

# On relative errors of floating-point operations: optimal bounds and applications

Claude-Pierre Jeannerod, Siegfried M. Rump

► **To cite this version:**

Claude-Pierre Jeannerod, Siegfried M. Rump. On relative errors of floating-point operations: optimal bounds and applications. *Mathematics of Computation*, American Mathematical Society, 2018, 87, pp.803-819. <10.1090/mcom/3234>. <hal-00934443v4>

**HAL Id: hal-00934443**

**<https://hal.inria.fr/hal-00934443v4>**

Submitted on 3 Nov 2016

**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 RELATIVE ERRORS OF FLOATING-POINT OPERATIONS: OPTIMAL BOUNDS AND APPLICATIONS

CLAUDE-PIERRE JEANNEROD AND SIEGFRIED M. RUMP

ABSTRACT. Rounding error analyses of numerical algorithms are most often carried out via repeated applications of the so-called standard models of floating-point arithmetic. Given a round-to-nearest function  $\text{fl}$  and barring underflow and overflow, such models bound the relative errors  $E_1(t) = |t - \text{fl}(t)|/|t|$  and  $E_2(t) = |t - \text{fl}(t)|/|\text{fl}(t)|$  by the unit roundoff  $u$ . This paper investigates the possibility and the usefulness of refining these bounds, both in the case of an arbitrary real  $t$  and in the case where  $t$  is the exact result of an arithmetic operation on some floating-point numbers. We show that  $E_1(t)$  and  $E_2(t)$  are optimally bounded by  $u/(1+u)$  and  $u$ , respectively, when  $t$  is real or, under mild assumptions on the base and the precision, when  $t = x \pm y$  or  $t = xy$  with  $x, y$  two floating-point numbers. We prove that while this remains true for division in base  $\beta > 2$ , smaller, attainable bounds can be derived for both division in base  $\beta = 2$  and square root. This set of optimal bounds is then applied to the rounding error analysis of various numerical algorithms: in all cases, we obtain significantly shorter proofs of the best-known error bounds for such algorithms, and/or improvements on these bounds themselves.

## 1. INTRODUCTION

Given two integers  $\beta, p \geq 2$ , let  $\mathbb{F}$  be the associated set of floating-point numbers having base  $\beta$ , precision  $p$ , and no restriction on the exponent range:

$$\mathbb{F} = \{0\} \cup \{M \cdot \beta^e : M, e \in \mathbb{Z}, \beta^{p-1} \leq |M| < \beta^p\}.$$

Let also  $\text{fl} : \mathbb{R} \rightarrow \mathbb{F}$  denote any round-to-nearest function, such that

$$(1.1) \quad |t - \text{fl}(t)| = \min_{f \in \mathbb{F}} |t - f|, \quad t \in \mathbb{R}.$$

In particular, no specific tie-breaking strategy is assumed for the function  $\text{fl}$ . Two *relative errors* can then be defined, depending on whether the exact value or the rounded value is used to divide the absolute error in (1.1): the error relative to  $t$  is

$$E_1(t) = \frac{|t - \text{fl}(t)|}{|t|} \quad \text{if } t \neq 0,$$

while the error relative to  $\text{fl}(t)$  is

$$E_2(t) = \frac{|t - \text{fl}(t)|}{|\text{fl}(t)|} \quad \text{if } \text{fl}(t) \neq 0.$$

(In each case the relative error may be defined to be zero if the denominator is zero: the exponent range being unbounded,  $\text{fl}(t) = 0$  implies  $t = 0$ , so in both cases a zero denominator means that no error occurs when rounding  $t$ .)

---

November 3, 2016.

2010 *Mathematics Subject Classification.* Primary 65G50.

For rounding error analysis purposes, the most commonly used bounds are  $E_1(t) \leq u$  and  $E_2(t) \leq u$ , where

$$u = \frac{1}{2}\beta^{1-p}$$

is the *unit roundoff* associated with  $\text{fl}$  and  $\mathbb{F}$ , and such bounds are typically handled via the so-called *standard models*  $\text{fl}(t) = t(1 + \delta_1) = t/(1 + \delta_2)$ ,  $|\delta_1|, |\delta_2| \leq u$ ; see [5, pp. 38–39]. It is also known that the bound on the first relative error can be refined slightly [12, p. 232],<sup>1</sup> so that

$$(1.2) \quad E_1(t) \leq \frac{u}{1+u} \quad \text{and} \quad E_2(t) \leq u.$$

These worst-case bounds hold for any real number  $t$  and any rounding function  $\text{fl}$ . Furthermore, they can be regarded as *optimal* in the sense that each of them is attained for some pair  $(t, \text{fl})$  with  $t$  expressed in terms of  $\beta$  and  $p$ : since  $1+u$  is exactly halfway between the two consecutive elements  $1$  and  $1+2u$  of  $\mathbb{F}$ , we have  $E_1(1+u) = u/(1+u)$  and, assuming further that  $\text{fl}$  rounds ties “to even,”  $E_2(1+u) = u$ .

In this paper, we investigate the possibility and the usefulness of refining the bounds in (1.2) when  $t$  is not just an arbitrary real number but the exact result of an operation on some floating-point number(s), that is, when

$$t = x \text{ op } y, \quad x, y \in \mathbb{F}, \quad \text{op} \in \{+, -, \times, /\} : \mathbb{F} \times \mathbb{F} \rightarrow \mathbb{R}$$

as well as  $t = \sqrt{x}$ .

Our first contribution is to establish optimal bounds on both  $E_1$  and  $E_2$  for each of these five basic operations, as shown in Table 1.

TABLE 1. Optimal relative error bounds for various inputs  $t$ .

$t$	bound on $E_1(t)$	bound on $E_2(t)$
real number	$\frac{u}{1+u}$	$u$
$x \pm y$	$\frac{u}{1+u}$	$u$
$xy$	$\frac{u}{1+u}$	$u$
$x/y$	$\begin{cases} u - 2u^2 & \text{if } \beta = 2, \\ \frac{u}{1+u} & \text{if } \beta > 2 \end{cases}$	$\begin{cases} \frac{u-2u^2}{1+u-2u^2} & \text{if } \beta = 2, \\ u & \text{if } \beta > 2 \end{cases}$
$\sqrt{x}$	$1 - \frac{1}{\sqrt{1+2u}}$	$\sqrt{1+2u} - 1$

As we shall see later in the paper, each of these bounds is attained for some explicit input values in  $\mathbb{F}$  and rounding functions  $\text{fl}$ , possibly under some mild (necessary and sufficient) conditions on  $\beta$  and  $p$ . Specifically, for addition, subtraction, and multiplication the condition for optimality is that  $\beta$  is *even*, and in the case of multiplication in base 2 it is that  $2^p + 1$  is *not* a Fermat prime. In most practical situations such conditions are satisfied and thus the general bounds in (1.2) remain

<sup>1</sup>The refined bound  $E_1(t) \leq u/(1+u)$  already appears in [18, p. 74] and, in some special cases, in [2] for  $\beta = 2$  and [6] for  $\beta$  even.

the best ones for floating-point addition, subtraction, and multiplication; as Table 1 shows, this is also the case of division in any base larger than 2. In contrast, for division in base 2 and for square root, the general bound  $u/(1+u) \approx u-u^2$  on  $E_1$  can be decreased further to  $u-2u^2$  and to  $1-(1+2u)^{-1/2} \approx u-\frac{3}{2}u^2$ , respectively; likewise, the general bound  $u$  on  $E_2$  can be decreased to  $(u-2u^2)/(1+u-2u^2) \approx u-3u^2$  and  $(1+2u)^{1/2}-1 \approx u-\frac{1}{2}u^2$ .

Our second contribution is to show that in the context of rounding error analysis of numerical algorithms, applying these optimal bounds in a systematic way leads to simpler and sharper bounds, and/or to more direct proofs of existing ones.

In particular, this allows us to establish the following three new error bounds:

- For the summation of  $n$  floating-point numbers  $x_1, \dots, x_n$  (using  $n-1$  floating-point additions and any ordering), we show that the resulting floating-point approximation  $\hat{s}$  satisfies

$$(1.3) \quad \left| \hat{s} - \sum_{i=1}^n x_i \right| \leq \sum_{i=1}^{n-1} |e_i| \leq \frac{(n-1)u}{1+u} \sum_{i=1}^n |x_i|,$$

where  $|e_i|$  denotes the absolute error of the  $i$ th floating-point addition.

