

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

Claude-Pierre Jeannerod, Nicolas Louvet, Jean-Michel Muller

► **To cite this version:**

Claude-Pierre Jeannerod, Nicolas Louvet, Jean-Michel Muller. On the componentwise accuracy of complex floating-point division with an FMA. 21st IEEE Symposium on Computer Arithmetic, Apr 2013, Austin, TX, United States. IEEE Computer Society, pp.83–90, 2013. <ensl-00734339v2>

**HAL Id: ensl-00734339**

**<https://hal-ens-lyon.archives-ouvertes.fr/ensl-00734339v2>**

Submitted on 7 Feb 2013

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

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

Claude-Pierre Jeannerod  
INRIA  
claude-pierre.jeannerod@ens-lyon.fr

Nicolas Louvet  
UCB Lyon 1  
nicolas.louvet@ens-lyon.fr

Jean-Michel Muller  
CNRS  
jean-michel.muller@ens-lyon.fr

AriC INRIA project, Laboratoire LIP (CNRS, ENS de Lyon, INRIA, UCB Lyon 1), Université de Lyon, France

**Abstract**—This paper deals with the accuracy of complex division in radix-two floating-point arithmetic. Assuming that a fused multiply-add (FMA) instruction is available and that no underflow/overflow occurs, we study how to ensure high relative accuracy in the componentwise sense. Since this essentially reduces to evaluating accurately three expressions of the form  $ac + bd$ , an obvious approach would be to perform three calls to Kahan’s compensated algorithm for 2 by 2 determinants. However, in the context of complex division, two of those expressions are such that  $ac$  and  $bd$  have the same sign, suggesting that cheaper schemes should be used here (since cancellation cannot occur). We first give a detailed accuracy analysis of such schemes for the sum of two nonnegative products, providing not only sharp bounds on both their absolute and relative errors, but also sufficient conditions for the output of one of them to coincide with the output of Kahan’s algorithm. By combining Kahan’s algorithm with this particular scheme, we then deduce two new division algorithms. Our first algorithm is a straight-line program whose componentwise relative error is always at most  $5u + 13u^2$  with  $u$  the unit roundoff; we also provide examples of inputs for which the error of this algorithm approaches  $5u$ , thus showing that our upper bound is essentially the best possible. When tests are allowed we show with a second algorithm that the bound above can be further reduced to  $4.5u + 9u^2$ , and that this improved bound is reasonably sharp.

*Keywords.* Floating-point arithmetic; complex division; fused multiply-add (FMA); rounding error analysis.

## I. INTRODUCTION

For two complex numbers  $x = a + ib$  and  $y = c + id$  given by their real and imaginary parts, the result  $z$  of the division of  $x$  by  $y$  is classically expressed as

$$z = \frac{ac + bd}{c^2 + d^2} + i \frac{bc - ad}{c^2 + d^2}. \quad (1)$$

When evaluating (1) in floating-point arithmetic and if neither underflow nor overflow occurs, then it is well known that high relative accuracy is ensured in the *normwise* sense [1, §3.6], [2, §4.5]. More precisely, assume standard binary floating-point arithmetic, in precision  $p$  and with rounding ‘to nearest even’ (RN), and let  $\hat{z}$  denote the approximation to  $z$  given by

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

for its real part and by a similar expression for its imaginary part  $\operatorname{Im} \hat{z}$ . Then the normwise relative error

$$E_n(\hat{z}) = |\hat{z} - z|/|z|$$

is bounded by  $\gamma u + O(u^2)$ , where  $u = 2^{-p}$  is the unit roundoff and  $\gamma$  is a modest constant. For example, one can take  $\gamma$  equal to  $3 + \sqrt{5} = 5.236\dots$ , as noted in [3, §3.6]. (A detailed proof of this fact is given in appendix.)

However, in the *componentwise* sense, the approximation  $\hat{z}$  can be highly inaccurate: there exist values for  $a, b, c, d$  such that the componentwise relative error

$$E_c(\hat{z}) = \max(|\operatorname{Re} \hat{z} - \operatorname{Re} z|/|\operatorname{Re} z|, |\operatorname{Im} \hat{z} - \operatorname{Im} z|/|\operatorname{Im} z|)$$

is (much) larger than one. For example, the input

$$(a, b, c, d) = (N, N + 1, -N, N - 1) \quad (3)$$

with  $N = 2^{p-1} + 2^{p-2} - 1$  leads to  $E_c(\hat{z}) \approx 9.0 \times 10^{15} \gg 1$  when  $p = 53$  and  $\hat{z}$  is computed according to (2). Here, the input consists of normal floating-point numbers and there is no underflow/overflow, but since  $ac$  and  $bd$  have opposite signs and similar magnitudes a cancellation occurs during the evaluation of  $ac + bd$  as

$$\hat{f}_o = \operatorname{RN}(\operatorname{RN}(ac) + \operatorname{RN}(bd)). \quad (4)$$

In this paper, we are interested in performing complex divisions very accurately in the componentwise sense, assuming that a fused multiply-add (FMA) instruction is available and that no underflow/overflow occurs. The FMA evaluates expressions of the form  $ab + c$  with one rounding error only and since it is required by the 2008 revision of the IEEE 754 standard [4], one can expect that it will soon belong to the instruction set of most general-purpose processors.

From (1) it is clear that in order to achieve high relative accuracy in the componentwise sense, it suffices to evaluate with high relative accuracy two-dimensional dot products, that is, expressions of the form

$$g = ac + bd.$$

With an FMA, a naive computation would produce, say,

$$\hat{f} = \operatorname{RN}(ac + \operatorname{RN}(bd)), \quad (5)$$

so that the expression in (2) can for example be replaced by

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

Implementing complex division by means of such schemes is, however, not enough for our purposes. Indeed, considering

again (3) and  $p = 53$  now gives  $E_c(\hat{z}) \approx 4.5 \times 10^{15}$ , but this componentwise relative error is still (much) larger than one.

To ensure high accuracy, the FMA is typically used further to recover the rounding errors committed when evaluating the products  $ac$  and  $bd$ , and then to compensate for them by addition; see for example [5, p. 273]. It turns out that if the FMA is already used to get  $\hat{f}$  as in (5), then only the rounding error for  $bd$  needs to be recovered. This is the basis of the algorithm below, attributed to Kahan in [1, p. 65]:

**algorithm** Kahan( $a, b, c, d$ )  
 $\hat{w} := \text{RN}(bd)$ ;  
 $e := \text{RN}(bd - \hat{w})$ ; // this operation is exact:  $e = bd - \hat{w}$ .  
 $\hat{f} := \text{RN}(ac + \hat{w})$ ;  
 $\hat{g} := \text{RN}(\hat{f} + e)$ ;  
**return**  $\hat{g}$ ;

The error  $e = bd - \text{RN}(bd)$  is computed exactly thanks to the FMA instruction, and then added to  $\hat{f}$  in order to yield the approximation  $\hat{g}$  to  $g = ac + bd$ . In [6] it was shown that in the absence of underflow and overflow

$$|\hat{g} - g| \leq 2u|g|, \quad (6)$$

so that Kahan's algorithm always approximates  $g$  to high relative accuracy; it was also shown that the bound  $2u$  is essentially best possible (even if  $ac$  and  $bd$  can be swapped) via the explicit construction of floating-point numbers  $a, b, c, d$  for which  $|\hat{g} - g|/|g| = 2u + O(u^2)$ .

