

Iterative Refinement in Three Precisions for Fast and Accurate Solution of Ill-Conditioned Sparse Linear Systems

Nick Higham
School of Mathematics
The University of Manchester

<http://www.maths.manchester.ac.uk/~higham/>

Joint work with Erin Carson (NYU)

**SIAM Conference on Parallel Processing for Scientific
Computing, Waseda University, Tokyo, March 7-10, 2018**



Type	Size	Range	$u = 2^{-t}$
half	16 bits	$10^{\pm 5}$	$2^{-11} \approx 4.9 \times 10^{-4}$
single	32 bits	$10^{\pm 38}$	$2^{-24} \approx 6.0 \times 10^{-8}$
double	64 bits	$10^{\pm 308}$	$2^{-53} \approx 1.1 \times 10^{-16}$
quadruple	128 bits	$10^{\pm 4932}$	$2^{-113} \approx 9.6 \times 10^{-35}$

- Arithmetic ops (+, -, *, /, $\sqrt{\quad}$) performed *as if* first calculated to infinite precision, then rounded.
- Default: round to nearest, round to even in case of tie.
- Half precision is a *storage format only*.

Overview

The NEON technology is a packed SIMD architecture. NEON registers are considered as vectors of elements of the same data type. Multiple data types are supported by the technology. The following table describes data types as supported by the architecture version.

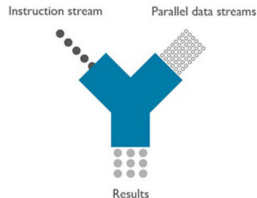
	ARMv7-A/R	ARMv8-A/R	ARMv8-A
		AArch32	AArch64
Floating-point	32-bit	16-bit*/32-bit	16-bit*/32-bit/64-bit
Integer	8-bit/16-bit/32-bit	8-bit/16-bit/32-bit/64-bit	8-bit/16-bit/32-bit/64-bit

The NEON instructions perform the same operations in all lanes of the vectors. The number of operations performed depends on the data types. NEON instructions allow up to:

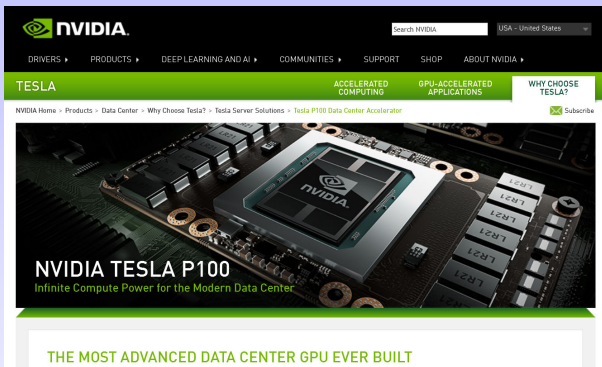
- 16x8-bit, 8x16-bit, 4x32-bit, 2x64-bit integer operations
- 8x16-bit*, 4x32-bit, 2x64-bit** floating-point operations

The implementation on NEON technology can also support issue of multiple instructions in parallel.

SIMD Architecture



NVIDIA Tesla P100 (2016)



The screenshot shows the NVIDIA website's product page for the Tesla P100. At the top, the NVIDIA logo is on the left, and a search bar and location dropdown (USA - United States) are on the right. Below the logo is a navigation menu with links for DRIVERS, PRODUCTS, DEEP LEARNING AND AI, COMMUNITIES, SUPPORT, SHOP, and ABOUT NVIDIA. A green banner highlights 'TESLA' with sub-links for 'ACCELERATED COMPUTING', 'GPU-ACCELERATED APPLICATIONS', and 'WHY CHOOSE TESLA?'. The main content area features a breadcrumb trail: 'NVIDIA Home > Products > Data Center > Why Choose Tesla? > Tesla Server Solutions > Tesla P100 Data Center Accelerator'. Below this is a large image of the Tesla P100 GPU. Overlaid on the image is the text 'NVIDIA TESLA P100' in white, with the tagline 'Infinite Compute Power for the Modern Data Center' in green below it. At the bottom of the image area, a green banner reads 'THE MOST ADVANCED DATA CENTER GPU EVER BUILT'.

“The Tesla P100 is the world’s first accelerator built for deep learning, and has native hardware ISA support for **FP16** arithmetic, delivering over 21 TeraFLOPS of **FP16** processing power.”

AMD Radeon Instinct MI25 GPU (2017)

RADEON INSTINCT
EXCITING. LIGHT. RELIABLE. LEARNING.