- For the summation of  $n$  real numbers  $x_1, \dots, x_n$ , we show that by first rounding each  $x_i$  into  $\text{fl}(x_i)$  and then summing the  $\text{fl}(x_i)$  in any order, the resulting approximation  $\hat{s}$  satisfies

$$(1.4) \quad \left| \hat{s} - \sum_{i=1}^n x_i \right| \leq \sum_{i=1}^n |d_i| + \sum_{i=1}^{n-1} |e_i| \leq \zeta_n \sum_{i=1}^n |x_i|, \quad \zeta_n < nu,$$

where the  $d_i$  are given by  $d_i = x_i - \text{fl}(x_i)$  and the  $e_i$  are as before.

- For the Euclidean norm of a vector of  $n$  floating-point numbers  $x_1, \dots, x_n$  we show that summing the squares in any order and then taking the square root of the result yields a floating-point number  $\hat{r}$  such that

$$(1.5) \quad \hat{r} = \left( \sum_{i=1}^n x_i^2 \right)^{1/2} \cdot (1 + \epsilon), \quad |\epsilon| \leq (n/2 + 1)u.$$

Note that each of these bounds holds without any restriction on  $n$ . The bounds in (1.3) and (1.4) improve upon the best previous ones, from [10], in two ways: they are sharper and apply not only to the absolute error of  $\hat{s}$  but also to the sum of the absolute *local* errors. Furthermore, we will see that the new bound in (1.3) implies the one in (1.4) almost immediately; in other words, using (1.3) allows us to recover the constant  $nu$  established in [10, Proposition 4.1] for sums of  $n$  reals (and thus  $n$ -dimensional inner products as well), but in a much more direct way. Finally, the bound in (1.5) nicely replaces the expression  $(n/2 + 1)u + \mathcal{O}(u^2)$  that would result from using the suboptimal bound  $E_1(\sqrt{x}) \leq u$ .

Besides the evaluation of sums and norms, and to illustrate further the benefits of applying the refined bounds in Table 1, we provide four other typical examples of rounding error analysis. These examples deal with small arithmetic expressions, minors of tridiagonal matrices, Cholesky factorization, and complex floating-point multiplication. We shall see that in each case the existing error bound can be either replaced by a simpler and sharper one, or recovered via a significantly shorter proof.

**Notation and assumptions.** All our results hold under the customary assumption that  $\beta, p \geq 2$ . Furthermore, following [5] (and unless stated otherwise) we

shall ignore underflow and overflow by assuming that the exponent range of  $\mathbb{F}$  is unbounded. Finally, the common tool used to establish all the error bounds in Table 1 is the function  $\text{ufp} : \mathbb{R} \rightarrow \mathbb{F}_{\geq 0}$  from [16], called *unit in the first place* (ufp) and defined as follows:  $\text{ufp}(0) = 0$  and, if  $t \in \mathbb{R} \setminus \{0\}$ ,  $\text{ufp}(t)$  is the largest integer power of  $\beta$  such that  $\text{ufp}(t) \leq |t|$ . Hence, in particular, for  $t$  nonzero,

$$(1.6a) \quad \text{ufp}(t) \leq |t| < \beta \text{ufp}(t),$$

and for any  $t$ ,

$$(1.6b) \quad \text{ufp}(t) \leq |\text{fl}(t)| \leq \beta \text{ufp}(t).$$

**Outline.** This paper is organized as follows. We begin in Section 2 by recalling how to derive the bounds in (1.2) and by completely characterizing their attainability. Section 3 then gives proofs for all the bounds announced in Table 1 together with explicit expressions of input values at which these bounds are attained. A first application of these results is described in Section 4, where we establish the new bounds for summation shown in (1.3) and (1.4). We conclude in Section 5 with the derivation of the bound (1.5) together with the analysis of four other examples.

## 2. OPTIMAL ERROR BOUNDS WHEN ROUNDING A REAL NUMBER

Using the function  $\text{ufp}$ , the bounds  $E_i(t) \leq u$  for  $i = 1, 2$  are easily derived as follows. First, recall from (1.6) that  $|t| \neq 0$  belongs to the right-open interval  $[\text{ufp}(t), \beta \text{ufp}(t))$ , which contains  $(\beta - 1)\beta^{p-1}$  equally-spaced elements of  $\mathbb{F}$ . The distance between two such consecutive elements is thus  $\frac{\beta \text{ufp}(t) - \text{ufp}(t)}{(\beta - 1)\beta^{p-1}} = 2u \text{ufp}(t)$ , and rounding to nearest implies that the absolute error is bounded as

$$(2.1) \quad |t - \text{fl}(t)| \leq u \text{ufp}(t).$$

This bound is sharp and the values at which it is attained are the *midpoints* of  $\mathbb{F}$ , that is, the rational numbers lying exactly halfway between two consecutive elements of  $\mathbb{F}$ .

Dividing both sides of (2.1) by either  $|t|$  or  $|\text{fl}(t)|$  gives

$$(2.2) \quad E_1(t) \leq u \frac{\text{ufp}(t)}{|t|} \quad \text{and} \quad E_2(t) \leq u \frac{\text{ufp}(t)}{|\text{fl}(t)|},$$

and by using the lower bounds in (1.6) we arrive at the classical bounds  $E_1(t) \leq u$  and  $E_2(t) \leq u$ . As noted in [5, Theorem 2.2], we have in fact the strict inequality  $E_1(t) < u$ , since  $t$  cannot be at the same time a midpoint and equal to its  $\text{ufp}$ ; on the other hand, the derivation above shows that the bound on  $E_2(t)$  is attained if and only if  $t$  is a midpoint such that  $|\text{fl}(t)| = \text{ufp}(t)$ .

Let us now refine the bound  $E_1(t) < u$ . As shown in [12, p. 232], all we need for this is a lower bound on  $|t|$  slightly sharper than the one in (1.6): by definition of rounding to nearest,  $|t - \text{fl}(t)| \leq |t - f|$  for all  $f \in \mathbb{F}$ , so taking in particular  $f = \text{sign}(t)\text{ufp}(t)$  gives  $|t - \text{fl}(t)| \leq |t - \text{ufp}(t)|$ ; in other words,

$$(2.3) \quad \text{ufp}(t) + |t - \text{fl}(t)| \leq |t|$$

and, because of (2.1), equality occurs if and only if  $|t| \leq (1 + u)\text{ufp}(t)$ . Thus, by applying (2.3) and then (2.1) to the definition of  $E_1$ , we find that

$$E_1(t) \leq \frac{|t - \text{fl}(t)|}{\text{ufp}(t) + |t - \text{fl}(t)|} \leq \frac{u}{1 + u}.$$

Furthermore, due to the conditions of attainability of (2.1) and (2.3) given above, this bound on  $E_1(t)$  is attained if and only if  $t$  is a midpoint such that  $|t| \leq (1+u)\text{ufp}(t)$ , that is, if and only if  $|t| = (1+u)\text{ufp}(t)$ .

We summarize our discussion in the theorem below. Although the bounds given there already appear in [18, p. 74] and [12, p. 232] for  $E_1$ , and in [5, Theorem 2.3] for  $E_2$ , the characterization of their attainability does not seem to have been reported elsewhere.

**Theorem 2.1.** *If  $t \in \mathbb{R} \setminus \{0\}$ , then*

$$E_1(t) \leq \frac{u}{1+u} \quad \text{and} \quad E_2(t) \leq u.$$

*Furthermore, the bound on  $E_1$  is attained if and only if  $|t| = (1+u)\text{ufp}(t)$ , and the bound on  $E_2$  is attained if and only if  $|t| = (1+u)\text{ufp}(t)$  and  $|\text{fl}(t)| = \text{ufp}(t)$ .*

### 3. OPTIMAL ERROR BOUNDS FOR FLOATING-POINT OPERATIONS

We establish here optimal bounds on both  $E_1$  and  $E_2$  for the operations of addition, subtraction, multiplication, fused multiply-add, division, and square root.

**3.1. Addition, subtraction, and fused multiply-add.** When  $t$  has the form  $t = x + y$  with  $x, y$  in  $\mathbb{F}$ , we show in the theorem below that the general bounds  $E_1(t) \leq \frac{u}{1+u}$  and  $E_2(t) \leq u$  given in (1.2) remain optimal unless the basis  $\beta$  is odd. This extends the analysis done by Holm [6], who considers only the first relative error and assumes implicitly that the basis is even.

