



HAL
open science

On the error of Computing $ab + cd$ using Cornea, Harrison and Tang's method

Jean-Michel Muller

► **To cite this version:**

Jean-Michel Muller. On the error of Computing $ab + cd$ using Cornea, Harrison and Tang's method. 2013. ensl-00862910v1

HAL Id: ensl-00862910

<https://ens-lyon.hal.science/ensl-00862910v1>

Preprint submitted on 17 Sep 2013 (v1), last revised 24 Sep 2013 (v2)

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 error of Computing $ab + cd$ using Cornea, Harrison and Tang's method

Jean-Michel Muller

September 2013

Abstract

In their book *Scientific Computing on The Itanium* [1], Cornea, Harrison and Tang introduce an accurate algorithm for evaluating expressions of the form $ab + cd$ in binary floating-point arithmetic, assuming an FMA instruction is available. They show that if p is the precision of the floating-point (FP) format and if $u = 2^{-p}$, the relative error of the result is of order u . We improve their proof to show that the relative error is bounded by $2u + 7u^2 + 6u^3$. Furthermore, by building an example for which the relative error is asymptotically (as $p \rightarrow \infty$ or, equivalently, as $u \rightarrow 0$) equivalent to $2u$, we show that our error bound is asymptotically optimal.

1 Introduction and notation

1.1 Computing $ab + cd$

Expressions of the form $ab + cd$, where a, b, c, d are floating-point numbers arise naturally in many numerical computations. Typical examples are complex multiplication and division; discriminant of quadratic equations; cross-products and 2D determinants. The naive way of computing $ab + cd$ may lead to very inaccurate results, due to catastrophic cancellations. Several algorithms have been introduced, to overcome this problem. An algorithm attributed to Kahan by Higham [2, p. 65] can be used when an FMA (fused multiply-add) instruction is available.

Algorithm 1 Kahan's algorithm for computing $x = ab + cd$ with fused multiply-adds. $\text{RN}(t)$ means t rounded to the nearest FP number.

```
 $\hat{w} \leftarrow \text{RN}(cd)$   
 $e \leftarrow \text{RN}(cd - \hat{w})$  // this operation is exact:  $e = cd - \hat{w}$ .  
 $\hat{f} \leftarrow \text{RN}(ab + \hat{w})$   
 $\hat{x} \leftarrow \text{RN}(\hat{f} + e)$   
return  $\hat{x}$ 
```

Jeannerod, Louvet and Muller [4] show that in radix- β floating-point arithmetic, the relative error of Kahan’s algorithm is bounded by $2u$, where $u = \frac{1}{2}\beta^{1-p}$ is the unit roundoff.

Another algorithm, that also uses an FMA instruction, was introduced by Cornea, Harrison and Tang in their book *Scientific Computing on The Itanium* [1]. Cornea et al’s algorithm is

Algorithm 2 Cornea, Harrison and Tang’s algorithm for computing $x = ab + cd$ with fused multiply-adds.

```

 $p_1 \leftarrow \text{RN}(ab)$ 
 $e_1 \leftarrow ab - p_1$            // exact with an FMA
 $p_2 \leftarrow \text{RN}(cd)$ 
 $e_2 \leftarrow cd - p_2$            // exact with an FMA
 $p \leftarrow \text{RN}(p_1 + p_2)$ 
 $e \leftarrow \text{RN}(e_1 + e_2)$ 
 $s \leftarrow \text{RN}(p + e)$ 
return  $s$ 

```

Cornea, Harrison and Tang provide a quick error analysis to show that the relative error of their algorithm is of the order of u . At the time of the publication of their algorithm, the bound on Kahan’s algorithm was not known, which made their algorithm a very attractive choice, although it requires slightly more computation than Kahan’s algorithm. Now, to compare these two algorithms, we need to evaluate the largest possible relative error of Cornea et al’s algorithm more accurately. This is the purpose of this paper.

1.2 Some notation and assumptions

Throughout the paper, we assume a binary floating-point system of precision $p \geq 2$, with unbounded exponent range (that is, our results will apply to real-life computations provided that no underflow or overflow occurs). In such a system, a floating-point number is a number x that can be expressed in the form

$$x = M_x \cdot 2^{e_x - p + 1},$$

where M_x and e_x are integers, and $2^{p-1} \leq |M_x| \leq 2^p - 1$. We denote $u = 2^{-p}$. If t is a nonzero real number, with $2^k \leq t < 2^{k+1}$, we define $\text{ulp}(t)$ as 2^{k-p+1} .

