



Accurate Floating Point Product and Exponentiation

Stef Graillat

► To cite this version:

Stef Graillat. Accurate Floating Point Product and Exponentiation. IEEE Transactions on Computers, 2009, 58 (7), pp.994-1000. 10.1109/TC.2008.215 . hal-00164607

HAL Id: hal-00164607

<https://hal.science/hal-00164607>

Submitted on 21 Jul 2007

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.

Accurate Floating Point Product and Exponentiation

Stef Graillat*

July 23, 2007

Abstract

Several different techniques and softwares intend to improve the accuracy of results computed in a fixed finite precision. Here we focus on a method to improve the accuracy of the product of floating point numbers. We show that the computed result is as accurate as if computed in twice the working precision. The algorithm is simple since it only requires addition, subtraction and multiplication of floating point numbers in the same working precision as the given data. Such an algorithm can be useful for example to compute the determinant of a triangular matrix and to evaluate a polynomial when represented by the root product form. It can also be used to compute the power of a floating point number.

Key words: accurate product, exponentiation, finite precision, floating point arithmetic, faithful rounding, error-free transformations

AMS Subject Classifications: 65-04, 65G20, 65G50

1 Introduction

In this paper, we present fast and accurate algorithms to compute the product of floating point numbers. Our aim is to increase the accuracy at a fixed precision. We show that the results have the same error estimates as if computed in twice the working precision and then rounded to working precision. Then we address the problem on how to compute a faithfully rounded result, that is to say one of the two adjacent floating point numbers of the exact result.

This paper was motivated by papers [15, 18, 6, 12] and [11] where similar approaches are used to compute summation, dot product, polynomial evaluation and power.

The applications of our algorithms are multiple. One of the examples frequently used in Sterbenz's book [20] is the computation of the product of some floating point numbers. Our algorithms can be used to compute the determinant of a triangle matrix

$$T = \begin{bmatrix} t_{11} & t_{12} & \cdots & t_{1n} \\ & t_{22} & & t_{2n} \\ & & \ddots & \vdots \\ & & & t_{nn} \end{bmatrix}.$$

Indeed, the determinant of T is

$$\det(T) = \prod_{i=1}^n t_{ii}.$$

*Laboratoire LIP6, Département Calcul Scientifique, Université Pierre et Marie Curie (Paris 6), 4 place Jussieu, F-75252, Paris cedex 05, France (stef.graillat@lip6.fr, <http://www-pequan.lip6.fr/~graillat>).

Another application is for evaluating a polynomial when represented by the root product form $p(x) = a_n \prod_{i=1}^n (x - x_i)$. We can also apply our algorithms to compute the power of a floating point number.

The rest of the paper is organized as follows. In Section 2, we recall notations and auxiliary results that will be needed in the sequel. We present the floating point arithmetic and the so-called error-free transformations. In Section 3, we present a classic algorithm to compute the product of floating point numbers. We give an error estimate as well as a validated error bound. We also present a new compensated algorithm together with an error estimate and a validated error bound. We show that under mild assumptions, our algorithm gives a faithfully rounded result. In Section 4, we apply our algorithm to compute the power of a floating point number. We propose two different algorithms: one with our compensated algorithm, the other one with the use of a double-double library.

2 Notation and auxiliary results

2.1 Floating point arithmetic

Throughout the paper, we assume to work with a floating point arithmetic adhering to IEEE 754 floating point standard in rounding to nearest [9]. We assume that no overflow nor underflow occurs. The set of floating point numbers is denoted by \mathbb{F} , the relative rounding error by \mathbf{eps} . For IEEE 754 double precision, we have $\mathbf{eps} = 2^{-53}$ and for single precision $\mathbf{eps} = 2^{-24}$.

We denote by $\mathbf{fl}(\cdot)$ the result of a floating point computation, where all operations inside parentheses are done in floating point working precision. Floating point operations in IEEE 754 satisfy [8]

$$\mathbf{fl}(a \circ b) = (a \circ b)(1 + \varepsilon_1) = (a \circ b)/(1 + \varepsilon_2) \text{ for } \circ = \{+, -, \cdot, /\} \text{ and } |\varepsilon_\nu| \leq \mathbf{eps}.$$

This implies that

$$|a \circ b - \mathbf{fl}(a \circ b)| \leq \mathbf{eps}|a \circ b| \text{ and } |a \circ b - \mathbf{fl}(a \circ b)| \leq \mathbf{eps}|\mathbf{fl}(a \circ b)| \text{ for } \circ = \{+, -, \cdot, /\} \quad (2.1)$$

2.2 Error-free transformations

One can notice that $a \circ b \in \mathbb{R}$ and $\mathbf{fl}(a \circ b) \in \mathbb{F}$ but in general we do not have $a \circ b \in \mathbb{F}$. It is known that for the basic operations $+, -, \cdot$, the approximation error of a floating point operation is still a floating point number (see for example [5]):

$$\begin{aligned} x = \mathbf{fl}(a \pm b) &\Rightarrow a \pm b = x + y && \text{with } y \in \mathbb{F}, \\ x = \mathbf{fl}(a \cdot b) &\Rightarrow a \cdot b = x + y && \text{with } y \in \mathbb{F}. \end{aligned} \quad (2.2)$$

These are *error-free* transformations of the pair (a, b) into the pair (x, y) .

Fortunately, the quantities x and y in (2.2) can be computed exactly in floating point arithmetic. For the algorithms, we use Matlab-like notations. For addition, we can use the following algorithm by Knuth [10, Thm B. p.236].