**Theorem 3.1.** *Let  $\beta, p \geq 2$  and  $t = x + y$  with  $x, y \in \mathbb{F}$ . The bounds in (1.2) are optimal if and only if  $\beta$  is even. Furthermore, when  $\beta$  is even, they are attained for  $(x, y) = (1, u)$  and rounding “to nearest even.”*

*Proof.* If  $\beta$  is odd, then  $u = \frac{1}{2}\beta^{1-p} = \sum_{i=0}^{\infty} \frac{\beta-1}{2}\beta^{-i-p}$  with  $\frac{\beta-1}{2} \in \{1, \dots, \beta-1\}$ , so  $u$  and  $1+u$  have infinite expansions in base  $\beta$ . Hence  $x + y$  cannot have the form  $\pm(1+u)\beta^e$  with  $e \in \mathbb{Z}$  and thus, by Theorem 2.1, equality never occurs in the bounds  $E_1(x+y) \leq u/(1+u)$  and  $E_2(x+y) \leq u$ .

If  $\beta$  is even, then  $u$  is in  $\mathbb{F}$ . Consequently, we can take  $(x, y) = (1, u)$ , which gives  $E_1(x+y) = u/(1+u)$  and, for ties rounded “to even”,  $E_2(x+y) = u$ .  $\square$

Since  $1+2u$  belongs to  $\mathbb{F}$  for any  $\beta$  and since  $1+u = (1+2u) - u = 1 \times 1 + u$ , the theorem above extends immediately to subtraction as well as to higher-level operations encompassing addition, like the fused multiply-add  $(x, y, z) \mapsto \text{fl}(xy+z)$ .

**3.2. Multiplication.** When  $t = xy$  with  $x, y \in \mathbb{F}$ , the theorem below shows that the situation is more subtle than for addition: although the condition that  $\beta$  is even remains necessary for the optimality of the bounds  $E_1(t) \leq u/(1+u)$  and  $E_2(t) \leq u$ , it is sufficient only when  $\beta \neq 2$ ; when  $\beta = 2$ , optimality turns out to be equivalent to  $2^p + 1$  being composite.

**Theorem 3.2.** *Let  $\beta, p \geq 2$  and  $t = xy$  with  $x, y \in \mathbb{F}$ . We then have the following, depending on the value of the base  $\beta$ :*

- *If  $\beta = 2$ , then the bounds in (1.2) are optimal if and only if  $2^p + 1$  is not prime; furthermore, if  $D$  is a non-trivial divisor of  $2^p + 1$ , then these bounds are attained for*

$$(x, y) = \left( \frac{(2+2u)\text{ufp}(D)}{D}, \frac{D}{\text{ufp}(D)} \right)$$

and rounding “to nearest even.”

- If  $\beta > 2$ , then the bounds in (1.2) are optimal if and only if  $\beta$  is even, and when the latter is true these bounds are attained for

$$(x, y) = (2 + 2u, 2^{-1})$$

and rounding “to nearest even.”

Before proving this result, note that the attainability of the bound  $E_1(xy) \leq u/(1+u)$  has been observed in [6, p. 10] in the particular case  $(\beta, p) = (2, 5)$ , by taking  $x$  and  $y$  in the form shown above with  $D = 11$ .

*Proof.* By Theorem 2.1, the optimality of the bounds in (1.2) when  $t = xy$  is equivalent to the existence of a pair  $(x, y) \in \mathbb{F} \times \mathbb{F}$  such that  $|xy| = (1+u)\text{ufp}(xy)$ .

Assume first that  $\beta > 2$ . Similarly to Theorem 3.1, a necessary condition for optimality is that  $\beta$  be even. Furthermore, if  $\beta$  is even, then both  $2 + 2u$  and  $2^{-1}$  are in  $\mathbb{F}$ , and since their product equals  $1+u$ , it suffices to take  $(x, y) = (2+2u, 2^{-1})$  to show that the bounds in (1.2) are optimal.

Let us now consider the case  $\beta = 2$ . Since  $\text{ufp}(t \cdot 2^e) = \text{ufp}(t) \cdot 2^e$  for all  $(t, e) \in \mathbb{R} \times \mathbb{Z}$ , we can assume with no loss of generality that  $1 \leq x, y < 2$ . This implies that  $\text{ufp}(xy) \in \{1, 2\}$  and, since  $x, y \in \{1, 1+2u, 1+4u, \dots\}$  and  $u > 0$ , that the product  $xy$  cannot be equal to  $1+u$ . Hence, optimality is equivalent to the existence of  $x, y \in \mathbb{F} \cap [1, 2)$  such that  $xy = 2 + 2u$ , that is, equivalent to the existence of integers  $X, Y$  such that

$$(3.1) \quad XY = (2^p + 1) \cdot 2^{p-1} \quad \text{and} \quad 2^{p-1} \leq X, Y < 2^p.$$

If  $2^p + 1$  is prime, then either  $X$  or  $Y$  must be larger than  $2^p$ , so (3.1) has no solution.

If  $2^p + 1$  is composite, one can construct a solution  $(X_0, Y_0)$  to (3.1) as follows. Let  $D$  denote a non-trivial divisor of  $2^p + 1$ , and let  $X_0 = \frac{2^p+1}{D} \text{ufp}(D)$  and  $Y_0 = D \frac{2^{p-1}}{\text{ufp}(D)}$ . Clearly,  $X_0$  is an integer and the product  $X_0 Y_0$  has the desired shape. Thus, it remains to check that  $Y_0 \in \mathbb{Z}$  and that both  $X_0$  and  $Y_0$  are in the range  $[2^{p-1}, 2^p)$ . Since  $2^p + 1$  is odd,  $D$  must be odd too, which implies that

$$(3.2) \quad \text{ufp}(D) + 1 \leq D < 2\text{ufp}(D) \quad \text{and} \quad D < 2^p.$$

Consequently,  $\text{ufp}(D) \leq 2^{p-1}$ , so that  $\text{ufp}(D)$  divides  $2^{p-1}$  and  $Y_0$  is an integer. Furthermore, (3.2) leads to  $2^{p-1} < \frac{2^p+1}{2} < X_0 \leq (2^p + 1)(1 - \frac{1}{D}) < 2^p$  and  $2^{p-1} < Y_0 < 2^p$ , so that  $X_0$  and  $Y_0$  satisfy the range constraint in (3.1).

Finally, multiplying  $X_0$  and  $Y_0$  by  $2u = 2^{1-p}$  gives  $x_0 = (2 + 2u)\text{ufp}(D)/D$  and  $y_0 = D/\text{ufp}(D)$  in  $\mathbb{F} \cap [1, 2)$  and such that  $x_0 y_0 = 2 + 2u = (1+u)\text{ufp}(x_0 y_0)$ .  $\square$

In the rest of this section, we show that the optimality condition “ $2^p + 1$  is not prime” arising in Theorem 3.2 for radix two is in fact satisfied in most practical situations.

First of all, if  $p$  is *not* a power of two, this condition is well known to hold [3, Exercise 18, §4], since  $p$  can be factored as  $p = mq$  with  $q$  odd and then

$$(3.3) \quad 2^p + 1 = (2^m + 1)(2^{p-m} - 2^{p-2m} + \dots - 2^m + 1).$$

The binary *basic* formats specified by the IEEE 754-2008 standard being such that  $p \in \{24, 53, 113\}$ , they fall into this category. More precisely, since here  $p$  is either odd or a multiple of three, we deduce from (3.3) explicit divisors of  $2^p + 1$  and thus explicit pairs  $(x, y) \in \mathbb{F}^2$  for which the bounds in Theorem 3.2 are attained:

- If  $p$  is odd, then 3 divides  $2^p + 1$  and we can take

$$(x, y) = \left(\frac{4+4u}{3}, \frac{3}{2}\right);$$

- If  $p \equiv 0 \pmod{3}$ , then  $2^p + 1$  can be factored as  $(2^{p/3} + 1)(2^{2p/3} - 2^{p/3} + 1)$ , so we can take

$$(x, y) = (2 - 2u^{1/3} + 2u^{2/3}, 1 + u^{1/3}).$$

In fact, the sufficient condition “ $p$  is not a power of two” is satisfied not only by those basic formats but also by all the binary *interchange* formats of IEEE 754-2008, for which either  $p \in \{11, 24, 53, 113\}$  or

$$(3.4) \quad p = k - d + 13 \quad \text{with} \quad d = \lceil 4 \log_2 k \rceil, \quad k = 32j, \quad j \in \mathbb{N}_{\geq 5},$$

and where  $\lceil \cdot \rceil$  denotes rounding to a nearest integer. (A proof of the fact that (3.4) implies that  $p$  is not a power of two is deferred to Appendix A.)

