

ARITH 21, Austin, TX, 2012

On the componentwise accuracy of complex floating point division with an FMA

Claude-Pierre Jeannerod
INRIA

Nicolas Louvet
UCB Lyon 1

Jean-Michel Muller
CNRS

Aric research team, LIP



Complex division

Given $x = a + ib$ and $y = c + id$, $z = x/y$ is can be expressed as

$$z = \frac{ac + bd}{c^2 + d^2} + i \frac{bc - ad}{c^2 + d^2}.$$

Several issues when performed in floating point arithmetic:

- **Loss of componentwise accuracy**, due to cancellations;
- **Spurious overflow/underflow**, e.g., $c^2 + d^2$ overflows, while the real and imaginary parts are representable.

In this talk:

- **We focus on accuracy problems.**
- We assume that a **Fused Multiply-Add (FMA)** instruction is available.

Assumptions and notation

- We assume a standard **radix-2**, **precision- p** ($p \geq 2$), floating point arithmetic.
- $\text{RN}(\cdot)$ denotes **rounding 'to the nearest even'**.
- An **FMA** instruction is available:
given three floating point numbers x , y and z , it computes $\text{RN}(xy + z)$;
- $u = 2^{-p}$ denotes the unit roundoff;
- ulp to denote the **unit in the last place** function:
 $\text{ulp}(0) = 0$ and for any nonzero real number t , $\text{ulp}(t) = 2^{\lfloor \log_2(t) \rfloor + 1 - p}$.
- We assume that **underflows and overflows never occur**, so that

$$|\text{RN}(t) - t| \leq \frac{1}{2} \text{ulp}(t) \leq u|t| \quad \text{for any real number } t.$$

High normwise accuracy

Let $x = a + ib$ and $y = c + id$ with floating point coefficients.

$z = x/y$ can be evaluated **without FMA** as

$$\operatorname{Re} \hat{z} = \operatorname{RN} \left(\frac{\operatorname{RN}(\operatorname{RN}(ac) + \operatorname{RN}(bd))}{\operatorname{RN}(\operatorname{RN}(c^2) + \operatorname{RN}(d^2))} \right),$$

and a similar expression for the imaginary part.

In the absence of underflow and overflow [Baudin, Smith, 2012],

$$\frac{|\hat{z} - z|}{|z|} < \underbrace{(3 + \sqrt{5})}_{=5.236\dots} u + 13u^2.$$

↪ The computed \hat{z} is always highly accurate in the normwise sense.

Componentwise relative error

$$E_c(\hat{z}) := \max\left(\frac{|\operatorname{Re} \hat{z} - \operatorname{Re} z|}{|\operatorname{Re} z|}, \frac{|\operatorname{Im} \hat{z} - \operatorname{Im} z|}{|\operatorname{Im} z|}\right).$$

The classic division formula can be **very inaccurate in the componentwise sense**.

We provide an example ($p = 53$) such that:

- $E_c(\hat{z}) \approx 9.0 \times 10^{15} \gg 1$ with the classic algorithm **without FMA**:

$$\operatorname{Re} \hat{z} = \operatorname{RN}\left(\frac{\operatorname{RN}(\operatorname{RN}(ac) + \operatorname{RN}(bd))}{\operatorname{RN}(\operatorname{RN}(c^2) + \operatorname{RN}(d^2))}\right).$$

- $E_c(\hat{z}) \approx 4.5 \times 10^{15} \gg 1$ with the classic algorithm **with FMA**:

$$\operatorname{Re} \hat{z} = \operatorname{RN}\left(\frac{\operatorname{RN}(ac + \operatorname{RN}(bd))}{\operatorname{RN}(c^2 + \operatorname{RN}(d^2))}\right).$$

↔ How to obtain high componentwise accuracy?

Kahan's algorithm for $g = ac + bd$

```

Kahan( $a, b, c, d$ )
   $\hat{w} := \text{RN}(bd)$ ;
   $e := \text{RN}(bd - \hat{w})$ ;           //  $e = bd - \hat{w}$  is computed exactly
   $\hat{f} := \text{RN}(ac + \hat{w})$ ;       //  $\hat{f}$  is a naive evaluation of  $ac + bd$  with FMA
   $\hat{g} := \text{RN}(\hat{f} + e)$ ;       // compensation step
  return  $\hat{g}$ ;

```

Kahan's algorithm always approximates $ac + bd$ to **high relative accuracy**.

