Mixed precision algorithms: an overview

Theo Mary Sorbonne Université, CNRS, LIP6

NHR PerfLab Seminar Series, 12 December 2023

Slides available at

https://bit.ly/NHRmixed





Floating-point arithmetics

	number of bits				
		signif. (t)	exp.	range	$u = 2^{-t}$
fp128	quadruple	113	15	$10^{\pm 4932}$	$1 imes 10^{-34}$
fp64	double	53	11	$10^{\pm 308}$	$1 imes 10^{-16}$
fp32	single	24	8	$10^{\pm 38}$	$6 imes 10^{-8}$
fp16	bolf	11	5	$10^{\pm 5}$	$5 imes 10^{-4}$
bfloat16	патт	8	8	$10^{\pm 38}$	$4 imes 10^{-3}$
fp8 (e4m3)	an e ret e re	4	4	$10^{\pm 2}$	$6 imes 10^{-2}$
fp8 (e5m2)	quarter	3	5	$10^{\pm 5}$	$1 imes 10^{-1}$
fp16 bfloat16 fp8 (e4m3) fp8 (e5m2)	half quarter	11 8 4 3	5 8 4 5	$egin{array}{c} 10^{\pm 5} \ 10^{\pm 38} \ 10^{\pm 2} \ 10^{\pm 5} \end{array}$	5×10^{-4} 4×10^{-3} 6×10^{-2} 1×10^{-1}

The unit roundoff $u = 2^{-t}$ determines the relative accuracy any number in the representable range can be approximated with:

If
$$x \in \mathbb{R}$$
 belongs to $[e_{\min}, e_{\max}]$, then $f(x) = x(1 + \delta), |\delta| \le u$

Moreover the standard model of arithmetic is

 $fl(x \text{ op } y) = (x \text{ op } y)(1 + \delta), \quad |\delta| \le u$, for $op \in \{+, -, \times, \div\}$

- Storage, data movement and communications are all proportional to total number of bits (mantissa + exponent)
 lower precision ⇒ lighter computations
- ③ Speed of computations also generally proportional
 - $\,\circ\,$ on most architectures, fp32 is 2× faster than fp64
 - $\,\circ\,$ on some architectures, fp16/bfloat16 up to $16\times$ faster than fp32

lower precision \Rightarrow faster computations

- \bigcirc Power consumption is proportional to the square of the number of mantissa bits
 - $\circ\,$ fp16 (11 bits) consumes 5× less energy than fp32 (24 bits)
 - $\,\circ\,$ bfloat16 (8 bits) consumes $9\times$ less energy than fp32

lower precision \Rightarrow greener computations

Errors are proportional to the unit roundoff lower precision => lower accuracy Mix several precisions in the same code with the goal of

- Getting the performance benefits of low precisions
- While preserving the accuracy and stability of the high precision

Terminology varies: Mixed precision, Multiprecision, Adaptive precision, Variable precision, Transprecision, Dynamic precision, ...

Mix several precisions in the same code with the goal of

- Getting the performance benefits of low precisions
- While preserving the accuracy and stability of the high precision

Terminology varies: Mixed precision, Multiprecision, Adaptive precision, Variable precision, Transprecision, Dynamic precision, ...

How to select the right precision for the right variable/operation

- Precision tuning: autotuning based on the source code
 - Does not need any understanding of what the code does
 - 😟 Does not have any understanding of what the code does
- In linear algebra: exploit as much as possible the knowledge we have about the code

Survey

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

Solution of Ax = b:

• Direct methods

- Robust, black box solvers
- $\circ~$ High time and memory cost for factorization of A

• Iterative methods

- $\circ~$ Low time and memory per-iteration cost
- $\circ~$ Convergence is application dependent

Solution of Ax = b:

• Direct methods

- Robust, black box solvers
- $\circ~$ High time and memory cost for factorization of A
- \Rightarrow Need fast factorization

• Iterative methods

- $\circ~$ Low time and memory per-iteration cost
- Convergence is application dependent
- \Rightarrow Need good preconditioner

Solution of Ax = b:

• Direct methods

- Robust, black box solvers
- High time and memory cost for factorization of A
- \Rightarrow Need fast factorization

• Iterative methods

- $\circ~$ Low time and memory per-iteration cost
- Convergence is application dependent
- \Rightarrow Need good preconditioner

\Rightarrow Mixed precision / approximate factorizations bridge the gap

- $\circ~$ as approximate fast direct methods
- $\circ~$ as high quality preconditioners

Standard method to solve Ax = b:

- 1. Factorize A = LU, where L and U are lower and upper triangular
- 2. Solve Ly = b and Ux = y

In uniform precision \boldsymbol{u} , the computed $\widehat{\boldsymbol{x}}$ satisfies

- Backward error $\frac{\|A\widehat{x}-b\|}{\|A\|\|\widehat{x}\|+\|b\|} \leq f(n)\rho_n u = O(u)$
- Forward error $\frac{\|\widehat{\kappa}-x\|}{\|x\|} \leq f(n)\rho_n\kappa(A)u = O(\kappa(A)u)$, with $\kappa(A) = \|A\| \|A^{-1}\|$

