中国科学院大学 夏季强化课程 20222

Fast Solvers for Large Algebraic Systems

Lecture 5. Mixed-precision methods

Chensong Zhang, AMSS

http://lsec.cc.ac.cn/~zhangcs



Table of Contents



- Lecture 1: Large-scale numerical simulation
- Lecture 2: Fast solvers for sparse linear systems
- Lecture 3: Methods for non-symmetric problems
- Lecture 4: Methods for nonlinear problems
- Lecture 5: Mixed-precision methods
- Lecture 6: Communication hiding and avoiding
- Lecture 7: Fault resilience and reliability
- Lecture 8: Robustness and adaptivity

Sources of Error in Simulation







Accuracy of Direct Solvers, Revisited

• Floating-point computation (IEEE compatible)

$$f(x \text{ op } y) = (x \text{ op } y)(1+\delta), \quad |\delta| \le u$$

• Forward error of matrix-multiplications (depends on num of inner multiplications)

Compute $C = AB, A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{n \times k} \Longrightarrow |\hat{C} - C| \le \gamma_n |A| |B| \qquad \gamma_n := \frac{1}{1 + 1}$

Backward perturbation analysis of the LU decomposition

$$(A + \delta A)\hat{x} = b, \quad |\delta A| \le \gamma_{3n} |\hat{L}| |\hat{U}|$$

Accuracy of direct solution methods

$$\frac{\|x - \hat{x}\|_{\infty}}{\|x\|_{\infty}} \le \gamma_{3n} \frac{\left\| |\hat{L}| |\hat{U}| \right\|_{\infty}}{\|A\|_{\infty}} \kappa_{\infty}(A)$$



N. Higham, 2002



- Multigrid methods: Using direct method as coarse grid solver
- Domain decomposition methods: Using direct method on local subdomains and "direct" preconditioner on interfaces
- Partial factorization: Incomplete factorizations as preconditioners (like ILU methods)
- Direct method used in part of a preconditioner (like MSP for the well part)
- Block iterative methods: Using direct method on sub-blocks (like Schwarz preconditioners)
- Low-precision direct method as a preconditioner (like GMRES-IR)
- Iterative refinement: improve precision of direct solvers (like LU-IR)
- Factorization of nearby problem as a preconditioner

Iain Duff. The SIAM Conference on Applied Linear Algebra. October 26-29, 2009. Monterey, California.

Some Historical Comments



- Long long ago ... computing machines are very different from now
 - Fixed-point arithmetic → Floating-point: SP, DP, QP, ...
 - Floating-point arithmetic: Different standards for different machines
 - IEEE 754 Standard, Intel 8087 coprocessor, ...
- Multi-precision arithmetic simulators available for a long time
 - Maple, Mathematica, ...
 - Julia, half, numpy, ...
- Solving an algebraic system to certain accuracy on a computer:
 - If the problem is ill-conditioned (sensitive to perturbation) ...; see Lecture 2
 - If the precision is low, quality of solution could be very bad
 - The way to improve quality is using higher precision (too costly) or IR

Why Mixed-Precision Now



	NVIDIA A100	NVIDIA Turing	NVIDIA Volta
Supported Tensor Core Precisions	FP64, TF32, bfloat16, FP16, INT8, INT4, INT1	FP16, INT8, INT4, INT1	FP16
Supported CUDA® Core Precisions	FP64, FP32, FP16, bfloat16, INT8	FP64, FP32, FP16, INT8	FP64, FP32, FP16, INT8

- *"Because it's there..."* George Mallory, 1923
- Fast hardware with mixed-precision capability is availab

- Make less data movement (see Lecture 6)
- Fill more data in the cache
- Traditional algorithms need to be modified
- Traditional algorithms must take advantage
- The main driving force is DL and it will continue

NVIDIA Volta		V100 PCle	V100 SXM2	V100S PCle				
ED14	GPU Architecture	NVIDIA Volta						
FFIO	NVIDIA Tensor Cores		640					
	NVIDIA CUDA [®] Cores		5,120					
4, FP32, FP16, IN18	Double-Precision Performance	7 TFLOPS	7.8 TFLOPS	8.2 TFLOPS				
	Single-Precision Performanc e	14 TFLOPS	15.7 TFLOPS	16.4 TFLOPS				
	Tensor Performance	112 TFLOPS	125 TFLOPS	130 TFLOPS				
available	GPU Memory	32 GB /10	GB HBM2	32 GB HBM2				
	Memory Bandwidth	900 0		1134 GB/sec				
	ECC	but not 2X faster						
Tensor Cores	Interconnect Bandwidth	32 GB/sec	300 GB/sec	32 GB/sec				
introduced	System Interface	PCle Gen3	NVIDIA NVLink™	PCle Gen3				
in Volta, 4X	Form Factor	PCIe Full Height/Length	SXM2	PCIe Full Height/Length				
regular FP <u>16</u>	Max Power Comsumption	250 W	300 W	250 W				
	Thermal Solution	Passive						
	Compute APIs	CUDA, DirectCompute, OpenCL™, OpenACC®						

Driving Force: Deep Learning

- They are very tolerant of reduced-precision computations
- The computations performed by most models are compositions of a relatively small handful of dense linear algebra calculations . . .
- Many of the mechanisms developed over the past 40 years to enable general-purpose programs to run with high performance on modern CPUs . . . are unnecessary for machine learning computations

Design." 20	020 IEEE International Solid-State Circuits Conference (ISSCC), 2020, pp. 8-14.	
Ref: J. Dea	in, "The Deep Learning Revolution and Its Implications for Computer Architecture a	nd Chip