Assume now that  $p$  is a power of two. In this case  $2^p + 1$  is not prime if and only if it is not a prime number of the form  $F_\ell = 2^{2^\ell} + 1$ , called a Fermat prime. Currently, the only known Fermat primes are  $F_0, F_1, F_2, F_3, F_4$  and, on the other hand,  $F_\ell$  is known to be composite for all  $5 \leq \ell \leq 32$ ; see for example [11, 17].

To summarize, the only values of  $p$  for which the optimality condition “ $2^p + 1$  is not prime” can fail to be satisfied are 2, 4, 8, 16 and  $p = 2^\ell \geq 2^{33} \approx 8.6 \times 10^9$ , and none of these values corresponds to an IEEE format.

**3.3. Division.** This section focuses on the largest possible relative errors committed when rounding  $x/y$  with  $x, y$  nonzero elements of  $\mathbb{F}$ . As the theorem below shows, the general bounds  $E_1(t) \leq \frac{u}{1+u}$  and  $E_2(t) \leq u$  given in (1.2) can be refined further for base 2, but remain optimal for all other bases. Note that unlike for addition or multiplication, optimality is achieved without any extra assumption on the parity of  $\beta$  or the number theoretic properties of  $p$ .

**Theorem 3.3.** *Let  $\beta, p \geq 2$  and let  $x, y \in \mathbb{F}$  be nonzero. Then*

$$E_1(x/y) \leq \begin{cases} u - 2u^2 & \text{if } \beta = 2, \\ \frac{u}{1+u} & \text{if } \beta > 2, \end{cases}$$

and

$$E_2(x/y) \leq \begin{cases} \frac{u-2u^2}{1+u-2u^2} & \text{if } \beta = 2, \\ u & \text{if } \beta > 2. \end{cases}$$

The bounds for  $\beta = 2$  are attained at  $(x, y) = (1, 1 - u)$  and, assuming ties are rounded “to even”, the bounds for  $\beta > 2$  are attained at  $(x, y) = (2 + 2u, 2)$ .

*Proof.* When  $\beta > 2$  the bounds are the general ones given in (1.2), and the fact they are attained for division follows immediately from  $2 + 2u \in \mathbb{F}$ . The rest of the proof is thus devoted to the case  $\beta = 2$ .

Let  $t = x/y$ . Since  $t$  cannot be a midpoint [13, 9], we have  $\text{fl}(-t) = -\text{fl}(t)$  and  $\text{fl}(t \cdot 2^e) = \text{fl}(t) \cdot 2^e$ ,  $e \in \mathbb{Z}$ , regardless of the tie-breaking rule of  $\text{fl}$ . Consequently, we can assume  $1 \leq t < 2$  and  $x, y > 0$ . When  $t = 1$  both  $E_1$  and  $E_2$  are zero, so we are left with handling  $t$  such that  $1 < t < 2$ .

The lower bound on  $t$  implies  $x > y$ , which for  $x$  and  $y$  in  $\mathbb{F}$  is equivalent to  $x \geq y + 2u \text{ufp}(y)$ . Hence, using  $y \leq (2 - 2u)\text{ufp}(y)$ ,

$$(3.5) \quad t \geq 1 + 2u \frac{\text{ufp}(y)}{y} \geq 1 + \frac{u}{1-u} = \frac{1}{1-u}.$$



Since  $1/(1-u)$  is strictly larger than the midpoint  $1+u$ , it follows that

$$(3.6) \quad \text{fl}(t) \in \{1+2u, 1+4u, \dots\}.$$

■ Assume for now that  $p \geq 3$ . (For simplicity, the case  $\beta = p = 2$  is handled separately at the end of the proof.)

If  $\text{fl}(t) \geq 1+4u$  then  $t \geq 1+3u$ , so that  $E_1(t) \leq u \text{ufp}(t)/t \leq u/(1+3u) < u-2u^2$  for  $p \geq 3$ . Similarly,  $E_2(t) \leq u/(1+4u) < (u-2u^2)/(1+u-2u^2)$  for  $p \geq 3$ .

If  $\text{fl}(t) = 1+2u$  then, recalling (3.5) and the fact that  $t$  is not a midpoint,

$$(3.7) \quad \frac{1}{1-u} \leq t < 1+3u.$$

We now distinguish between the following two sub-cases, depending on how  $t$  compares to  $\text{fl}(t) = 1+2u$ :

- If  $t \leq 1+2u$ , then  $E_1(t) = (1+2u)/t - 1$  and  $E_2(t) = 1 - t/(1+2u)$ , so that the first inequality in (3.7) gives immediately the desired bounds, which are attained only when  $t = 1/(1-u)$ ; since both 1 and  $1-u$  are in  $\mathbb{F}$  when  $\beta = 2$ , this value of  $t$  is obtained for  $(x, y) = (1, 1-u)$ .
- If  $t > 1+2u$ , then  $E_1(t) = 1 - (1+2u)/t$  and  $E_2(t) = t/(1+2u) - 1$ . The bound  $t < 1+3u$  from (3.7) then gives immediately  $E_1(t) < u/(1+3u)$ , which is less than  $u-2u^2$  for  $p \geq 3$ .

For  $E_2(t)$ , however, using  $t < 1+3u$  is not enough and we show first how to replace this bound by the slightly sharper one

$$(3.8) \quad t \leq \frac{2+7u}{2+u} = 1+3u - \frac{3}{2}u^2 + \mathcal{O}(u^3).$$

The range of  $t$  implies  $y+2u \text{ufp}(y) < x < y+6u \text{ufp}(y)$ . Note that  $y$  cannot be equal to  $(2-2u)\text{ufp}(y)$ , for otherwise  $2\text{ufp}(y) < x < (2+4u)\text{ufp}(y)$ , thus contradicting the fact that  $x \in \mathbb{F}$ . Therefore,  $y \leq (2-4u)\text{ufp}(y)$  and

$$x = y + 4u \text{ufp}(y).$$

Writing  $y = (1+2ku)\text{ufp}(y)$  with  $k$  a nonnegative integer, we deduce that

$$t = 1 + \frac{4u}{1+2ku},$$

so  $t < 1+3u$  is equivalent to  $k > 2^{p-1}/3$ , that is,  $k \geq (2^{p-1}+1)/3$  for  $k$  is an integer. Hence  $2ku \geq (1+2u)/3$  and (3.8) follows.

Recalling that  $E_2(t) = t/(1+2u) - 1$  and applying (3.8), we arrive at  $E_2(t) \leq \frac{2u(1-u)}{(2+u)(1+2u)}$ , which is less than  $\frac{u-2u^2}{1+u-2u^2}$  for  $p \geq 3$ .

■ Assume now that  $p = 2$ . We have  $u = 1/4$  and, assuming with no loss of generality that  $\text{ufp}(x) = 1$ , we see that  $x \in \{1, 3/2\}$  and that  $y$  is either  $\text{ufp}(y)$  or  $3/2 \cdot \text{ufp}(y)$ . This yields four possibilities for  $t = x/y$ , all leading to  $E_1(t) = E_2(t) = 0$  except for  $x = 1$  and  $y = 3/2 \cdot \text{ufp}(y)$ . In this case, the constraint  $1 < t < 2$  implies further  $y = 3/4 = 1-u$ . It follows that  $t = 1/(1-u)$ , from which we deduce  $\text{fl}(t) = 1+2u$ .

The announced bounds thus hold also in the case  $\beta = p = 2$ , and since 1 and  $1-u$  are in  $\mathbb{F}$ , they are attained for  $(x, y) = (1, 1-u)$ .  $\square$

**3.4. Square root.** Finally, we show how to refine further the bounds  $E_1(t) \leq u/(1+u)$  and  $E_2(t) \leq u$  in the special case where  $t = \sqrt{x}$  for some positive floating-point number  $x$ , thereby establishing the optimal bounds in the last row of Table 1. This result is independent of any specific property of the base and the precision, and holds for any tie-breaking strategy.

**Theorem 3.4.** *Let  $\beta, p \geq 2$  and let  $x \in \mathbb{F}$  be positive. Then*

$$E_1(\sqrt{x}) \leq 1 - \frac{1}{\sqrt{1+2u}} \quad \text{and} \quad E_2(\sqrt{x}) \leq \sqrt{1+2u} - 1,$$

and these bounds are attained only for  $x = (1+2u)\beta^{2e}$  with  $e \in \mathbb{Z}$ .