Consequently, by using Kahan's algorithm to evaluate  $ac + bd$ ,  $bc - ad$ ,  $c^2 + d^2$  in (1) and then performing two floating-point divisions, we obviously get high componentwise relative accuracy: considering for example the real part,

$$\text{Re } \hat{z} = \frac{g(1 + \epsilon)}{(c^2 + d^2)(1 + \epsilon')} (1 + \epsilon''),$$

where  $|\epsilon|, |\epsilon'| \leq 2u$  by (6) and where the relative error  $|\epsilon''|$  of floating-point division is bounded by  $u$ . Hence the real part of  $\hat{z}$  has the form  $(1 + \theta) \cdot \text{Re } z$  with  $|\theta| \leq \frac{1+2u}{1-2u}(1+u) - 1$  and since the same holds for the imaginary part, it follows that the componentwise relative error satisfies

$$E_c(\hat{z}) \leq 5u + O(u^2). \quad (7)$$

Yet this approach raises two issues, which we address in this paper. First, among the three dot products  $ac + bd$ ,  $bc - ad$ ,  $c^2 + d^2$  appearing in (1) only one of them can cancel: if  $abcd \geq 0$  cancellation can occur only when evaluating  $bc - ad$ , while if  $abcd < 0$  cancellation can occur only when evaluating  $ac + bd$ . Hence, there is always exactly one potentially 'difficult' case (one of the numerators) and two 'easy' cases (the denominator  $c^2 + d^2$  and the other numerator). Although it seems obvious not to use Kahan's compensated algorithm for such easy cases where cancellation cannot occur, the behavior of simpler algorithms like those producing  $\hat{f}_o$  in (4) or  $\hat{f}$  in (5) deserves further investigation: Is the relative error bound  $2u$  given in (6) for  $\hat{g}$  also a sharp bound for  $\hat{f}_o$  and  $\hat{f}$ ? How are  $\hat{f}$  and  $\hat{g}$  related, and when do they coincide?

We address this first issue in Section II by presenting a thorough study of the case ' $ac + bd$ ' with  $ac \geq 0$  and  $bd \geq 0$ '. In particular,

- we show that without any further assumption on  $a, b, c, d$  the bound  $2u$  is indeed already sharp for  $\hat{f}_o$  and  $\hat{f}$ ;
- we also show that either  $\hat{f}$  and  $\hat{g}$  coincide, or one of them is a faithfully rounded value of the exact result  $g$ ;
- third, we show that if  $0 \leq bd \leq ac$  then  $\hat{f} = \hat{g}$  and a sharp bound is  $\frac{3}{2}u$ .

Along with these relative error bounds, we also systematically provide bounds on absolute errors, which we express in terms of ulps (units in the last place) of  $g$ . In particular, this provides a good illustration of the impact the FMA can have on error bounds: for  $|\hat{f} - g|$  and  $|\hat{f}_o - g|$ , that is, for the naive evaluation of  $g$  with or without FMA, sharp bounds are proven to be  $\text{ulp}(g)$  and  $\frac{5}{4}\text{ulp}(g)$ , respectively. Finally, in the more constrained case  $g = c^2 + d^2$  of a sum of two squares, all our error bounds are shown to remain either sharp or reasonably tight for the basic IEEE formats. By combining the results in Section II with the bound in (6) we immediately deduce two other FMA-based, highly accurate algorithms for complex floating-point division:

- a straight-line program that evaluates the denominator  $c^2 + d^2$  using (5) and the two numerators using Kahan's algorithm, and whose relative error is again clearly bounded as in (7);
- an algorithm using tests in order to always reduce the error bound on  $c^2 + d^2$  to  $\frac{3}{2}u$ , thus improving (7) to

$$E_c(\hat{z}) \leq \frac{9}{2}u + O(u^2). \quad (8)$$

The second issue concerns the sharpness of the bounds in (7) and (8). Indeed, although it is straightforward to derive the two algorithms just mentioned and to perform their rounding error analysis, it is a priori not clear that sharp bounds for the numerators and for the denominator will lead to sharp bounds for the whole complex division algorithms.

Section III is devoted to this second issue, and we show in particular through numerical examples that for all the basic binary IEEE formats

- the error bound in  $5u + O(u^2)$  is sharp for the straight-line program;
- the error bound in  $\frac{9}{2}u + O(u^2)$  is reasonably sharp for the algorithm using tests.

Furthermore, for the straight-line program and when the precision  $p$  is even, we show that the bound  $5u + O(u^2)$  is asymptotically optimal by explicitly constructing an input  $(a, b, c, d)$  for which  $E_c(\hat{z}) = 5u - O(u^{3/2})$ .

**Assumptions and notation.** Throughout this paper we make the customary assumption that  $p \geq 2$ . Also, we write  $\text{ulp}$  to denote the *unit in the last place* function:  $\text{ulp}(0) = 0$  and for any nonzero real number  $t$ ,  $\text{ulp}(t)$  is the unique integer power of two such that  $2^{p-1} \leq |t|/\text{ulp}(t) < 2^p$ . Furthermore, we assume that underflows and overflows never occur, so that

$$|\text{RN}(t) - t| \leq \frac{1}{2}\text{ulp}(t) \leq u|t| \quad \text{for any real number } t. \quad (9)$$

In particular, the rightmost inequality leads to the so-called *standard model* of floating-point arithmetic [1, p. 40]. For some of our proofs we will also use the following classical modified version of (9): for any real number  $t$ ,

$$|\text{RN}(t) - t| \leq \frac{1}{2} \text{ulp}(\text{RN}(t)) \leq u |\text{RN}(t)|.$$

## II. ON THE SUM OF TWO NONNEGATIVE PRODUCTS

In this section we focus on the approximation of  $g = ac + bd$  by  $\widehat{f}_o$ ,  $\widehat{f}$ , or  $\widehat{g}$  in the case where

$$ac \geq 0 \quad \text{and} \quad bd \geq 0. \quad (10)$$

From [6] and [7] we already know that  $\text{ulp}(g)$  and  $2u$  are sharp absolute/relative error bounds for  $\widehat{g}$ , and that  $2u$  is a relative error bound for  $\widehat{f}_o$ . As a first extension of these results, we show in Section II-A that this bound  $2u$  on  $|\widehat{f}_o - g|/g$  is sharp, and also that a sharp bound on  $|\widehat{f}_o - g|$  is  $\frac{5}{4} \text{ulp}(g)$ . Then, we show in Section II-B that similarly to  $\widehat{g}$ , sharp absolute and relative error bounds for  $\widehat{f}$  are  $\text{ulp}(g)$  and  $2u$ . Since  $\widehat{f}$  and  $\widehat{g}$  share the same error bounds, we study further how they relate to each other. Specifically, we prove in Section II-C that either  $\widehat{f}$  and  $\widehat{g}$  are equal, or one of them is a faithfully rounded value of  $g$ . We conclude in Section II-D by assuming further that

$$\text{ulp}(bd) \leq \text{ulp}(ac). \quad (11)$$

In this case, which may be of interest when tests are allowed and which is implied by  $0 \leq bd \leq ac$ , we prove that  $\widehat{f}$  and  $\widehat{g}$  always coincide, thus showing that the compensation step in Kahan's algorithm simply has no effect on the quality of the result. Furthermore, the error bounds  $\text{ulp}(g)$  and  $2u$  given in the general case can now be improved to  $\frac{3}{4} \text{ulp}(g)$  and  $\frac{3}{2}u$ , which are shown to be sharp when (11) holds.