8

NVIDIA A100 TE Unprecedented Accele	INSOF	R CORE GF Every Scale	U	100	
D	=	С	+	A	В
$[\times \times \times \times]$	[×	(] [$ \times \times \times \times] $	$[\times \times \times \times]$
$ \times \times \times \times $	_ ×	$\times \times \times \times$		$\times \times \times \times$	$ \times \times \times \times $
$\times \times \times \times$	- ×	$\times \times \times \times$		$\times \times \times \times$	$ \times \times \times \times $
$\left[\times\times\times\times\right]$	Ĺ×	$\times \times \times \times_{}$		$[\times \times \times \times]$	$[\times \times \times \times]$
fp16 or fp32	fp	16 or fp32		fp16	fp16

Fused multiply-add (FMA) used in Tensor Core



NVIDIA

Range of FP Representation



	signif.	exp.	u	$x_{\min}^{(s)}$	X min	X _{max}
bfloat16	8	8	3.91×10^{-3}	$9.18 imes 10^{-41}$	1.18×10^{-38}	3.39 × 10 ³⁸
fp16	11	5	4.88×10^{-4}	5.96 × 10 ⁻⁸	6.10×10^{-5}	$6.55 imes 10^4$
fp32	24	8	5.96×10^{-8}	$1.40 imes 10^{-45}$	1.18×10^{-38}	3.40 × 10 ³⁸
fp64	53	11	1.11×10^{-16}	$4.94 imes 10^{-324}$	2.22×10^{-308}	1.80×10^{308}
fp128	113	15	9.63 × 10 ⁻³⁵	$6.48 imes 10^{-4966}$	$3.36 imes 10^{-4932}$	1.19 × 10 ⁴⁹³²

Number of bits + 1

- For fp16, it is very likely to overflow / underflow; for bfloat16, low precision
- We can ignore the small numbers (underflow)
- But the real trouble is: how to handle the large numbers?
- Need to scale the problem first

Taking A Detour



- Mixed-precision or multi-precision: More accurate result -> Better performance
- How to measure performance? We can compare wall time (actual computing time)!
- But we might solve different problems with different solvers on different computers ...



Source: https://www.geeksforgeeks.org/computer-organization-von-neumann-architecture/



Memory Hierarchy





Memory Hierarchy



A Simplified Cache Mechanism





C.-S. Zhang, AMSS



- Program usually access a relatively small portion of the memory
 - **Temporal Locality:** If an item is referenced, it might be referenced again soon
 - Spatial Locality: If an item is referenced, "nearby" items tend to be referenced
- Misses at start up: cold start
- Hit time: 1 4 cycles ← → Miss penalty: 8 32 cycles (extra work)





Hennessy, John L., and David A. Patterson. Computer architecture: a quantitative approach. Elsevier, 2011

A Sample CPU Specific





2019 Production

Benchmarks https://askgeek.io/en/cpus/Intel/Core-i9-9980HK



Intel[®] Core i9-9980HK

- Total num of cores = 8, num of threads = 16
- Maximum CPU clock speed = 5.00 GHz
- Intel Smart Cache:
 - ✓ L1 = 512 KB (cache line: 64 Bytes)
 - ✓ L2 = 2 MB
 - ✓ L3 = 16 MB
- Maximum memory size = 128 GB
- Power consumption (TDP) = 45 Watt
- Maximum operating temperature = 100 °C
- Manufacturing process technology = 14 nm

Memory Address Spaces





Communication, Revisited



Communication → Data Movement

- Physics: Longer the distance, higher the cost
- Technology: Increase speed with different methods
- Economy: Manufacturing cost is also important
- Combined: Computing ability (FLOPS) improves faster than communication ability

	Today's Systems	Predicted Exascale Systems*	Factor Improvement
System Peak	10 ¹⁶ flops/s	10 ¹⁸ flops/s	100
Node Memory Bandwidth	10 ² GB/s	10 ³ GB/s	10
Interconnect Bandwidth	10 ¹ GB/s	10 ² GB/s	10
Memory Latency	$10^{-7} { m s}$	$5\cdot 10^{-8}$ s	2
Interconnect Latency	10 ⁻⁶ s	$5\cdot 10^{-7}$ s	2
	()		(

*Sources: from P. Beckman (ANL), J. Shalf (LBL), and D. Unat (LBL)

Annual improvements									
Time/flop		Bandwidth	Latency						
59%	Network	26%	15%						
	DRAM	23%	5%						

Machine Balance Over Time





The gap between the computing power and the memory bandwidth keeps increasing in the last 45 years

- Computing power is (relatively) cheap
- Data movement is slow
- Data movement will be slower
- Take-home message:

Make less data movement!

Source: A. Abdelfattah et al. "A survey of numerical linear algebra methods utilizing mixed-precision arithmetic". The International Journal of High Performance Computing Applications. 2021;35(4):344-369.

Measuring Performance





Source: Introduction to the Roofline Model, by Samuel Williams, Lawrence Berkeley National Lab

CPU-Bound vs Memory-Bound





Source: Samuel Williams, LBNL https://crd.lbl.gov/divisions/amcr/computerscience-amcr/par/research/roofline/introduction/ a[i] = b[i] + scalar * c[i]: 2 FLOP, 2 Read, 1 Write, CI = 2/(3*8) = 0.083 7-point stencil: 7 Flop, 1 Read (cache-hit), 1 Write, CI= 7/(2*8) = 0.438