Theorem ([Jeannerod, Louvet, Muller, 2012])

In the absence of underflow/overflow,

$$\frac{|\hat{g} - g|}{|g|} \leq 2u.$$

*This bound is proved to be **asymptotically optimal**, via the explicit construction of inputs for which $\frac{|\hat{g} - g|}{|g|} = 2u - O(u^2)$.*

Complex division with Kahan's algorithm

Kahan's algorithm to evaluate the numerators and the denominator in

$$z = \frac{ac + bd}{c^2 + d^2} + i \frac{bc - ad}{c^2 + d^2},$$

then performing two floating point divisions.

↪ **High componentwise accuracy** can easily be obtained: $E_c(\hat{z}) \leq 5u + O(u^2)$.

- ① However, among the three terms in the division formula **only one** can cancel. In particular, a simpler scheme should be used for the denominator.
 → Accuracy of simpler algorithms for the **sum of two nonnegative products**.
- ② We deduce two accurate FMA-based algorithms for complex division:
 - ▶ a straight-line algorithm,
 - ▶ an algorithm using tests to reduce error bound.
 → The **sharpness of the error bounds** was investigated.

Outline

- 1 Introduction
- 2 On the sum of two nonnegative products
 - Summary of the error bounds
 - The case $\text{ulp}(bd) \leq \text{ulp}(ac)$
- 3 Two complex division algorithms
 - A straight-line algorithm
 - If tests are allowed
- 4 Conclusion and future work

Given a, b, c, d four floating point numbers, we focus in this part on the approximation of $g = ac + bd$ under the assumption

$$ac \geq 0 \quad \text{and} \quad bd \geq 0.$$

We consider the following ways of evaluating g :

- $\hat{f}_o = \text{RN}(\text{RN}(ac) + \text{RN}(bd))$,
- $\hat{f} = \text{RN}(ac + \text{RN}(bd))$,
- $\hat{g} = \text{Kahan}(a, b, c, d)$.

Known results:

- From [Brent, Percival, Zimmermann, 2007]:
 $\hookrightarrow 2u$ is a relative error bound for \hat{f}_o .
- From [Jeannerod, Louvet, Muller, 2012]:
 $\hookrightarrow \text{ulp}(g)$ and $2u$ are sharp absolute/relative error bounds for \hat{g} .

Summary of the error bounds

\widehat{r}	bound on $ \widehat{r} - g $	bound on $\frac{ \widehat{r} - g }{g}$
$\widehat{f}_o = \text{RN}(\text{RN}(ac) + \text{RN}(bd))$	$\frac{5}{4} \text{ulp}(g)$	$2u$ [BrePerZim07]
$\widehat{f} = \text{RN}(ac + \text{RN}(bd))$	$\text{ulp}(g)$	$2u$
$\widehat{g} = \text{Kahan}(a, b, c, d)$	$\text{ulp}(g)$ [JeLoMu12]	$2u$ [JeLoMu12]
$\text{ulp}(bd) \leq \text{ulp}(ac), \widehat{f} = \widehat{g}$	$\frac{3}{4} \text{ulp}(g)$	$\frac{3}{2}u$

- The **optimality or the asymptotic optimality** (as $p \rightarrow \infty$) of these error bounds is established using examples of the form $ac + bd$ parametrized by p .
- We investigated the sharpness of these bounds for the evaluation of **sums of squares** $c^2 + d^2$: in most cases the bounds remain asymptotically optimal.
- When no example parametrized by p were found (p odd), the quality of the error bounds is illustrated with **numerical examples in precisions** $p = 53, 113$.

The case $\text{ulp}(bd) \leq \text{ulp}(ac)$

The error bounds $\text{ulp}(g)$ and $2u$ for the general case can be improved:

Theorem

If $\text{ulp}(bd) \leq \text{ulp}(ac)$ then $|\hat{f} - g| \leq \frac{3}{4}\text{ulp}(g)$ and $\frac{|\hat{f}-g|}{g} \leq \frac{3}{2}u$.

