# An attempt of exploiting low precision computing in the GMRES(m) method

March 2, 2023
2023 SIAM Conference on Computational Science and Engineering (CSE23)
@RAI Congress Centre, Amsterdam, The Netherlands

## Takeshi Fukaya  (Hokkaido University)

Collaborator:
Yingqi Zhao and Takeshi Iwashita (Hokkaido University)

# Introduction

## ◆Problem setting

Solving a linear system with a sparse coefficient matrix:

$$A\boldsymbol{x} = \boldsymbol{b}$$

$A$: $n$-dimensional sparse matrix (regular, real and not symmetric)

## ◆Target algorithm

GMRES(m) method: restarted GMRES method

## ◆Goal

Through numerical experiments, to investigate possibilities of exploiting low precision computing in the GMRES(m) method without the loss of accuracy of obtained solutions (baseline: GMRES(m) with FP64 only).

# Key idea

◆ **Iterative refinement (IR) & mixed precision**

> Step 1.  computing residual:  $r_k := b - A x_k$
>
> Step 2.  solving error equation:  $e_k = A^{-1} r_k$  (solving a linear system)
>
> Step 3.  updating the solution:  $x_{k+1} := x_k + e_k$

**In Step 2, low precision computing can be acceptable.**

◆ **Relation between IR and GMRES(m)**

> 1: **repeat**
>
> 2:  Solve $A x = b$ by $m$-iteration GMRES with the initial guess $x_0$, and find the solution $x_m$.
>
> 3:  $x_0 \leftarrow x_m$  (update the initial guess)
>
> 4: **until** satisfy required accuracy condition or attain maximum iteration number

**GMRES(m) has the structure of IR (IR using m-step GMRES for Step 2).**

**Input**: An initial guess $\boldsymbol{x}_0$

Check convergence here (using FP64)

1: **repeat**
2: $\quad \boldsymbol{r}_0 \leftarrow \boldsymbol{b} - A\boldsymbol{x}_0, \quad \beta \leftarrow \|\boldsymbol{r}_0\|_2$
3: $\quad \boldsymbol{v}_0 \leftarrow \boldsymbol{r}_0 / \beta$
4: $\quad$ Compute $m$-step Arnoldi process with $A$ and $\boldsymbol{v}_0$, and get $V_m$ and $\overline{H}_m$.
5: $\quad$ Compute $\boldsymbol{y}_m$ from $\beta$ and $\overline{H}_m$.
6: $\quad \boldsymbol{e}_m \leftarrow V_m \boldsymbol{y}_m$
7: $\quad \boldsymbol{x}_0 \leftarrow \boldsymbol{x}_0 + \boldsymbol{e}_m$
8: **until** satisfy required accuracy condition or attain maximum iteration number

corresponds to Step 2 in IR (solving error equation)

Low precision computing can be acceptable.

## ◆ What we present in this talk

Numerical results of two attempts of introducing low precision computing:

- GMRES(m) using FP32 and FP64

- GMRES(m) using low precision data including those lower than FP32.

# GMRES(m) using FP32 & FP64

**Input**: An initial guess $\boldsymbol{x}_0$

1: $A^{(\text{FP32})} \leftarrow \text{ToFP32}(A)$ — Prepare matrix data in FP32

2: **repeat**

3:    $\boldsymbol{r}_0 \leftarrow \boldsymbol{b} - A\boldsymbol{x}_0, \quad \beta \leftarrow \|\boldsymbol{r}_0\|_2$

convert to FP32 data

4:    $\boldsymbol{v}_0^{(\text{FP32})} \leftarrow \text{ToFP32}(\boldsymbol{r}_0/\beta), \quad \beta^{(\text{FP32})} \leftarrow \text{ToFP32}(\beta)$

5:    Compute $m$-step Arnoldi process in low precision with $A^{(\text{FP32})}$ and $\boldsymbol{v}_0^{(\text{FP32})}$, and get $V_m^{(\text{FP32})}$ and $\overline{H}^{(\text{FP32})}{}_m$.

6:    Compute in low precision $\boldsymbol{y}_m^{(\text{FP32})}$ from $\beta^{(\text{FP32})}$ and $\overline{H}^{(\text{FP32})}{}_m$.