A Sample TRIAD Code



11_			
10	Shot (Purchase unlo	nck	to remove watermark) ports for "/mnt/a/u/staff/arnoldg/stream/triad.c"
11	ionor (Porchase unic	JUK	to remove watermarky
//IN	LINING OPTION VALUE	s:	
11	-inline-factor: 100		
11	-inline-min-size: 3	0	
11	-inline-max-size: 2	20	
10	deline new total a		. 2000
11	-inline-max-total-s	1 Ze	2000
11	-inline-max-per-rou	tin	2: 16698
11	-inline-max-per-com	p11	e: 580000
11			
1 :	# include <stdio.h></stdio.h>		
2 :	# include <unistd.h< td=""><td>></td><td></td></unistd.h<>	>	
3 :	# include <math.h></math.h>		
4	# include cfloat.h>		
5	# include climits.h		
6	t include csys/time		
7	 dofino STREAM A 	ap.	V 517E 10000000
	define STREAM	MAR	
•	# define STREAM_T	TPE	00015
9	# define OFFSET 0		
10			
11	static STREAM_TYPE		a[STREAM_ARRAY_SIZE+OFFSET],
12			b[STREAM_ARRAY_SIZE+OFFSET],
13			c[STREAM_ARRAY_SIZE+OFFSET];
14			
15	void tuned STREAM T	ria	d(double scalar)
16	(
//TN	LINE REPORT: (tuned	ST	REAM Triad(double)) [1] /mmt/a/u/staff/armolds/stream/triad.c(16.1)
11	same merontri (cuneu	-	real in sources to the function of sources and an end of the real
111-	ak (a lu lakatt lar	-	Annual Asiad a /16 1) annual \$24051, BESISTED ALL MATTON . [Aunual CTREAM Taine] (and (a
///m	nc/a/u/starr/arnoid	g/s	creamycriau.clio,ij.remark wowo51: KEOISIEK ALLOCATION : [tuned_SIKEAM_Irlad]/mmt/a/u/starf/armoldg/stream/triad.clib
11			
11	Hardware register	s	
11	Reserved	:	2[rsp rip]
11	Available	:	39[rax rdx rcx rbx rbp rsi rdi r8-r15 mm0-mm7 zmm0-zmm15]
11	Callee-save	:	6[rbx rbp r12-r15]
11	Assigned	:	12[rax rdx rcx rbx rsi rdi r8-r10 zmm0-zmm2]
11			
11	Routine temperart	es	
11	Total	:	192
11	Global		10
11	loco?		87
11	Local	1	65
11	Regenerable	:	35
11	Spilled	:	1
11			
11	Routine stack		
11	Variables	:	48 bytes*
11	Reads	:	4 [8.080+00 ~ 8.0%]
11	Writes	:	6 [5.08e+00 ~ 0.0%]
11	Spills		48 bytes*
11	Reads		13 [5.699-81 ~ 8.8%]
11	Weiter		11 [1 10440] - 0 00]
	WILLES	•	TE FETODATE - 01001
11			
11	NOTES		
11			
11	*Non-overlapp	ing	variables and spills may share stack space,
11	so the total	. st	ack size might be less than this.
11			
11			
17	ssize_t i:		
18	#pragma omp paralle	1 f	10
//00	enMP Construct at /	mnt	/a/u/staff/accolde/stream/triad.c(18.1)
11			
ine	OD REGIN at /met/-		taff/annolde/stoarm/triad s/10 1)
11.0	or ordin at /mit/a/	3/5	Let i / et troategy als tenty (1 and 10 (and 1)
// <p< td=""><td>eeted 100p for vect</td><td>001</td><td>Lacron></td></p<>	eeted 100p for vect	001	Lacron>
//LO	OF END		
11			
//LO	OP BEGIN at /mnt/a/	u/s	taff/arnoldg/stream/triad.c(18,1)
11	remark #15300: LOO	PW	AS VECTORIZED
//10	OP END		
11			
1110	OP REGIN at /met/a/	u/e	taff/arnolde/stream/triad.c(18.1)
11/1	amainden loon for	art	anizations
// <r< td=""><td>emainder toop for w</td><td>ect</td><td>N. TSACTON A</td></r<>	emainder toop for w	ect	N. TSACTON A
//10	UP END		
19	for (j=0; j	<st< td=""><td>(LAPI_ANNAY_SIZE;]++)</td></st<>	(LAPI_ANNAY_SIZE;]++)
28	a[j] =	b[j]+scalar*c[j];
21	}		
22			

```
18 #pragma omp parallel for
//OpenMP Construct at /mnt/a/u/staff/arnoldg/stream/triad.c(18,1)
//remark #16200: OpenMP DEFINED LOOP WAS PARALLELIZED
11
//LOOP BEGIN at /mnt/a/u/staff/arnoldg/stream/triad.c(18,1)
//<Peeled loop for vectorization>
//LOOP END
11
//LOOP BEGIN at /mnt/a/u/staff/arnoldg/stream/triad.c(18,1)
11
    remark #15300: LOOP WAS VECTORIZED
//LOOP END
11
//LOOP BEGIN at /mnt/a/u/staff/arnoldg/stream/triad.c(18,1)
//<Remainder loop for vectorization>
//LOOP END
            for (j=0; j<STREAM ARRAY SIZE; j++)</pre>
19
                a[i] = b[i]+scalar*c[i];
20
21 }
```

Source: https://wiki.ncsa.illinois.edu/display/AAG/Copy+of+triad.c+and+vector+optimization

Roofline Model Papers





performance of software and hardware.

Source: https://www.connectedpapers.com/



Mixed-Precision Algorithms



• Research on mixed-precision algorithms nowadays:

Try to accelerate the applications by using lower-precision formats while maintaining high accuracy results

- Main goals:
 - Efficient memory access (also better cache performance)
 - Less data movement (also less energy consumption)
 - Shift gears: Optimize number of flops in good old days

→ → → Fast Computation + Fast Communication

• Your generation will face "slower" communication (relative to computation)!