The table below summarizes the bounds presented in this section:

$\widehat{r}$	bound on $ \widehat{r} - g $	bound on $ \widehat{r} - g /g$
$\widehat{f}_o$	$\frac{5}{4} \text{ulp}(g)$	$2u$ [7]
$\widehat{f}$	$\text{ulp}(g)$	$2u$
$\widehat{g}$	$\text{ulp}(g)$ [6]	$2u$ [6]
$\widehat{f} = \widehat{g}$ for (11)	$\frac{3}{4} \text{ulp}(g)$	$\frac{3}{2}u$

The sharpness of these error bounds is established using examples of the form  $ac + bd$  parametrized by the precision  $p$ . We also investigated the sharpness of these error bounds for the evaluation of sums of squares  $c^2 + d^2$ : in most cases we show that the bounds remain asymptotically optimal as  $p \rightarrow \infty$ , but for some of the algorithms considered we have not found examples parametrized by  $p$  to show this. In those cases, however, we illustrate the quality of the error bounds with examples in precisions  $p = 53, 113$ , which correspond to the binary64 and binary128 formats of the IEEE 754-2008 standard [4].

### A. Error bounds for the approximation $\widehat{f}_o$ without FMA

From [7] it is known that when approximating  $g$  as in (10) using  $\widehat{f}_o$  as in (4), the relative error is always at most of  $2u$ :

$$|\widehat{f}_o - g| \leq 2ug.$$

Since here  $g = |g|$ , this bound thus is the same as the bound obtained in (6) using Kahan's algorithm. Furthermore, just like for Kahan's algorithm, it is asymptotically optimal, as the example below shows.

**Example 1** ( $2ug$  is an asymptotically optimal bound on  $|\widehat{f}_o - g|$ ). For any  $p \geq 2$ , consider

$$a = b = 2^p - 1, \quad c = 2^{p-1} + 2, \quad d = 2^{p-1} + 1.$$

One easily shows that  $g = 2^{2p} + 2^{p+1} - 3$  and  $\widehat{f}_o = 2^{2p}$ , so that we have

$$\frac{|\widehat{f}_o - g|}{g} = \frac{2^{2p} - 3 \cdot 2^{p-1}}{2^{2p} + 2^{p+1} - 3} \cdot 2u.$$

Thus,  $|\widehat{f}_o - g|/(2ug)$  is in  $1 - O(2^{-p})$  as  $p \rightarrow \infty$ .

Let us now give an error bound in ulps for  $\widehat{f}_o$  and an example showing the asymptotic optimality of this bound.

**Property 1.**  $|\widehat{f}_o - g| \leq \frac{5}{4} \text{ulp}(g)$ .

*Proof:* Let  $\widehat{v} = \text{RN}(ac)$ ,  $\widehat{w} = \text{RN}(bd)$ , and  $f_o = \widehat{v} + \widehat{w}$ . Then  $|\widehat{f}_o - g| \leq |\widehat{f}_o - f_o| + |f_o - g|$  and we shall prove first that  $|f_o - g|$  is at most  $\frac{3}{4} \text{ulp}(g)$ . With  $m = \min\{ac, bd\}$  and  $m' = \max\{ac, bd\}$ , we have

$$\begin{aligned} |f_o - g| &\leq |\widehat{v} - ac| + |\widehat{w} - bd| \\ &\leq \frac{1}{2} \text{ulp}(ac) + \frac{1}{2} \text{ulp}(bd) = \frac{1}{2} \text{ulp}(m) + \frac{1}{2} \text{ulp}(m'). \end{aligned}$$

From  $0 \leq m \leq \frac{g}{2}$  and  $0 \leq m' \leq g$  we get  $\text{ulp}(m) \leq \frac{1}{2} \text{ulp}(g)$  and  $\text{ulp}(m') \leq \text{ulp}(g)$ , and then

$$|f_o - g| \leq \frac{3}{4} \text{ulp}(g). \quad (12)$$

We conclude the proof by using the same case analysis as in [7, p. 1470]:

- If  $\text{ulp}(f_o) \leq \text{ulp}(g)$  then  $|\widehat{f}_o - f_o| \leq \frac{1}{2} \text{ulp}(f_o) \leq \frac{1}{2} \text{ulp}(g)$ , which can be combined with (12) to give the result;
- If  $\text{ulp}(g) < \text{ulp}(f_o)$  then  $0 \leq g < 2^k \leq f_o$  for some integer  $k$ . Hence, using (12) and  $\text{ulp}(g) \leq \frac{1}{2} \text{ulp}(f_o)$ ,

$$0 \leq f_o - 2^k < f_o - g \leq \frac{3}{4} \text{ulp}(g) \leq \frac{3}{8} \text{ulp}(f_o).$$

Consequently, the floating-point number  $2^k$  must be equal to  $\widehat{f}_o = \text{RN}(f_o)$ , and then  $0 < \widehat{f}_o - g \leq f_o - g \leq \frac{3}{4} \text{ulp}(g)$ . Thus, in this case  $|\widehat{f}_o - g|$  is also at most  $\frac{5}{4} \text{ulp}(g)$ . ■

**Example 2** ( $\frac{5}{4} \text{ulp}(g)$  is an asymptotically optimal bound on  $|\widehat{f}_o - g|$ ). Assuming  $p \geq 7$ , let

$$a = c = 2^{p-1} + 2^{p-3} + 1, \quad b = d = 2^{p-1} + 2^{p-2} + 1.$$

One easily shows that  $g = 61 \cdot 2^{2p-6} + 2^{p+1} + 3 \cdot 2^{p-2} + 2$ ,  $\widehat{f}_o = 61 \cdot 2^{2p-6} + 2^{p+2}$ , and  $\text{ulp}(g) = 2^p$ , so that

$$|\widehat{f}_o - g| = \left(\frac{5}{4} - 2^{1-p}\right) \text{ulp}(g).$$

**Special case of the sum of two squares.** When computing an expression of the form  $c^2 + d^2$  using (4), Example 2 shows that the absolute error bound  $\frac{5}{4}\text{ulp}(g)$  is asymptotically optimal. Still in the particular case of an expression of the form  $c^2 + d^2$ , Example 3 below shows that the relative bound  $2u$  for  $\hat{f}_o$  remains asymptotically optimal when  $p$  is even. When  $p$  is odd, we did not manage to find an example parametrized by  $p$  to prove that the error bound  $2u$  is asymptotically optimal; nevertheless, when  $p$  is odd, the quality of the bound is illustrated in Table I for  $p = 53, 113$ .

**Example 3.** Assuming  $p \geq 2$  is even, let us consider

$$c = 2^{p/2-1} \quad \text{and} \quad d = 2^{p-1} + 2^{p/2-1}.$$

We have  $g = 2^{2p-2} + 2^{\frac{3p}{2}-1} + 2^{p-1}$  and  $\hat{f}_o = 2^{2p-2} + 2^{\frac{3p}{2}-1}$ , so that

$$\frac{|\hat{f}_o - g|}{g} = \frac{2^{2p-2}}{2^{2p-2} + 2^{\frac{3p}{2}-1} + 2^{p-1}} \cdot 2u.$$

Hence  $|\hat{f}_o - g|/(2ug)$  is in  $1 - O(2^{-\frac{p}{2}})$  as  $p \rightarrow \infty$ .

