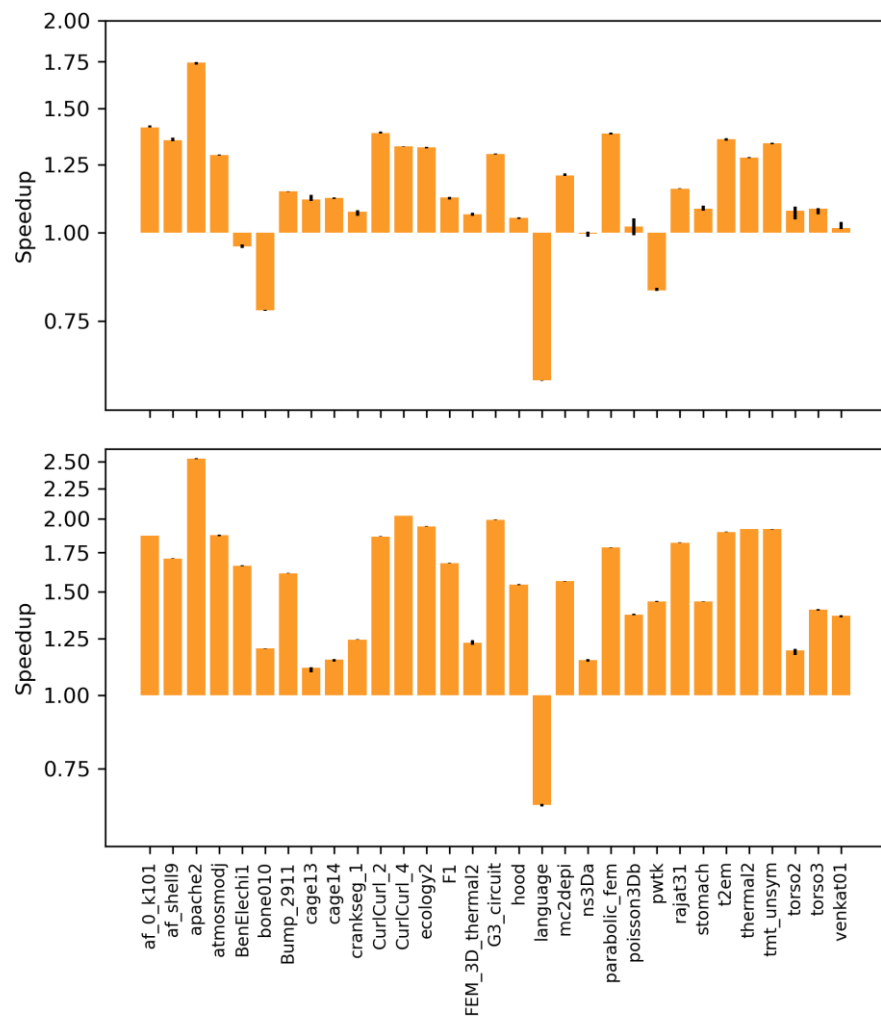
$$\text{for } k = 0, 1, 2, \ldots$$

Double: $\quad r_k \leftarrow b - A x_k$

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

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

# GMRES Simplified Algorithm

$$\text{GMRES}_{res}(\boldsymbol{A}, \boldsymbol{x_0}, \boldsymbol{b}, \boldsymbol{M}^{-1})$$
$$\text{for } k = 0, 1, 2, \dots$$

Double: $\quad \boldsymbol{r_k} \leftarrow \boldsymbol{b} - \boldsymbol{A}\boldsymbol{x_k}$