7:    $\boldsymbol{e}_m^{(\text{FP32})} \leftarrow V_m^{(\text{FP32})} \boldsymbol{y}_m^{(\text{FP32})}$

convert to FP64 data

8:    $\boldsymbol{x}_0 \leftarrow \boldsymbol{x}_0 + \text{ToFP64}(\boldsymbol{e}_m^{(\text{FP32})})$

9: **until** satisfy required accuracy condition or attain maximum iteration number
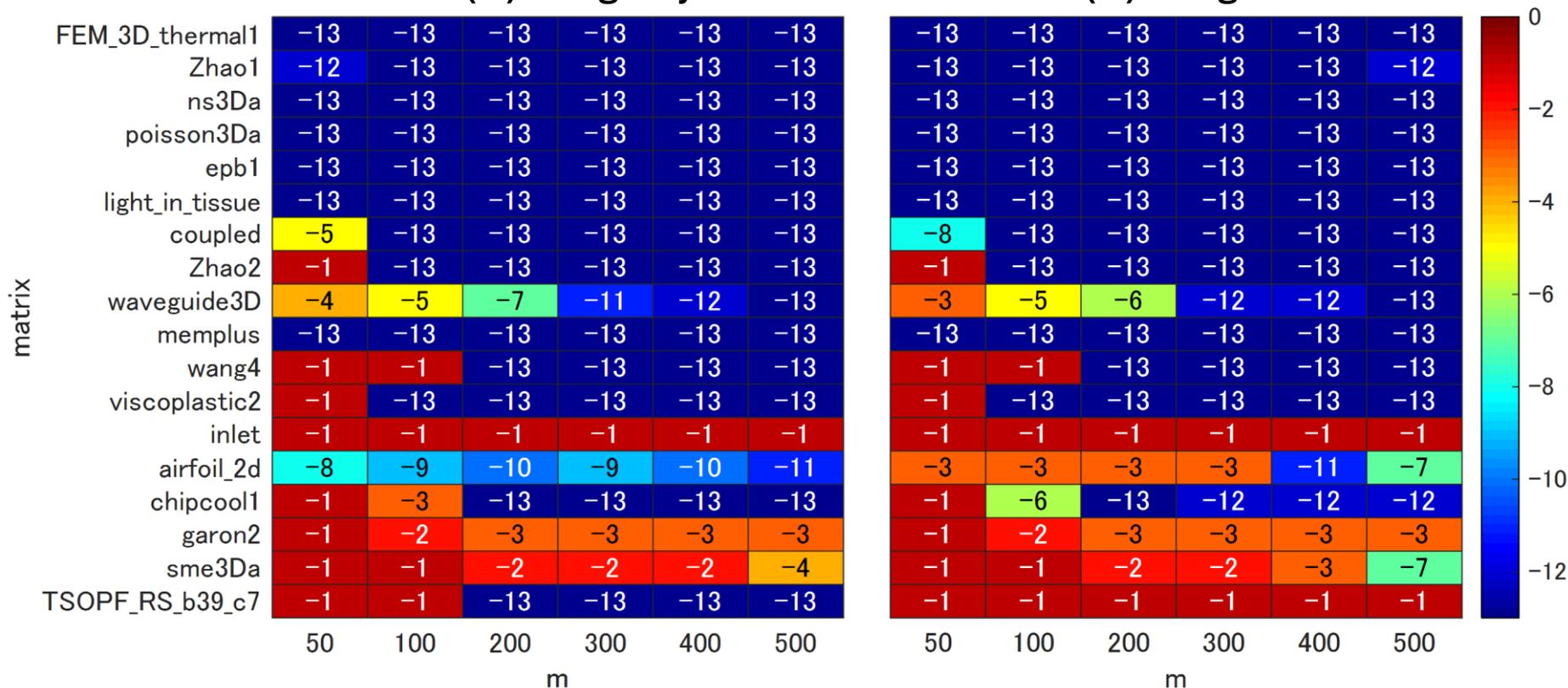
We focus on the numerical behavior (convergence property) of the mixed-precision GMRES(m) method using FP32 and FP64 compared with that of GMRES(m) using only FP64. (For some problems, its effectiveness in execution time has been already reported.)

$$\log_{10} \frac{\|b - Ax\|_2}{\|b\|_2}$$ **at the maximum iterations (or convergence condition)**