- Relative error in  $\hat{f}_o$  close to  $2u$ :

$p$	example
53	$c = 4503608217436160$ $d = 4503599711256576$ $ \hat{f}_o - g /(ug) = 1.9374 \dots$
113	$c = 5192296858544272361496373058600960$ $d = 5192296858534827718602488876630016$ $ \hat{f}_o - g /(ug) = 1.9374 \dots$

- Relative error in  $\hat{f} = \hat{g}$  close to  $2u$ :

$p$	example
53	$c = 8426657115275263$ $d = 4503608217436160 \cdot 2^{26}$ $ \hat{f} - g /(ug) = 1.9980 \dots$
113	$c = 9715274200149150133070733366001663$ $d = 5192296858544272361496373058600960 \cdot 2^{56}$ $ \hat{f} - g /(ug) = 1.9980 \dots$

- For  $|d| \leq |c|$ , errors in  $\hat{f} = \hat{g}$  close to  $\frac{3}{4}\text{ulp}(g)$  and  $\frac{3}{2}u$ :

$p$	example
53	$c = 4503600164241409$ $d = 4503599900000256$ $ \hat{f} - g /\text{ulp}(g) = 0.7480 \dots$ $ \hat{f} - g /(ug) = 1.4960 \dots$
113	$c = 5192296858534828204991248632643585$ $d = 5192296858534827921264472108302336$ $ \hat{f} - g /\text{ulp}(g) = 0.7480 \dots$ $ \hat{f} - g /(ug) = 1.4960 \dots$

TABLE I  
EXAMPLES OF THE FORM  $c^2 + d^2$ , FOR  $p = 53, 113$ .

### B. Error bounds for the naive approximation $\hat{f}$ with an FMA

We now turn to the computation of  $\hat{f}$  as in (5). The property and example below show that the relative error bound  $2u$

obtained for  $\hat{f}_o$  in Section II-A is still sharp, while the absolute error is now tightly bounded by one ulp of the exact result.

**Property 2.**  $|\hat{f} - g| \leq \text{ulp}(g) \leq 2ug$ .

*Proof:* Since  $g = f + e$ , we have  $|\hat{f} - g| \leq |\hat{f} - f| + |e|$ . By definition,  $|e| \leq \frac{1}{2}\text{ulp}(bd)$  and using  $0 \leq bd \leq g$  leads to

$$|e| \leq \frac{1}{2}\text{ulp}(g). \quad (13)$$

On the other hand, we can check that  $|\hat{f} - f| \leq \frac{1}{2}\text{ulp}(g)$  as follows:

- If  $\text{ulp}(f) \leq \text{ulp}(g)$  then  $|\hat{f} - f| \leq \frac{1}{2}\text{ulp}(f) \leq \frac{1}{2}\text{ulp}(g)$ ;
- If  $\text{ulp}(g) < \text{ulp}(f)$  then  $0 \leq g < 2^k \leq f$  for some integer  $k$ . Since  $\hat{f} = \text{RN}(f)$  and  $2^k$  is a floating-point number, this implies  $|\hat{f} - f| \leq f - 2^k < f - g = -e \leq \frac{1}{2}\text{ulp}(g)$ .

Thus, in both cases we have  $|\hat{f} - g| \leq \text{ulp}(g)$  as wanted, and (9) implies further that  $\text{ulp}(g) \leq 2ug$ . ■

**Example 4** ( $\text{ulp}(g)$  and  $2ug$  are asymptotically optimal bounds on  $|\hat{f} - g|$ ; Example 6.2 in [6]). For any  $p \geq 2$ , consider

$$a = 2^p - 2, \quad b = (2^p - 1)2^p, \quad c = d = 2^{p-1} + 1.$$

One easily shows that  $g = 2^{3p-1} + 2^{2p} - 2^p - 2$ ,  $\text{ulp}(g) = 2^{2p}$ , and  $\hat{f} = 2^{3p-1}$ . Hence  $|\hat{f} - g| = (1 - 2^{-p} - 2^{1-2p})\text{ulp}(g)$  and

$$\frac{|\hat{f} - g|}{g} = \frac{2^{3p-1} - 2^{2p} - 2^p}{2^{3p-1} + 2^{2p} - 2^p - 2} \cdot 2u.$$

Consequently, both  $|\hat{f} - g|/\text{ulp}(g)$  and  $|\hat{f} - g|/(2ug)$  are in  $1 - O(2^{-p})$  as  $p \rightarrow \infty$ . In this example, since  $e = 2^{2p-1} - 2^p$  and  $\text{ulp}(\hat{f}) = 2^{2p}$ , it may also be noticed that the result  $\hat{g}$  returned by Kahan's algorithm satisfies  $\hat{g} = \hat{f}$ .

**Special case of the sum of two squares.** Examples 6.3 and 6.4 in [6] are of the form  $c^2 + d^2$ , with  $\hat{f} = \hat{g}$ . Hence we already know from that paper that:

- the absolute error bound  $|\hat{f} - g| \leq \text{ulp}(g)$  is asymptotically optimal for  $p$  even, and optimal for  $p$  odd;
- the relative error bound  $|\hat{f} - g| \leq 2ug$  is asymptotically optimal when  $p$  is even.

When  $p$  is odd, numerical examples of the form  $c^2 + d^2$  in precisions 53 and 113 are given in Table I that illustrate the quality of the relative error bound.

### C. Relationship between $\hat{f}$ and $\hat{g}$

The sharp bounds we have so far for  $\hat{f}$  and  $\hat{g}$  are the same: one ulp of  $g$  for absolute errors, and  $2u$  for relative errors. In this section, we go one step further by investigating more precisely how  $\hat{f}$  and  $\hat{g}$  are related. We show in Theorem 1 below that if  $\hat{f}$  and  $\hat{g}$  differ then they must be two consecutive floating-point numbers surrounding the exact value  $g$ . For this, we start with a lemma giving a necessary condition for  $\hat{f}$  to be different from  $\hat{g}$ .

**Lemma 1.** If  $\hat{f} \neq \hat{g}$  then

$$2^k \leq \hat{w} < \hat{f} < 2^{k+1} \text{ for some } k \in \mathbb{Z}, \quad (14a)$$

and

$$|e| = \frac{1}{2}\text{ulp}(\widehat{w}) = \frac{1}{2}\text{ulp}(\widehat{f}). \quad (14b)$$

*Proof:* Using  $0 \leq \widehat{w} \leq ac + \widehat{w}$  together with the monotonicity of rounding, we have

$$0 \leq \widehat{w} \leq \widehat{f}.$$

If  $\widehat{f}$  is zero, this implies  $\widehat{w} = 0$  and thus  $bd = e = 0$ , so that  $\widehat{g} = 0 = \widehat{f}$ . Hence  $\widehat{f} \neq \widehat{g}$  implies  $\widehat{f}$  is positive and, therefore,  $2^k \leq \widehat{f} < 2^{k+1}$  for some integer  $k$ .

(i) If  $\widehat{w} < 2^k$  then  $\text{ulp}(\widehat{w}) \leq \frac{1}{2}\text{ulp}(\widehat{f})$  and then  $|e| \leq \frac{1}{2}\text{ulp}(\widehat{w}) \leq \frac{1}{4}\text{ulp}(\widehat{f})$ , so that  $\widehat{g} = \text{RN}(\widehat{f} + e)$  equals  $\widehat{f}$ .

(ii) If  $\widehat{w} = \widehat{f}$  then  $\widehat{g} = \text{RN}(\widehat{w} + e) = \text{RN}(bd) = \widehat{w}$  and thus  $\widehat{f} = \widehat{g}$ .