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

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

# Performance

- Target accuracy $10^{-10} = \dfrac{\|b - Ax\|_2}{\|A\|_F \|x\|_2 + \|b\|_2}$

- 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
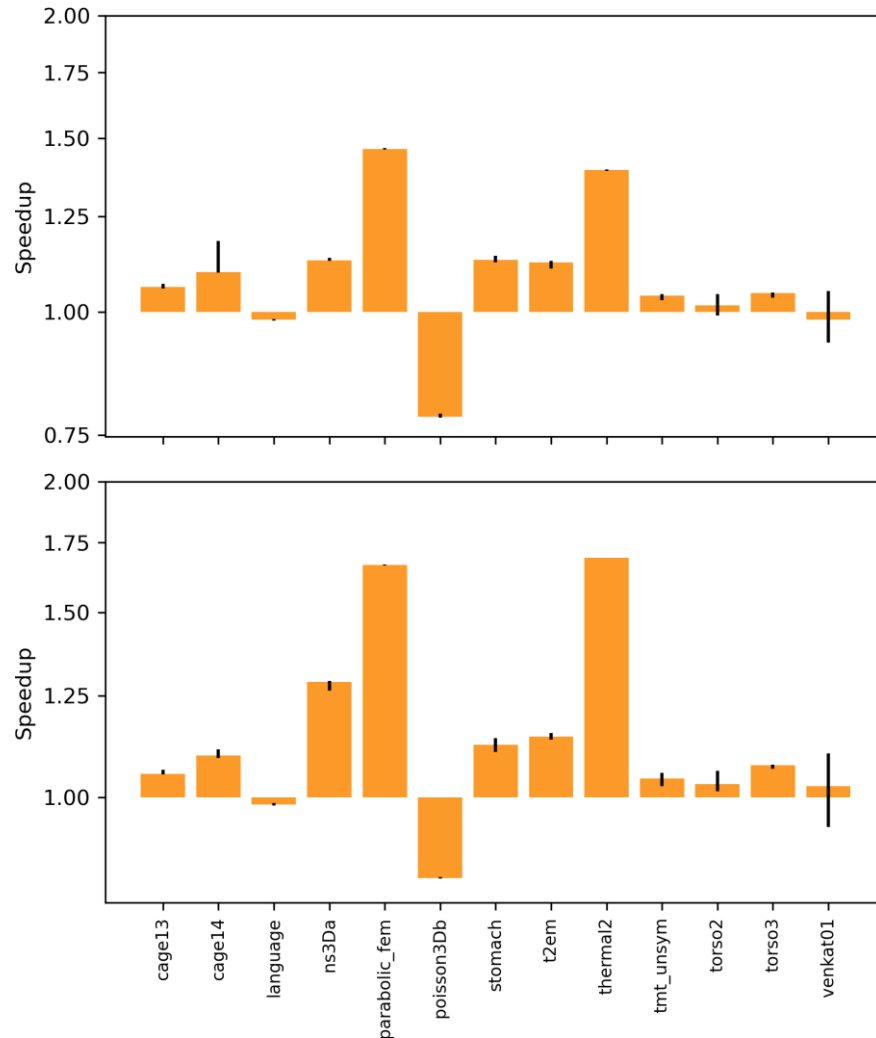
# Performance – Scalar Jacobi



- Speedups
  - Median time of 3 run
  - 3 runs
  - Error bars: mins and maxes

- Geometric mean of speedup
  - MGS: 14%
  - CGSR: 54%

# Performance – ILU(0)



- Speedups
  - Median time of 3 run
  - 3 runs
  - Error bars: mins and maxes

- Geometric mean of speedup
  - MGS: -7%
  - CGSR: -4%

# Performance – ILU(0) with Jacobi Solves



- ILU(0) w/ 5 Jacobi iterations for each triangular solve
- Speedups
  - Median time of 3 run
  - 3 runs
  - Error bars: mins and maxes
- Geometric mean of speedup
  - MGS: 8%
  - CGSR: 14%

# Future Directions

- Choice of low-precision
  - Half, Bfloat16
  - Compression

- Distributed systems

- Other Krylov methods

- Applications

# Conclusions

- When restarted, mixed-precision GMRES often outperforms double-precision GMRES

# 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-perf

# Publications

- N. Lindquist, P. Luszczek, 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](#)

- [Submitted] N. Lindquist, P. Luszczek, and J. Dongarra, "Accelerating restarted GMRES with mixed precision arithmetic," in Transactions on Parallel and Distributed Systems.