Algorithm 2.1 (Knuth [10]). Error-free transformation of the sum of two floating point numbers

```
function [x, y] = TwoSum(a, b)
    x = fl(a + b)
    z = fl(x - a)
    y = fl((a - (x - z)) + (b - z))
```

Another algorithm to compute an error-free transformation is the following algorithm from Dekker [5]. The drawback of this algorithm is that we have $x + y = a + b$ provided that $|a| \geq |b|$. Generally, on modern computers, a comparison followed by a branching and 3 operations costs more than 6 operations. As a consequence, **TwoSum** is generally more efficient than **FastTwoSum**.

Algorithm 2.2 (Dekker [5]). Error-free transformation of the sum of two floating point numbers.

```
function  $[x, y] = \text{FastTwoSum}(a, b)$ 
   $x = \text{fl}(a + b)$ 
   $y = \text{fl}((a - x) + b)$ 
```

For the error-free transformation of a product, we first need to split the input argument into two parts. Let p be given by $\text{eps} = 2^{-p}$ and define $s = \lceil p/2 \rceil$. For example, if the working precision is IEEE 754 double precision, then $p = 53$ and $s = 27$. The following algorithm by Dekker [5] splits a floating point number $a \in \mathbb{F}$ into two parts x and y such that

$$a = x + y \quad \text{and} \quad x \text{ and } y \text{ nonoverlapping with } |y| \leq |x|.$$

Algorithm 2.3 (Dekker [5]). Error-free split of a floating point number into two parts

```
function  $[x, y] = \text{Split}(a, b)$ 
   $\text{factor} = \text{fl}(2^s + 1)$ 
   $c = \text{fl}(\text{factor} \cdot a)$ 
   $x = \text{fl}(c - (c - a))$ 
   $y = \text{fl}(a - x)$ 
```

With this function, an algorithm from Veltpkamp (see [5]) enables to compute an error-free transformation for the product of two floating point numbers. This algorithm returns two floating point numbers x and y such that

$$a \cdot b = x + y \quad \text{with } x = \text{fl}(a \cdot b).$$

Algorithm 2.4 (Veltpkamp [5]). Error-free transformation of the product of two floating point numbers

```
function  $[x, y] = \text{TwoProduct}(a, b)$ 
   $x = \text{fl}(a \cdot b)$ 
   $[a_1, a_2] = \text{Split}(a)$ 
   $[b_1, b_2] = \text{Split}(b)$ 
   $y = \text{fl}(a_2 \cdot b_2 - (((x - a_1 \cdot b_1) - a_2 \cdot b_1) - a_1 \cdot b_2))$ 
```

The following theorem summarizes the properties of algorithms **TwoSum** and **TwoProduct**.

Theorem 2.1 (Ogita, Rump and Oishi [15]). Let $a, b \in \mathbb{F}$ and let $x, y \in \mathbb{F}$ such that $[x, y] = \text{TwoSum}(a, b)$ (Algorithm 2.1). Then,

$$a + b = x + y, \quad x = \text{fl}(a + b), \quad |y| \leq \text{eps}|x|, \quad |y| \leq \text{eps}|a + b|. \quad (2.3)$$

The algorithm **TwoSum** requires 6 flops.

Let $a, b \in \mathbb{F}$ and let $x, y \in \mathbb{F}$ such that $[x, y] = \text{TwoProduct}(a, b)$ (Algorithm 2.4). Then,

$$a \cdot b = x + y, \quad x = \text{fl}(a \cdot b), \quad |y| \leq \text{eps}|x|, \quad |y| \leq \text{eps}|a \cdot b|. \quad (2.4)$$

The algorithm **TwoProduct** requires 17 flops.

The `TwoProduct` algorithm can be re-written in a very simple way if a Fused-Multiply-and-Add (FMA) operator is available on the targeted architecture [14]. Some computers have a *Fused-Multiply-and-Add* (FMA) operation that enables a floating point multiplication followed by an addition to be performed as a single floating point operation. The Intel IA-64 architecture, implemented in the Intel Itanium processor, has an FMA instruction as well as the IBM RS/6000 and the PowerPC before it. On the Itanium processor, the FMA instruction enables a multiplication and an addition to be performed in the same number of cycles than one multiplication or one addition. As a result, it seems to be advantageous for speed as well as for accuracy.

Theoretically, this means that for $a, b, c \in \mathbb{F}$, the result of $\text{FMA}(a, b, c)$ is the nearest floating point number of $a \cdot b + c \in \mathbb{R}$. The FMA satisfies

$$\text{FMA}(a, b, c) = (a \cdot b + c)(1 + \varepsilon_1) = (a \cdot b + c)/(1 + \varepsilon_2) \text{ with } |\varepsilon_\nu| \leq \text{eps}.$$

Thanks to the FMA, the `TwoProduct` algorithm can be re-written as follows which costs only 2 flops.

Algorithm 2.5 (Ogita, Rump and Oishi [15]). Error-free transformation of the product of two floating point numbers using an FMA.

```
function [x, y] = TwoProductFMA(a, b)
    x = a · b
    y = FMA(a, b, -x)
```

3 Accurate floating point product

In this section, we present a new accurate algorithm to compute the product of floating point numbers. In Subsection 3.1, we recall the classic method and we give a theoretical error bound as well as a validated computable error bound. In Subsection 3.2, we present our new algorithm based on a compensated scheme together with a theoretical error bound. In Subsection 3.3, we give sufficient conditions on the number of floating point numbers so as to get a faithfully rounded result. Finally, in Subsection 3.4, we give a validated computable error bound for our new algorithm.