LU-based refinement (LU-IR)

```
Factorize A = LU
Solve Ax_1 = b via x_1 = U^{-1}(L^{-1}b)
repeat
r_i = b - Ax_i
Solve Ad_i = r_i via d_i = U^{-1}(L^{-1}r_i)
x_{i+1} = x_i + d_i
until converged
```

LU-based refinement (LU-IR)

Factorize A = LU in precision **u** Solve $Ax_1 = b$ via $x_1 = U^{-1}(L^{-1}b)$ in precision **u** repeat $r_i = b - Ax_i$ in precision \mathbf{u}^2 Solve $Ad_i = r_i$ via $d_i = U^{-1}(L^{-1}r_i)$ in precision **u** $x_{i+1} = x_i + d_i$ in precision **u until** converged

🖹 Wilkinson (1948) 🛛 🖹 Moler (1967)

Assuming $\kappa(A)\mathbf{u} < 1$:

- Backward error $\frac{\|A\widehat{\mathbf{x}}-b\|}{\|A\|\|\widehat{\mathbf{x}}\|+\|b\|} = O(\mathbf{u})$
- Forward error $\frac{\|\widehat{x}-x\|}{\|x\|} = O(\mathbf{u})$

LU-IR with fp32 LU

```
Factorize A = LU in precision \mathbf{u}_{\mathbf{f}}
Solve Ax_1 = b via x_1 = U^{-1}(L^{-1}b) in precision \mathbf{u}_{\mathbf{f}}
repeat
r_i = b - Ax_i in precision \mathbf{u}
Solve Ad_i = r_i via d_i = U^{-1}(L^{-1}r_i) in precision \mathbf{u}_{\mathbf{f}}
x_{i+1} = x_i + d_i in precision \mathbf{u}
until converged
with \mathbf{u}_{\mathbf{f}} \equiv \text{fp32} and \mathbf{u} \equiv \text{fp64}
```

```
🖹 Langou et al (2006) 📑 Buttari et al (2007) 📑 Baboulin et al (2009)
```

LU-IR with fp32 LU

```
Factorize A = LU in precision \mathbf{u}_{\mathbf{f}}
Solve Ax_1 = b via x_1 = U^{-1}(L^{-1}b) in precision \mathbf{u}_{\mathbf{f}}
repeat
r_i = b - Ax_i in precision \mathbf{u}
Solve Ad_i = r_i via d_i = U^{-1}(L^{-1}r_i) in precision \mathbf{u}_{\mathbf{f}}
x_{i+1} = x_i + d_i in precision \mathbf{u}
until converged
```

with $\mathbf{u_f} \equiv \text{fp32}$ and $\mathbf{u} \equiv \text{fp64}$

Langou et al (2006) E Buttari et al (2007) Baboulin et al (2009)

- For *n* × *n* matrices:
 - $O(n^3)$ flops in fp32
 - $O(n^2)$ flops per iteration in fp64
- Assuming $\kappa(A)\mathbf{u_f} < 1$:
 - Backward error $\frac{\|A\hat{x}-b\|}{\|A\|\|\hat{x}\|+\|b\|} = O(\mathbf{u})$ • Forward error $\frac{\|\hat{x}-x\|}{\|x\|} = O(\kappa(A)\mathbf{u})$

LU-IR with CELL processor







CELL processor (2006–2008) fp64 peak: 21 GFLOPS fp32 peak: 205 GFLOPS $\Rightarrow 10 \times$ speedup!



NVIDIA GPUs



NVIDIA Hopper (H100) GPU

Peak performance (TFLOPS)							
	P100		A100	H100			
	2016	2018	2020	2022			
fp64	5	8	$10 \rightarrow 20$	$33 \rightarrow 67$			
fp32	10	16	20	67			
tfloat32			160	495			
fp16	20	125	$40 \rightarrow 320$	$134 \rightarrow 990$			
bfloat16			$40 \rightarrow 320$	$134 \rightarrow 990$			
fp8				1979			
with tensor cores							

 $^{11/46}\text{Since}$ A100, 16-bit arithmetic is 16× faster than 32-bit

NVIDIA GPU tensor cores

Tensor cores units available on NVIDIA GPUs carry out a fixed size (e.g., 4 \times 4) matrix multiplication :



• Performance boost vs fp32: 8–16× speedup vs fp32

NVIDIA GPU tensor cores

Tensor cores units available on NVIDIA GPUs carry out a fixed size (e.g., 4×4) matrix multiplication :



- Performance boost vs fp32: 8–16× speedup vs fp32
- Accuracy boost vs fp16: let C = AB, with $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times p}$, the computed \widehat{C} satisfies $\begin{pmatrix} nu_{16} & (\text{fp16}) \end{pmatrix}$

$$|\widehat{C} - C| \lesssim c_n |A| |B|, \quad c_n = \begin{cases} 2u_{16} + nu_{32} & \text{(tensor cores)} \\ nu_{32} & \text{(fp32)} \end{cases}$$

12/46 📑 Blanchard, Higham, Lopez, M., Pranesh (2020)