Floating-Point Number Storage Format



IEEE single-precision floating-point number

```
31
                  23
                         22
       30
                                                            0
  Χ
       XXXXXXXX
                                      significand
 sign
         exponent
                                                  With leading
Aorma.
                Number = (-1)^{s} 2^{e-127} (1 \cdot f)
                                                  1, normalized
 where s = 0 for positive numbers, 1 for negative numbers,
          e = exponent (between 0 and 255), and
                         f = \text{mantissa.}
```

- Fixed-size FP format
- Representation range
- Significand / Mantissa
- Unit round-off / precision
- Machine epsilon
- Little-endian or big-endian
- IEEE 754-1985 (fp32, fp64)
- IEEE 754-2008 (fp16, fp128)

Ref: Higham, Nicholas J. and Mary, Theo. "Mixed Precision Algorithms in Numerical Linear Algebra", Acta Numerica, to appear

IEEE754 Single-Precision Numbers



	Exponent	Decimal Value	Shift-127 Value
Reserved INF/NAN -	→ 1111,1111	255	128 —
1.0 * 2 0111,1111	→ 0111,1111	127	0
Reserved for	0000,0001	1	-126
0.0 (all bits	→ 0000,0000	0	???
are 0)			
	Exponent	Mantissa	Purpose
1.f * 2 0000,0000 ?	0	zero	0.0
No! Need special —	→ 0	nonzero	subnormal
treatment	1-254	anything	normal
Max number	255	zero	INF
1.11 * 2 1111,1110	255	nonzero	NaN

IEEE754 2008 Version



Table 1. Parameters for floating-point arithmetics: number of bits in significand (mantissa), including the implicit most significant bit; number of bits in exponent; and, to three significant figures, unit roundoff *u*, smallest positive (subnormal) number $x_{min}^{(s)}$, smallest normalized positive number x_{min} and largest finite number x_{max} . In Intel's bfloat16 specification, subnormal numbers are not supported [10], so any number less than x_{min} in magnitude is flushed to zero; the value shown holds if subnormal numbers are supported.

	signif.	exp.	и	$x_{\min}^{(s)}$	X_{min}	X _{max}
bfloat16	8	8	$3.91 imes 10^{-3}$	$9.18 imes 10^{-41}$	1.18×10^{-38}	3.39×10^{38}
fp16	11	5	4.88×10^{-4}	5.96 × 10 ⁻⁸	6.10 × 10 ⁻⁵	6.55 × 10 ⁴
fp32	24	8	5.96 × 10 ⁻⁸	$1.40 imes 10^{-45}$	1.18×10^{-38}	3.40×10^{38}
fp64	53	11	1.11×10^{-16}	$4.94 imes 10^{-324}$	2.22×10^{-308}	1.80×10^{308}
fp128	113	15	9.63 × 10 ⁻³⁵	$6.48 imes 10^{-4966}$	3.36×10^{-4932}	1.19 × 10 ⁴⁹³²

Source: Dongarra J, Grigori L, Higham NJ. "Numerical algorithms for high-performance computational science". Phil. Trans. R. Soc. A 378: 20190066, 2020



$(-1)^{s} \times 2^{0-127+1} \times 0.f$

- Subnormal numbers have all zero bits for the exponent
- Numbers that are smaller than the smallest normal positive number (consider positive only)
- Subnormal numbers are useful in filling gaps of floating point scale near zero
- Shift by -126 instead of -127! Why +1? Remember the exponent of the smallest normal number?
- Cannot add leading ONE any more: less bit → less accurate
- The largest possible subnormal number is $0.111...11 \times 2^{-126} < 1.0...0 \times 2^{-126}$
- The smallest possible subnormal number is $0.000...01 \times 2^{-126} \approx 1.40E-45$
- Floating-point operations in the subnormal region can be very very slow
- Flush subnormals (denormals) to zero using compiler arguments

Simple Examples of FP Arithmetic

Example using decimal

A = 9.999 × 10⁻¹, B = 1.610 × 10⁻¹, A+B =?

Step 1. Align the smaller exponent with the larger one.

 $B = 0.0161 \times 10^1 = 0.016 \times 10^1$ (round off)

Step 2. Add significands

9.999 + 0.016 = 10.015, so $A+B = 10.015 \times 10^{10}$

Step 3. Normalize

 $A+B = 1.0015 \times 10^2$

Step 4. Round off

 $A+B = 1.002 \times 10^2$

Example using decimal

 $A = 1.110 \times 10^{10}$, $B = 9.200 \times 10^{-5}$ $A \times B = ?$

Step 1. Exponent of $A \times B = 10 + (-5) = 5$ **Step 2**. Multiply significands 1.110× 9.200 = 10.212000 **Step 3**. Normalize the product 10.212 × 10⁵ = 1.0212 × 10⁶ **Step 4**. Round off $A \times B = 1.021 \times 10^{6}$ **Step 5**. Decide the sign of $A \times B (+ x + = +)$

So, A x B = + 1.021 x 10⁶

Source: https://homepage.divms.uiowa.edu/~ghosh/

Unit round-off



Working With Mixed Precisions



N. Higham, 2002

• Compute matrix-vector multiplications in finite precision (forward error):

 $fl(Ax) = Ax + \delta, \quad |\delta| \le \gamma_n |A| |x|$

• Compute residual in finite precision (forward error):

$$fl(b - Ax) = b - Ax + \delta, \quad |\delta| \le \gamma_{n+1} \left(|b| + |A| |x| \right)$$

- Pessimistic bounds! Improvement using blocking, probabilistic error bounds, ...
- How about working on two-precision computation?