3.1 Classic method

The classic method for evaluating a product of n numbers $a = (a_1, a_2, \dots, a_n)$

$$p = \prod_{i=1}^n a_i$$

is the following algorithm.

Algorithm 3.1. Product evaluation

```
function res = Prod(a)
    p1 = a1
    for i = 2 : n
        pi = fl(pi-1 · ai)
    end
    res = pn
```

This algorithm requires $n - 1$ flops. Let us now analyse its accuracy.

We will use standard notations and standard results for the following error estimations (see [8]). The quantities γ_n are defined as usual [8] by

$$\gamma_n := \frac{n\mathbf{eps}}{1 - n\mathbf{eps}} \quad \text{for } n \in \mathbb{N}.$$

When using γ_n , we implicitly assume that $n\mathbf{eps} \leq 1$. A forward error bound is

$$|a_1 a_2 \cdots a_n - \mathbf{res}| = |a_1 a_2 \cdots a_n - \mathbf{fl}(a_1 a_2 \cdots a_n)| \leq \gamma_{n-1} |a_1 a_2 \cdots a_n|. \quad (3.5)$$

Indeed, by induction,

$$\mathbf{res} = \mathbf{fl}(a_1 a_2 \cdots a_n) = a_1 a_2 \cdots a_n (1 + \varepsilon_2)(1 + \varepsilon_3) \cdots (1 + \varepsilon_n), \quad (3.6)$$

with $\varepsilon_i \leq \mathbf{eps}$ for $i = 2 : n$. It follows from Lemma 3.1 of [8, p.63] that $(1 + \varepsilon_2)(1 + \varepsilon_3) \cdots (1 + \varepsilon_n) = 1 + \theta_n$ where $|\theta_n| \leq \gamma_{n-1}$.

A convenient device for keeping track of power of $1 + \varepsilon$ term is describe in [8, p.68]. The relative error counter $\langle k \rangle$ denotes the product

$$\langle k \rangle = \prod_{i=1}^k (1 + \varepsilon_i), \quad |\varepsilon_i| \leq \mathbf{eps}.$$

A useful rule for the counter is $\langle j \rangle \langle k \rangle = \langle j + k \rangle$. Using this notation, Equation (3.6) can be written $\mathbf{res} = \mathbf{fl}(a_1 a_2 \cdots a_n) = a_1 a_2 \cdots a_n \langle n - 1 \rangle$.

It is shown in [16] that for $a \in \mathbb{F}$, we have

$$\begin{aligned} (1 + \mathbf{eps})^n &\leq \frac{1}{(1 - \mathbf{eps})^n} \leq \frac{1}{1 - n\mathbf{eps}}, \\ \frac{|a|}{1 - n\mathbf{eps}} &\leq \mathbf{fl}\left(\frac{|a|}{1 - (n + 1)\mathbf{eps}}\right). \end{aligned} \quad (3.7)$$

From Equation (3.6), it follows that

$$|a_1 a_2 \cdots a_n - \mathbf{res}| \leq (1 + \mathbf{eps})^{n-1} \gamma_{n-1} |\mathbf{res}|.$$

If $m\mathbf{eps} \leq 1$ for $m \in \mathbb{N}$, $\mathbf{fl}(m\mathbf{eps}) = m\mathbf{eps}$ and $\mathbf{fl}(1 - m\mathbf{eps}) = 1 - m\mathbf{eps}$. Therefore,

$$\gamma_m \leq (1 + \mathbf{eps}) \mathbf{fl}(\gamma_m). \quad (3.8)$$

Hence,

$$\begin{aligned} |a_1 a_2 \cdots a_n - \mathbf{res}| &\leq (1 + \mathbf{eps})^n \mathbf{fl}(\gamma_{n-1}) |\mathbf{res}| \\ &\leq (1 + \mathbf{eps})^{n+1} \mathbf{fl}(\gamma_{n-1} |\mathbf{res}|), \end{aligned}$$

and so

$$|a_1 a_2 \cdots a_n - \mathbf{res}| \leq \mathbf{fl}\left(\frac{\gamma_{n-1} |\mathbf{res}|}{1 - (n + 2)\mathbf{eps}}\right).$$

The previous inequality gives us a validated error bound that can be computed in pure floating point arithmetic in rounding to nearest.

3.2 Compensated method

We present hereafter a compensated scheme to evaluate the product of floating point numbers, i.e. the error of individual multiplication is somehow corrected. The technique used here is based on the paper [15].

Algorithm 3.2. Product evaluation with a compensated scheme

```
function res = CompProd(a)
    p1 = a1
    e1 = 0
    for i = 2 : n
        [p_i, pi] = TwoProduct(p_{i-1}, a_i)
        e_i = fl(e_{i-1}a_i + pi)
    end
    res = fl(p_n + e_n)
```

This algorithm requires $19n - 18$ flops if we use **TwoProduct**. It only requires $4n - 3$ flops if we use **TwoProductFMA** instead of **TwoProduct** (if, of course, an FMA is available).