PRODUCTS APPLICATIONS SOFTWARE SUPPORT STORE

RADEON INSTINCT MI25

World's **Fastest** Training Accelerator for Machine Intelligence and Deep Learning ¹

NOTIFY ME

Introducing the **Radeon Instinct™ MI25**

PERFORMANCE FEATURES SPECIFICATION

FACEBOOK TWITTER YOUTUBE

World's Most Advanced GPU Memory Architecture and Next-Generation Compute Engine
Powered by the Vega™ Architecture.

“24.6 TFLOPS FP16 or 12.3 TFLOPS FP32 peak GPU compute performance on a single board . . . Up to 82 GFLOPS/watt FP16 or 41 GFLOPS/watt FP32 peak GPU compute performance”

Accelerating the Solution of $Ax = b$

$A \in \mathbb{R}^{n \times n}$ nonsingular.

Standard method for solving $Ax = b$: factorize $A = LU$, solve $LUx = b$, all at working precision.

Can we solve $Ax = b$ **faster** and/or **more accurately** by exploiting multiprecision arithmetic?

Iterative Refinement for $Ax = b$ (classic)

Solve $Ax_0 = b$ by LU factorization in **double precision**.

- $r = b - Ax_0$ **quad precision**
- Solve $Ad = r$ **double precision**
- $x_1 = x_0 + d$ **double precision**

($x_0 \leftarrow x_1$ and iterate as necessary.)

- Programmed in **J. H. Wilkinson**, *Progress Report on the Automatic Computing Engine* (1948).
- Popular up to 1970s, exploiting cheap accumulation of inner products.

Iterative Refinement (1970s, 1980s)

Solve $Ax_0 = b$ by LU factorization.

- $r = b - Ax_0$
- Solve $Ad = r$
- $x_1 = x_0 + d$

Everything in **double precision**.

- **Skeel** (1980).
- **Jankowski & Woźniakowski** (1977) for a general solver.

Iterative Refinement (2000s)

Solve $Ax_0 = b$ by LU factorization in **single precision**.

- $r = b - Ax_0$ **double precision**
 - Solve $Ad = r$ **single precision**
 - $x_1 = x_0 + d$ **double precision**
-
- **Dongarra et al.** (2006).
 - Motivated by single precision being at least twice as fast as double.

Iterative Refinement in Three Precisions

A, b given in **precision u** .

Solve $Ax_0 = b$ by LU factorization in **precision u_f** .

■ $r = b - Ax_0$ **precision u_r**

■ Solve $Ad = r$ **precision u_f**

■ $x_1 = \text{fl}(x_0 + d)$ **precision u**

■ Three previous usages are special cases.

■ Choose precisions from half, single, double, quadruple subject to $u_r \leq u \leq u_f$.

■ **Can we compute *more accurate* solutions *faster*?**

Existing Rounding Error Analysis

- **Wilkinson** (1963): fixed-point arithmetic.
- **Moler** (1967): floating-point arithmetic.
- **Higham** (1997, 2002): more general analysis for arbitrary solver.
- **Dongarra et al.** (2006): lower precision LU.

At most two precisions and require $\kappa(A)u < 1$.

New Analysis

- Applies to any solver.
- Covers b'err and f'err. Focus mainly on f'err.
- Allows $\kappa(A)u \gtrsim 1$.

New Analysis

Assume computed solution to $Ad_i = r_i$ has normwise relative error $O(u_f)$ and satisfies

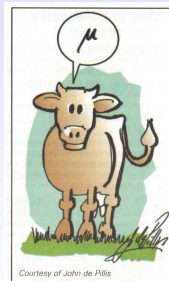
$$\frac{\|d_i - \hat{d}_i\|_\infty}{\|d_i\|} \leq \mu_i < 1.$$

Define μ_i by

$$\|A(x - \hat{x}_i)\|_\infty = \mu_i \|A\|_\infty \|x - \hat{x}_i\|_\infty,$$

and note that

$$\kappa_\infty(A)^{-1} \leq \mu_i \leq 1.$$



Condition Numbers

$$|A| = (|a_{ij}|).$$

$$\text{cond}(A, x) = \frac{\| |A^{-1}| |A| |x| \|_{\infty}}{\|x\|_{\infty}},$$

$$\text{cond}(A) = \text{cond}(A, e) = \| |A^{-1}| |A| \|_{\infty},$$

$$\kappa_{\infty}(A) = \|A\|_{\infty} \|A^{-1}\|_{\infty}.$$

$$1 \leq \text{cond}(A, x) \leq \text{cond}(A) \leq \kappa_{\infty}(A).$$

Convergence Result

Theorem (Carson & H, 2017)