First compute residual in precision \overline{u} and then store it in precision u

$$fl(b - Ax) = b - Ax + \delta, \quad |\delta| \le u|b - Ax| + (1 + u(\overline{\gamma}_{n+1})|b| + |A||x|)$$

Accuracy of Direct Solvers, Revisited

• Floating-point computation (IEEE compatible)

$$f(x \text{ op } y) = (x \text{ op } y)(1+\delta), \quad |\delta| \le u$$

Forward error of matrix-multiplications (depends on num of inner multiplications)

Compute
$$C = AB, A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{n \times k} \Longrightarrow |\hat{C} - C| \le \gamma_n |A| |B| \qquad \gamma_n := \frac{nu}{1 - nu}$$

Backward perturbation analysis of the LU decomposition

$$(A + \delta A)\hat{x} = b, \quad |\delta A| \le \gamma_{3n} |\hat{L}| |\hat{U}|$$

Accuracy of direct solution methods

$$\frac{\|x-\hat{x}\|_{\infty}}{\|x\|_{\infty}} \leq \gamma_{3n} \frac{\left\||\hat{L}||\hat{U}|\right\|_{\infty}}{\|A\|_{\infty}} \kappa_{\infty}(A)$$



N. Higham, 2002

Iterative Refinement Methods



- Sometimes, the condition number of a linear system could be very large
- Direct solvers might not be able to give a solution as accurate as we wanted (or thought)
- This has been noticed long ago (Wilkinson 1948), but using high-precision everywhere is costly
- We can improve precision iteratively using the iterative refinement

Algorithm 9: General iterative refinement algorithm

```
1%% Given a nonsingular matrix A \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n, and an initial guess x_0 \in \mathbb{R}^n;2%% Set three precisions u_r \leq u \leq u_\ell for residual, work, and solve, respectively;3for i = 0: MaxIter or converged4Compute r_i \leftarrow b - Ax_i in precision u_r;5Solve the error equation Ae_i = r_i in precision u_\ell;6Update x_{i+1} \leftarrow x_i + e_i in precision u;7end
```

Higham, N. and Mary, T. " Mixed Precision Algorithms in Numerical Linear Algebra", Acta Numerica

Accuracy of IR Methods







Algorithm 10: LU-IR method

- 1 %% Given a nonsingular matrix $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^n$;
- 2 %% Set three precisions $u_r \leqslant u \leqslant u_\ell$ for residual, work, and solve, respectively;
- 3 Compute the factorization A = LU in precision u_ℓ ;
- 4 Solve $LUx_0 = b$ by substitution in precision u_ℓ ;
- 5 for i = 0 : MaxIter or converged
- 6 Compute $r_i \leftarrow b Ax_i$ in precision u_r ;
- 7 Solve $LUe_i = r_i$ by substitution in precision u_ℓ ;

```
8 Update x_{i+1} \leftarrow x_i + e_i, in precision u;
```

- 9 end
- LU factorization only need to be computed once for each coefficient matrix
- The Carson-Higham theorem can be applied to this method
- Less stable but faster factorization methods can be used to improve efficiency (e.g. in SuperLU and PARDISO) [Li and Demmel 1998; Arioli, Duff, Gratton, and Pralet 2007]

GMRES-IR



Algorithm 11: Preconditioned GMRES-IR method

- 1 %% Given a nonsingular matrix $A \in \mathbb{R}^{n \times n}$ and $b \in \mathbb{R}^{n}$; 2 %% Set five precisions u_r, u_g, u_p, u, u_ℓ for different purposes; 3 Compute the factorization A = LU in precision u_ℓ ;
- 4 Solve $LUx_0 = b$ by substitution in precision u_ℓ ;
- 5 for i = 0 : MaxIter or converged

6 Compute
$$r_i \leftarrow b - Ax_i$$
 in precision u_r ;

7 Solve
$$Ae_i = r_i$$
 by PGMRES in precision u_g ;

8 In PGMRES, preconditioner $U^{-1}L^{-1}$ applied in precision u_p ;

```
9 Update x_{i+1} \leftarrow x_i + e_i in precision u;
```

10 end

If $\kappa(A)^2 u_\ell^2 (u_g + \kappa(A)u_p)$ is small, GMRES-IR will converge till

$$\frac{\|\hat{x} - x\|_{\infty}}{\|x\|_{\infty}} \lesssim u + 4p \operatorname{cond}(A, x)u_r$$

Accuracy of IR Methods





- GMRES-IR shows advantages compared with LU-IR in practice
- GMRES-IR is used in HPL-AI
- One can also use GMRES for

the original equation instead of

the error equation

• One can use low precision for

GMRES iterations

• Carson and Higham, SISC, 2017;

Haidar et al. 2020; ...

Different Versions of IR



- Traditional IR methods for ancient hardware (some computers can compute residual using higher precision than working precision without extra work); since 1970s these methods die due to such hardware does not exist any more.
- Fixed-precision IR methods (improve stability)
- IR methods with low-precision solvers: emphasize on efficiency and energy consumption

Advances in Mixed Precision Algorithms: 2021 Edition

by the ECP Multiprecision Effort Team (Lead: Hartwig Anzt)

Ahmad Abdelfattah, Hartwig Anzt, Alan Ayala, Erik G. Boman, Erin Carson, Sebastien Cayrols, Terry Cojean, Jack Dongarra, Rob Falgout, Mark Gates, Thomas Grützmacher, Nicholas J. Higham, Scott E. Kruger, Sherry Li, Neil Lindquist, Yang Liu, Jennifer Loe, Piotr Luszczek, Pratik Nayak, Daniel Osei-Kuffuor, Sri Pranesh, Sivasankaran Rajamanickam, Tobias Ribizel, Barry Smith, Kasia Swirydowicz, Stephen Thomas, Stanimire Tomov, Yaohung M. Tsai, Ichi Yamazaki, Urike Meier Yang

Mixed precision algorithms in numerical linear algebra

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

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

August 21, 2021

Mixed-Precision Algorithms

A few selected topics on mixed-precision numerical linear algebra

Blocked Algorithms



 $\gamma_n = nu + O(u^2)$

• Simple summation (forward error analysis)

$$s = x^T y = \sum_{i=1}^n x_i y_i \quad \Longrightarrow \quad |s - \hat{s}| \leq n u |x|^T |y|$$

Algorithm 12: Blocked summation method



Further Improvement on Summation



- So the basic idea is to reduce the number of additions in the summation
- A more advanced algorithm (FABsum) can give error bound independent of *n* to first order
- Using hardware: Fused Multiply-Add (FMA)

No FMA:
$$fl(x + yz) = (x + yz(1 + \delta_1))(1 + \delta_2), \quad |\delta_1|, |\delta_2| \le u$$

With FMA: $fl(x+yz) = (x+yz)(1+\delta), \quad |\delta| \le u$

Not only faster in computation

But also more accurate (by a factor of 2)!