Block version to use matrix-matrix operations

```
for k = 1: n/b do
    Factorize L_{kk}U_{kk} = A_{kk} (with unblocked alg.)
    for i = k + 1: n/b do
         Solve L_{ik} U_{kk} = A_{ik} and L_{kk} U_{ki} = A_{ki} for L_{ik} and U_{ki}
    end for
    for i = k + 1: n/b do
         for i = k + 1: n/b do
             A_{ii} \leftarrow A_{ii} - \widetilde{L}_{ik}\widetilde{U}_{ki}
         end for
    end for
end for
```

Block LU factorization with tensor cores

- Block version to use matrix-matrix operations
- $O(n^3)$ part of the flops done with tensor cores

```
for k = 1: n/b do
     Factorize L_{kk}U_{kk} = A_{kk} (with unblocked alg.)
    for i = k + 1: n/b do
          Solve L_{ik} U_{kk} = A_{ik} and L_{kk} U_{ki} = A_{ki} for L_{ik} and U_{ki}
     end for
     for i = k + 1: n/b do
         for j = k + 1: n/b do
               \widetilde{L}_{ik} \leftarrow \mathsf{fl}_{16}(L_{ik}) and \widetilde{U}_{ki} \leftarrow \mathsf{fl}_{16}(U_{ki})
              A_{ii} \leftarrow A_{ii} - L_{ik} U_{ki} using tensor cores
          end for
     end for
end for
```

LU factorization with tensor cores

Error analysis for LU follows from matrix multiplication analysis and gives same bounds to first order 🖹 Blanchard et al. (2020)



Impact on iterative refinement





- TC accuracy boost can be critical!
- TC performance suboptimal here \Rightarrow why?

LU factorization is memory bound

- LU factorization is traditionally a compute-bound operation...
- $\bullet\,$ With Tensor Cores, flops are 8–16 $\times\,$ faster
- Matrix is stored in fp32 \Rightarrow data movement is unchanged
- \Rightarrow LU with tensor cores becomes memory-bound !





LU factorization is memory bound

- LU factorization is traditionally a compute-bound operation...
- With Tensor Cores, flops are 8–16 \times faster
- Matrix is stored in fp32 \Rightarrow data movement is unchanged
- \Rightarrow LU with tensor cores becomes memory-bound !



- Idea: store matrix in fp16
- ^{16/46} Problem: huge accuracy loss, tensor cores accuracy boost completely negated

Two ingredients to reduce data movement with no accuracy loss:

Reducing data movement

Two ingredients to reduce data movement with no accuracy loss:

1. Mixed fp16/fp32 representation



Matrix after 2 steps:

Reducing data movement

Two ingredients to reduce data movement with no accuracy loss:

1. Mixed fp16/fp32 representation



Matrix after 2 steps:

Reducing data movement

Two ingredients to reduce data movement with no accuracy loss:

- 1. Mixed fp16/fp32 representation
- 2. Right-looking \rightarrow left-looking factorization



Matrix after 2 steps:

 $O(n^3)$ fp32 + $O(n^2)$ fp16 $\rightarrow O(n^2)$ fp32 + $O(n^3)$ fp16

Experimental results



- Nearly 50 TFLOPS without significantly impacting accuracy
 Lopez and M. (2020)
- Even more critical on A100: 50 TFLOPS (A in fp32) \rightarrow **175 TFLOPS** (A in fp16+left-looking)

Use of fp16 presents two risks:

- Overflow/underflow in the LU factors
 - $|||L||U||| \le f(n)\rho_n ||A|| \Rightarrow$ even if A fits in the range, its LU factors may not
 - Image: Higham, Pranesh, Zounon (2019) : two-sided diagonal scaling $A' \leftarrow D_r A D_c$ so that $\|A\| \le c$
 - $\circ~$ To minimize underflow and better utilize the range of fp16, helpful to take c as close as possible to maximum safe value
 - Zounon et al. (2020) : appearance of subnormal numbers (in fp32) can lead to slowdowns if they are not flushed to zero
- Loss of positive definiteness
 - $\circ~$ Rounding a posdef A to fp16 might make it indefinite \Rightarrow Cholesky factorization breaks down
 - ▷ 🖹 Higham & Pranesh (2021) : factorize $A + \sigma D$ instead $(D = \text{diag}(A), \sigma = O(u_{16}))$

Three-precision LU-IR

```
Factorize A = LU in precision \mathbf{u}_{\mathbf{f}}
Solve Ax_1 = b via x_1 = U^{-1}(L^{-1}b) in precision \mathbf{u}_{\mathbf{f}}
repeat
r_i = b - Ax_i in precision \mathbf{u}_r
Solve Ad_i = r_i via d_i = U^{-1}(L^{-1}r_i) in precision \mathbf{u}_{\mathbf{f}}
x_{i+1} = x_i + d_i in precision \mathbf{u}
until converged
e.g., with \mathbf{u}_{\mathbf{f}} \equiv \text{fp16}, \mathbf{u} \equiv \text{fp32}, and \mathbf{u}_r \equiv \text{fp64}
```