Algorithm 3.3. Product evaluation with a compensated scheme with **TwoProductFMA**

```
function res = CompProdFMA1(a)
    p1 = a1
    e1 = 0
    for i = 2 : n
        [p_i, pi] = TwoProductFMA(p_{i-1}, a_i)
        e_i = fl(e_{i-1}a_i + pi)
    end
    res = fl(p_n + e_n)
```

It finally costs $3n - 2$ flops if we also use $e_i = \text{FMA}(e_{i-1}, a_i, \pi_i)$ instead of $e_i = \text{fl}(e_{i-1}a_i + \pi_i)$.

Algorithm 3.4. Product evaluation with a compensated scheme with **TwoProductFMA** and FMA

```
function res = CompProdFMA2(a)
    p1 = a1
    e1 = 0
    for i = 2 : n
        [p_i, pi] = TwoProductFMA(p_{i-1}, a_i)
        e_i = FMA(e_{i-1}, a_i, pi)
    end
    res = fl(p_n + e_n)
```

We will provide an error analysis only for Algorithm **CompProd**. The error analysis for **CompProdFMA1** is the same as **CompProd** since they share the same operations. For the error analysis of **CompProdFMA2**, little changes must be done to take into account the operation $e_i = \text{FMA}(e_{i-1}, a_i, \pi_i)$. This changes nearly nothing so it is straightforward to modify the analysis to deal with it.

For error analysis, we note that

$$p_n = \text{fl}(a_1 a_2 \cdots a_n) \quad \text{and} \quad e_n = \text{fl} \left(\sum_{i=2}^n \pi_i a_{i+1} \cdots a_n \right).$$

We also have

$$p = a_1 a_2 \dots a_n = \text{fl}(a_1 a_2 \dots a_n) + \sum_{i=2}^n \pi_i a_{i+1} \dots a_n = p_n + e, \quad (3.9)$$

where $e = \sum_{i=2}^n \pi_i a_{i+1} \dots a_n$.

Before proving the main theorem, we will need two technical lemmas. The next lemma makes it possible to obtain a bound on the individual error of the multiplication namely π_i in function of the initial data a_i .

Lemma 3.1. *Suppose floating point numbers $\pi_i \in \mathbb{F}$, $2 \leq i \leq n$ are computed by the following algorithm*

```
p1 = a1
for i = 2 : n
    [pi, pi_i] = TwoProduct(p_{i-1}, a_i)
end
```

Then,

$$|\pi_i| \leq \mathbf{eps}(1 + \gamma_{i-1})|a_1 \dots a_i| \quad \text{for } i = 2 : n.$$

Proof. From Equation (2.1), it follows that

$$|\pi_i| \leq \mathbf{eps}|p_i|.$$

Moreover, $p_i = \text{fl}(a_1 \dots a_i)$ so that from (3.5),

$$|p_i| \leq (1 + \gamma_{i-1})|a_1 \dots a_i|.$$

Hence, $|\pi_i| \leq \mathbf{eps}(1 + \gamma_{i-1})|a_1 \dots a_i|$. □

The following lemma enables us to bound the rounding errors during the computation of the error during the full product.

Lemma 3.2. *Suppose floating point numbers $e_i \in \mathbb{F}$, $1 \leq i \leq n$ are computed by the following algorithm*

```
e1 = 0
for i = 2 : n
    [pi, pi_i] = TwoProduct(p_{i-1}, a_i)
    e_i = fl(e_{i-1} a_i + pi_i)
end
```

Then,

$$|e_n - \sum_{i=2}^n \pi_i a_{i+1} \dots a_n| \leq \gamma_{n-1} \gamma_{2n} |a_1 a_2 \dots a_n|.$$

Proof. First, one notices that $e_n = \text{fl}(\sum_{i=2}^n (\pi_i a_{i+1} \dots a_n))$. We will use the error counters described above. For n floating point numbers x_i , it is easy to see that [8, chap.4]

$$\text{fl}(x_1 + x_2 + \dots + x_n) = x_1 \langle n-1 \rangle + x_2 \langle n-1 \rangle + x_3 \langle n-2 \rangle + \dots + x_n \langle 1 \rangle.$$

This implies that

$$e_n = \text{fl}\left(\sum_{i=2}^n (\pi_i a_{i+1} \dots a_n)\right) = \text{fl}(\pi_2 a_3 \dots a_n) \langle n-2 \rangle + \text{fl}(\pi_3 a_4 \dots a_n) \langle n-2 \rangle + \dots + \text{fl}(\pi_n) \langle 1 \rangle.$$

Furthermore, we have shown before that $\text{fl}(a_1 a_2 \cdots a_n) = a_1 a_2 \cdots a_n \langle n-1 \rangle$. Consequently,

$$e_n = \pi_2 a_3 \cdots a_n \langle n-2 \rangle \langle n-1 \rangle + \pi_3 a_4 \cdots a_n \langle n-3 \rangle \langle n-1 \rangle + \cdots + \pi_n \langle 1 \rangle.$$

A straightforward computation yields

$$|e_n - \sum_{i=2}^n \pi_i a_{i+1} \cdots a_n| \leq \gamma_{2n-3} \sum_{i=2}^n |\pi_i a_{i+1} \cdots a_n|.$$

From Lemma 3.1, we have $|\pi_i| \leq \text{eps}(1 + \gamma_{i-1})|a_1 \cdots a_i|$ and hence

$$|e_n - \sum_{i=2}^n \pi_i a_{i+1} \cdots a_n| \leq (n-1)\text{eps}(1 + \gamma_{n-1})\gamma_{2n-3}|a_1 a_2 \cdots a_n|.$$

Since $\text{eps}(1 + \gamma_{n-1}) = \gamma_{n-1}/(n-1)$ and $\gamma_{2n-3} \leq \gamma_{2n}$, we obtain the desired result. \square

One may notice that the computation of e_n is similar to the Horner scheme. One could have directly applied a result on the error of the Horner scheme [8, Eq.(5.3),p.95].