Blanchard, Pierre, Nicholas John Higham and Théo Mary. "A Class of Fast and Accurate Summation Algorithms." SIAM J. Sci. Comput. 42 (2020): A1541-A1557.

Considerations for Sparse Problems



- Much harder to predict savings of MP iterative methods for sparse matrices
- LU factorization of a sparse matrix
 - Reordering and analysis phase of the algorithm does not involve floating-point arithmetic and so does not benefit from reducing the precision
 - LU can generate cascading fill-ins (small multipliers combine to produce subnormal numbers)
- Even if matrix entries can be scaled from double to normalized single precision numbers, subnormal numbers may still be generated during the factorization
- Floating-point operations on subnormal numbers can be very slow → This can cause significant performance loss

Ref: Zounon M, Higham NJ, Lucas C, Tisseur F. "Performance impact of precision reduction in sparse linear systems solvers". Peer J Computer Science 2022, 8:e778





- If using same precision, then this is the restarted GMRES method (see Lecture 3)
- FGMRES as outer iteration and GMRES as inner solver (Buttari et al. 2008)
- FGMRES as outer iteration and low precision LU as inner solver is stable (Arioli and Duff 2009)

A. Abdelfattah et al., "Advances in Mixed Precision Algorithms: 2021 Edition". doi:10.2172/1814677.

Performance of MP-GMRES





- MP-GMRES has almost same convergence behavior as DP-GMRES, much better than SP-GMRES
- GEMV can be speed-up by 1.28 and 1.57 times
- SpMV can be speed-up by 2.48 times!!!
- But ... such information is not enough! How about Preconditioned GMRES?

Source: A. Abdelfattah et al., "Advances in Mixed Precision Algorithms: 2021 Edition". doi:10.2172/1814677.

Performance of Preconditioned MP-GMRES





Figure 7: Convergence (left) and solve times (right) for a Laplacian on a stretched grid with a degree 40 polynomial preconditioner. We tested i) fp64 GMRES with fp64 polynomial [Double Prec], ii) fp64 GMRES with fp32 polynomial [Single

Source: Abdelfattan, A. et al, "Advances in Mixed Precision Algorithms: 2021 Edition". doi:10.2172/1814677.

Robustness of Scaling Algorithms



Listing 12: Scaling method

- 1 %% Given a nonsingular matrix $A \in \mathbb{R}^{n imes n}$ and a parameter $heta \in (0,1]$;
- $2 \quad | R = \operatorname{diag}(\|(A(i,:))\|_\infty^{-1})$
- 3 $A \leftarrow RA$ %% row equilibration
- 4 $P = diag(||(A(:,j))||_{\infty}^{-1})$
- 5 $A \leftarrow AP$ %% column equilibration
- 6 $|A \leftarrow \operatorname{fl}(\theta x_{\max} \beta^{-1} A)$ %% β is the maximum magnitude of |A|

Table 3. Number of iterations required by the proposed iterative refinement solver having the FP32 as lower precision (dsgesv) using the GM, IR and IRGM refinement methods. Also, we show the effect of the different scaling techniques in each block-column.