*Proof.* Let  $t = \sqrt{x}$ . Writing  $x = \mu\beta^{2e}$  with  $e \in \mathbb{Z}$  and  $\mu \in \mathbb{F} \cap [1, \beta^2)$ , we see that  $t = \sqrt{\mu}\beta^e$  with  $\mu \in \{1, 1+2u, 1+4u, \dots\}$ . If  $\mu = 1$  then  $E_1(t) = E_2(t) = 0$ , so we are left with the following two cases.

If  $\mu = 1+2u$  then we deduce from  $1 < \sqrt{1+2u} < 1+u$  that  $\text{fl}(t) = \beta^e$  and, consequently, that  $E_1(t) = 1 - 1/\sqrt{1+2u}$  and  $E_2(t) = \sqrt{1+2u} - 1$ .

If  $\mu \geq 1+4u$  then, recalling that  $\mu < \beta^2$ , we have  $\sqrt{1+4u}\beta^e \leq t < \beta^{e+1}$  and  $\text{ufp}(t) = \beta^e$ . This implies that  $E_1(t) \leq u\text{ufp}(t)/t \leq u/\sqrt{1+4u} =: \varphi$  and it can be checked that  $\varphi < 1 - 1/\sqrt{1+2u}$  for  $u \leq 1/2$ . Furthermore, using  $\sqrt{1+4u} > 1+u$  gives  $\text{fl}(t) \geq (1+2u)\beta^e$  and then  $E_2(t) \leq u/(1+2u) < \sqrt{1+2u} - 1$ .  $\square$

#### 4. APPLICATION TO SUMMATION

**4.1. Sums of floating-point numbers.** Consider first the evaluation of the sum of  $n$  floating-point numbers. Here  $x_1, \dots, x_n \in \mathbb{F}$  are given and we assume that an approximation  $\hat{s} \in \mathbb{F}$  to the exact sum

$$s = x_1 + \dots + x_n$$

is produced after  $n-1$  floating-point additions, using any evaluation order.<sup>2</sup> Each of these additions takes a pair of floating-point numbers, say  $(a, b) \in \mathbb{F}^2$ , and returns  $\text{fl}(a+b)$ , thus committing the local error  $e := a+b - \text{fl}(a+b)$ . When evaluating the sum as above,  $n-1$  such errors can occur and, denoting their set as  $\{e_1, \dots, e_{n-1}\}$ , it is easy to see that

$$(4.1) \quad s = e_1 + \dots + e_{n-1} + \hat{s};$$

see for example [5, p. 81].

The theorem below shows how to bound the sum of the  $|e_i|$  and, therefore,  $|\hat{s} - s|$  as well. This bound is slightly sharper than the one given in [10, Proposition 3.1] and can be established in the same way, using  $u/(1+u)$  instead of  $u$  to bound the relative rounding error committed by each floating-point addition. (For completeness a detailed proof is presented in Appendix B; note that as in [10] this bound holds even if underflow occurs.)

**Theorem 4.1.** *For  $x_1, \dots, x_n \in \mathbb{F}$ , any order of evaluation of the sum  $s = \sum_{i=1}^n x_i$  produces an approximation  $\hat{s}$  such that the rounding errors  $e_1, \dots, e_{n-1}$  satisfy*

$$|\hat{s} - s| \leq \sum_{i=1}^{n-1} |e_i| \leq \sigma_{n-1} \sum_{i=1}^n |x_i|, \quad \sigma_n := \frac{nu}{1+u}.$$

<sup>2</sup>For example, for  $n = 4$  possible evaluation orders include  $((x_1 + x_2) + x_3) + x_4$  and  $((x_4 + x_3) + x_2) + x_1$  and  $(x_1 + x_2) + (x_3 + x_4)$ .

A direct consequence of this result is a sharper and simpler bound for the following compensated summation scheme (see [14] and the references therein):

$$\widehat{s}_{\text{comp}} = \text{fl}(\widehat{s} + \widehat{e}), \quad \widehat{e} := \text{a floating-point evaluation of } e_1 + \cdots + e_{n-1}.$$

Since each  $e_i$  belongs to  $\mathbb{F}$  and can be computed exactly—using for example Knuth's *TwoSum* algorithm [12, p. 236], we can apply Theorem 4.1 twice and, writing  $e = \sum_{i=1}^{n-1} e_i$ , we deduce that

$$(4.2) \quad |\widehat{e} - e| \leq \sigma_{n-1} \sigma_{n-2} \sum_{i=1}^n |x_i|.$$

Since  $s = \widehat{s} + e$ , we have also

$$(4.3) \quad \begin{aligned} |\widehat{s}_{\text{comp}} - s| &\leq |\text{fl}(\widehat{s} + \widehat{e}) - (\widehat{s} + \widehat{e})| + |\widehat{e} - e| \\ &\leq \frac{u}{1+u} |\widehat{s} + \widehat{e}| + |\widehat{e} - e| \\ &\leq \frac{u}{1+u} |s| + \left(1 + \frac{u}{1+u}\right) |\widehat{e} - e|. \end{aligned}$$

Then, combining (4.2) and (4.3) and using the fact that  $(1 + \frac{u}{1+u})(\frac{u}{1+u})^2 \leq \frac{u^2}{1+u^2}$  we arrive at

$$(4.4) \quad |\widehat{s}_{\text{comp}} - s| \leq \frac{u}{1+u} |s| + (n-1)(n-2) \frac{u^2}{1+u^2} \sum_{i=1}^n |x_i|.$$

The bound in (4.4) holds without any restriction on  $n$  and regardless of the orderings used for adding the  $x_i$  and adding the  $e_i$ . Furthermore, it is slightly sharper than the bound  $u|s| + (n-1)^2 \frac{u^2}{(1-(n-1)u)^2} \sum_{i=1}^n |x_i|$  given in [14, Proposition 4.5] in the special case of recursive compensated summation.

**4.2. Sums of real numbers.** We now turn to the case where  $x_1, \dots, x_n$  are in  $\mathbb{R}$  instead of  $\mathbb{F}$ . An approximation  $\widehat{s} \in \mathbb{F}$  to  $s = x_1 + \cdots + x_n \in \mathbb{R}$  is obtained in two steps, by first rounding each  $x_i$  into  $\text{fl}(x_i)$  and then evaluating  $\text{fl}(x_1) + \cdots + \text{fl}(x_n)$  as above. A typical example is the computation of inner products, where each  $x_i$  is the exact product of two given floating-point numbers.

Writing  $d_i = x_i - \text{fl}(x_i)$  for the errors due to rounding the data and, as before,  $e_i$  for the errors due to the  $n-1$  additions, we deduce that the exact and computed sums are now related as

$$(4.5) \quad s = d_1 + \cdots + d_n + e_1 + \cdots + e_{n-1} + \widehat{s}.$$

By combining (1.2) and Theorem 4.1 we obtain almost immediately the following bound on the sum of the  $|d_i|$  and the  $|e_i|$ .

**Theorem 4.2.** *For  $x_1, \dots, x_n \in \mathbb{R}$ , any order of evaluation of the sum  $\sum_{i=1}^n \text{fl}(x_i)$  produces an approximation  $\widehat{s}$  to the sum  $s = \sum_{i=1}^n x_i$  such that the rounding errors  $d_1, \dots, d_n$  and  $e_1, \dots, e_{n-1}$  satisfy*

$$|\widehat{s} - s| \leq \sum_{i=1}^n |d_i| + \sum_{i=1}^{n-1} |e_i| \leq \zeta_n \sum_{i=1}^n |x_i|, \quad \zeta_n := \frac{(1+2u)nu - u^2}{(1+u)^2} < nu.$$

*Proof.* The lower bound follows from (4.5). To establish the upper bound, let  $v = u/(1+u)$ . Theorem 4.1 gives  $\sum_{i < n} |e_i| \leq (n-1)v \sum_{i \leq n} |\text{fl}(x_i)|$  and, on the other hand, (1.2) implies that  $|d_i| \leq v|x_i|$  and  $|\text{fl}(x_i)| \leq (1+v)|x_i|$ . Hence

$\sum_{i \leq n} |d_i| + \sum_{i < n} |e_i| \leq (v + (n-1)v(1+v)) \sum_{i \leq n} |x_i|$ , and it is easily checked that  $v + (n-1)v(1+v)$  simplifies to  $((1+2u)nu - u^2)/(1+u)^2$ , which is less than  $nu$ .  $\square$

Thus, the theorem above gives in particular the bound

$$\left| \widehat{s} - \sum_{i=1}^n x_i \right| \leq \zeta_n \sum_{i=1}^n |x_i|.$$

This bound has a slightly smaller constant than the one given in [10, Proposition 4.1] and, perhaps more importantly, its proof is significantly shorter and avoids a tedious induction and ufp-based case distinctions.