We can finally state the main theorem.

Theorem 3.3. *Suppose Algorithm 3.2 is applied to floating point number $a_i \in \mathbb{F}$, $1 \leq i \leq n$, and set $p = \prod_{i=1}^n a_i$. Then,*

$$|\text{res} - p| \leq \text{eps}|p| + \gamma_n \gamma_{2n}|p|.$$

Proof. The fact that $\text{res} = \text{fl}(p_n + e_n)$ implies that $\text{res} = (1 + \varepsilon)(p_n + e_n)$ with $\varepsilon \leq \text{eps}$. So it follows

$$\begin{aligned} |\text{res} - p| &= |\text{fl}(p_n + e_n) - p| = |(1 + \varepsilon)(p_n + e_n - p) + \varepsilon p| \\ &= |(1 + \varepsilon)(p_n + \sum_{i=2}^n \pi_i a_{i+1} \cdots a_n - p) + (1 + \varepsilon)(e_n - \sum_{i=2}^n \pi_i a_{i+1} \cdots a_n) + \varepsilon p| \\ &= |(1 + \varepsilon)(e_n - \sum_{i=2}^n \pi_i a_{i+1} \cdots a_n) + \varepsilon p| \quad \text{by (3.9)} \\ &\leq \text{eps}|p| + (1 + \text{eps})|e_n - \sum_{i=2}^n \pi_i a_{i+1} \cdots a_n| \\ &\leq \text{eps}|p| + (1 + \text{eps})\gamma_{n-1}\gamma_{2n}|a_1 a_2 \cdots a_n|. \end{aligned}$$

Since $(1 + \text{eps})\gamma_{n-1} \leq \gamma_n$, it follows that $|\text{res} - p| \leq \text{eps}|p| + \gamma_n \gamma_{2n}|p|$. \square

It may be interesting to study the condition number of the product evaluation. One defines

$$\text{cond}(a) = \limsup_{\varepsilon \rightarrow 0} \left\{ \frac{|(a_1 + \Delta a_1)(a_2 + \Delta a_2) \cdots (a_n + \Delta a_n) - a_1 a_2 \cdots a_n|}{\varepsilon |a_1 a_2 \cdots a_n|} : |\Delta a_i| \leq \varepsilon |a_i| \right\}.$$

A standard computation yields

$$\text{cond}(a) = n.$$

Corollary 3.4. *Suppose Algorithm 3.2 is applied to floating point number $a_i \in \mathbb{F}$, $1 \leq i \leq n$, and set $p = \prod_{i=1}^n a_i$. Then,*

$$\frac{|\text{res} - p|}{|p|} \leq \text{eps} + \frac{\gamma_n \gamma_{2n}}{n} \text{cond}(a).$$

3.3 Faithful rounding

We define the floating point predecessor and successor of a real number r satisfying $\min\{f : f \in \mathbb{R}\} < r < \max\{f : f \in \mathbb{F}\}$ by

$$\text{pred}(r) := \max\{f \in \mathbb{F} : f < r\} \quad \text{and} \quad \text{succ}(r) := \min\{f \in \mathbb{F} : r < f\}.$$

Definition 3.1. A floating point number $f \in \mathbb{F}$ is called a *faithful rounding* of a real number $r \in \mathbb{R}$ if

$$\text{pred}(f) < r < \text{succ}(f).$$

We denote this by $f \in \square(r)$. For $r \in \mathbb{F}$, this implies that $f = r$.

A faithful rounding is then one of the two adjacent floating point numbers of the exact result.

Lemma 3.5 (Rump, Ogita and Oishi [18, lem. 2.5]). *Let $r, \delta \in \mathbb{R}$ and $\tilde{r} := \text{fl}(r)$. Suppose that $2|\delta| < \text{eps}|\tilde{r}|$. Then $\tilde{r} \in \square(r + \delta)$, that means \tilde{r} is a faithful rounding of $r + \delta$.*

Let **res** be the result of **CompProd**. Then we have $p = p_n + e$ and **res** = $\text{fl}(p_n + e_n)$ with $e = \sum_{i=2}^n \pi_i a_{i+1} \cdots a_n$. It follows that $p = (p_n + e_n) + (e - e_n)$. This leads to the following lemma which gives a criterion to ensure that the result of **CompProd** is faithfully rounded.

Lemma 3.6. *With the previous notations, if $2|e - e_n| < \text{eps}|\text{res}|$ then **res** is a faithful rounding of p .*

Since we have $|e - e_n| \leq \gamma_n \gamma_{2n} |p|$ and $(1 - \text{eps})|p| - \gamma_n \gamma_{2n} |p| \leq |\text{res}|$, a sufficient condition to ensure a faithful rounding is

$$2\gamma_n \gamma_{2n} |p| < \text{eps}((1 - \text{eps})|p| - \gamma_n \gamma_{2n} |p|)$$

that is

$$\gamma_n \gamma_{2n} < \frac{1 - \text{eps}}{2 + \text{eps}} \text{eps}.$$

Since $\gamma_n \gamma_{2n} \leq 2(n\text{eps})^2 / (1 - 2n\text{eps})^2$, a sufficient condition is

$$2 \frac{(n\text{eps})^2}{(1 - 2n\text{eps})^2} < \frac{1 - \text{eps}}{2 + \text{eps}} \text{eps}$$

which is equivalent to

$$\frac{n\text{eps}}{1 - 2n\text{eps}} < \sqrt{\frac{(1 - \text{eps})\text{eps}}{2(2 + \text{eps})}}$$

and then to

$$n < \frac{\sqrt{1 - \text{eps}}}{\sqrt{2}\sqrt{2 + \text{eps}} + 2\sqrt{(1 - \text{eps})\text{eps}}} \text{eps}^{-1/2}.$$