Carson and Higham (2018)

```
Factorize A = LU in precision \mathbf{u}_{\mathbf{f}}
Solve Ax_1 = b via x_1 = U^{-1}(L^{-1}b) in precision \mathbf{u}_{\mathbf{f}}
repeat
r_i = b - Ax_i in precision \mathbf{u}_r
Solve Ad_i = r_i via d_i = U^{-1}(L^{-1}r_i) in precision \mathbf{u}_{\mathbf{f}}
x_{i+1} = x_i + d_i in precision \mathbf{u}
until converged
e.g., with \mathbf{u}_{\mathbf{f}} \equiv \mathbf{fp}\mathbf{16}, \mathbf{u} \equiv \mathbf{fp}\mathbf{32}, and \mathbf{u}_r \equiv \mathbf{fp}\mathbf{64}
```

Carson and Higham (2018)

Assuming $\kappa(A)\mathbf{u_f} < 1$:

• Backward error
$$rac{\|A\widehat{\mathbf{x}}-b\|}{\|A\|\|\widehat{\mathbf{x}}\|+\|b\|}=O(\mathbf{u})$$

• Forward error
$$\frac{\|\widehat{\mathbf{x}}-\mathbf{x}\|}{\|\mathbf{x}\|} = O(\mathbf{u} + \kappa(A)\mathbf{u}_{\mathbf{r}})$$

Three-precision LU-IR is as general (as modular) as possible $_{20/46}$

GMRES-based IR (GMRES-IR)

Factorize A = LU in precision \mathbf{u}_{f} Solve $Ax_{1} = b$ via $x_{1} = U^{-1}(L^{-1}b)$ in precision \mathbf{u}_{f} repeat $r_{i} = b - Ax_{i}$ in precision \mathbf{u}_{r} Solve $U^{-1}L^{-1}Ad_{i} = U^{-1}L^{-1}r_{i}$ with GMRES in precision \mathbf{u} with products with $U^{-1}L^{-1}A$ in precision \mathbf{u}^{2} $x_{i+1} = x_{i} + d_{i}$ in precision \mathbf{u} until converged

Carson and Higham (2017)

• Replace LU solver by preconditioned GMRES:

- $\,\circ\,$ GMRES can be asked to converge to accuracy $u \ll u_f$
- $\kappa(\widetilde{A})$ often smaller than $\kappa(A)$
- $\Rightarrow \widetilde{A}d_i = \widetilde{r_i}$ is solved with accuracy $\phi_i = \kappa(\widetilde{A})\mathbf{u}$
 - Convergence condition improved from $\kappa(A)\mathbf{u_f} < 1$ to $\kappa(\widetilde{A})\mathbf{u} < 1$

• **The catch**: the matrix-vector products are with $\tilde{A} = U^{-1}L^{-1}A$, introduce an extra

 $_{21/46}$ $\kappa(A)$ unless performed in higher precision

	u _f	u	u _r	$\max \kappa(A)$	Forward error
LU-IR	fp32	fp64	fp128	10 ⁸	10^{-16}
GMRES-IR	fp32	fp64	fp128	10 ¹⁶	10^{-16}
LU-IR	fp16	fp64	fp128	10 ³	10^{-16}
GMRES-IR	fp16	fp64	fp128	10^{11}	10^{-16}

GMRES-IR can handle much more ill-conditioned matrices.

	u _f	u	u _r	$\max \ \kappa(A)$	Forward error
LU-IR	fp32	fp64	fp128	10 ⁸	10^{-16}
GMRES-IR	fp32	fp64	fp128	10^{16}	10^{-16}
LU-IR	fp16	fp64	fp128	10 ³	10^{-16}
GMRES-IR	fp16	fp64	fp128	10^{11}	10^{-16}

GMRES-IR can handle much more ill-conditioned matrices. However: LU solves are performed in precision u^2 instead of $u_f \Rightarrow$ practical limitation

Rethinking GMRES-IR

- Goal: solve $Ad_i = r_i$ with GMRES and bound $\phi_i = \|\widehat{d}_i d_i\|/\|d_i\|$
 - In what precision do we really need to run GMRES?
 - · How much extra precision is really needed in the matvec products?

```
Solve Ax_1 = b by LU factorization in precision \mathbf{u}_{\mathbf{f}}

repeat

r_i = b - Ax_i in precision \mathbf{u}_{\mathbf{r}}

Solve U^{-1}L^{-1}Ad_i = U^{-1}L^{-1}r_i with GMRES in precision \mathbf{u}

except products with U^{-1}L^{-1}A in precision \mathbf{u}^2

x_{i+1} = x_i + d_i in precision \mathbf{u}

until converged
```
Rethinking GMRES-IR