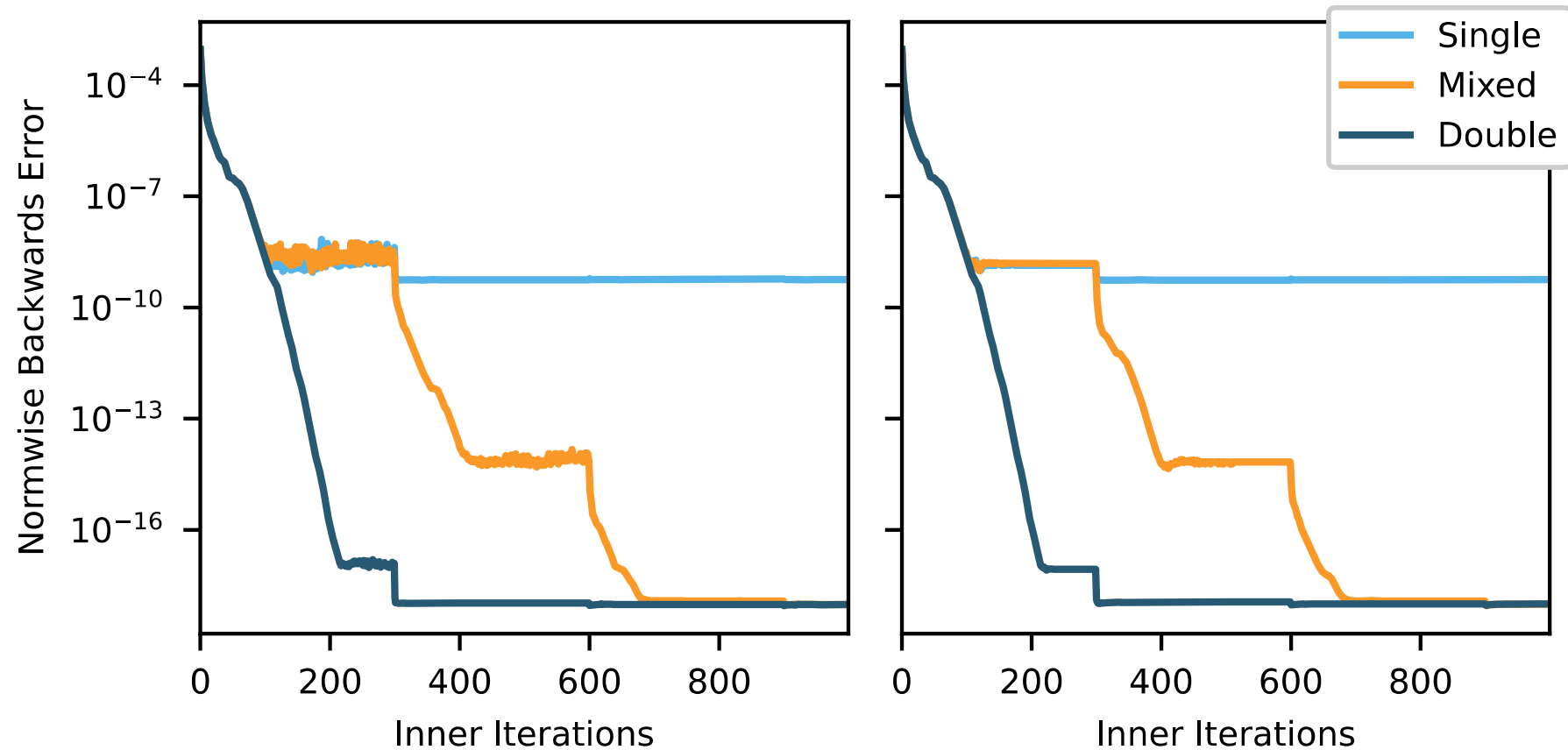# Effect on Convergence: Configuration

- 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)