We assume that an FMA instruction is available. 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 [3], one can expect that it will soon belong to the instruction set of most general-purpose processors. In the following we assume that the rounding mode is *round to nearest even*, and we denote RN the rounding function, so that the result returned when computing $\text{FMA}(a, b, c)$ is $\text{RN}(ab + c)$.

We will frequently use the following property:

for any real number t ,

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

and

$$|\text{RN}(t) - t| \leq u|\text{RN}(t)|.$$

2 Preliminary properties of Algorithm 2

Remark 2.1. *If $ab = -cd$ then $ab+cd = 0$ is exactly computed by the algorithm.*

Remark 2.2. *Let cd be the product of two binary floating-point numbers of precision p . Define $p_2 = \text{RN}(cd)$ and $e_2 = cd - p_2$. We have:*

- *either e_2 is a multiple of $2^{-p+1} \text{ulp}(p_2)$ (which implies that it fits in $p - 2$ bits);*
- *or $|cd| \leq (2^p - 2 + 2^{-p}) \text{ulp}(p_2)$.*

Proof. Since c and d are two precision- p binary floating-point numbers, one has

$$c = M_c 2^{e_c - p + 1} \quad \text{and} \quad d = M_d 2^{e_d - p + 1},$$

where M_c , M_d , e_c , and e_d are integers, with $2^{p-1} \leq |M_c|, |M_d| \leq 2^p - 1$. The number cd is a multiple of $2^{e_c + e_d - 2p + 2}$, hence $p_2 = \text{RN}(cd)$ and $e_2 = cd - p_2$ are multiple of $2^{e_c + e_d - 2p + 2}$ too.

- if $p_2 < 2^{e_c + e_d + 1}$ then $\text{ulp}(p_2) \leq 2^{e_c + e_d - p + 1}$, so that (since $\text{ulp}(p_2)$ is a power of 2) e_2 is a multiple of $2^{-p+1} \text{ulp}(p_2)$;
- if $p_2 \geq 2^{e_c + e_d + 1}$ then $\text{ulp}(p_2) = 2^{e_c + e_d - p + 2}$, therefore

$$|cd| = |M_c M_d| \cdot 2^{e_c + e_d - 2p + 2} \leq (2^p - 1)^2 \cdot 2^{e_c + e_d - 2p + 2} = (2^p - 2 + 2^{-p}) \text{ulp}(p_2).$$

□

Remark 2.3. *Denote $u = 2^{-p}$. We have,*

- $p_1 + e_1 = ab$, $|e_1| \leq u \cdot |p_1|$, and $|e_1| \leq u \cdot |ab|$;
- $p_2 + e_2 = cd$, $|e_2| \leq u \cdot |p_2|$, and $|e_2| \leq u \cdot |cd|$;
- $p = (p_1 + p_2) \cdot (1 + \epsilon_1)$, with $|\epsilon_1| \leq u$;
- $e = (e_1 + e_2) \cdot (1 + \epsilon_2)$, with $|\epsilon_2| \leq u$.

We have,

$$p + e = (ab + cd)(1 + \epsilon_1) + \gamma,$$

with

$$\gamma = (e_1 + e_2) \cdot (\epsilon_2 - \epsilon_1),$$

which implies

$$|\gamma| = 2u^2 \cdot (|ab| + |cd|).$$

3 Discussion on the various cases that occur in Algorithm 2

3.1 If ab and cd have the same sign

In that case, $|\gamma| \leq 2u^2 \cdot |ab + cd|$, so that the final relative error is bounded by $2u + 3u^2 + 2u^3$.

3.2 If ab and cd have different signs

Without loss of generality, we assume $|ab| \geq |cd|$, $ab > 0$ and $cd < 0$ (notice that if $ab = 0$ or $cd = 0$ the analysis becomes straightforward).

3.2.1 If $|cd| \leq \frac{1}{2}ab$

In that case,

$$|ab + cd| \geq \frac{1}{3}(|ab| + |cd|),$$

so that $|\gamma| \leq 6u^2 \cdot |ab + cd|$, which implies that the final relative error is bounded by $2u + 7u^2 + 6u^3$.

