





MUMPS workshop @ CINES

25 June 2025

Mixed precision algorithms in numerical linear algebra

Theo Mary Sorbonne Université, CNRS, LIP6

Floating-point landscape

	Signif. bits	Exp. bits	Range ($f_{\rm max}/f_{\rm min}$)	Unit roundoff <i>u</i>
fp128	113	15	$2^{32766}pprox 10^{9863}$	$2^{-114}\approx 1\times 10^{-34}$
fp64	52	11	$2^{2046}pprox 10^{616}$	$2^{-53}pprox 1 imes 10^{-16}$
fp32	23	8	$2^{254}pprox 10^{76}$	$2^{-24}pprox 6 imes 10^{-8}$
tfloat32	10	8	$2^{254}pprox 10^{76}$	$2^{-11}\approx 5\times 10^{-4}$
fp16	10	5	$2^{30}pprox 10^9$	$2^{-11}pprox 5 imes 10^{-4}$
bfloat16	7	8	$2^{254}pprox 10^{76}$	$2^{-8}pprox 4 imes 10^{-3}$
fp8 (E4M3)	3	4	$2^{15}pprox 3 imes 10^4$	$2^{-4}pprox 6 imes 10^{-2}$
fp8 (E5M2)	2	5	$2^{30}pprox 10^9$	$2^{-3}pprox 1 imes 10^{-1}$
fp6 (E2M3)	3	2	$2^3 \approx 8$	$2^{-4}pprox 6 imes 10^{-2}$
fp6 (E3M2)	2	3	$2^7 pprox 128$	$2^{-3}pprox 0.125$
fp4 (E2M1)	1	2	$2^3 \approx 8$	$2^{-2} pprox 0.25$

Lower precisions:

- Faster, consume less memory and energy
- Eower accuracy and narrower range
- \Rightarrow Mixed precision algorithms

Standard model of FPA:

$$\begin{array}{l} \text{For any } x \text{ such that } |x| \in [f_{\min}, f_{\max}], \\ \quad \mathsf{fl}(x) = x(1+\delta), \quad |\delta| \leq u \end{array}$$



Acta Numerica (2022), pp. 347-414 doi:10.1017/S0962492922000022

Mixed precision algorithms in numerical linear algebra

Nicholas J. Higham Department of Mathematics, University of Manchester, Manchester, M13 9PL, UK E-mail: nick.higham@manchester.ac.uk

> Theo Mary Sorbonne Université, CNRS, LIP6, Paris, F-75005, France E-mail: theo.mary@lip6.fr

https://bit.ly/mixed-survey



CONTENTS

1	Introduction	2
2	Floating-point arithmetics	6
3	Rounding error analysis model	14
4	Matrix multiplication	15
5	Nonlinear equations	18
6	Iterative refinement for $Ax = b$	22
7	Direct methods for $Ax = b$	25
8	Iterative methods for $Ax = b$	35
9	Mixed precision orthogonalization and QR factoriza-	
	tion	39
10	Least squares problems	42
11	Eigenvalue decomposition	43
12	Singular value decomposition	46
13	Multiword arithmetic	47
14	Adaptive precision algorithms	50
15	Miscellany	52

Iterative refinement

Run baseline algorithm in low precision, refine result to high accuracy

Ø Multiword arithmetic

Emulate high precision with low precision

Memory accessors

Decouple the storage (low) precision and the compute (high) precision

Adaptive precision

Adapt the precision of each instruction to the problem/input at hand

O Discussion

Comparison of strengths and weaknesses



2 Multiword arithmetic





Iterative refinement for Ax = b

- 1: Compute an initial approximation x
- 2: repeat
- 3: r = b Ax
- 4: Solve $Ac \approx r$
- 5: x = x + c
- 6: until convergence
- If $Ac \approx r$ is solved with precision ε :
 - $\bullet\,$ Attainable accuracy independent of ε
 - Convergence rate $\propto \kappa(A) \varepsilon$
- Can use direct, iterative, or any kind of approximate solvers
 E. Carson and N. J. Higham. Accelerating the solution of linear systems by iterative refinement in three precisions. SISC 2018.

Iterative refinement for Ax = b

- 1: Compute an initial approximation x
- 2: repeat
- 3: r = b Ax
- 4: Solve $Ac \approx r$
- 5: x = x + c
- 6: until convergence
- If $Ac \approx r$ is solved with precision ε :
 - $\bullet\,$ Attainable accuracy independent of $\varepsilon\,$
 - Convergence rate $\propto \kappa(A) \varepsilon$
- Can use direct, iterative, or any kind of approximate solvers
 E. Carson and N. J. Higham. Accelerating the solution of linear systems by iterative refinement in three precisions. SISC 2018.