name size $\kappa_{\infty}(A)$ crystk03 24 696 2.09×10^2 crystm03 24 696 8.24×10^1 TEM27623 27 623 4.70×10^3 Goodwin054 32 510 4.49×10^6 thermal1 17 880 1.60×10^3 Zhao1 33 861 4.45×10^2 nd6k 18 000 3.82×10^2 nd12k 36 000 3.51×10^2 epb1 14 734 2.62×10^4 appu 14 000 1.02×10^4				scalar scaling			diagonal scaling				diag + scalar scaling		
crystk03 24 696 2.09 × 10 ² crystm03 24 696 8.24 × 10 ¹ TEM27623 27 623 4.70 × 10 ³ Goodwin054 32 510 4.49 × 10 ⁶ thermal1 17 880 1.60 × 10 ³ Zhao1 33 861 4.45 × 10 ² nd6k 18 000 3.82 × 10 ² nd12k 36 000 3.51 × 10 ² epb1 14 734 2.62 × 10 ⁴ appu 14 000 1.02 × 10 ⁴ ns3Da 20 414 3.79 × 10 ³	IR	IRGM	GM	IR	IRGM	GM	$\kappa_{\infty}(A^{s})$	IR	IRGM	GM	IR	IRGM	GM
crystm03 24 696 8.24 × 10 ¹ TEM27623 27 623 4.70 × 10 ³ Goodwin054 32 510 4.49 × 10 ⁶ thermal1 17 880 1.60 × 10 ³ Zhao1 33 861 4.45 × 10 ² nd6k 18 000 3.82 × 10 ² nd12k 36 000 3.51 × 10 ² epb1 14 734 2.62 × 10 ⁴ appu 14 000 1.02 × 10 ⁴ ns3Da 20 414 3.79 × 10 ³	1	1	2	1	1	2	2.5×10^{1}	1	1	2	1	1	2
TEM27623 27 623 4.70 × 10 ³ Goodwin054 32 510 4.49 × 10 ⁶ thermal1 17 880 1.60 × 10 ³ Zhao1 33 861 4.45 × 10 ² nd6k 18 000 3.82 × 10 ² nd12k 36 000 3.51 × 10 ² epb1 14 734 2.62 × 10 ⁴ appu 14 000 1.02 × 10 ⁴ ns3Da 20 414 3.79 × 10 ³	1	1	2	1	1	2	7.7×10^{0}	1	1	2	1	1	2
Goodwin054 32 510 4.49 × 10 ⁶ thermal1 17 880 1.60 × 10 ³ Zhao1 33 861 4.45 × 10 ² nd6k 18 000 3.82 × 10 ² nd12k 36 000 3.51 × 10 ² epb1 14 734 2.62 × 10 ⁴ appu 14 000 1.02 × 10 ⁴ ns3Da 20 414 3.79 × 10 ³	1	1	1	1	1	2	2.1 × 10 ¹	1	1	2	1	1	2
thermal1 17 880 1.60 × 10 ³ Zhao1 33 861 4.45 × 10 ² nd6k 18 000 3.82 × 10 ² nd12k 36 000 3.51 × 10 ² epb1 14 734 2.62 × 10 ⁴ appu 14 000 1.02 × 10 ⁴ ns3Da 20 414 3.79 × 10 ³	2	2	3	2	2	3	2.2×10^{4}	2	2	3	2	2	3
Zhao1 33 861 4.45 × 10 ² nd6k 18 000 3.82 × 10 ² nd12k 36 000 3.51 × 10 ² epb1 14 734 2.62 × 10 ⁴ appu 14 000 1.02 × 10 ⁴ ns3Da 20 414 3.79 × 10 ³	1	1	3	1	1	2	8.0×10^{1}	1	2	3	2	2	2
nd6k 18 000 3.82 × 10 ² nd12k 36 000 3.51 × 10 ² epb1 14 734 2.62 × 10 ⁴ appu 14 000 1.02 × 10 ⁴ ns3Da 20 414 3.79 × 10 ³	1	1	3	2	1	2	1.6 × 10 ¹	2	2	2	2	2	2
nd12k 36 000 3.51 × 10 ² epb1 14 734 2.62 × 10 ⁴ appu 14 000 1.02 × 10 ⁴ ns3Da 20 414 3.79 × 10 ³	1	1	2	1	1	2	1.0×10^{2}	1	1	2	1	1	2
epb1 14 734 2.62 × 10 ⁴ appu 14 000 1.02 × 10 ⁴ ns3Da 20 414 3.79 × 10 ³	1	1	2	1	1	2	1.0×10^{2}	1	1	2	1	1	2
$\begin{array}{c} appu & 14\ 000 & 1.02 \times 10^4 \\ ns3Da & 20\ 414 & 3.79 \times 10^3 \end{array}$	2	2	4	2	2	3	$4.4 imes10^4$	2	2	3	2	2	3
ns3Da 20 414 3.79×10^3	2	2	3	2	2	2	5.8×10^{3}	2	2	3	2	2	2
	2	2	3	2	2	2	3.4×10^{3}	2	2	3	2	2	2
mixtank new 29 957 4.69 × 10 ¹¹	2	2	4	2	2	3	2.1 × 10 ⁵	2	2	4	2	2	3
e40r0100 17 281 2.23 × 10 ⁸	2	3	4	3	4	4	1.1 × 10 ⁴	2	2	3	2	2	2
bcsstk25 15 439 1.02 × 10 ¹⁰	1	1	1	1	1	2	6.2×10^{2}	1	1	2	1	1	2
bcsstk37 25 503 2.27 × 10 ⁹	1	1	2	1	1	2	1.5×10^{2}	1	1	2	1	1	2
raefsky3 21 200 5.18 × 10 ¹¹	1	1	3	1	1	2	1.2×10^{4}	2	2	3	2	2	2
thread 29 736 1.78×10^8	1	1	1	1	1	2	1.3×10^{3}	1	1	2	1	1	2
sme3Db 29 067 1.05 × 10 ⁸	15	10	9	15	11	9	3.8 × 10 ⁷	16	11	12	16	11	9
mult dcop 01 25 187 9.55 × 10 ¹¹	3	2	6	3	4	4	1.5×10^{8}	2	2	5	2	2	3
ramage02 16 830 4.52 × 10 ⁸	—	60	39	—	58	38	4.5 × 10 ⁷	11	12	15	37	28	22

accelerate solution of linear systems. Proc. R. Soc. 2020

Source: Haidar A, et al.,

refinement using

Mixed-precision iterative

tensor cores on GPUs to

Table 4. Number of iterations required by the proposed iterative refinement solver having the FP16-TC as lower precision (dhgesv-TC) using the GM, IR and IRGM refinement methods. Also, we show the effect of the different scaling techniques in each block-column.