- Goal: solve $Ad_i = r_i$ with GMRES and bound $\phi_i = \|\widehat{d}_i d_i\|/\|d_i\|$
 - In what precision do we really need to run GMRES?
 - · How much extra precision is really needed in the matvec products?

```
Solve Ax_1 = b by LU factorization in precision \mathbf{u}_{\mathbf{f}}

repeat

r_i = b - Ax_i in precision \mathbf{u}_{\mathbf{r}}

Solve U^{-1}L^{-1}Ad_i = U^{-1}L^{-1}r_i with GMRES in precision \mathbf{u}

except products with U^{-1}L^{-1}A in precision \mathbf{u}^2

x_{i+1} = x_i + d_i in precision \mathbf{u}

until converged
```

Rethinking GMRES-IR

- Goal: solve $Ad_i = r_i$ with GMRES and bound $\phi_i = \|\widehat{d}_i d_i\|/\|d_i\|$
 - In what precision do we really need to run GMRES?
 - · How much extra precision is really needed in the matvec products?

```
Solve Ax_1 = b by LU factorization in precision \mathbf{u}_{\mathbf{f}}

repeat

r_i = b - Ax_i in precision \mathbf{u}_{\mathbf{r}}

Solve U^{-1}L^{-1}Ad_i = U^{-1}L^{-1}r_i with GMRES in precision \mathbf{u}_{\mathbf{g}}

except products with U^{-1}L^{-1}A in precision \mathbf{u}_{\mathbf{p}}

x_{i+1} = x_i + d_i in precision \mathbf{u}

until converged
```

Relax the requirements on the GMRES precisions: run at precision $u_g \leq u$ with matvecs at precision $u_p \leq u^2$

 \Rightarrow **FIVE precisions** in total!

23/46 What can we say about the convergence of this GMRES-IR5?

- **Unpreconditioned** GMRES in precision u_g for Ax = b:
 - Backward error of order ug 📄 Paige, Rozloznik, Strakos (2006)
 - Forward error of order $\kappa(A)u_g$
- **Two-precision preconditioned** GMRES for $\widetilde{A}x = b$:
 - Backward error of order $\kappa(A)\mathbf{u_p} + \mathbf{u_g}$
 - Forward error of order $\kappa(\widetilde{A})(\kappa(A)\mathbf{u}_{p} + \mathbf{u}_{g})$
 - $\circ \ \kappa(\widetilde{A}) \leq (1+\kappa(A)\mathsf{u}_{\mathsf{f}})^2$

Side-result: generalization of the backward stability of GMRES to a preconditioned two-precision GMRES

Amestoy, Buttari, Higham, L'Excellent, M., Vieublé (2021)

Five-precision GMRES-IR

Solve $Ax_1 = b$ by LU factorization in precision $\mathbf{u}_{\mathbf{f}}$ **repeat** $r_i = b - Ax_i$ in precision $\mathbf{u}_{\mathbf{r}}$ Solve $U^{-1}L^{-1}Ad_i = U^{-1}L^{-1}r_i$ with GMRES in precision $\mathbf{u}_{\mathbf{g}}$ except products with $U^{-1}L^{-1}A$ in precision $\mathbf{u}_{\mathbf{p}}$ $x_{i+1} = x_i + d_i$ in precision \mathbf{u} **until** converged

Theorem (convergence of GMRES-IR5)

Under the condition $(\mathbf{u}_{g} + \kappa(A)\mathbf{u}_{p})\kappa(A)^{2}\mathbf{u}_{f}^{2} < 1$, the forward error converges to its limiting accuracy $\frac{\|\widehat{x} - x\|}{\|x\|} \leq \mathbf{u}_{r}\kappa(A) + \mathbf{u}$

Amestoy, Buttari, Higham, L'Excellent, M., Vieublé (2021)

With five arithmetics (fp16, bfloat16, fp32, fp64, fp128) there are over **3000 different** combinations of GMRES-IR5!

They are not all relevant !

Meaningful combinations: those where none of the precisions can be lowered without worsening either the limiting accuracy or the convergence condition.

Filtering rules	
• $u^2 \le u_r \le u \le u_f$	• $\mathbf{u_p} < \mathbf{u}, \ \mathbf{u_p} = \mathbf{u}, \ \mathbf{u_p} > \mathbf{u}$ all possible
• $u_p \leq u_g$	• $u_g \ge u$
• $u_p < u_f$	• $u_g < u_f$, $u_g = u_f$, $u_g > u_f$ all possible

ug	u _p	Convergence Condition $\max(\kappa(A))$
LU-I	R	$2 imes 10^3$
bfloat16	fp32	$3 imes 10^4$
fp16	fp32	$4 imes 10^4$
fp16	fp64	$9 imes 10^4$
fp32	fp64	$8 imes 10^6$
fp64	fp64	$3 imes 10^7$
fp64	fp128	$2 imes 10^{11}$

Meaningful combinations of GMRES-IR5 for $u_f \equiv$ fp16 and $u \equiv$ fp64

Five combinations between LU-IR and Carson & Higham's GMRES-IR \Rightarrow More flexible precisions choice to fit at best the hardware constraints and the problem difficulty.