For IR in precisions $u_r \leq u \leq u_f$ if

$$\phi_i = 2u_f \min(\text{cond}(\mathbf{A}), \kappa_\infty(\mathbf{A})\mu_i) + u_f \theta_i$$

is sufficiently less than 1, the forward error is reduced on the i th iteration by a factor $\approx \phi_i$ until an iterate \hat{x} satisfies

$$\frac{\|x - \hat{x}\|_\infty}{\|x\|_\infty} \lesssim 4nu_r \text{cond}(\mathbf{A}, x) + u.$$

Analogous standard bound would have

- $\mu_i = 1$,
- $u_f \theta_i = \kappa(\mathbf{A})u_f$.

Backward Error Result

Theorem (Carson & H, 2017)

For IR in precisions $u_r \leq u \leq u_f$, if

$$\phi = (c_1 \kappa_\infty(A) + c_2) u_f$$

is sufficiently less than 1 then the residual is reduced on each iteration by a factor approximately ϕ until

$$\|b - A\hat{x}_i\|_\infty \lesssim nu(\|b\|_\infty + \|A\|_\infty \|\hat{x}_i\|_\infty).$$

Precision Combinations

H = half, S = single, D = double, Q = quad. “ $u_f u u_r$ ”:

Traditional:

SSD
DDQ
HHS
HHD
HHQ
SSQ

1970s/1980s:

SSS
DDD
HHH
QQQ

2000s:

SDD
HSS
DQQ
HDD
HQQ
SQQ

3 precisions:

HSD
HSQ
HDQ
SDQ

Results for LU Factorization (1)

u_f	u	u_r	$\kappa_\infty(A)$	Backward error		Forward error
				norm	comp	
H	S	S	10^4	S	S	$\text{cond}(A, x) \cdot S$
H	S	D	10^4	S	S	S
H	D	D	10^4	D	D	$\text{cond}(A, x) \cdot D$
H	D	Q	10^4	D	D	D
S	S	S	10^8	S	S	$\text{cond}(A, x) \cdot S$
S	S	D	10^8	S	S	S
S	D	D	10^8	D	D	$\text{cond}(A, x) \cdot D$
S	D	Q	10^8	D	D	D

Results (2): HSD vs. SSD

u_f	u	u_r	$\kappa_\infty(A)$	Backward error		Forward error
				norm	comp	
H	S	S	10^4	S	S	$\text{cond}(A, x) \cdot S$
H	S	D	10^4	S	S	S
H	D	D	10^4	D	D	$\text{cond}(A, x) \cdot D$
H	D	Q	10^4	D	D	D
S	S	S	10^8	S	S	$\text{cond}(A, x) \cdot S$
S	S	D	10^8	S	S	S
S	D	D	10^8	D	D	$\text{cond}(A, x) \cdot D$
S	D	Q	10^8	D	D	D

Results (2): HSD vs. SSD

u_f	u	u_r	$\kappa_\infty(A)$	Backward error		Forward error
				norm	comp	
H	S	S	10^4	S	S	$\text{cond}(A, x) \cdot S$
H	S	D	10^4	S	S	S
H	D	D	10^4	D	D	$\text{cond}(A, x) \cdot D$
H	D	Q	10^4	D	D	D
S	S	S	10^8	S	S	$\text{cond}(A, x) \cdot S$
S	S	D	10^8	S	S	S
S	D	D	10^8	D	D	$\text{cond}(A, x) \cdot D$
S	D	Q	10^8	D	D	D

Can we get the benefit of “HSD” while allowing a larger range of $\kappa_\infty(A)$?

Extending the Range of Applicability

Recall that the convergence condition is

$$\phi_i = 2u_f \min(\text{cond}(A), \kappa_\infty(A)\mu_i) + u_f\theta_i \ll 1.$$

- We need **both terms** to be smaller than $\kappa_\infty(A)u_f$.

Recall that

$$\frac{\|d_i - \hat{d}_i\|_\infty}{\|d_i\|} \leq u_f\theta_i,$$

$$\mu_i \|A\|_\infty \|x - x_i\|_\infty = \|A(x - x_i)\|_\infty = \|b - Ax_i\|_\infty = \|r_i\|_\infty.$$

Bounding μ_i

For a stable solver, in the **early stages** we expect

$$\frac{\|r_i\|}{\|A\|\|x_i\|} \approx u \ll \frac{\|x - x_i\|}{\|x\|},$$

or equivalently $\mu_i \ll 1$. But **close to convergence**

$$\frac{\|r_i\|}{\|A\|\|x_i\|} \approx u \approx \frac{\|x - x_i\|}{\|x\|} \quad \text{or} \quad \mu_i \approx 1.$$

Conclude

$\mu_i \ll 1$ initially and $\mu_i \rightarrow 1$ as the iteration converges.