We have just shown that if $n < \alpha \text{eps}^{-1/2}$ where $\alpha \approx 1/2$ then the result is faithfully rounded. More precisely, in double precision where $\text{eps} = 2^{-53}$, if $n < 2^{25} \approx 5 \cdot 10^7$, we get a faithfully rounded result.

3.4 Validated error bound

We present here how to compute a valid error bound in pure floating point arithmetic in rounding to nearest. It holds that

$$\begin{aligned} |\mathbf{res} - p| &= |\mathbf{fl}(p_n + e_n) - p| = |\mathbf{fl}(p_n + e_n) - (p_n + e_n) + (p_n + e_n) - p| \\ &\leq \mathbf{eps}|\mathbf{res}| + |p_n + e_n - p| \\ &\leq \mathbf{eps}|\mathbf{res}| + |e_n - e|. \end{aligned}$$

Since $|e_n - e| \leq \gamma_{n-1}\gamma_{2n}|p|$ and $|p| \leq (1 + \mathbf{eps})^{n-1} \mathbf{fl}(|a_1 a_2 \cdots a_n|)$ we obtain

$$\begin{aligned} |\mathbf{res} - p| &\leq \mathbf{eps}|\mathbf{res}| + \gamma_{n-1}\gamma_{2n}|p| \\ &\leq \mathbf{eps}|\mathbf{res}| + \gamma_{n-1}\gamma_{2n}(1 + \mathbf{eps})^{n-1} \mathbf{fl}(|a_1 a_2 \cdots a_n|). \end{aligned}$$

Using (3.7) and (3.8), we get

$$\begin{aligned} |\mathbf{res} - p| &\leq \mathbf{fl}(\mathbf{eps}|\mathbf{res}|) + (1 + \mathbf{eps})^n \mathbf{fl}(\gamma_n) \mathbf{fl}(\gamma_{2n}) \mathbf{fl}(|a_1 a_2 \cdots a_n|) \\ &\leq \mathbf{fl}(\mathbf{eps}|\mathbf{res}|) + (1 + \mathbf{eps})^{n+2} \mathbf{fl}(\gamma_n \gamma_{2n} |a_1 a_2 \cdots a_n|) \\ &\leq \mathbf{fl}(\mathbf{eps}|\mathbf{res}|) + \mathbf{fl}\left(\frac{\gamma_n \gamma_{2n} |a_1 a_2 \cdots a_n|}{1 - (n+3)\mathbf{eps}}\right) \\ &\leq (1 + \mathbf{eps}) \mathbf{fl}\left(\mathbf{eps}|\mathbf{res}| + \frac{\gamma_n \gamma_{2n} |a_1 a_2 \cdots a_n|}{1 - (n+3)\mathbf{eps}}\right) \\ &\leq \mathbf{fl}\left(\left(\mathbf{eps}|\mathbf{res}| + \frac{\gamma_n \gamma_{2n} |a_1 a_2 \cdots a_n|}{1 - (n+3)\mathbf{eps}}\right) / (1 - 2\mathbf{eps})\right). \end{aligned}$$

We can summarize this as follows.

Lemma 3.7. *Suppose Algorithm 3.2 is applied to floating point numbers $a_i \in \mathbb{F}$, $1 \leq i \leq n$ and set $p = \prod_{i=1}^n a_i$. Then, the absolute forward error affecting the product is bounded according to*

$$|\mathbf{res} - p| \leq \mathbf{fl}\left(\left(\mathbf{eps}|\mathbf{res}| + \frac{\gamma_n \gamma_{2n} |a_1 a_2 \cdots a_n|}{1 - (n+3)\mathbf{eps}}\right) / (1 - 2\mathbf{eps})\right).$$

3.5 Validated error bound and faithful rounding

In the previous subsection, we have shown that

$$|e_n - e| \leq \mathbf{fl}\left(\frac{\gamma_n \gamma_{2n} |a_1 a_2 \cdots a_n|}{1 - (n+3)\mathbf{eps}}\right). \quad (3.10)$$

Lemma 3.6 tells us that if $2|e - e_n| < \mathbf{eps}|\mathbf{res}|$ then \mathbf{res} is a faithful rounding of p (where \mathbf{res} is the result of `CompProd`).

As a consequence, if

$$\mathbf{fl}\left(2\frac{\gamma_n \gamma_{2n} |a_1 a_2 \cdots a_n|}{1 - (n+3)\mathbf{eps}}\right) < \mathbf{fl}(\mathbf{eps}|\mathbf{res}|)$$

then we got a faithfully rounded result. This makes it possible to check *a posteriori* if the result is faithfully rounded.

4 Exponentiation

In this section, we study two exponentiation algorithms (for computing x^n with $x \in \mathbb{F}$ and $n \in \mathbb{N}$). The first one is linear (in $\mathcal{O}(n)$) whereas the second one is logarithmic (in $\mathcal{O}(\log n)$).