- What about solving $AU^{-1}L^{-1}y = b$ with flexible GMRES ?
- Stability not provable as easily/unconditionally as in the left-preconditioned case
 Arioli and Duff (2009) Carson and Daužickaitė (2023)
 but works very well in practice
- $\Rightarrow\,$ Can relax the need to apply preconditioner in higher precision, at the cost of using flexible GMRES

Sparse matrices

Fill-in: $nnz(A) \ll nnz(LU)$



- Memory footprint of high-precision copy of A negligible ⇒ modern IR saves memory (not the case for dense systems!)
- \bigcirc Matvecs with A cheap compared with LU solves \Rightarrow can afford to compute residual very accurately!
- Example 2 Relative weight of refinement is higher: $O(n^2)$ vs $O(n^3)$ for dense $\Rightarrow O(n^{4/3})$ vs $O(n^2)$ for sparse (at best, even worse for 2D problems) \Rightarrow less room to amortize iterations

Higher weight of symbolic operations

Sparse LU-IR vs GMRES-IR

Results with the MUMPS solver

Amestoy, Buttari, Higham, L'Excellent, M., Vieublé (2023)

Matrix	time (s)			memory (GB)			
	fp64	fp32	fp32	fp64	fp32	fp32	
		+LU-IR	+GMRES-IR		+LU-IR	+GMRES-IR	
ElectroPhys10M	265.2	154.0	166.5	272.0	138.0	171.3	
Bump_2911	205.4	129.3	144.5	135.7	68.4	77.8	
DrivAer6M	91.8	67.6	77.9	81.6	41.7	52.9	
Queen_4147	284.2	165.2	184.7	178.0	89.8	114.5	
tminlet3M	294.5	136.2	157.9	241.1	121.0	169.9	
perf009ar	46.1	57.5	52.0	55.6	28.9	38.1	
elasticity-3d	156.7	—	118.6	153.0	—	103.6	
lfm_aug5M	536.2	254.5	269.3	312.0	157.0	187.5	
Long_Coup_dt0	67.2	46.6	49.0	52.9	26.7	33.1	
CarBody25M	62.9	_	109.8	77.6		54.3	
thmgaz	97.6	65.4	79.8	192.0	97.7	141.7	

- Up to $2 \times$ time and memory reduction, even for ill-conditioned problems
- 30/46 GMRES-IR usually more expensive than LU-IR, but more robust

```
Compute M^{-1} \approx A^{-1} and x_1 = M^{-1}b in precision \mathbf{u}_{\mathbf{f}}

repeat

r_i = b - Ax_i in precision \mathbf{u}_{\mathbf{r}}

Solve M^{-1}Ad_i = M^{-1}r_i with GMRES in precision \mathbf{u}_{\mathbf{g}}

except products with M^{-1}A in precision \mathbf{u}_{\mathbf{p}}

x_{i+1} = x_i + d_i in precision \mathbf{u}

until converged
```

- Replace $U^{-1}L^{-1}$ with general M^{-1}
- Equivalent to restarted GMRES

Cheaper preconditioners than LU

```
Initialize x_1

repeat

r_i = b - Ax_i in precision \mathbf{u}_r \equiv \mathbf{u}_{high}

Solve Ad_i = r_i with GMRES in precision \mathbf{u}_g \equiv \mathbf{u}_{low}

x_{i+1} = x_i + d_i in precision \mathbf{u} \equiv \mathbf{u}_{high}

until converged
```

- Replace $U^{-1}L^{-1}$ with general M^{-1}
- Equivalent to restarted GMRES
- No preconditioner (M = I): mixed precision inner-outer scheme
 Turner and Walker (1992)
 Buttari et al. (2008)
 Lindquist et al. (2020)
 Loe et al. (2021)

Sparsification (dropping)

Dropping: replace with zero any value sufficiently small

$$|a_{ij}| \leq \epsilon ||A|| \quad \Rightarrow \quad a_{ij} \leftarrow 0$$



Sparsification (dropping)

Dropping: replace with zero any value sufficiently small

$$|a_{ij}| \le \epsilon \|A\| \Rightarrow a_{ij} \leftarrow 0$$



32/46

Data sparsification (low-rank approximations)

Low-rank compression: given $A = U\Sigma V^T$, if we truncate singular vectors associated with $\sigma_i \leq \epsilon$, we obtain \widetilde{A} such that $\|\widetilde{A} - A\| \leq \epsilon$



Compress A_{ij} such that $\|\widetilde{A}_{ij} - A_{ij}\| \le \epsilon \|A\|$:

• If
$$||A_{ij}|| \le \epsilon ||A|| \Rightarrow A_{ij} \leftarrow 0$$
 (drop block)

$$_{/46}$$
 • otherwise replace A_{ij} with $\widetilde{A}_{ij}=X_{ij}Y_{ij}^{T}$

33



Block Low Rank

Example on tminlet3M matrix fp64 MUMPS reference: time \rightarrow 295.5 memory \rightarrow 241.1

	time (s)		memory (GB)	
	LU-IR	GMRES-IR	LU-IR	GMRES-IR
fp32 MUMPS	136.2	157.9	121.0	169.9

Example on tminlet3M matrix fp64 MUMPS reference: time \rightarrow 295.5 memory \rightarrow 241.1

	time (s)		merr	nory (GB)
	LU-IR GMRES-IR		LU-IR	GMRES-IR
fp32 MUMPS	136.2	157.9	121.0	169.9
$\epsilon = 10^{-8}$	149.7	165.3	114.0	161.9