From (i) and (ii) we deduce that  $\widehat{f} \neq \widehat{g}$  implies the condition in (14a). This condition gives in particular  $\text{ulp}(\widehat{w}) = \text{ulp}(\widehat{f})$  and, since  $\widehat{f}$  is a floating-point number,

$$2^k + \text{ulp}(\widehat{f}) \leq \widehat{f} \leq 2^{k+1} - \text{ulp}(\widehat{f}). \quad (15)$$

Furthermore, by definition  $|e| \leq \frac{1}{2}\text{ulp}(\widehat{w})$  and, on the other hand,  $|e| < \frac{1}{2}\text{ulp}(\widehat{w})$  implies  $\text{RN}(\widehat{f} + e) = \widehat{f}$ . Thus,  $\widehat{f} \neq \widehat{g}$  also implies the condition in (14b). ■

**Theorem 1.** *Let  $m = \min(\widehat{f}, \widehat{g})$  and  $m' = \max(\widehat{f}, \widehat{g})$ . If  $\widehat{f} \neq \widehat{g}$  then  $g \in [m, m']$  and  $m' = m + \text{ulp}(m)$ .*

*Proof:* First, by Lemma 1 we have  $|e| = \frac{1}{2}\text{ulp}(\widehat{f})$ . Hence, if  $e \geq 0$  then  $\widehat{g} = \text{RN}(\widehat{f} + \frac{1}{2}\text{ulp}(\widehat{f}))$  must be either  $\widehat{f}$  or  $\widehat{f} + \text{ulp}(\widehat{f})$ , while if  $e < 0$  it must be either  $\widehat{f} - \text{ulp}(\widehat{f})$  or  $\widehat{f}$ . Since  $\widehat{f} \neq \widehat{g}$ , we deduce that in both cases

$$\widehat{g} = \widehat{f} + 2e. \quad (16)$$

Note also that  $|\widehat{f} - \widehat{g}| = \text{ulp}(m)$ . Indeed, if  $m = \widehat{f}$  then this follows immediately from (16); else, we have  $\widehat{g} < \widehat{f}$  and, since  $\widehat{g}$  is a floating-point number, using (15) gives  $\text{ulp}(\widehat{g}) = \text{ulp}(\widehat{f})$ . Together with (16), this implies  $|\widehat{f} - \widehat{g}| = \text{ulp}(\widehat{g}) = \text{ulp}(m)$ .

Writing  $\epsilon$  for the error committed when rounding  $ac + \widehat{w}$  to nearest, we have  $\widehat{f} = ac + \widehat{w} + \epsilon$  and  $|e| \leq \frac{1}{2}\text{ulp}(\widehat{f}) = |e|$ . Hence, using the identity  $e = bd - \widehat{w}$ ,

$$g = \widehat{f} + e - \epsilon.$$

If  $e > 0$  then  $-\epsilon \in [-e, e]$ , so that  $g$  belongs to  $[\widehat{f}, \widehat{f} + 2e] = [\widehat{f}, \widehat{g}]$ ; if  $e \leq 0$  then  $-\epsilon \in [e, -e]$  and in this case  $g$  belongs to  $[\widehat{f} + 2e, \widehat{f}] = [\widehat{g}, \widehat{f}]$ . We conclude that  $g$  ranges in  $[m, m']$  with  $m' - m = |\widehat{f} - \widehat{g}| = \text{ulp}(m)$ . ■

Recall from [8, Definition 2.3] that a *faithful rounding* of a real number  $t$  is any floating-point number  $s$  whose floating-point predecessor and successor satisfy the strict inequalities  $\text{pred}(s) < t < \text{succ}(s)$ . Consequently, Theorem 1 shows in particular that if  $\widehat{f}$  and  $\widehat{g}$  differ then at least one of them must be a faithful rounding of the exact result  $g$ .

Since  $|\widehat{f} - \widehat{g}| = \text{ulp}(m)$  and  $0 \leq m \leq g$ , another consequence of the theorem above is that

$$|\widehat{f} - \widehat{g}| \leq \text{ulp}(g).$$

This bound improves upon the bound  $|\widehat{f} - \widehat{g}| \leq 2\text{ulp}(g)$  immediately implied by  $|\widehat{f} - g| \leq \text{ulp}(g)$  and  $|\widehat{g} - g| \leq \text{ulp}(g)$ ; furthermore, it is shown to be optimal in the example below.

**Example 5** ( $\text{ulp}(g)$  is an optimal bound on  $|\widehat{f} - \widehat{g}|$ ). *Assuming  $p \geq 4$ , let us consider*

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

*Then one can check that  $|\widehat{f} - \widehat{g}| = \text{ulp}(g)$ . It may also be checked that in this example the result computed by the naive algorithm with FMA is always more accurate than the one computed with Kahan's algorithm, since we have*

$$|\widehat{f} - g| = 2^{-p-1}\text{ulp}(g) \quad \text{and} \quad |\widehat{g} - g| = (1 - 2^{-p-1})\text{ulp}(g).$$

*D. Special case where  $\text{ulp}(bd) \leq \text{ulp}(ac)$*

The goal of this section is to show that if our input  $a, b, c, d$  satisfy (10) and (11), then  $\widehat{f}$  is always equal to  $\widehat{g}$  and the error bounds  $\text{ulp}(g)$  and  $2u$  can be decreased further.

**Lemma 2.** *If  $\text{ulp}(bd) \leq \text{ulp}(ac)$  then  $|e| \leq \frac{1}{4}\text{ulp}(g)$ .*

*Proof:* Since  $|e| \leq \frac{1}{2}\text{ulp}(bd)$ , it suffices to show that  $\text{ulp}(bd) \leq \frac{1}{2}\text{ulp}(g)$ . We do this by considering two cases.

If  $bd \leq ac$  then  $bd \leq g/2$  and thus  $\text{ulp}(bd) \leq \frac{1}{2}\text{ulp}(g)$ .

Assume now that  $ac < bd$ . This implies  $\text{ulp}(ac) \leq \text{ulp}(bd)$  as well as  $g > 0$  and  $ac < g/2$ . Consequently,  $\text{ulp}(bd) = \text{ulp}(ac) \leq 2^{1-p}ac < 2^{-p}g < \text{ulp}(g)$ , and the conclusion follows from  $\text{ulp}(bd) < \text{ulp}(g)$  being equivalent to  $\text{ulp}(bd) \leq \frac{1}{2}\text{ulp}(g)$ . ■

**Theorem 2.** *If  $\text{ulp}(bd) \leq \text{ulp}(ac)$  then  $\widehat{f} = \widehat{g}$ .*

*Proof:* If  $\widehat{f} = 0$  then it has already been seen in the proof of Lemma 1 that  $\widehat{g} = 0$  as well. Let us now assume that  $\widehat{f} > 0$ . Since  $\widehat{f} = \text{RN}(f)$  is a floating-point number, we have  $\widehat{f} \leq (2^p - 1)\text{ulp}(\widehat{f})$  and it follows that for RN 'to nearest even'

$$f < (2^p - \frac{1}{2})\text{ulp}(\widehat{f}).$$

On the other hand,  $0 \leq \widehat{w} \leq f = ac + \widehat{w}$ , which implies  $\text{ulp}(\widehat{w}) \leq \text{ulp}(f) \leq \text{ulp}(\widehat{f})$  and then

$$|e| \leq \frac{1}{2}\text{ulp}(\widehat{f}).$$

Hence  $g = f + e \leq f + |e| < 2^p\text{ulp}(\widehat{f})$  and, using  $g \geq 0$ ,

$$\text{ulp}(g) \leq \text{ulp}(\widehat{f}).$$

By Lemma 2 we thus obtain  $|e| \leq \frac{1}{4}\text{ulp}(\widehat{f})$ , from which we conclude that  $\widehat{g} = \text{RN}(\widehat{f} + e) = \widehat{f}$ . ■

The condition  $\text{ulp}(bd) \leq \text{ulp}(ac)$  is sufficient but *not* necessary to have  $\widehat{f}$  equal to  $\widehat{g}$ ; see Example 4. In addition, this condition is equivalent to

$$\text{exponent of } bd \leq \text{exponent of } ac,$$

and, for  $ac$  and  $bd$  nonnegative, it is implied by

$$bd \leq ac.$$

Furthermore, in this special case where both (10) and (11) hold, Property 3 below shows that we can also improve the bounds we had given for the general case where only (10) holds. To get this, all we have to do is to replace, in the proof of Property 2, the bound in (13) by the bound in Lemma 2.