3.2.2 If $|cd| > \frac{1}{2}ab$

In that case, since function $t \rightarrow \text{RN}(t)$ is an increasing function, we easily find

$$\frac{1}{2}p_1 \leq |p_2| \leq p_1.$$

Applying Sterbenz Lemma, we find that $p = p_1 + p_2$ exactly, so that $\epsilon_1 = 0$, which gives

$$p + e = ab + cd + \gamma,$$

with

$$\gamma = (e_2 + e_1)\epsilon_2,$$

which implies

$$|\gamma| \leq u^2 \cdot (|ab| + |cd|).$$

1. **if $|ab + cd| \geq u \cdot (|ab| + |cd|)$** , then $|\gamma| \leq u \cdot |ab + cd|$, so that the final relative error is bounded by $2u + u^2$.
2. **if $|ab + cd| < u \cdot (|ab| + |cd|)$ and p_1 and p_2 have the same floating-point exponent e** . In that case, we have,

- $|e_1| \leq (1/2)\text{ulp}(p_1) = 2^{e-p}$,
- $|e_2| \leq (1/2)\text{ulp}(p_2) = 2^{e-p}$,
- e_1 and e_2 are multiple of 2^{e-2p+1} ,

Hence, $e_1 + e_2$ is a multiple of 2^{e-2p+1} , say $e_1 + e_2 = K \cdot 2^{e-2p+1}$, $k \in \mathbb{Z}$, that satisfies

$$|K \cdot 2^{e-2p+1}| \leq 2^{e-p+1},$$

i.e., $|K| \leq 2^p$. This implies that $e_1 + e_2$ is a floating-point number. Hence, $e = \text{RN}(e_1 + e_2) = e_1 + e_2$, so that $\epsilon_2 = 0$. As a consequence, $p + e = ab + cd$ exactly, and the final relative error is bounded by u .

3. **if $|ab + cd| < u \cdot (|ab| + |cd|)$ and p_1 and p_2 do not have the same floating-point exponent.** In such a case, $\frac{1}{2}p_1 \leq |p_2| \leq p_1$ implies that the exponent of p_2 is the exponent of p_1 minus one, so that $\text{ulp}(p_2) = \frac{1}{2}\text{ulp}(p_1)$.

Remark 3.1. Notice that we necessarily have $(p_1 + p_2) \leq 4\text{ulp}(p_2)$: p_1 and p_2 are obviously multiples of $\text{ulp}(p_2)$, and if we had $(p_1 + p_2) \leq 4\text{ulp}(p_2)$, that would imply

$$|ab + cd| = |p_1 + p_2 + e_1 + e_2| \geq 5\text{ulp}(p_2) - \text{ulp}(p_2) - \frac{1}{2}\text{ulp}(e_2) = 7/2\text{ulp}(p_2),$$

whereas

$$|ab| + |cd| < 2^p\text{ulp}(p_1) + 2^p\text{ulp}(p_2) = 3 \cdot 2^p\text{ulp}(p_2),$$

so that

$$\frac{|ab| + |cd|}{|ab + cd|} \leq \frac{6}{7} \cdot 2^p = \frac{6}{7u},$$

which contradicts the assumption $|ab + cd| < u \cdot (|ab| + |cd|)$.

The fact that p_1 and p_2 do not have the same floating-point exponent (so that there is a power of 2 between them), and that $(p_1 + p_2) \leq 4\text{ulp}(p_2)$ implies that there remain only a very few cases to examine. Define e_{p_1} as the floating-point exponent of p_1 :

- either p_1 is the floating-point number immediately above $2^{e_{p_1}}$. In such a case $-p_2$ is either $2^{e_{p_1}} - \text{ulp}(p_2)$ or $2^{e_{p_1}} - 2\text{ulp}(p_2)$;
- or $p_1 = 2^{e_{p_1}}$. In such a case, $p_2 = 2^{e_{p_1}} - i \cdot \text{ulp}(p_2)$, with $i = 1, 2, 3$, or 4 .

We can even reduce further the number of cases to be considered:

- First, one can apply Remark 2.2. If e_2 is a multiple of $2^{-p+1}\text{ulp}(p_2)$, then $e_1 + e_2$ is a multiple of $2^{-p+1}\text{ulp}(p_2)$, say $e_1 + e_2 = K \cdot 2^{-p+1} \cdot \text{ulp}(p_2)$. Since $|e_1 + e_2| \leq \frac{1}{2}(\text{ulp}(p_1) + \text{ulp}(p_2)) = \frac{3}{2}\text{ulp}(p_2)$, we deduce that $|K| \leq 3 \cdot 2^{p-2} < 2^p$. This shows that $e_1 + e_2$ is a precision- p floating-point number. Hence, $e = \text{RN}(e_1 + e_2) = e_1 + e_2$, so that $\epsilon_2 = 0$. As a consequence, $p + e = ab + cd$ exactly, and the final relative error is bounded by u . Now, Remark 2.2 tells us that if e_2 is not a multiple of $2^{-p+1}\text{ulp}(p_2)$, then $|cd| \leq (2^p - 2 + 2^{-p})\text{ulp}(p_2)$, so that $|p_2| = |\text{RN}(cd)| \leq 2^{e_{p_1}} - 2\text{ulp}(p_2)$. Hence the case $p_2 = 2^{e_{p_1}} - \text{ulp}(p_2)$ need not be considered.