Example on tminlet3M matrix fp64 MUMPS reference: time \rightarrow 295.5 memory \rightarrow 241.1

	time (s)		memory (GB)	
	LU-IR GMRES-IR		LU-IR	GMRES-IR
fp32 MUMPS	136.2	157.9	121.0	169.9
$\epsilon = 10^{-8}$	149.7	165.3	114.0	161.9
$\epsilon = 10^{-6}$	88.3	98.8	82.4	93.8

Example on tminlet3M matrix fp64 MUMPS reference: time \rightarrow 295.5 memory \rightarrow 241.1

	time (s)		memory (GB)	
	LU-IR GMRES-IR		LU-IR	GMRES-IR
fp32 MUMPS	136.2	157.9	121.0	169.9
$\epsilon = 10^{-8}$	149.7	165.3	114.0	161.9
$\epsilon = 10^{-6}$	88.3	98.8	82.4	93.8
$\epsilon = 10^{-4}$	_	105.6	_	70.9

• GMRES-IR allows to push BLR further!

- Sparsification only deals in absolutes: either we keep the data at full accuracy, or we discard it completely!
- We need a new paradigm that uses multiple, gradual levels of approximation
- \Rightarrow Adaptive precision sparsification



Adapt the precisions to the data at hand by storing and computing "less important" (usually meaning smaller) data in lower precision

Adaptive precision SpMV

- **Goal:** compute the SpMV y = Ax with accuracy ϵ using q precisions $u_1 \le \epsilon < u_2 < \ldots < u_q$
- Split elements a_{ij} on each row *i* into *q* buckets B_{i1}, \ldots, B_{iq} , where bucket B_{ik} uses precision u_k

Adaptive precision SpMV

- **Goal:** compute the SpMV y = Ax with accuracy ϵ using q precisions $u_1 \le \epsilon < u_2 < \ldots < u_q$
- Split elements a_{ij} on each row *i* into *q* buckets B_{i1}, \ldots, B_{iq} , where bucket B_{ik} uses precision u_k
- How should we build the buckets?

$$\begin{cases} |a_{ij}| \leq \epsilon ||A|| & \Rightarrow \text{ drop} \\ |a_{ij}| \in [\epsilon||A||/u_{k+1}, \epsilon ||A||/u_k) & \Rightarrow \text{ place in } B_{ik} \\ |a_{ij}| > \epsilon ||A||/u_2 & \Rightarrow \text{ place in } B_{i1} \end{cases}$$



Adaptive precision SpMV

36/46

- **Goal:** compute the SpMV y = Ax with accuracy ϵ using q precisions $u_1 \le \epsilon < u_2 < \ldots < u_q$
- Split elements a_{ij} on each row *i* into *q* buckets B_{i1}, \ldots, B_{iq} , where bucket B_{ik} uses precision u_k
- How should we build the buckets?

$$\begin{cases} |a_{ij}| \leq \epsilon ||A|| & \Rightarrow \text{ drop} \\ |a_{ij}| \in [\epsilon||A||/u_{k+1}, \epsilon||A||/u_k) & \Rightarrow \text{ place in } B_{ik} \\ |a_{ij}| > \epsilon ||A||/u_2 & \Rightarrow \text{ place in } B_{i1} \end{cases}$$



• **Theorem**: the computed \hat{y} satisfies $\|\hat{y} - y\| \le c\epsilon \|A\| \|x\|$ Graillat, Jézéquel, M., Molina (2022)

Adaptive precision SpMV: results

- Experimental results on 36-core computer
- Adaptive algorithm uses fp64, fp32, and dropping and $\epsilon = 2^{-53}$
- Comparison vs uniform fp64

Matrix	adaptive storage	storage adaptive time		backward error	
	(% of fp64)	(% of fp64)	fp64	adaptive	
Hook_1498	97%	99%	3e-16	7e-16	
Flan_1565	78%	82%	3e-16	3e-16	
Long_Coup_dt0	75%	83%	5e-16	2e-15	
imagesensor	15%	56%	2e-16	9e-16	
power9	16%	37%	1e-16	4e-16	
nv2	19%	26%	2e-16	2e-15	

Seven-precision SpMV

Emulated formats						
	Bits					
Format	Signif.(t)	Exponent	Range	$u = 2^{-t}$		
bf16	8	8	10 ^{±38}	$4 imes 10^{-3}$		
fp24	16	8	$10^{\pm 38}$	$2 imes 10^{-5}$		
fp32	24	8	$10^{\pm 38}$	$6 imes 10^{-8}$		
fp40	29	11	$10^{\pm 308}$	$2 imes 10^{-9}$		
fp48	37	11	$10^{\pm 308}$	$8 imes 10^{-12}$		
fp56	45	11	$10^{\pm 308}$	$3 imes 10^{-14}$		
fp64	53	11	$10^{\pm308}$	$1 imes 10^{-16}$		