Furthermore, the applications mentioned in [10] obviously benefit directly from this new bound: for inner products, matrix-vector products, and matrix-matrix products, the constants  $nu$  obtained in [10, Theorem 4.2 and p. 343] can all be replaced by  $\zeta_n$ .

## 5. OTHER APPLICATION EXAMPLES

**5.1. Example 1: Small arithmetic expressions.** Rounding error analyses as those done in [5] typically involve bounds on  $|\theta_n|$ , where  $\theta_n$  is an expression of the form  $\theta_n = \prod_{i=1}^n (1 + \delta_i) - 1$  with  $\delta_i$  a relative error term associated with a single floating-point operation. Using the classical bound  $|\delta_i| \leq u$  it is easily checked that for all  $n$ ,

$$(5.1) \quad -nu \leq \theta_n \leq (1+u)^n - 1.$$

Note that only the upper bound has the form  $nu + \mathcal{O}(u^2)$ , that is, contains terms nonlinear in  $u$ . From (5.1) it follows that

$$|\theta_n| \leq (1+u)^n - 1$$

and, assuming  $nu < 1$ , this bound itself is usually bounded above by the classical fraction  $\gamma_n = nu/(1 - nu)$ .

Using the refined bound  $E_1(t) \leq u/(1+u)$  from (1.2), we can replace  $|\delta_i| \leq u$  by  $|\delta_i| \leq u/(1+u)$ , which immediately leads to the refined enclosure

$$(5.2) \quad -\frac{nu}{1+u} \leq \theta_n \leq \left(1 + \frac{u}{1+u}\right)^n - 1.$$

Although the upper bound in (5.2) still has  $\mathcal{O}(u^2)$  terms in general, it is bounded by  $nu$  as long as  $n \leq 3$ . Consequently, by just systematically using  $E_1(t) \leq u/(1+u)$  instead of  $E_1(t) \leq u$ , we can replace  $\gamma_n = nu + \mathcal{O}(u^2)$  by  $nu$  in every error analysis where  $\theta_n$  appears with  $n \leq 3$ .

For example, when evaluating  $ab + cd$  in the usual way as  $\widehat{r} = \text{fl}(\text{fl}(ab) + \text{fl}(cd))$ , we have  $\widehat{r} = (ab(1 + \delta_1) + cd(1 + \delta_2))(1 + \delta_3)$  with  $|\delta_i| \leq u/(1+u)$ . Hence

$$\widehat{r} = ab(1 + \theta_2) + cd(1 + \theta'_2), \quad |\theta_2|, |\theta'_2| \leq 2u,$$

so the usual  $\gamma_2$  is indeed replaced by  $2u$ . Note that once we have replaced  $E_1(t) \leq u$  by  $E_1(t) \leq u/(1+u)$ , we obtain that term  $2u$  immediately, without having to resort to a sophisticated ufp-based argument as the one introduced in [1, pp. 1470–1471].

Similar examples include the evaluation of small arithmetic expressions like the product  $x_1 x_2 x_3 x_4$  (using any parenthesization) or the sums  $((x_1 + x_2) + x_3) + x_4$  and  $((x_1 + x_2) + (x_3 + x_4)) + ((x_5 + x_6) + (x_7 + x_8))$  (using these specific parenthesizations); in each case the forward error bound classically involves  $\gamma_3$ , which we now replace by  $3u$ .

**5.2. Example 2: Leading principal minors of a tridiagonal matrix.** Consider the tridiagonal matrix

$$A = \begin{bmatrix} d_1 & e_1 & & & & \\ c_2 & d_2 & e_2 & & & \\ & \ddots & \ddots & \ddots & & \\ & & \ddots & \ddots & e_{n-1} & \\ & & & c_n & d_n & \end{bmatrix} \in \mathbb{F}^{n \times n},$$

and let  $\mu_1, \mu_2, \dots, \mu_n$  be the sequence of its  $n$  leading principal minors. Writing  $\mu_{-1} = 0$  and  $\mu_0 = 1$ , those minors are thus defined by the linear recurrence

$$\mu_k = d_k \mu_{k-1} - c_k e_{k-1} \mu_{k-2}, \quad 1 \leq k \leq n.$$

Using the usual bound  $E_1(t) \leq u$  and barring underflow and overflow, Wilkinson shows in [19, §3] that the evaluation of this recurrence produces floating-point numbers  $\hat{\mu}_1, \dots, \hat{\mu}_n$  such that

$$\hat{\mu}_k = d_k(1 + \epsilon_k) \hat{\mu}_{k-1} - c_k(1 + \epsilon'_k) e_{k-1}(1 + \epsilon''_k) \hat{\mu}_{k-2},$$

where  $(1 - u)^2 - 1 \leq \epsilon_k \leq (1 + u)^2 - 1$  and  $(1 - u)^{3/2} - 1 \leq \epsilon'_k, \epsilon''_k \leq (1 + u)^{3/2} - 1$ . In other words, the computed  $\hat{\mu}_k$  are the leading principal minors of a nearby tridiagonal matrix  $A + \Delta A = [a_{ij}(1 + \delta_{ij})]$  that satisfies

$$(5.3) \quad -2u < \delta_{ii} \leq 2u + u^2 \quad \text{and} \quad -\frac{3}{2}u < \delta_{ij} \leq \frac{3}{2}u + \mathcal{O}(u^2) \quad \text{if } i \neq j.$$

Notice that the terms  $u^2$  and  $\mathcal{O}(u^2)$  come exclusively from the upper bounds on  $\epsilon_k, \epsilon'_k, \epsilon''_k$ . By using the refined bound  $E_1(t) \leq u/(1 + u)$  from (1.2) instead of just  $E_1(t) \leq u$ , these upper bounds are straightforwardly improved to

$$\epsilon_k \leq \left(1 + \frac{u}{1+u}\right)^2 - 1 < 2u \quad \text{and} \quad \epsilon'_k, \epsilon''_k \leq \left(1 + \frac{u}{1+u}\right)^{3/2} - 1 < \frac{3}{2}u.$$

Consequently, Wilkinson's bounds in (5.3) can be replaced by the following more concise and slightly sharper ones:

$$|\delta_{ii}| < 2u \quad \text{and} \quad |\delta_{ij}| < \frac{3}{2}u \quad \text{if } i \neq j.$$

**5.3. Example 3: Euclidean norm of an  $n$ -dimensional vector.** Given a vector  $[x_1, \dots, x_n]^T \in \mathbb{F}^n$ , let its norm

$$r = \sqrt{x_1^2 + \dots + x_n^2}$$

be evaluated in floating point in the usual way: form the squares  $\text{fl}(x_i^2)$ , sum them up in any order into  $\hat{s}$ , and return  $\hat{r} = \text{fl}(\sqrt{\hat{s}})$ .

By applying the usual bound  $E_1(t) \leq u$ , all we can say is  $\hat{s} = (\sum_{i=1}^n x_i^2)(1 + \theta_n)$  with  $\theta_n$  as in (5.1), and  $\hat{r} = \sqrt{\hat{s}}(1 + \delta)$  with  $|\delta| \leq u$ . Consequently,

$$\hat{r} = r(1 + \epsilon),$$

where  $\epsilon = \sqrt{1 + \theta_n} \cdot (1 + \delta) - 1$  satisfies  $(1 - u)^{n/2+1} - 1 \leq \epsilon \leq (1 + u)^{n/2+1} - 1$ . Although the lower bound has absolute value at most  $(n/2+1)u$ —see for example [4, p. 42], the upper bound is strictly larger than this, so that

$$(5.4) \quad -(n/2 + 1)u \leq \epsilon \leq (n/2 + 1)u + \mathcal{O}(u^2).$$

To avoid the  $\mathcal{O}(u^2)$  term above, we can use the refined bound  $E_1(t) \leq u/(1 + u)$ , which says  $|\delta| \leq u/(1 + u)$ , together with the improved bound for inner products from [10], which says  $|\theta_n| \leq nu$ . Indeed, from these two bounds we deduce that  $\epsilon$

is upper bounded by  $\sqrt{1 + nu} \cdot (1 + u/(1 + u)) - 1$ , and the latter quantity is easily checked to be at most  $(n/2 + 1)u$ . Thus, recalling the lower bound in (5.4), we conclude that

$$(5.5) \quad |\epsilon| \leq (n/2 + 1)u.$$

In particular, evaluating the hypotenuse  $\sqrt{x_1^2 + x_2^2}$  in floating-point produces a relative error of at most  $2u$ ; this improves over the classical bound  $2u + \mathcal{O}(u^2)$ , stated for example in [7, p. 225].