General technique also applicable to

 Nonlinear solvers (Newton's method)
 F. Tisseur. Newton's method in floating point arithmetic and iterative refinement of generalized eigenvalue problems. SIMAX 2001.

Least-squares

E. Carson, N. J. Higham, and S. Pranesh. Three-Precision GMRES-Based Iterative Refinement for Least Squares Problems. SISC 2020.

• Eigenvalue decomposition

■ J. J. Dongarra, C. B. Moler, and J. H. Wilkinson. Improving the accuracy of computed eigenvalues and eigenvectors. SINUM 1983.

• Singular value decomposition

J. J. Dongarra. Improving the accuracy of computed singular values. SISSC 1983.

Low-rank approximations

M. Baboulin, O. Kaya, T. M., and M. Robeyns. Mixed precision iterative refinement for low-rank matrix and tensor approximations. SISC 2025.

Iterative refinement: success stories

LU-based IR

1: $A \approx LU$ in low precision 2: $x = U^{-1}L^{-1}b$ in low precision

3: repeat

- 4: r = b Ax in high precision
- 5: $c = U^{-1}L^{-1}r$ in low precision
- 6: x = x + c in high precision
- 7: **until** convergence

P. R. Amestoy, A. Buttari, N. J. Higham, J.-Y. L'Excellent, T. M., and B. Vieublé. Combining sparse approximate factorizations with mixed-precision iterative refinement. TOMS 2023.

Results with fp32 BLR($arepsilon$)-MUMPS solver $+$ fp64 refinement						
Matrix	n	$\kappa(A)$	ε	lts	Time	Memory
ElectroPhys	10M	10 ¹	10^{-6}	3	5.2×	3.7×
tminlet	3M	107	10^{-6}	7	4.2×	2.9 imes
CarBody	25M	10^{13}	F	F	F	F

Iterative refinement: success stories

GMRES-based IR

- 1: $A \approx LU$ in low precision
- 2: $x = U^{-1}L^{-1}b$ in low precision

3: repeat

- r = b Ax in high precision 4:
- Solve $Ac \approx r$ with LU-preconditioned 5: **GMRES** in mixed precision
- 6: x = x + c in high precision
- 7: **until** convergence

P. R. Amestov, A. Buttari, N. J. Higham. J.-Y. L'Excellent, T. M., and B. Vieublé, Combining sparse approximate factorizations with mixed-precision iterative refinement. TOMS 2023.

P. R. Amestov, A. Buttari, N. J. Higham, J.-Y. L'Excellent, T. M., and B. Vieublé, Five-precision GMRES-based iterative refinement SIMAX 2024

Results with fp32 BLR($arepsilon$)-MUMPS solver $+$ fp64 refinement						
Matrix	n	$\kappa(A)$	ε	lts	Time	Memory
ElectroPhys tminlet CarBody	10M 3M 25M	10 ¹ 10 ⁷ 10 ¹³	$\begin{array}{c} 10^{-6} \rightarrow \text{same} \\ 10^{-6} \rightarrow 10^{-4} \\ \text{F} \rightarrow 10^{-8} \end{array}$	$\begin{array}{c} 3 ightarrow { m same} \ 7 ightarrow { m 69} \ { m F} ightarrow { m 23} \end{array}$	$\begin{array}{l} 5.2\times\rightarrow \text{ same}\\ 4.2\times\rightarrow 3.3\times\\ F\rightarrow 0.7\times\end{array}$	$\begin{array}{l} 3.7\times\rightarrow\text{same}\\ 2.9\times\rightarrow3.4\times\\ \text{F}\rightarrow1.9\times\end{array}$



2 Multiword arithmetic





Goal: compute C = AB to high accuracy

• Step 1: Compute the multiword decompositions

$$A \approx \sum_{i=1}^{w} \alpha_i A_i$$
 and $B \approx \sum_{j=1}^{w} \beta_j B_j$

where the words A_i and B_j are stored in low precision \Rightarrow decomposition error