Bounding θ_j

- $u_f \theta_j$ bounds rel error in solution of $Ad_j = r_j$.
- We need $u_f \theta_j \ll 1$.

Standard solvers cannot achieve this for very ill conditioned A !

Empirically observed by **Rump (1990)** that if \hat{L} and \hat{U} are computed LU factors of A from GEPP then

$$\kappa(\hat{L}^{-1}A\hat{U}^{-1}) \approx 1 + \kappa(A)u,$$

even for $\kappa(A) \gg u^{-1}$.

Preconditioned GMRES

To compute the updates d_i we apply **GMRES** to

$$\tilde{A}d_i \equiv \hat{U}^{-1}\hat{L}^{-1}Ad_i = \hat{U}^{-1}\hat{L}^{-1}r_i.$$

- \tilde{A} is applied in twice the working precision.
- $\kappa(\tilde{A}) \ll \kappa(A)$ typically.
- Rounding error analysis shows we get an accurate \hat{d}_i even for numerically singular A .
- Call the overall alg **GMRES-IR**.
- GMRES cgce rate not directly related to $\kappa(\tilde{A})$.

Cf. **Kobayashi & Ogita** (2015), who explicitly form \tilde{A} .

Benefits of GMRES-IR

Recall $H = 10^{-4}$, $S = 10^{-8}$, $D = 10^{-16}$, $Q = 10^{-34}$.

	u_f	u	u_r	$\kappa_\infty(A)$	Backward error		
					nrm	cmp	F'error
LU	H	D	Q	10^4	D	D	D
LU	S	D	Q	10^8	D	D	D

Benefits of GMRES-IR

Recall $H = 10^{-4}$, $S = 10^{-8}$, $D = 10^{-16}$, $Q = 10^{-34}$.

	u_f	u	u_r	$\kappa_\infty(A)$	Backward error nrm cmp F'error		
LU	H	D	Q	10^4	D	D	D
LU	S	D	Q	10^8	D	D	D
GMRES-IR	H	D	Q	10^{12}	D	D	D
GMRES-IR	S	D	Q	10^{16}	D	D	D

Benefits of GMRES-IR

Recall $H = 10^{-4}$, $S = 10^{-8}$, $D = 10^{-16}$, $Q = 10^{-34}$.

	u_f	u	u_r	$\kappa_\infty(A)$	Backward error nrm cmp F'error		
LU	H	D	Q	10^4	D	D	D
LU	S	D	Q	10^8	D	D	D
GMRES-IR	H	D	Q	10^{12}	D	D	D
GMRES-IR	S	D	Q	10^{16}	D	D	D

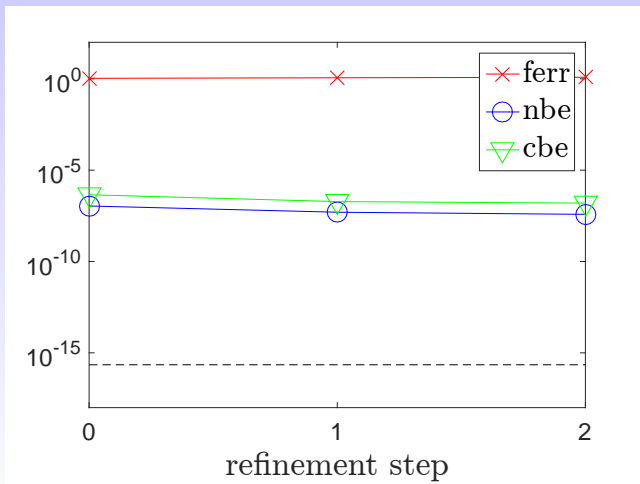
How many GMRES iterations are required?

Some tests with 100×100 matrices ...