Graillat, Jézéquel, M., Molina, Mukunoki (2023)

GMRES-based iterative refinement

GMRES $r = b - Ax_0$ $\beta = \|\boldsymbol{r}\|_2$ $q_1 = r/\beta$ for k = 1, 2, ... do $y = Aq_k$ for i = 1: k do $h_{ik} = q_i^T y$ $y = y - h_{ik}q_i$ end for $h_{k+1,k} = \|y\|_2$ $q_{k+1} = y/h_{k+1,k}$ Solve min_{c_k} $||Hc_k - \beta e_1||_2$. $x_{k} = x_{0} + Q_{k}c_{k}$ end for

GMRES-IR for i = 1, 2, ... do $r_i = b - Ax_{i-1}$ Solve $Ad_i = r_i$ by GMRES $x_i = x_{i-1} + d_i$ end for

GMRES-based iterative refinement

GMRES $r = b - Ax_0$ $\beta = \|\boldsymbol{r}\|_2$ $q_1 = r/\beta$ for k = 1, 2, ... do $y = Aq_k \rightarrow \epsilon_{low}$ for i = 1: k do $h_{ik} = q_i^T y$ $y = y - h_{ik}q_i$ end for $h_{k+1,k} = \|y\|_2$ $q_{k+1} = y/h_{k+1,k}$ Solve min_{c_k} $||Hc_k - \beta e_1||_2$. $x_{k} = x_{0} + Q_{k}c_{k}$ end for

GMRES-IRfor i = 1, 2, ... do $r_i = b - Ax_{i-1} \rightarrow \epsilon_{high}$ Solve $Ad_i = r_i$ by GMRES $x_i = x_{i-1} + d_i$ end for

GMRES-IR with adaptive precision SpMV



GMRES-IR with adaptive precision SpMV



Adaptive precision low rank compression



How to increase low-rank compression?

Adaptive precision low rank compression



How to increase low-rank compression?

• Standard approach: increase ϵ to discard more vectors

Adaptive precision low rank compression



How to increase low-rank compression?

- Standard approach: increase ϵ to discard more vectors
- Adaptive precision compression: partition U and V into q groups of decreasing precisions u₁ ≤ ε < u₂ < ... < u_q
- Why does it work? $B = B_1 + B_2 + B_3$ with $|B_i| \le O(||\Sigma_i||)$ Amestoy, Boiteau, Buttari, Gerest, Jézéquel, L'Excellent, M. (2021)

Adaptive precision BLR LU factorization

Stability of LU factorization: $\widehat{L}\widehat{U}=A+\Delta A$

- Standard LU (Wilkinson) : $\|\Delta A\| \lesssim 3n^3 \rho_n u_1 \|A\|$
- **BLR LU** (Higham & M.) : $\|\Delta A\| \lesssim (c_1 \epsilon + c_2 \rho_n u_1) \|A\|$ \blacksquare Higham and M. (2021)
- Adaptive prec. BLR LU (this work) : $\|\Delta A\| \lesssim (c'_1 \epsilon + c'_2 \rho_n u_1) \|A\|$

Example of kernel: LR \times matrix multiplication:



Step k:

- Compute $L_{kk}U_{kk} = A_{kk}$
- Update $A_{ij} \leftarrow A_{ij} - (A_{ik}U_{kk}^{-1}) \times (L_{kk}^{-1}A_{kj})$



Adaptive precision BLR implementation in MUMPS

1) Toward storage gains

- A large number of precisions (current version: **7 formats**)
- For storage only



- A preliminary version has been completed

- 2) Toward time gains
 - Small number of precisions
 - Chosen according to availability in hardware
 - For computations



- Ongoing development

Matrix		LU factors size	Total memory	Backward error
	fp64 MUMPS	141	194	
thmgaz	BLR double	95	120	6.4e-14
-	BLR mixed	59	86	6.5e-14
	fp64 MUMPS	235	547	
knuckle8M	BLR double	117	281	1.6e-10
	BLR mixed	71	236	7.7e-09
	fp64 MUMPS	38	58	
perf009ar	BLR double	26	36	1.3e-10
	BLR mixed	20	25	1.4e-10

Memory consumption reduced by up to

- 1.7× (LU factors size)
- $1.4 \times$ (total memory, including working arrays)

A plausible scenario for solving Ax = b in mixed precision:

```
Compute adaptive precision BLR LU factorization A \approx LU at accuracy \epsilon_{\rm f}
Solve Ax_1 = b with adaptive precision BLR LU solves at accuracy \epsilon_f
repeat
   Compute r_i = b - Ax_i with adaptive precision SpMV at accuracy \epsilon_r
   Solve Ad_i = r_i with adaptive precision GMRES preconditioned by BLR LU factors, using
           adaptive precision SpMV at accuracy \epsilon_a,
           adaptive precision BLR LU solves at accuracy \epsilon_{\rm n}.
          and precision ug for the rest
   x_{i+1} = x_i + d_i (in uniform precision u!)
until converged
```



Modularity



Slides https://bit.ly/NHRmixed



Survey https://bit.ly/mixed-survey



Thanks! Questions?