• Step 2: compute

$$C = \sum_{i+j \le k} \alpha_i \beta_j A_i B_j$$

for some k (e.g., k = 2w or k = w) \Rightarrow accumulation error

fp32 emulation

- NVIDIA Tensor Cores fp16 matrix multiplication with fp32 accumulation
- fp16-TC/fp32 speed ratio: Volta Ampere Hopper Blackwell $8 \times 16 \times 15 \times 56 \times$
- Compute $A \approx A_1 + A_2$ and $B \approx B_1 + B_2$ \Rightarrow decomposition error = 2⁻²²
- Compute $C \approx A_1B_1 + A_1B_2 + A_2B_1$ with fp32 accumulation \Rightarrow accumulation error $\propto 2^{-24}$
- Three products \Rightarrow large speedups
- bfloat16×6 or bfloat16×9 also possible. See
 M. Fasi, N. J. Higham, F. Lopez, T. M. and M. Mikaitis. Matrix Multiplication in Multiword Arithmetic: Error Analysis and Application to GPU Tensor Cores. SISC 2023.

and references therein

fp64 emulation

- fp16 and int8 words both possible, int8 usually more efficient
 H. Ootomo, K. Ozaki, and R. Yokota. DGEMM on integer matrix multiplication unit. IJHPCA 2024.
- Number of products quite large (${\sim}35$), but int8/fp64 speed ratio is $112{\times}$
- Number of products must depend on dynamical range of values; otherwise componentwise error can be large for badly scaled matrices.
 A. Abdelfattah, J. Dongarra, M. Fasi, M. Mikaitis, and F. Tisseur. Analysis of floating-point matrix multiplication computed via integer arithmetic. Preprint 2025.



2 Multiword arithmetic





Memory accessors: main principle

 Decouple storage and compute precisions: data is stored (compressed) in low precision and accessed (decompressed) back to high precision to be computed on/with
 Description of the precision of the preci

B H. Anzt, G. Flegar, T. Grützmacher, and E. S. Quintana-Ortí. Toward a modular precision ecosystem for high-performance computing. IJHPCA 2019.

- Lowering storage precision reduces memory consumption and volume of data transfers ⇒ faster memory-bound computations
 T. Grützmacher, H. Anzt, and E. S. Quintana-Ortí. Using Ginkgo's memory accessor for improving the accuracy of memory-bound low precision BLAS. Software: Practice and Experience 2023.
- Compression may use custom formats without hardware support, e.g., floating-point numbers with truncated mantissa
 D. Mukunoki, M. Kawai, and T. Imamura. Sparse Matrix-Vector Multiplication with Reduced-Precision Memory Accessor MCSoC 2023.

Memory accessors: success stories

- Tensor Cores implement mixed precision matrix multiply–accumulate
 C ← C + AB, where A and B are stored in, and C is accumulated in, fp32 or fp16
- Accumulating the update operations A_{ij} ← A_{ij} L_{ik} U_{kj} of LU factorization in fp32 reduces the error bound from nu₁₆ to 2u₁₆ + nu₃₂

 P. Blanchard, N. J. Higham, F. Lopez, T. M., and S. Pranesh. Mixed Precision Block Fused Multiply-Add: Error Analysis and Application to GPU Tensor Cores. SISC 2020.
- Significant accuracy boost, but performance limited by storage/data transfers
 ⇒ store matrix in fp16, and preserve accuracy by accumulating in fp32 buffers

$$B_{ij} = \sum_{k} L_{ik} U_{kj}, \quad A_{ij} \leftarrow A_{ij} - B_{ij}$$

F. Lopez and T. M.. Mixed Precision LU Factorization on GPU Tensor Cores: Reducing Data Movement and Memory Footprint. IJHPCA 2023.

	fp16	fp16/fp32 Tensor Cores	
		fp32 storage	fp16 storage
Accuracy	$1 imes 10^{-3}$	$1 imes 10^{-5}$	$3 imes 10^{-5}$
Speed (TFLOPS)	—	50	140

At what data granularity should we use memory accessors?

- Too small (e.g., variable-wise) \Rightarrow need to rewrite all the code $\ensuremath{\mathfrak{S}}$
- Too large (e.g., matrix-wise) ⇒ accessed data does not fit into fast memory, inefficient ☺
- Just right (e.g., block-wise) ⇒ blocks fit into fast memory and computations can use BLAS ! ☺

P. R. Amestoy, A. Jego, J.-Y. L'Excellent, T. M., and G. Pichon. BLAS-based Block Memory Accessor with Applications to Mixed Precision Sparse Direct Solvers. Preprint 2025.