Test 1: LU-IR, $(u_f, u, u_r) = (S, D, D)$

$$\kappa_\infty(A) \approx 10^{10}, \quad \sigma_j = \alpha^j$$

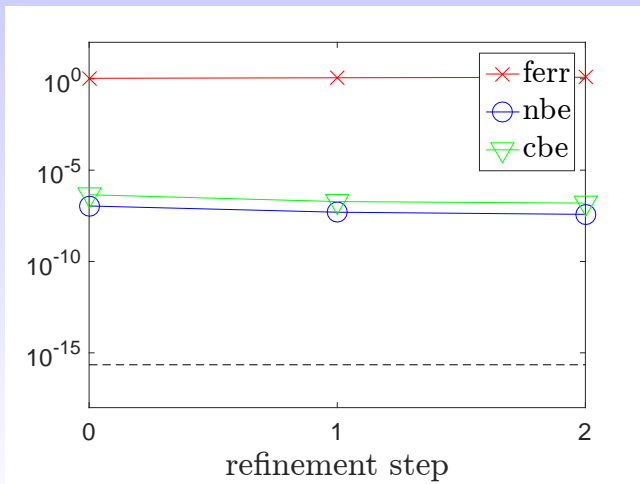
Divergence!



Test 1: LU-IR, $(u_f, u, u_r) = (S, D, Q)$

$$\kappa_\infty(\mathbf{A}) \approx 10^{10}, \quad \sigma_j = \alpha^j$$

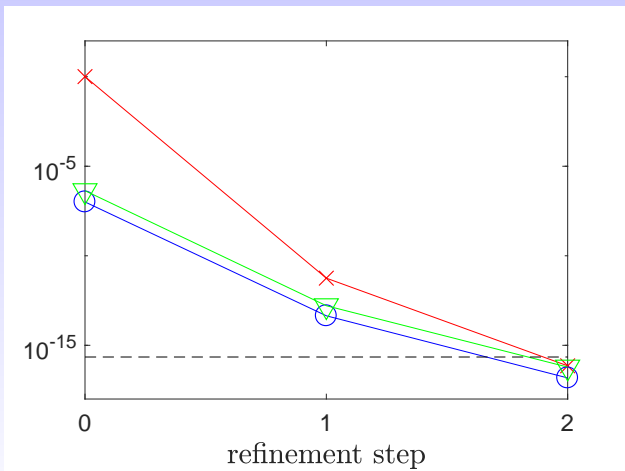
Divergence!



Test 1: LU-IR, $(u_f, u, u_r) = (S, D, Q)$

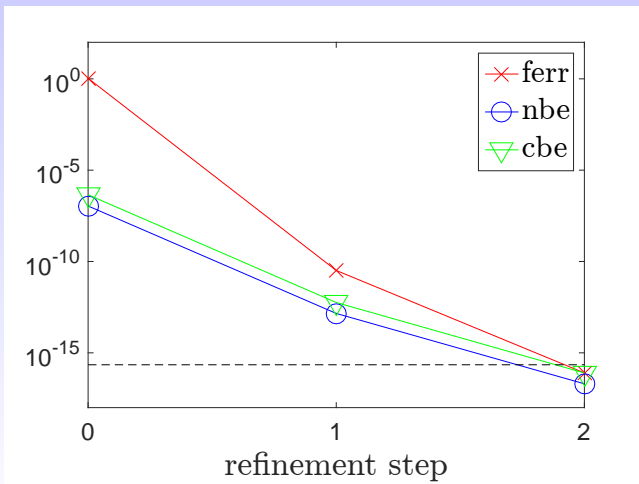
$$\kappa_\infty(A) \approx 10^4, \quad \sigma_j = \alpha^j$$

Convergence



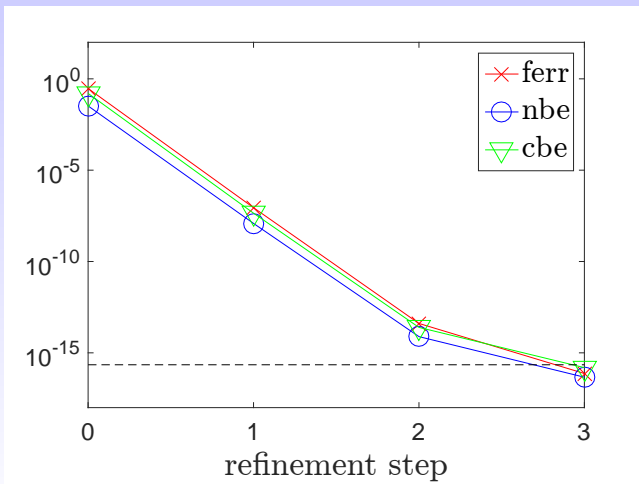
Test 1: GMRES-IR, $(u_f, u, u_r) = (S, D, Q)$