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

The next example shows that these new bounds are sharp.

**Example 6** (In the special case where  $\text{ulp}(bd) \leq \text{ulp}(ac)$  the bound  $\frac{3}{4}\text{ulp}(g)$  is optimal and the bound  $\frac{3}{2}ug$  is asymptotically optimal). *Assuming  $p \geq 6$ , consider*

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

*In this example, we have  $\text{ulp}(ac) = \text{ulp}(bd)$ ,  $|\hat{f} - g|$  equals  $\frac{3}{4}\text{ulp}(g)$ , and  $|\hat{f} - g|/(\frac{3}{2}ug)$  is in  $1 - O(2^{-\frac{p}{2}})$  as  $p \rightarrow \infty$ .*

**Special case of the sum of two squares.** Example 7 below shows that the error bounds of Property 3 remain asymptotically optimal when evaluating an expression of the form  $c^2 + d^2$  with  $p$  even. For  $p$  odd, we give examples for  $p = 53, 113$  in Table I showing that our bounds are reasonably sharp.

**Example 7.** *Assuming  $p \geq 6$  even, consider*

$$c = 2^{p-1} + 2^{p/2+1} \quad \text{and} \quad d = 2^{p-1} + 2^{p/2-1} + 1.$$

*With these inputs,  $0 \leq d \leq c$  and  $\text{ulp}(d^2) = \text{ulp}(c^2)$ . Also, both  $|\hat{f} - g|/(\frac{3}{4}\text{ulp}(g))$  and  $|\hat{f} - g|/(\frac{3}{2}ug)$  are in  $1 - O(2^{-\frac{p}{2}})$  as  $p \rightarrow \infty$ .*

### III. TWO COMPLEX DIVISION ALGORITHMS AND THEIR COMPONENTWISE RELATIVE ACCURACY

When the denominator  $c^2 + d^2$  in (1) is evaluated as  $\hat{\delta} = \text{RN}(c^2 + \text{RN}(d^2))$  with no comparison between  $|c|$  and  $|d|$ , Property 2 shows that  $\hat{\delta} = (c^2 + d^2)(1 + \epsilon_1)$  with  $|\epsilon_1| \leq 2u$ . On the other hand, Property 3 suggests that one should choose to compare  $|c|$  and  $|d|$ , and to compute  $\hat{\delta} = \text{RN}(c^2 + \text{RN}(d^2))$  if  $|d| \leq |c|$  and  $\hat{\delta} = \text{RN}(d^2 + \text{RN}(c^2))$  otherwise: this ensures the smaller bound  $|\epsilon_1| \leq \frac{3}{2}u$ .

It is well known that conditional branches can limit the performance of programs on pipelined architectures, essentially because each wrong branch prediction causes the pipeline to stall (see Chapters 3 and 4, and Appendix A in [9]). Hence, using a conditional branch to choose how to compute  $\hat{\delta}$  depending on a comparison between  $|d|$  and  $|c|$  may have a bad impact on the performance of the program.

However, some architectures support predicated instructions [9, p. 340], [5, Chap. 2]. Each predicated instruction in a program refers to a condition: if the condition is true then the instruction is executed normally, otherwise it is executed as a 'no-operation' causing no particular hazard in the pipeline. The IEEE 754-2008 standard [4] also describes new operations with two floating-point numbers as operands:

- `minNumMag`, which returns the one of smallest magnitude (or, in the case of equal magnitudes, the minimum);

- `maxNumMag`, which returns the one of largest magnitude (or, in the case of equal magnitudes, the maximum).

The operations `minNumMag` and `maxNumMag` can be used to sort two floating-point numbers by order of magnitude: they are already available for example on the Itanium processor through the instructions `famin` and `famax` [5, p. 291]. Either predicated instructions, or instructions implementing the `minNumMag` and `maxNumMag` operations, could be used in our case to benefit from Property 3 while avoiding conditional branches.

In this section, we consider two complex division algorithms. We first analyze a straight-line algorithm, with a componentwise relative error bound of  $5u + O(u^2)$ ; this bound is proved to be asymptotically optimal when  $p$  is even, and its quality is illustrated with numerical examples when  $p = 53, 113$ . We next describe a second division algorithm in which a comparison between  $|c|$  and  $|d|$  is used to ensure a componentwise relative error bound of  $\frac{9}{2}u + O(u^2)$ , and whose quality is again illustrated with numerical examples for all the basic binary IEEE formats.

#### A. Straight-line version: no tests allowed

We assume here that we can neither compare the magnitudes of  $c$  and  $d$ , nor determine which one of the two numerators  $ac + bd$  and  $bc - ad$  is a sum of two products of the same sign. Hence, we shall compute the numerators using Kahan's algorithm and the denominator using (5).

**algorithm** `CompDivS`( $a + ib, c + id$ )

```

 $\hat{\delta} := \text{RN}(c^2 + \text{RN}(d^2));$ 
 $\hat{g}_{\text{re}} := \text{Kahan}(a, b, c, d);$  // evaluates  $ac + bd$ 
 $\hat{g}_{\text{im}} := \text{Kahan}(b, -a, c, d);$  // evaluates  $bc - ad$ 
 $\hat{z}_{\text{re}} := \text{RN}(\hat{g}_{\text{re}}/\hat{\delta});$   $\hat{z}_{\text{im}} := \text{RN}(\hat{g}_{\text{im}}/\hat{\delta});$ 
return  $\hat{z}_{\text{re}} + i\hat{z}_{\text{im}};$ 

```

For reasons of symmetry, for bounding the componentwise relative error of that algorithm it suffices to consider the relative error in the real part of the quotient. Using Property 2 together with (6) and (9) leads to

$$\hat{z}_{\text{re}} = \frac{ac + bd}{c^2 + d^2} \cdot \frac{(1 + \epsilon_2)}{(1 + \epsilon_1)} (1 + \epsilon_3),$$

with  $|\epsilon_1|, |\epsilon_2| \leq 2u$  and  $|\epsilon_3| \leq u$ . Hence  $\hat{z}_{\text{re}} = (1 + \epsilon) \cdot \text{Re}(z)$  with  $|\epsilon|$  bounded by  $\frac{(1+2u)}{(1-2u)}(1+u) - 1$ , which is less than  $5u + 13u^2$  if and only if  $p \geq 5$ .