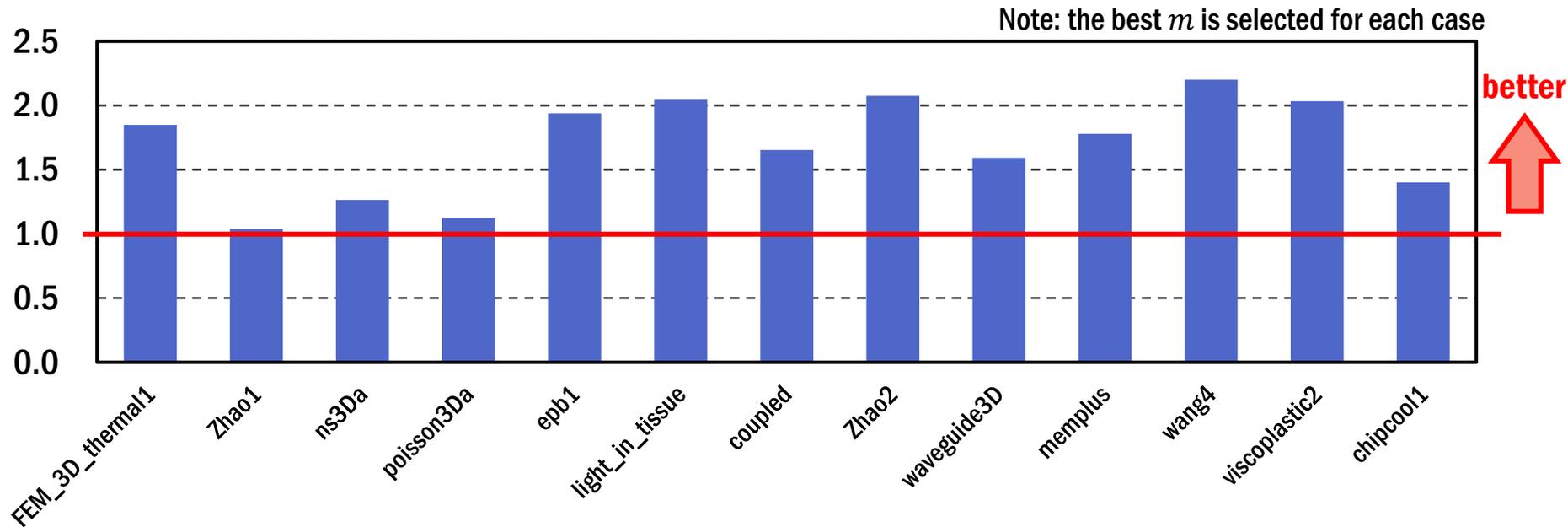GMRES(m) using only FP64     MP-GMRES(m) using FP32 & FP64

**If a problem can be solved by GMRES(m) using only FP64, the problem is expected to be solved also by MP-GMRES(m) using FP32 & FP64 (with the almost same # of iterations).**

# Evaluation on execution time

| | Both converged | Only FP64 converged | Both not converged |
|---|---|---|---|
| # of matrices | 13 | 1 | 4 |

## Speedup of MP-GMRES(m) using PF32 & FP64 over FP64 GMRES(m)
### (by a standard thread parallel implementation, on a system with 2x 20-core Skylake Xeon)



Note: the best $m$ is selected for each case

better

For details, please see our paper: Y. Zhao et al., Numerical Investigation into the Mixed Precision GMRES(m) Method Using FP64 and FP32, JIP, 30 (2022), 525-537 (Open access).

**Unpublished results**

# Conclusion

# Conclusion

## ◆ Summary

Through numerical experiments, we investigated possibilities of introducing low precision computing into the GMRES(m) method.

- The MP-GMRES(m) using FP32 and FP64 shows the similar convergence property as that of GMRES(m) using only FP64.

- There is a considerable possibility of introducing lower precision data than FP32 into the GMRES(m) method if a problem is not difficult.

- The impact of reducing the precision of $A$ and $V$ is different; more aggressive reduction for $A$ will be acceptable than for $V$.

## ◆ Future work

- Further numerical experiments

- Theoretical analysis

- Discussion on expected speedup (e.g., performance modeling)

- Study on the case of using preconditioners