			no scali	ng		scalar scaling			diagonal scal	diagonal scaling				diag + scalar scaling		
name	size	$\kappa_{\infty}(A)$	IR	IRGM	GM	IR	IRGM	GM	$\kappa_{\infty}(A^{s})$	IR	IRGM	GM	IR	IRGM	GM	
crystk03	24 696	2.09×10^{2}	1	1	2	1	1	2	$2.5 imes 10^{1}$	1	1	2	1	1	2	
crystm03	24 696	8.24×10^{1}	1	1	2	1	1	2	7.7×10^{0}	1	1	2	1	1	2	
TEM27623	27 623	4.70×10^{3}	1	1	1	1	1	2	2.1×10^{1}	1	1	2	1	1	2	
Goodwin054	32 510	$4.49 imes 10^{6}$	8	5	6	8	5	5	$2.2 imes 10^4$	5	5	6	9	5	4	
thermal1	17 880	1.60×10^{3}	3	4	4	3	4	3	$8.0 imes 10^1$	3	3	4	3	3	3	
Zhao1	33 861	4.45×10^{2}	3	4	4	3	4	3	1.6×10^{1}	3	4	4	3	4	3	
nd6k	18 000	3.82×10^{2}	3	4	4	3	4	3	1.0×10^{2}	2	2	3	2	2	2	
nd12k	36 000	3.51 × 10 ²	3	4	4	3	4	3	1.0 × 10 ²	1	1	2	1	1	2	
epb1	14 734	2.62×10^{4}	3	4	4	3	4	3	$4.4 imes 10^4$	3	4	4	4	4	3	
арри	14 000	1.02×10^{4}	4	4	5	3	4	3	5.8×10^3	4	4	5	3	4	3	
ns3Da	20 414	3.79 × 10 ³	3	4	4	3	4	3	3.4×10^{3}	3	4	4	—	_	—	
mixtank new	29 957	4.69 × 10 ¹¹	—	36	31	13	11	10	2.1 × 10 ⁵	8	9	8	—	_	—	
e40r0100	17 281	2.23×10^{8}	44	9	11	—	9	8	1.1×10^{4}	5	6	6	—	—	—	
bcsstk25	15 439	1.02×10^{10}	—	—	—	1	2	4	6.2 × 10 ²	3	4	4	3	4	3	
bcsstk37	25 503	2.27 × 10 ⁹	—	_	—	2	2	4	1.5 × 10 ²	3	4	4	3	4	3	
raefsky3	21 200	5.18 × 10 ¹¹	—	—	—	4	4	11	1.2×10^{4}	7	6	8	6	6	5	
thread	29 736	1.78×10^{8}	—	—	—	2	2	3	1.3 × 10 ³	3	4	4	3	3	3	
sme3Db	29 067	1.05×10^{8}	—	15	9	74	13	10	3.8 × 10 ⁷	—	20	11	—	18	9	
mult dcop 01	25 187	9.55 × 10 ¹¹	—	_	—	—	108	62	1.5 × 10 ⁸	—	47	34	—	63	26	
ramage02	16 830	4.52 × 10 ⁸	—	_	—	—	—	—	4.5 × 10 ⁷	—	179	83	—	—	—	



Algorithm 15: MPV (multi-precision multigrid V-cycle)



Ref: McCormick, Stephen F., Joseph Benzaken, and Rasmus Tamstorf. "Algebraic error analysis for mixed-precision multigrid solvers". SIAM Journal on Scientific Computing 43.5 (2021): S392-S419.

Sources of Error in Simulation, Revisited



Floating-point error

 $u(x) = U_h(x) + \mathcal{E}_{dis} + \mathcal{E}_{alg} + \mathcal{E}_{guant} + \mathcal{E}_{round}$

- Quantization error: The error caused by representing algebraic system in floating-point storage formats
- Rounding error: The error caused by storing an exact solution to the algebraic system in floatingpoint storage formats
- Quantization error might be the dominating term!?

Ref: Rasmus Tamstorf, Joseph Benzaken, and Stephen F. McCormick. "Discretization-Error-Accurate Mixed-Precision Multigrid Solvers". SIAM Journal on Scientific Computing 2021 43:5, S420-S447

Some Software with MP Support



Software	Description	MP Support
GINKGO	A modern linear algebra library engineered towards performance portability, and productivity	 ✓ memory accessor that encapsulates on-the-fly compression for decoupling the memory precision from the arithmetic precision ✓ mix & match precisions in the sense of combining linear operators and vectors stored in different precision formats without explicit conversion ✓ mixed precision SAI preconditioning or CB-GMRES
Trilinos, Kokkos	An object-oriented software framework for the solution of large-scale, complex multi-physics engineering and scientific problems on new and emerging high-performance computing (HPC) architectures.	 ✓ half precision capabilities in Kokkos Core ✓ single precision MueLu (multigrid) and IfPack2 (factorization- based) preconditioners ✓ GMRES-IR solvers
MAGMA	dense (BLAS and LAPACK), sparse, and batched linear algebra routines for single nodes with GPUs	 mixed-precision iterative refinement solvers based on LU, Cholesky, and QR factorizations

Mix-Precision NLA Papers





A survey of numerical linear algebra methods utilizing mixed-precision arithmetic

A. Abdelfattah, H. Anzt, E. Boman, E. Carson, T. Cojean, J. Dongarra, Alyson Fox, M. Gates, N. Higham, Xiaoye S. Li, J. Loe, P. Luszczek, S. Pranesh, S. Rajamanickam, T. Ribizel, Barry Smith, K. Swirydowicz, Stephen J. Thomas, S. Tomov, Y. Tsai, U. Yang

2021, Int. J. High Perform. Comput. Appl.

•••

13 Citations

🛛 Save

Open in: 🔥 🏷 💩 🎖

The efficient utilization of mixed-precision numerical linear algebra algorithms can offer attractive acceleration to scientific computing applications. Especially with the hardware integration of low-precision special-function units designed for machine learning applications, the traditional numerical algorithms community urgently needs to reconsider the floating point formats used in the distinct operations to efficiently leverage the available compute power. In this work, we provide a comprehensive survey of mixed-precision numerical linear algebra routines, including the underlying concepts, theoretical background, and experimental results for both dense and sparse linear algebra problems.

Source: https://www.connectedpapers.com/



Acta Numerica (2021), pp. 1–69 doi:10.1017/S09624929XXXXXXX © Cambridge University Press, 2021 Printed in the United Kingdom

Mixed precision algorithms in numerical linear algebra

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

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

- Did you have trouble with floatingpoint precision related problems?
- Did you experience trouble to reach
 - desired precision in computation?
- What was the reason? How did you handle them?
- Is the solvers or other algorithms you use numerically stable?
- Can you use mixed precision some where in your code?

Contact Me

NCMIS

- Office hours: Mon 14:00—15:00
- Walk-in or online with appointment
- <u>zhangcs@lsec.cc.ac.cn</u>
- http://lsec.cc.ac.cn/~zhangcs



My sincere gratitude to:

Xiaowen Xu, Ningxi Tian, Ting Lai, Bin Dai

中国科学院大学 夏季强化课程 20222

Fast Solvers for Large Algebraic Systems

THANKS

Chensong Zhang, AMSS

http://lsec.cc.ac.cn/~zhangcs

Release version 2022.07.8