**Property 4.** *For  $p \geq 5$  and in the absence of underflow and overflow, the quotient  $\hat{z}$  computed by `CompDivS` satisfies*

$$E_c(\hat{z}) \leq 5u + 13u^2.$$

The following example shows that, when  $p$  is even, the bound given by Property 4 is asymptotically optimal: here, the relative error in the real part of the computed quotient is  $5u - O(u^{3/2})$ . Since this example is more involved than the previous ones, we give a detailed analysis.

**Example 8** (For  $p$  even, the bound in Property (4) is asymptotically optimal). *Let*

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

and assume  $p$  is even. Then, the quotient  $\hat{z}$  computed by *CompDivS* satisfies

$$\frac{|Re \hat{z} - Re z|}{|Re z|} = 5u - O(u^{3/2}). \quad (17)$$

**Proof of (17).** Let us define  $R = 2^{p/2}$ . We have

$$\begin{aligned} Re z &= -\frac{2R^3 + 5R^2 - 4R}{2R^6 + 4R^5 + 4R^4 - 8R^2 + 8} \\ &= -\frac{1}{R^3} - \frac{1}{2R^4} + \frac{5}{R^5} + O\left(\frac{1}{R^6}\right). \end{aligned}$$

Elementary manipulations show that the value of  $\hat{\delta}$  in Algorithm *CompDivS* is  $R^6 + 2R^5$  and that the value of  $\hat{g}_{re}$  is  $-R^3 - \frac{5}{2}R^2$ . Hence

$$\frac{\hat{g}_{re}}{\hat{\delta}} = \frac{-R^3 - \frac{5}{2}R^2}{R^6 + 2R^5} = -\frac{1}{R^3} - \frac{1}{2R^4} + \frac{1}{R^5} + O\left(\frac{1}{R^6}\right),$$

and we will show that  $\hat{z}_{re} = \tau$ , where  $\tau$  is the floating-point number defined by

$$\tau := -\frac{1}{R^3} - \frac{1}{2R^4} = -2^{-3p/2} - 2^{2p-1}.$$

Since  $\hat{z}_{re} = RN(\hat{g}_{re}/\hat{\delta})$ , it suffices to show that

$$\tau \leq \frac{\hat{g}_{re}}{\hat{\delta}} < \tau + \frac{1}{2}ulp(\tau),$$

and since  $ulp(\tau) = 2^{-5p/2+1} = 2/R^5$ , this is equivalent to

$$\underbrace{-\frac{1}{R^3} - \frac{1}{2R^4}}_{Left(R)} \leq \underbrace{\frac{-R^3 - \frac{5}{2}R^2}{R^6 + 2R^5}}_{Med(R)} < \underbrace{-\frac{1}{R^3} - \frac{1}{2R^4} + \frac{1}{R^5}}_{Right(R)}. \quad (18)$$

An elementary analysis of the functions  $Med(R) - Left(R)$  and  $Right(R) - Med(R)$  allows us to check that (18) is satisfied when  $R > 0$ , which is true by definition of  $R$ . Hence we deduce that *CompDivS* returns  $\hat{z}_{re} = \tau$ . It follows that

$$\begin{aligned} \left| \frac{\hat{z}_{re}}{Re z} - 1 \right| &= \frac{10R^5 + 2R^4 - 8R^3 - 4R^2 + 8R + 4}{2R^7 + 5R^6 - 4R^5} \\ &= \frac{5}{R^2} - \frac{23}{2R^3} + \frac{139}{4R^4} + O\left(\frac{1}{R^5}\right), \end{aligned}$$

which leads to (17).  $\blacksquare$

For example, when  $p = 24$ , which corresponds to the binary32 IEEE format, the value in (17) is about 4.997. When  $p$  is odd, we have not yet constructed an example parametrized by  $p$ . Nevertheless, the examples given in Table II already illustrate well the quality of the bound in Property 4 for other formats of practical interest like binary64 ( $p = 53$ ) and binary128 ( $p = 113$ ).

$p$	example
53	$a = 2^{52} + 1$ $b = -142398041 = -\left\lceil \frac{2^{105} + 2^{104}}{d} \right\rceil - 38644$ $c = 2^{52}$ $d = 94906267 \cdot 2^{52} = (1 + \lceil 2^{53}/2 \rceil) \cdot 2^{52}$ $ \hat{z}_{re} - Re z /(u Re z ) = 4.9987\dots$
113	$a = 2^{112} + 1$ $b = -152857240142482713 = -\left\lceil \frac{2^{225} + 2^{224}}{d} \right\rceil - 1864174$ $c = 2^{112}$ $d = 101904826760412363 \cdot 2^{112} = (1 + \lceil 2^{113}/2 \rceil) \cdot 2^{112}$ $ \hat{z}_{re} - Re z /(u Re z ) = 4.9999\dots$

TABLE II  
RELATIVE ERROR IN  $\hat{z}_{re}$  COMPUTED USING *COMPDIVS* CLOSE TO THE UPPER BOUND  $5u + 13u^2$ .

*B. If tests are allowed*

If tests are allowed, then one can compare  $|c|$  and  $|d|$  to select between  $RN(c^2 + RN(d^2))$  and  $RN(d^2 + RN(c^2))$  the computation that ensures the smallest relative error bound according to Property 3, and still use Kahan's algorithm for the evaluation of both numerators  $ac + bd$  and  $bc - ad$ . This gives the following complex division algorithm:

```

algorithm CompDivT( $a + ib, c + id$ )
  if  $|d| \leq |c|$  then  $\hat{\delta} := RN(c^2 + RN(d^2))$ ;
  else  $\hat{\delta} := RN(d^2 + RN(c^2))$ ;
   $\hat{g}_{re} := \text{Kahan}(a, b, c, d)$ ; // evaluates  $ac + bd$ 
   $\hat{g}_{im} := \text{Kahan}(b, -a, c, d)$ ; // evaluates  $bc - ad$ 
   $\hat{z}_{re} := RN(\hat{g}_{re}/\hat{\delta})$ ;  $\hat{z}_{im} := RN(\hat{g}_{im}/\hat{\delta})$ ;
  return  $\hat{z}_{re} + i\hat{z}_{im}$ ;