Example

For $p \geq 6$, let

$$a = 2^{p-1} + 2^{\lfloor \frac{p}{2} \rfloor}, \quad b = c = 2^{p-1} + 2^{\lceil \frac{p}{2} \rceil - 1}, \quad d = 2^{p-1} + 2^{\lfloor \frac{p}{2} \rfloor - 1}.$$

In this example, we have $\text{ulp}(bd) = \text{ulp}(ac)$, and

- $|\hat{f} - g| = \frac{3}{4}\text{ulp}(g) \Rightarrow$ optimal absolute error bound
- $\frac{|\hat{f}-g|}{g} = \frac{3}{2}u - O(u^{3/2}) \Rightarrow$ asymptotically optimal relative error bound

Outline

- 1 Introduction
- 2 On the sum of two nonnegative products
 - Summary of the error bounds
 - The case $\text{ulp}(bd) \leq \text{ulp}(ac)$
- 3 Two complex division algorithms
 - A straight-line algorithm
 - If tests are allowed
- 4 Conclusion and future work

A straight-line algorithm

CompDivS($a + ib, c + id$)

$\hat{\delta} := \text{RN}(c^2 + \text{RN}(d^2))$ // $c^2 + d^2$ $\rightsquigarrow 2u$

$\hat{g}_{\text{re}} := \text{Kahan}(a, b, c, d)$

$\hat{g}_{\text{im}} := \text{Kahan}(b, -a, c, d)$ // $ac + bd$ and $bc - ad$ $\rightsquigarrow 2u$

$\hat{z}_{\text{re}} := \text{RN}(\hat{g}_{\text{re}}/\hat{\delta})$

$\hat{z}_{\text{im}} := \text{RN}(\hat{g}_{\text{im}}/\hat{\delta})$ // $\text{Re}(\hat{z})$ and $\text{Im}(\hat{z})$ $\rightsquigarrow u$

return $\hat{z}_{\text{re}} + i\hat{z}_{\text{im}}$

Property

For $p \geq 5$ and in the absence of underflow and overflow, $E_c(\hat{z}) \leq 5u + 13u^2$.

For p odd, the sharpness of the bound illustrated through numerical examples:

p	$ \hat{z}_{\text{re}} - \text{Re } z /(u \text{Re } z)$
53	4.9987...
113	4.9987...

For p even, the bound is proved to be asymptotically optimal.

Assuming p is even, let

$$\begin{aligned} a &= 2^p - 5 \cdot 2^{\frac{p}{2}-1}, & b &= -2^{-\frac{p}{2}} \cdot (2^p - 5 \cdot 2^{\frac{p}{2}-1} + 3), \\ c &= 2^p - 2, & d &= 2^{\frac{p}{2}+1} \cdot (2^{p-1} + 2^{\frac{p}{2}-1}). \end{aligned}$$

Defining $R = 2^{p/2}$, it can be check that

$$\begin{aligned} \operatorname{Re} z &= -\frac{2R^3 + 5R^2 - 4R}{2R^6 + 4R^5 + 4R^4 - 8R^2 + 8}, \\ \frac{\widehat{g}_{\operatorname{re}}}{\widehat{\delta}} &= \frac{-R^3 - \frac{5}{2}R^2}{R^6 + 2R^5} = \underbrace{-\frac{1}{R^3} - \frac{1}{2R^4}}_{=\operatorname{Re} \widehat{z}} + O\left(\frac{1}{R^5}\right). \end{aligned}$$

Then we deduce

$$\frac{|\operatorname{Re} \widehat{z} - \operatorname{Re} z|}{|\operatorname{Re} z|} = \frac{5}{R^2} - O\left(\frac{1}{R^3}\right) = 5u - O(u^{3/2}),$$

which is equivalent to the componentwise bound $5u + 13u^2$ as $p \rightarrow \infty$.

If tests are allowed

CompDivT($a + ib, c + id$)

if $|d| \leq |c|$ **then** $\hat{\delta} := \text{RN}(c^2 + \text{RN}(d^2))$ // $|d| \leq |c| \implies \text{ulp}(c^2) \leq \text{ulp}(d^2)$
else $\hat{\delta} := \text{RN}(d^2 + \text{RN}(c^2))$ // $c^2 + d^2 \rightsquigarrow 1.5u$
 $\hat{g}_{\text{re}} := \text{Kahan}(a, b, c, d)$
 $\hat{g}_{\text{im}} := \text{Kahan}(b, -a, c, d)$ // $ac + bd$ and $bc - ad \rightsquigarrow 2u$
 $\hat{z}_{\text{re}} := \text{RN}(\hat{g}_{\text{re}}/\hat{\delta})$
 $\hat{z}_{\text{im}} := \text{RN}(\hat{g}_{\text{im}}/\hat{\delta})$ // $\text{Re}(\hat{z})$ and $\text{Im}(\hat{z}) \rightsquigarrow u$
return $\hat{z}_{\text{re}} + i\hat{z}_{\text{im}}$