- If $p_1 = 2^{e_{p_1}}$, then, since $p_1 = \text{RN}(ab)$, $2^{e_{p_1}} - \frac{1}{4}\text{ulp}(p_1) \leq ab \leq 2^{e_{p_1}} + \frac{1}{2}\text{ulp}(p_1)$. However the case $ab \leq 2^{e_{p_1}}$ is easily dealt with: in that case, we have $|e_1| \leq \frac{1}{2}\text{ulp}(p_2)$, so that it is very similar to a case already met: $e_1 + e_2$ is a floating-point number. Hence, $e = \text{RN}(e_1 + e_2) = e_1 + e_2$, so that $e_2 = 0$. As a consequence, $p + e = ab + cd$ exactly, and the final relative error is bounded by u .

Therefore, we only need to consider two cases:

- **Case 1** p_1 is the floating-point number immediately above $2^{e_{p_1}}$, and $2^{e_{p_1}} - 2\text{ulp}(p_2)$. When reasoning on the consequences of Remark 2.2, we have seen that we can further assume that $|cd| \leq (2^p - 2 + 2^{-p})\text{ulp}(p_2) = 2^{e_{p_1}} - (2 - 2^{-p})\text{ulp}(p_2)$. This case is exemplified by Figure 1. In that case,

$$|ab + cd| > (3 - 2^{-p})\text{ulp}(p_2),$$

and

$$|ab| + |cd| < \left(2^{p-1} + \frac{3}{2}\right)\text{ulp}(p_1) + (2^{p+1} - 2 + 2^{-p})\text{ulp}(p_2) = (2^{p+1} + 1 + 2^{-p})\text{ulp}(p_2),$$

so that

$$\gamma < u^2 \frac{2^{p+1} + 1 + 2^{-p}}{3 - 2^{-p}} \cdot |ab + cd|.$$

Elementary manipulations show that as soon as $u = 2^{-p}$ is less than $1/2$ (i.e., $p \geq 1$, which always holds), the ratio

$$\frac{2^{p+1} + 1 + 2^{-p}}{3 - 2^{-p}} = \frac{2}{3u} + \frac{5}{9} + \frac{14u}{27} + \frac{14u^2}{81} + \dots$$

is less than

$$\frac{2}{3u} + 1.$$

As a consequence, $\gamma \leq \left(\frac{2u}{3} + u^2\right) |ab + cd|$, so that the final relative error is less than $\frac{5}{3}u + \frac{5}{3}u^2 + u^3$.

- **Case 2** $p_1 = 2^{e_{p_1}}$ and $-p_2$ is $p_1 - 2\text{ulp}(p_2)$, $p_1 - 3\text{ulp}(p_2)$, or $p_1 - 4\text{ulp}(p_2)$. We have seen that we can further assume $|cd| \leq 2^{e_{p_1}} - (2 - 2^{-p})\text{ulp}(p_2)$, and $ab > 2^{e_{p_1}}$. This case is exemplified by Figure 2. In that case,

$$|ab + cd| > (2 - 2^{-p})\text{ulp}(p_2),$$

and

$$|ab| + |cd| < [(2^p - 1) + (2^p - 2 - 2^{-p})]\text{ulp}(p_2) = (2^{p+1} - 1 + 2^{-p})\text{ulp}(p_2).$$

We deduce

$$\gamma \leq u^2 \frac{2^{p+1} - 1 + 2^{-p}}{2 - 2^{-p}} |ab + cd|.$$

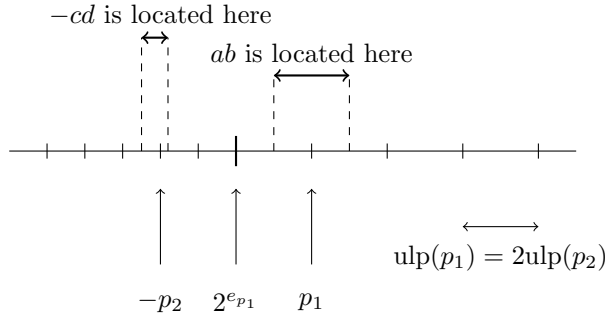


Figure 1: Case $p_1 = 2^{e_{p_1}} \cdot (1 + 2^{-p+1})$.