```

When *CompDivT* is implemented on an architecture supporting the operations *minNumMag* and *maxNumMag* defined in the IEEE 754-2008 standard [4], its 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  $|\underline{d}| \leq |c|$ .
 $\hat{\delta} := RN(\underline{c}^2 + RN(\underline{d}^2))$ ;

```

Similar techniques avoiding branches can also be used more generally on architectures supporting predicated instructions [9, Chap. 4, p. 340], [5, Chap. 2].

According to Property 3 we now have  $\hat{\delta} = (c^2 + d^2)(1 + \epsilon_1)$  with  $|\epsilon_1| \leq \frac{3}{2}u$ , which leads to the following error bound.

**Property 5.** *For  $p \geq 6$  and in the absence of underflow and overflow, the quotient  $\hat{z}$  computed by *CompDivT* satisfies*

$$E_c(\hat{z}) \leq \frac{9}{2}u + 9u^2.$$

The examples given in Table III show that the bound  $\frac{9}{2}u + 9u^2$  is reasonably sharp for the basic binary IEEE formats ( $p = 24, 53, 113$ ).



$p$	example
24	$a = 8391768$ $b = -8392368$ $c = 8391504$ $d = 8390648$ $ \widehat{z}_{\text{re}} - \text{Re } z /(u \text{Re } z ) = 4.4932\dots$
53	$a = 4503599627378010$ $b = -4503599627377047$ $c = 6369051672541039$ $d = 6369051672534109$ $ \widehat{z}_{\text{re}} - \text{Re } z /(u \text{Re } z ) = 4.4421\dots$
113	$a = 7343016637207168931428032607349357$ $b = -7343016637207168931428032607349412$ $c = 7343016637207168931428032607356045$ $d = 7343016637207168931428032607355264$ $ \widehat{z}_{\text{re}} - \text{Re } z /(u \text{Re } z ) = 4.4620\dots$

TABLE III  
RELATIVE ERROR IN  $\widehat{z}_{\text{re}}$  COMPUTED USING COMPDIVT (GIVEN WITH  $|d| \leq |c|$ ) CLOSE TO THE UPPER BOUND  $\frac{9}{2}u + 9u^2$ .

Concerning the computation of the numerators  $ac + bd$  and  $bc - ad$ , Property 3 could also be used: instead of using Kahan’s algorithm for computing both numerators, one could determine which one is the sum of two products with the same sign, and then use Property 3 as in CompDivT to evaluate the corresponding numerator with a relative error bounded by  $\frac{3}{2}u$ . However, this would only improve the error bound for one of the components of the computed quotient, and the componentwise relative error bound would remain  $\frac{9}{2}u + 9u^2$  as with CompDivT: for this reason, we have not considered this algorithm here.

#### IV. CONCLUDING REMARKS AND FUTURE WORK

By combining Kahan’s algorithm with a cheaper scheme, we have introduced two complex division algorithms that take advantage of the availability of an FMA instruction. We have also given sharp componentwise relative error bounds for these algorithms. This work should be pursued further in the following three directions. First, we should implement our algorithms in order to compare their running times and measure the overhead induced by the use of Kahan’s algorithm to improve the accuracy. Second, extensions to non binary radices and other rounding modes should be proposed. Third, it would be interesting to study how our division algorithms can be combined with the scaling techniques from [10], [11], [12], [13] in order to avoid spurious under/overflows while preserving high componentwise accuracy.

#### ACKNOWLEDGMENTS

This research was partly supported by the French *Agence Nationale de la Recherche (ANR)*, under grants ANR 11 BS02 013 (HPAC project) and ANR 2010 BLAN 0203 01 (TaMaDi project).

#### REFERENCES

- [1] N. J. Higham, *Accuracy and Stability of Numerical Algorithms*, 2nd ed. Philadelphia, PA, USA: SIAM, 2002.
- [2] J.-M. Muller, N. Brisebarre, F. de Dinechin, C.-P. Jeannerod, V. Lefèvre, G. Melquiond, N. Revol, D. Stehlé, and S. Torres, *Handbook of Floating-Point Arithmetic*. Birkhäuser Boston, 2010.
- [3] M. Baudin, “Error bounds of complex arithmetic,” June 2011, available at [http://forge.scilab.org/upload/compdiv/files/complexerrorbounds\\_v0.2.pdf](http://forge.scilab.org/upload/compdiv/files/complexerrorbounds_v0.2.pdf).
- [4] IEEE Computer Society, *IEEE Standard for Floating-Point Arithmetic*. IEEE Standard 754-2008, Aug. 2008, available at <http://ieeexplore.ieee.org/servlet/opac?punumber=4610933>.
- [5] M. Cornea, J. Harrison, and P. T. P. Tang, *Scientific Computing on Itanium<sup>®</sup>-based Systems*. Intel Press, Hillsboro, OR, 2002.
- [6] C.-P. Jeannerod, N. Louvet, and J.-M. Muller, “Further analysis of Kahan’s algorithm for the accurate computation of  $2 \times 2$  determinants,” *Mathematics of Computation*, 2012, to appear. Preliminary version available at <http://hal-ens-lyon.archives-ouvertes.fr/ensl-00649347/en/>.
- [7] R. P. Brent, C. Percival, and P. Zimmermann, “Error bounds on complex floating-point multiplication,” *Mathematics of Computation*, vol. 76, pp. 1469–1481, 2007.
- [8] S. M. Rump, T. Ogita, and S. Oishi, “Accurate floating-point summation part I: Faithful rounding,” *SIAM Journal on Scientific Computing*, vol. 31, no. 1, pp. 189–224, 2008.
- [9] J. L. Hennessy and D. A. Patterson, *Computer Architecture: A Quantitative Approach*, 3rd ed. Morgan Kaufmann, 2003.
- [10] R. L. Smith, “Algorithm 116: Complex division,” *Comm. ACM*, vol. 5, no. 8, p. 435, 1962.
- [11] G. W. Stewart, “A note on complex division,” *ACM Trans. Math. Software*, vol. 11, no. 3, pp. 238–241, 1985.
- [12] D. M. Priest, “Efficient scaling for complex division,” *ACM Trans. Math. Software*, vol. 30, no. 4, Dec. 2004.
- [13] M. Baudin and R. L. Smith, “A robust complex division in Scilab,” October 2012, available at <http://arxiv.org/abs/1210.4539>.

#### APPENDIX

**Property 6** (§3.5 and §3.6 in [3]). *Given  $x = a + ib$  and  $y = c + id$  two complex numbers with floating-point coefficients, let  $\widehat{z}$  be the evaluation of  $z = x/y$  computed according to (2). In the absence of underflow and overflow,*

$$E_n(\widehat{z}) \leq \frac{(3 + \sqrt{5})u + \sqrt{5}u^2}{1 - 2u} < (3 + \sqrt{5})u + 13u^2,$$

where the latter inequality holds assuming  $p \geq 7$ .

*Proof:* Writing  $\nu = (ac + bd) + i(bc - ad)$  and  $\delta = c^2 + d^2$ , we have  $z = \nu/\delta$ . Let also  $\varphi = \widehat{\nu}/\widehat{\delta}$  with

$$\widehat{\nu} = \text{RN}(\text{RN}(ac) + \text{RN}(bd)) + i \text{RN}(\text{RN}(bc) - \text{RN}(ad))$$

and

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

Then  $\widehat{z}$  satisfies  $\widehat{z} = \text{RN}(\text{Re } \varphi) + i \text{RN}(\text{Im } \varphi)$  and we have

$$\begin{aligned} |\widehat{z} - \varphi|^2 &= (\text{RN}(\text{Re } \varphi) - \text{Re } \varphi)^2 + (\text{RN}(\text{Im } \varphi) - \text{Im } \varphi)^2 \\ &\leq u^2(\text{Re } \varphi)^2 + u^2(\text{Im } \varphi)^2. \end{aligned}$$

Hence  $|\widehat{z} - \varphi| \leq u|\varphi|$ , which implies  $\widehat{z} = \varphi(1 + \epsilon_1)$  for some  $\epsilon_1 \in \mathbb{C}$  such that  $|\epsilon_1| \leq u$ . Now, by [7] we know that  $\widehat{\nu} = \nu(1 + \epsilon_2)$  for some  $\epsilon_2 \in \mathbb{C}$  such that  $|\epsilon_2| \leq \sqrt{5}u$ , and also that  $\widehat{\delta} = \delta(1 + \epsilon_3)$  for some  $\epsilon_3 \in \mathbb{R}$  such that  $|\epsilon_3| \leq 2u$ . Therefore,  $\widehat{z} = z(1 + \epsilon)$  with  $\epsilon$  given by

$$\epsilon = \frac{(1 + \epsilon_1)(1 + \epsilon_2)}{1 + \epsilon_3} - 1.$$

The bounds on  $|\epsilon|$  then follow immediately from the bounds on the  $|\epsilon_i|$ ’s and from the fact that  $u = 2^{-p}$ . ■