Of course, the bound in (5.5) also applies when scaling by integer powers of the base is introduced to avoid underflow and overflow.

**5.4. Example 4: Cholesky factorization.** We consider  $A \in \mathbb{F}^{n \times n}$  symmetric and its triangularization in floating-point arithmetic using the classical Cholesky algorithm. If the algorithm runs to completion, then by using the bounds  $E_i(t) \leq u$ ,  $i = 1, 2$ , the traditional rounding error analysis concludes that the computed factor  $\widehat{R}$  satisfies  $\widehat{R}^T \widehat{R} = A + \Delta A$  with

$$|\Delta A| \leq \gamma_{n+1} |\widehat{R}^T| |\widehat{R}|;$$

see for example [5, Theorem 10.3]. Here  $\gamma_{n+1} = \frac{(n+1)u}{1-(n+1)u}$  has the form  $(n+1)u + \mathcal{O}(u^2)$  and requires  $n+1 < u^{-1}$ . It was shown in [15] that both the quadratic term in  $u$  and the restriction on  $n$  can be removed, resulting in the improved backward error bound

$$|\Delta A| \leq (n+1)u |\widehat{R}^T| |\widehat{R}|.$$

In the proof of [15, Theorem 4.4], one of the ingredients used to suppress the  $\mathcal{O}(u^2)$  term is the following property:

$$(5.6) \quad \left( a \in \mathbb{F}_{\geq 0} \quad \text{and} \quad b = \text{fl}(\sqrt{a}) \right) \quad \Rightarrow \quad |b^2 - a| \leq 2ub^2.$$

In [15] it is shown that this property may not hold if only the bound  $E_2(t) \leq u$  is assumed, and that in this case all we can say is  $-(2u + u^2)b^2 \leq b^2 - a \leq 2ub^2$ . Furthermore, a proof of (5.6) is given, which is about 10 lines long and based on a ufp-based case analysis; see [15, p. 692].

Instead, our optimal bound on  $E_2$  for square root provides a direct proof: Theorem 3.4 gives  $b(1 + \delta) = \sqrt{a}$  with  $|\delta| \leq \sqrt{1 + 2u} - 1$ ; hence  $|b^2 - a| = b^2 |2\delta + \delta^2|$  and  $|2\delta + \delta^2| \leq (2 + |\delta|)|\delta| \leq (\sqrt{1 + 2u} + 1)(\sqrt{1 + 2u} - 1) = 2u$ , from which (5.6) follows immediately.

**5.5. Example 5: Complex multiplication with an FMA.** Given  $a, b, c, d \in \mathbb{F}$ , consider the complex product

$$z = (a + ib)(c + id).$$

Various approximations  $\widehat{z} = \widehat{R} + i\widehat{I}$  to  $z$  can be obtained, depending on how  $R = ac - bd$  and  $I = ad + bc$  are evaluated in floating-point. It was shown in [1] that the conventional way, which uses 4 multiplications and 2 additions, gives  $\widehat{z} = z(1 + \epsilon)$  with  $\epsilon \in \mathbb{C}$  such that  $|\epsilon| < \sqrt{5}u$ , and that the constant  $\sqrt{5}$  is, at least in base 2, best possible. Assume now that an FMA is available, so that we compute, say,

$$\widehat{R} = \text{fl}(ac - \text{fl}(bd)) \quad \text{and} \quad \widehat{I} = \text{fl}(ad + \text{fl}(bc)).$$

For this algorithm and its variants<sup>3</sup> it was shown in [8] that the bound  $\sqrt{5}u$  can be reduced further to  $2u$ , and that the latter is essentially optimal. The fact that  $2u$  is an *upper* bound is established in [8, Theorem 3.1] with a rather long proof. As we shall see in the paragraph below, a much more direct proof follows from simply applying the refined bound  $E_1(t) \leq u/(1+u)$  in a systematic way.

Denoting by  $\delta_1, \dots, \delta_4$  the four rounding errors involved, we have

$$\begin{aligned}\widehat{R} &= (ac - bd(1 + \delta_1))(1 + \delta_2) \\ &= R + R\delta_2 - bd\delta_1(1 + \delta_2)\end{aligned}$$

and, similarly,  $\widehat{I} = I + I\delta_4 + bc\delta_3(1 + \delta_4)$ . Now let  $\lambda, \mu \in \mathbb{R}_{\geq 0}$  be such that

$$|\delta_2|, |\delta_4| \leq \lambda \quad \text{and} \quad |\delta_1(1 + \delta_2)|, |\delta_3(1 + \delta_4)| \leq \mu.$$

This implies that  $|R - \widehat{R}| \leq \lambda|R| + \mu|bd|$  and  $|I - \widehat{I}| \leq \lambda|I| + \mu|bc|$ , from which we deduce

$$\begin{aligned}(5.7) \quad |z - \widehat{z}|^2 &= (R - \widehat{R})^2 + (I - \widehat{I})^2 \\ &\leq \lambda^2|z|^2 + 2\lambda\mu A + \mu^2 B,\end{aligned}$$

where  $A = |R||bd| + |I||bc|$  and  $B = (bd)^2 + (bc)^2$ . It turns out that

$$(5.8) \quad A, B \leq |z|^2.$$

For  $B$ , this bound simply follows from the equality  $|z|^2 = (ac)^2 + (bd)^2 + (ad)^2 + (bc)^2$ . For  $A$ , define  $\pi = abcd$  and notice that  $A = |\pi - (bd)^2| + |\pi + (bc)^2|$  is equal to either  $B$  or  $\pm(2\pi + (bc)^2 - (bd)^2)$ ; furthermore, in the latter case we have

$$\begin{aligned}A &\leq 2|\pi| + |(bc)^2 - (bd)^2| \\ &\leq \min \{(ac)^2 + (bd)^2, (ad)^2 + (bc)^2\} + \max \{(bc)^2, (bd)^2\} \\ &\leq |z|^2.\end{aligned}$$

Thus, combining (5.7) and (5.8),  $|z - \widehat{z}| \leq (\lambda + \mu)|z|$ . Since the refined bound  $E_1(t) \leq u/(1+u)$  implies  $|\delta_i| \leq u/(1+u)$  for all  $i$ , we can take  $\lambda = u/(1+u)$  and  $\mu = u/(1+u) \cdot (1 + u/(1+u))$ , which are both less than  $u$ . Hence, barring underflow and overflow and since  $z = 0$  implies  $\widehat{z} = 0$ , we conclude that

$$\widehat{z} = z(1 + \epsilon), \quad |\epsilon| \leq \frac{2u+3u^2}{(1+u)^2} < 2u.$$

Note that  $\frac{2u+3u^2}{(1+u)^2}$  has the form  $2u - u^2 + \mathcal{O}(u^4)$  as  $u \rightarrow 0$ . Thus, our approach not only yields a shorter and more direct proof of the bound  $2u$  of [8], but it also improves on that bound.

#### APPENDIX A. PROOF THAT $p$ AS IN (3.4) IS NOT A POWER OF TWO

If  $j \in \{5, 6, 7\}$ , then  $p \in \{144, 175, 206\}$  and is not a power of two. Assume now that  $j \geq 8$ . Writing  $d = 4m + i$  for integers  $m, i$  with  $0 \leq i \leq 3$ , we have

$$4m + i - \frac{1}{2} \leq 4 \log_2 k \leq 4m + i + \frac{1}{2}.$$

If  $i \neq 0$ , this implies  $2^{m+1/8} \leq k \leq 2^{m+7/8}$  and then

$$2^{m+1/8} - 4m + 10 \leq p \leq 2^{m+7/8} - 4m + 12.$$

---

<sup>3</sup>There are three other ways to insert the innermost rounding fl, all giving the same error as the one developed here.

Since the assumption  $j \geq 8$  implies  $k \geq 2^8$  and thus  $m \geq 8$ , it follows that  $2^m < p < 2^{m+1}$ . Consequently,  $p$  cannot be a power of two when  $i \neq 0$ . On the other hand, when  $i = 0$ , we see that  $p = 32j - 4m + 13$  must be odd, and thus cannot be a power of two neither.