If p_2 is the largest possible,
 $-cd$ is located here

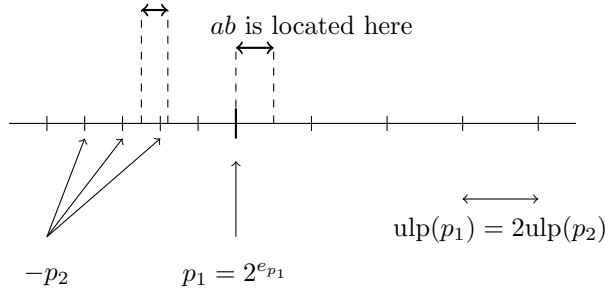


Figure 2: Case $p_1 = 2^{e_{p_1}}$.

We easily find

$$\frac{2^{p+1} - 1 + 2^{-p}}{2 - 2^{-p}} \leq \frac{1}{u} + u,$$

Hence $\gamma \leq (u + u^3)|ab + cd|$, from which we deduce that the final relative error is bounded by $2u + u^2 + u^3 + u^4$.

4 General result

The results obtained in the various cases considered in Section 3 can be summarized as follows

Theorem 4.1. *Provided no underflow/overflow occurs, and assuming radix-2, precision- p floating-point arithmetic, the relative error of Cornea et al's algorithm is bounded by $2u + 7u^2 + 6u^3$.*

Now, interestingly enough, we are going to see that the bound given by Theorem 4.1 is asymptotically optimal (as $p \rightarrow \infty$ or, equivalently, as $u \rightarrow 0$). To

show this, it suffices to consider, in radix-2, precision- p floating-point arithmetic:

$$\begin{cases} a &= 2^p - 1, \\ b &= 2^{p-3} + \frac{1}{2}, \\ c &= 2^p - 1, \\ d &= 2^{p-3} + \frac{1}{4}, \end{cases}$$

for which we find:

$$\begin{aligned} ab + cd &= 2^{2p-2} + 2^{p-1} - \frac{3}{4}, \\ p_1 &= 2^{2p-3} + 2^{p-2}, \\ e_1 &= 2^{p-3} - \frac{1}{2}, \\ p_2 &= 2^{2p-3}, \\ e_2 &= 2^{p-3} - \frac{1}{4}, \\ p &= 2^{2p-2}, \\ e &= 2^{p-2} - \frac{3}{4}, \\ s &= 2^{2p-2}. \end{aligned}$$

The relative error $|s - (ab + cd)|/|ab + cd|$ is equal to

$$\frac{2^{p-1} - \frac{3}{4}}{2^{2p-2} + 2^{p-1} - \frac{3}{4}} = \frac{2u - 3u^2}{1 + 2u - 3u^2} = 2u - 7u^2 + 20u^3 + \dots$$

which is asymptotically equivalent to $2u$. This shows that our relative error bound is asymptotically optimal.

In the frequent case where the considered floating-point format is the binary64/double precision format of the IEEE 754 Standard, the relative error bound provided by Theorem 4.1 is

$$u \times 2.000000000000000777156 \dots,$$

and the relative error attained with our example is

$$u \times 1.99999999999999922284 \dots$$

Conclusion

We have provided a relative error bound for Cornea, Harrison and Tang's algorithm (Algorithm 2), and we have shown that our bound is asymptotically optimal. Since that bound is not better than the (also asymptotically optimal) error bound for Kahan's algorithm (Algorithm 1), it is in general preferable to use Algorithm 1. A possible exception is when one wants to always get the same result when computing $ab + cd$ and $cd + ab$ (for instance to implement a commutative complex multiplication): in this case, the natural symmetry of Algorithm 2 will guarantee the required property, whereas it is easy to build examples for which Algorithm 1 does not satisfy it.

References

- [1] M. Cornea, J. Harrison, and P. T. P. Tang. *Scientific Computing on Itanium[®]-based Systems*. Intel Press, Hillsboro, OR, 2002.
- [2] N. J. Higham. *Accuracy and Stability of Numerical Algorithms*. SIAM, Philadelphia, 1996.
- [3] IEEE Computer Society. *IEEE Standard for Floating-Point Arithmetic*. IEEE Standard 754-2008, August 2008. available at <http://ieeexplore.ieee.org/servlet/opac?punumber=4610933>.
- [4] 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*, 82, October 2013.