4.1 A linear algorithm

The natural method to compute x^n is to apply algorithm **CompProd** for $a_i \in \mathbb{F}$ with $a_i = x$ for $1 \leq i \leq n$. This leads to the following algorithm.

Algorithm 4.1. Power evaluation with a compensated scheme

```
function res = CompLinPower(x, n)
    p1 = x
    e1 = 0
    for i = 2 : n
        [pi, πi] = TwoProduct(pi-1, x)
        ei = fl(ei-1x + πi)
    end
    res = fl(pn + en)
```

All that has been done concerning **CompProd** is still valid with **CompLinPower**. For example, if $n < 2^{25}$ then the result is faithfully rounded.

This algorithm is similar to the one of [11]. It is also the same as the Compensated Horner scheme [6] applied to the polynomial $p(x) = x^n$. Results concerning faithful polynomial evaluation can be found in [12].

4.2 A double-double library

Compensated methods are a possible way to improve the accuracy. Another possibility is to increase the working precision. For this purpose, some multiprecision libraries have been developed. One can divide those libraries into three categories.

- Arbitrary precision library using a *multiple-digit* format where a number is expressed as a sequence of digits coupled with a single exponent. Examples of this format are Bailey's MPFUN [2], Brent's MP [4] or MPFR [1].
- Arbitrary precision library using a *multiple-component* format where a number is expressed as unevaluated sums of ordinary floating point words. Examples of this format are Priest [17] and Shewchuk [19].
- Extended fixed precision library using the *multiple-component* format but with a limited number of components. Examples of this format are Bailey's *double-double* [3] (double-double numbers are represented as an unevaluated sum of a leading double and a trailing double) and *quad-double* [7] (quad-double numbers are represented as an unevaluated sum of four IEEE doubles).

In the sequel, we present two algorithms to compute product of two double-double or a double times a double-double. Those algorithms are taken from [13].

Algorithm 4.2. Multiplication of two double-double numbers

```
function [rh, rl] = prod_dd_dd(ah, al, bh, bl)
    [t1, t2] = TwoProduct(ah, bh)
    t3 = fl(((ah · bl) + (al · bh)) + t2)
    [rh, rl] = TwoProduct(t1, t3)
```

Algorithm 4.3. Multiplication of double-double number by a double number

```

function  $[r_h, r_l] = \text{prod\_dd\_d}(a, b_h, b_l)$ 
 $[t_1, t_2] = \text{TwoProduct}(a, b_h)$ 
 $t_3 = \text{fl}((a \cdot b_l) + t_2)$ 
 $[r_h, r_l] = \text{TwoProduct}(t_1, t_3)$ 

```

Theorem 4.1 (Lauter [13, thm. 4.7]). *Let $a_h + a_l$ and $b_h + b_l$ be the double-double arguments of Algorithm 4.2. Then the returned values r_h and r_l satisfy*

$$r_h + r_l = ((a_h + a_l) \cdot (b_h + b_l))(1 + \varepsilon)$$

where ε is bounded as follows : $|\varepsilon| \leq 16\text{eps}^2$. Furthermore, we have $|r_l| \leq \text{eps}|r_h|$.

Results for Algorithm 4.3 are very similar with $a = a_h$ and $a_l = 0$.

4.3 A logarithmic algorithm

A logarithmic algorithm was introduced in [11] using the classic right-to-left binary exponentiation algorithm and the double-double library. Hereafter, we propose a variant of this algorithm with the left-to-right binary exponentiation algorithm together with the double-double library. Contrary to the right-to-left binary exponentiation algorithm which needs two multiplications of two double-double numbers, the left-to-right binary exponentiation algorithm only needs a multiplication of two double-double numbers and a multiplication of a double number by a double-double numbers. Moreover, the multiplication of a double-double by a double-double is actually a square so that it can be a little bit optimized.

Algorithm 4.4. Power evaluation with a compensated scheme

```

function res = CompLogPower( $x, n$ )                                %  $n = (n_t n_{t-1} \dots n_1 n_0)_2$ 
 $[h, l] = [1, 0]$ 
for  $i = t : -1 : 0$ 
     $[h, l] = \text{prod\_dd\_dd}(h, l, h, l)$ 
    if  $n_i = 1$ 
         $[h, l] = \text{prod\_dd\_d}(x, h, l)$ 
    end
end
res =  $\text{fl}(h + l)$ 

```

Theorem 4.2. *The two values h and l returned by Algorithm 4.4 satisfy*

$$h + l = x^n(1 + \varepsilon)$$

with

$$(1 - 16\text{eps}^2)^{n-1} \leq 1 + \varepsilon \leq (1 + 16\text{eps}^2)^{n-1}.$$

Proof. The proof is very similar to the one of [11, Thm 4]. It comes from the fact that by induction one can show that the approximation of x^k is of the form $x^k(1 + \varepsilon_k)$ with $(1 - 16\text{eps}^2)^{k-1} \leq 1 + \varepsilon_k \leq (1 + 16\text{eps}^2)^{k-1}$. \square