Property

For $p \geq 6$ and in the absence of underflow and overflow, $E_c(\hat{z}) \leq 4.5u + 9u^2$.

p	$ \hat{z}_{\text{re}} - \text{Re } z /(u \text{Re } z)$
24	4.4932...
53	4.4421...
113	4.4620...

If tests are allowed

CompDivT($a + ib, c + id$)

if $|d| \leq |c|$ **then** $\widehat{\delta} := \text{RN}(c^2 + \text{RN}(d^2))$ // $|d| \leq |c| \implies \text{ulp}(c^2) \leq \text{ulp}(d^2)$

else $\widehat{\delta} := \text{RN}(d^2 + \text{RN}(c^2))$ // $c^2 + d^2 \rightsquigarrow 1.5u$

$\widehat{g}_{\text{re}} := \text{Kahan}(a, b, c, d)$

$\widehat{g}_{\text{im}} := \text{Kahan}(b, -a, c, d)$ // $ac + bd$ and $bc - ad \rightsquigarrow 2u$

$\widehat{z}_{\text{re}} := \text{RN}(\widehat{g}_{\text{re}}/\widehat{\delta})$

$\widehat{z}_{\text{im}} := \text{RN}(\widehat{g}_{\text{im}}/\widehat{\delta})$ // $\text{Re}(\widehat{z})$ and $\text{Im}(\widehat{z}) \rightsquigarrow u$

return $\widehat{z}_{\text{re}} + i\widehat{z}_{\text{im}}$

On architectures supporting the `minNumMag` and `maxNumMag` operations, the first two lines may be replaced by the following straight-line code:

$\underline{c} := \text{maxNumMag}(c, d);$

$\underline{d} := \text{minNumMag}(c, d);$

// Here, $\underline{c}^2 + \underline{d}^2 = c^2 + d^2$ and $|\underline{d}| \leq |\underline{c}|$.

$\widehat{\delta} := \text{RN}(\underline{c}^2 + \text{RN}(\underline{d}^2));$

Conclusion







By combining Kahan's algorithm with cheaper schemes, we have proposed two complex division algorithms that take advantage of the FMA instruction :

- 1 A straight-line algorithm, with $E_c(\hat{z}) \leq 5u + O(u^2)$.
 \hookrightarrow bound proved to be asymptotically optimal when p is even.
- 2 A version in which $|c|$ and $|d|$ are compared to ensure $E_c(\hat{z}) \leq 4.5u + O(u^2)$.

Future work:

- Overhead induced by the use of Kahan's algorithm?
- Extensions to non binary radices and other rounding modes.
- Using scaling techniques to avoid spurious under/overflows while preserving high componentwise accuracy [Smith, 1962], [Stewart, 1985], [Priest, 2004], [Baudin, Smith, 2012].

References

-  R. P. Brent, C. Percival, and P. Zimmermann, *Error bounds on complex floating-point multiplication*, *Mathematics of Computation* **76** (2007), 1469–1481.
-  Michael Baudin and Robert L. Smith, *A robust complex division in Scilab*, October 2012, Available at <http://arxiv.org/abs/1210.4539>.
-  C.-P. Jeannerod, N. Louvet, and J.-M. Muller, *Further analysis of Kahan's algorithm for the accurate computation of 2×2 determinants*, *Mathematics of Computation* (2012), to appear. Preliminary version available at <http://hal-ens-lyon.archives-ouvertes.fr/ensl-00649347/en/>.
-  D. M. Priest, *Efficient scaling for complex division*, *ACM Trans. Math. Software* **30** (2004), no. 4.
-  Robert L. Smith, *Algorithm 116: Complex division*, *Comm. ACM* **5** (1962), no. 8, 435.
-  G. W. Stewart, *A note on complex division*, *ACM Trans. Math. Software* **11** (1985), no. 3, 238–241.