#### APPENDIX B. PROOF OF THEOREM 4.1

The lower bound follows from (4.1) and the triangle inequality. For the upper bound, the proof is by induction on  $n$ , the case  $n = 1$  being trivial (since then there is no rounding error at all). For  $n \geq 2$ , we assume that the result is true up to  $n - 1$ , and we fix one evaluation order in dimension  $n$ . The approximation  $\widehat{s}$  obtained with this order has the form  $\widehat{s} = \text{fl}(\widehat{s}_1 + \widehat{s}_2)$ , where  $\widehat{s}_j$  is the result of a floating-point evaluation of  $s_j = \sum_{i \in I_j} x_i$  for  $j = 1, 2$  and with  $\{I_1, I_2\}$  a partition of the set  $I = \{1, 2, \dots, n\}$ . For  $j = 1, 2$  let  $n_j$  be the cardinality of  $I_j$  and let  $e_j^{(1)}, \dots, e_j^{(n_j-1)}$  be the rounding errors committed when evaluating  $s_j$ , so that  $s_j = \sum_{i < n_j} e_j^{(i)} + \widehat{s}_j$ . Consequently,

$$\sum_{i < n} e_i = \delta + \sum_{i < n_1} e_1^{(i)} + \sum_{i < n_2} e_2^{(i)}, \quad \delta = \widehat{s}_1 + \widehat{s}_2 - \text{fl}(\widehat{s}_1 + \widehat{s}_2).$$

Since  $1 \leq n_j < n$ , the inductive assumption leads to

$$\sum_{i < n} |e_i| \leq |\delta| + \frac{u}{1+u} \left( (n_1 - 1)\widetilde{s}_1 + (n_2 - 1)\widetilde{s}_2 \right),$$

where  $\widetilde{s}_j = \sum_{i \in I_j} |x_i|$  for  $j = 1, 2$ . Since  $n = n_1 + n_2$  and  $\sum_{i=1}^n |x_i| = \widetilde{s}_1 + \widetilde{s}_2$ , it remains to check that

$$(B.1) \quad |\delta| \leq \frac{u}{1+u} (n_2 \widetilde{s}_1 + n_1 \widetilde{s}_2).$$

To do so, we note that (1.2) and [10, Lemma 2.2] imply that

$$|\delta| \leq \min \left\{ |\widehat{s}_1|, |\widehat{s}_2|, \frac{u}{1+u} |\widehat{s}_1 + \widehat{s}_2| \right\},$$

and we consider the following three cases. Assume first that  $\widetilde{s}_2 \leq u/(1+u) \cdot \widetilde{s}_1$ . Then  $\widetilde{s}_2 \leq \widetilde{s}_1$  and, using  $|\delta| \leq |\widehat{s}_2|$ , we obtain

$$|\delta| \leq |\widehat{s}_2 - s_2| + \widetilde{s}_2 \leq \sum_{i < n_2} |e_2^{(i)}| + \widetilde{s}_2 \leq (n_2 - 1) \frac{u}{1+u} \widetilde{s}_2 + \frac{u}{1+u} \widetilde{s}_1 \leq n_2 \frac{u}{1+u} \widetilde{s}_1$$

and (B.1) thus follows. Second, when  $\widetilde{s}_1 \leq u/(1+u) \cdot \widetilde{s}_2$  we proceed similarly, simply swapping the indices 1 and 2. Third, when  $u/(1+u) \cdot \widetilde{s}_1 < \widetilde{s}_2$  and  $u/(1+u) \cdot \widetilde{s}_2 < \widetilde{s}_1$ , we have  $|\delta| \leq u/(1+u) \cdot |\widehat{s}_1 + \widehat{s}_2|$  with

$$|\widehat{s}_1 + \widehat{s}_2| \leq |\widehat{s}_1 - s_1| + \widetilde{s}_1 + |\widehat{s}_2 - s_2| + \widetilde{s}_2$$

and  $|\widehat{s}_j - s_j| \leq (n_j - 1)u/(1+u) \cdot \widetilde{s}_j \leq (n_j - 1)\widetilde{s}_k$  for  $(j, k) \in \{(1, 2), (2, 1)\}$ . Hence (B.1) follows in this third case as well, thus completing the proof of Theorem 4.1.



## REFERENCES

- [1] R. Brent, C. Percival, and P. Zimmermann, *Error bounds on complex floating-point multiplication*, Math. Comp. **76** (2007), 1469–1481.
- [2] T. J. Dekker, *A floating-point technique for extending the available precision*, Numer. Math. **18** (1971/72), 224–242.
- [3] R. L. Graham, D. E. Knuth, and O. Patashnik, *Concrete Mathematics: A Foundation for Computer Science*, 2nd ed., Addison-Wesley, Reading, MA, 1994.
- [4] G. H. Hardy, J. E. Littlewood, and G. Pólya, *Inequalities*, 2nd ed., Cambridge University Press, Cambridge, UK, 1952.
- [5] N. J. Higham, *Accuracy and Stability of Numerical Algorithms*, 2nd ed., Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002.
- [6] J. E. Holm, *Floating-Point Arithmetic and Program Correctness Proofs*, Ph.D. thesis, Cornell University, Ithaca, NY, August 1980, pp. vii+133.
- [7] T. E. Hull, T. F. Fairgrieve, and P. T. P. Tang, *Implementing complex elementary functions using exception handling*, ACM Trans. Math. Software **20** (1994), no. 2, 215–244.
- [8] C.-P. Jeannerod, P. Kornerup, N. Louvet, and J.-M. Muller, *Error bounds on complex floating-point multiplication with an FMA*, Math. Comp. (published electronically), 2016.
- [9] C.-P. Jeannerod, N. Louvet, J.-M. Muller, and A. Panhaleux, *Midpoints and exact points of some algebraic functions in floating-point arithmetic*, IEEE Trans. Comput. **60** (2011), no. 2, 228–241.
- [10] C.-P. Jeannerod and S. M. Rump, *Improved error bounds for inner products in floating-point arithmetic*, SIAM J. Matrix Anal. Appl. **34** (2013), no. 2, 338–344.
- [11] W. Keller, *Prime factors  $k \cdot 2^n + 1$  of Fermat numbers  $F_m$  and complete factoring status*, August 2016, web page available at <http://www.prothsearch.net/fermat.html>.
- [12] D. E. Knuth, *The Art of Computer Programming, Volume 2, Seminumerical Algorithms*, 3rd ed., Addison-Wesley, Reading, MA, 1998.
- [13] P. W. Markstein, *Computation of elementary functions on the IBM RISC System/6000 processor*, IBM Journal of Research and Development **34** (1990), no. 1, 111–119.
- [14] T. Ogita, S. M. Rump, and S. Oishi, *Accurate sum and dot product*, SIAM J. Sci. Comput. **26** (2005), no. 6, 1955–1988.
- [15] S. M. Rump and C.-P. Jeannerod, *Improved backward error bounds for LU and Cholesky factorizations*, SIAM J. Matrix Anal. Appl. **35** (2014), no. 2, 684–698.
- [16] S. M. Rump, T. Ogita, and S. Oishi, *Accurate floating-point summation part I: Faithful rounding*, SIAM J. Sci. Comput. **31** (2008), no. 1, 189–224.
- [17] N. J. A. Sloane and D. W. Wilson, *Sequence A019434*, The On-Line Encyclopedia of Integer Sequences, published electronically at <http://oeis.org>.
- [18] P. H. Sterbenz, *Floating-Point Computation*, Prentice-Hall, 1974.
- [19] J. H. Wilkinson, *Error analysis of floating-point computation*, Numer. Math. **2** (1960), 319–340.

INRIA AND UNIVERSITÉ DE LYON, LABORATOIRE LIP (CNRS, ENS DE LYON, INRIA, UCBL),  
 46 ALLÉE D’ITALIE 69364 LYON CEDEX 07, FRANCE  
*E-mail address:* `claude-pierre.jeannerod@inria.fr`

HAMBURG UNIVERSITY OF TECHNOLOGY, SCHWARZENBERGSTRASSE 95, HAMBURG 21071, GER-  
 MANY, AND FACULTY OF SCIENCE AND ENGINEERING, WASEDA UNIVERSITY, 3-4-1 OKUBO, SHINJUKU-  
 KU, TOKYO 169-8555, JAPAN  
*E-mail address:* `rump@tuhh.de`