$\kappa_\infty(\mathbf{A}) \approx 10^{10}$, $\sigma_i = \alpha^i$, GMRES its (2,3) **Convergence**



Test 2: GMRES-IR, $(u_f, u, u_r) = (H, D, Q)$

$\kappa_\infty(\mathbf{A}) \approx 10^2$, 1 small σ_i , GMRES its (8,8,8) **Convergence**

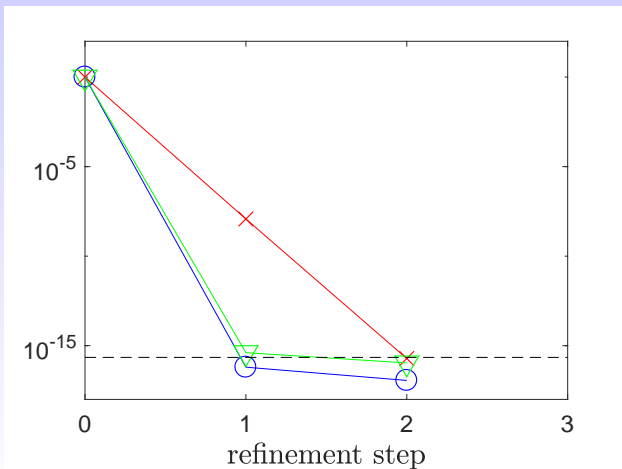


Test 3: GMRES-IR, $(u_f, u, u_r) = (H, D, Q)$

$\kappa_\infty(A) \approx 10^{12}$, $\sigma_i = \alpha^i$, GMRES (100,100)

Take $x_0 = 0$ because of overflow!

Convergence




Conclusions



- Both **low and high precision** floating-point arithmetic becoming more prevalent, in hardware and software.
- Showed that using **three precisions** in iter ref brings major benefits.
 - LU in lower precision \Rightarrow twice as fast as trad. IR.
 - GMRES-IR cges when trad. IR doesn't thanks to preconditioned GMRES solved of $Ad_i = r_i$.
- **GMRES-IR** handles $\kappa_\infty(A) \approx u^{-1}$.
- Dongarra et al. (2017) have achieved 1.8x speedups on GPUs.

Erin Carson & H, **A New Analysis of Iterative Refinement and its Application to Accurate Solution of Ill-Conditioned Sparse Linear Systems**,
SIAM J. Sci. Comput., 39(6):A2834–A2856, 2017.




Erin Carson & H, **Accelerating the Solution of Linear Systems by Iterative Refinement in Three Precisions**,
MIMS EPrint 2017.24, The University of Manchester, July 2017; to appear in *SIAM J. Sci. Comput.*

-  E. Carson and N. J. Higham.
Accelerating the solution of linear systems by iterative refinement in three precisions.
MIMS EPrint 2017.24, Manchester Institute for
Mathematical Sciences, The University of Manchester,
UK, July 2017.
30 pp.
Revised January 2018. To appear in SIAM J. Sci.
Comput.


References II

-  E. Carson and N. J. Higham.
A new analysis of iterative refinement and its application to accurate solution of ill-conditioned sparse linear systems.
SIAM J. Sci. Comput., 39(6):A2834–A2856, 2017.
-  A. Haidar, P. Wu, S. Tomov, and J. Dongarra.
Investigating half precision arithmetic to accelerate dense linear system solvers.
In *Proceedings of the 8th Workshop on Latest Advances in Scalable Algorithms for Large-Scale Systems*, *ScalA '17*, pages 10:1–10:8, Nov. 2017.

References III

-  N. J. Higham.
Iterative refinement for linear systems and LAPACK.
IMA J. Numer. Anal., 17(4):495–509, 1997.
-  N. J. Higham.
Accuracy and Stability of Numerical Algorithms.
Society for Industrial and Applied Mathematics,
Philadelphia, PA, USA, second edition, 2002.
ISBN 0-89871-521-0.
xxx+680 pp.
-  Y. Kobayashi and T. Ogita.
A fast and efficient algorithm for solving ill-conditioned
linear systems.
JSIAM Lett., 7:1–4, 2015.

References IV

-  J. Langou, J. Langou, P. Luszczek, J. Kurzak, A. Buttari, and J. Dongarra.

Exploiting the performance of 32 bit floating point arithmetic in obtaining 64 bit accuracy (revisiting iterative refinement for linear systems).

In Proceedings of the 2006 ACM/IEEE Conference on Supercomputing, Nov. 2006.