2 Multiword arithmetic





Adaptive precision: main principle

• Not all variables/operations need the same precision! Example:



- \Rightarrow Here, *b* can be stored and computed in low precision
 - Adaptive precision algorithms exploit this observation by dynamically selecting the minimal precision for each variable/operation, depending on the data and on the prescribed accuracy ε

Adaptive precision: success stories

- Goal: compute y = Ax, where A is a sparse matrix, with a prescribed accuracy ε
- Given p available precisions $u_1 < \varepsilon < u_2 < \ldots < u_p$, define partition

$$A = \sum_{k=1}^{p} A^{(k)}, \qquad a_{ij}^{(k)} = \begin{cases} \mathsf{fl}_k(a_{ij}) & \text{if } |a_{ij}| \in (\varepsilon ||A|| / u_k, \varepsilon ||A|| / u_{k+1}] \\ 0 & \text{otherwise} \end{cases}$$

 \Rightarrow the precision of each element is chosen inversely proportional to its magnitude \blacksquare S. Graillat, F. Jézéquel, T. M., and R. Molina. Adaptive precision sparse matrix-vector product and its application to Krylov solvers. SISC 2024.

• Example on Long_Coup_dt6 matrix with $\varepsilon = 2^{-53}$ and 7 precisions:

fp64	fp56	fp48	fp40	fp32	fp24	fp16	drop
0.05%	2%	25%	25%	4%	20%	14%	10%

Performance impact: 1.5× storage reduction \Rightarrow 1.4× time reduction

Adaptive precision: success stories



P. R. Amestoy, O. Boiteau, A. Buttari, M. Gerest, F. Jézéquel, J.-Y. L'Excellent, and T. M.. Mixed Precision Low Rank Approximations and their Application to Block Low Rank LU Factorization. IMAJNA 2022.

A. Buttari, T. M., and A. Pacteau. Truncated QR factorization with pivoting in mixed precision. SISC 2025.

- Image of size 1057×1600 and of rank 191 (with $\varepsilon = 0.04$)
- In fp32/bf16, only 13 columns in fp32



 $\begin{array}{c} \text{precision } u_1 \\ \text{precision } u_2 \\ \text{precision } u_3 \end{array}$



 \Rightarrow the precision of each singular vector is chosen inversely proportional to its singular value







Adaptive precision: success stories

- Adastra MUMPS4FWI project led by WIND team
- Application: Gorgon Model, reservoir 23km × 11km × 6.5km, grid size 15m, Helmholtz equation, 25-Hz
- Complex matrix, 531 Million dofs, storage(A)=220 GBytes;
- FR cost: flops for one LU factorization = 2.6 × 10¹⁸; Estimated storage for LU factors = 73 TBytes



(25-Hz Gorgon FWI velocity model)

FR (Full-Rank); BLR with $\varepsilon = 10^{-5}$;					48 000	cores (5	00 MPI \times	96 threads/MPI)
FR: fp32; Adaptive precision BLR: 3 precisions (32bits, 24bits, 16bits) for storage) for storage
L	U size (TBytes)	F	lops	Time BL	R + Mix	ed (sec)	Scaled Resid.
FR	BLR	+ adapt.	FR	BLR+adapt.	Analysis	Facto	Solve	BLR+adapt.
73	34	26	$2.6 imes10^{18}$	$0.5 imes10^{18}$	446	5500	27	$7 imes 10^{-4}$

1 Iterative refinement

2 Multiword arithmetic





Discussion

- Accuracy: how far is the computed result from the exact one?
- Performance: how much does it cost to reach this result?
- Robustness: for what range of inputs is the method guaranteed to work?

	Iterative	Multiword	Memory	Adaptive
	Refinement	Arithmetic	Accessor	Precision
Accuracy	Rigorously	Rigorously	Rigorously	Rigorously
Accuracy	controlled	controlled	controlled	controlled
Porformanco	High if it	High on GPUs and	High if	Data-dependent
Ferformance	converges quickly	if compute-bound	memory-bound	Data-dependent
Robustness	$\kappa(A)$ not too large	Black box	Application- specific	Black box

Discussion

- Accuracy: how far is the computed result from the exact one?
- Performance: how much does it cost to reach this result?
- Robustness: for what range of inputs is the method guaranteed to work?

	Iterative	Multiword	Memory	Adaptive
	Refinement	Arithmetic	Accessor	Precision
Accuracy	Rigorously	Rigorously	Rigorously	Rigorously
Accuracy	controlled	controlled	controlled	controlled
Performance	High if it	High on GPUs and	High if	Data-dependent
Feriormance	converges quickly	if compute-bound	memory-bound	Data-dependent
Pobustnoss	r(A) not too large	Black box	Application-	Black box
NUDUSTILESS	R(A) not too large		specific	DIACK DOA

Thanks! Questions?