Let us denote $\varphi = \text{eps}(1 - \text{eps})$ and

$$\overline{\gamma}_n = \frac{16n\text{eps}^2}{1 - 16n\text{eps}^2}.$$

It follows that $1 + \varepsilon \leq (1 + 16\text{eps}^2)^{n-1} \leq 1 + \overline{\gamma}_n$ [8, p.63]. Since $h + l = x^n(1 + \varepsilon)$ then it holds $x^n = h + l - \varepsilon x^n$. Furthermore, it also holds $\text{res} = \text{fl}(h + l)$. From Lemma 3.5, if we prove that $2\varepsilon|x^n| < \text{eps}|\text{res}|$ then res is a faithful rounding of x^n . From the definition of res , it follows that $(1 - \text{eps})|h + l| \leq |\text{res}|$ and so $(1 - \text{eps})(1 + \varepsilon)|x^n| \leq |\text{res}|$. A sufficient condition to ensure a faithful rounding is then

$$2\varepsilon|x^n| < \text{eps}(1 - \text{eps})(1 + \varepsilon)|x^n| \quad \text{and so} \quad 2\varepsilon < \text{eps}(1 - \text{eps})(1 + \varepsilon),$$

which is equivalent to

$$\varepsilon < \frac{\varphi}{2 - \varphi}.$$

Since $\varepsilon \leq \overline{\gamma}_n$, a sufficient condition is

$$\overline{\gamma}_n < \frac{\varphi}{2 - \varphi},$$

which is equivalent to

$$n < \frac{\varphi}{16\text{eps}^2} = \frac{(1 - \text{eps})\text{eps}^{-1}}{16}.$$

For example, in double precision where $\text{eps} = 2^{-53}$, if $n < 2^{49} \approx 5 \cdot 10^{14}$, then we get a faithfully rounded result.

5 Conclusion

In this paper, we provided an accurate algorithm for computing product of floating point numbers. We gave some sufficient conditions to obtain a faithfully rounded result as well as validated error bounds. We applied this algorithm to compute exponentiation of floating point numbers. We improved this algorithm by using a double-double library.

References

- [1] *MPFR, the Multiprecision Precision Floating Point Reliable library*. Available at URL = <http://www.mpfr.org>.
- [2] David H. Bailey. A Fortran 90-based multiprecision system. *ACM Trans. Math. Softw.*, 21(4):379–387, 1995.
- [3] David H. Bailey. *A Fortran-90 double-double library*, 2001. Available at URL = <http://crd.lbl.gov/~dhbailey/mpdist/index.html>.
- [4] Richard P. Brent. A Fortran multiple-precision arithmetic package. *ACM Trans. Math. Softw.*, 4(1):57–70, 1978.
- [5] T. J. Dekker. A floating-point technique for extending the available precision. *Numer. Math.*, 18:224–242, 1971.
- [6] Stef Graillat, Nicolas Louvet, and Philippe Langlois. Compensated Horner scheme. Research Report 04, Équipe de recherche DALI, Laboratoire LP2A, Université de Perpignan Via Domitia, France, 52 avenue Paul Alduy, 66860 Perpignan cedex, France, July 2005.
- [7] Yozo Hida, Xiaoye S. Li, and David H. Bailey. Algorithms for quad-double precision floating point arithmetic. In *Proc. 15th IEEE Symposium on Computer Arithmetic*, pages 155–162. IEEE Computer Society Press, Los Alamitos, CA, USA, 2001.

- [8] Nicholas J. Higham. *Accuracy and stability of numerical algorithms*. Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, second edition, 2002.
- [9] *IEEE Standard for Binary Floating-Point Arithmetic, ANSI/IEEE Standard 754-1985*. Institute of Electrical and Electronics Engineers, New York, 1985. Reprinted in SIGPLAN Notices, 22(2):9–25, 1987.
- [10] Donald E. Knuth. *The Art of Computer Programming, Volume 2, Seminumerical Algorithms*. Addison-Wesley, Reading, MA, USA, third edition, 1998.
- [11] Peter Kornerup, Vincent Lefevre, and Jean-Michel Muller. Computing integer powers in floating-point arithmetic. arXiv:0705.4369v1 [cs.NA].
- [12] Philippe Langlois and Nicolas Louvet. How to ensure a faithful polynomial evaluation with the compensated horner algorithm. In *Proceedings of the 18th IEEE Symposium on Computer Arithmetic (ARITH '07), Montpellier, France*, pages 141–149. IEEE Computer Society, Los Alamitos, CA, USA, 2007.
- [13] Christoph Quirin Lauter. Basic building blocks for a triple-double intermediate format. Research Report RR-5702, INRIA, September 2005.
- [14] Yves Nievergelt. Scalar fused multiply-add instructions produce floating-point matrix arithmetic provably accurate to the penultimate digit. *ACM Trans. Math. Software*, 29(1):27–48, 2003.
- [15] Takeshi Ogita, Siegfried M. Rump, and Shin’ichi Oishi. Accurate sum and dot product. *SIAM J. Sci. Comput.*, 26(6):1955–1988, 2005.
- [16] Takeshi Ogita, Siegfried M. Rump, and Shin’ichi Oishi. Verified solution of linear systems without directed rounding. Technical Report No. 2005-04, Advanced Research Institute for Science and Engineering, Waseda University, 2005.
- [17] Douglas M. Priest. *On Properties of Floating Point Arithmetics: Numerical Stability and the Cost of Accurate Computations*. PhD thesis, Mathematics Department, University of California, Berkeley, CA, USA, November 1992. <ftp://ftp.icsi.berkeley.edu/pub/theory/priest-thesis.ps.Z>.
- [18] Siegfried M. Rump, Takeshi Ogita, and Shin’ichi Oishi. Accurate floating-point summation. Technical Report 05.12, Faculty for Information and Communication Sciences, Hamburg University of Technology, nov 2005.
- [19] Jonathan Richard Shewchuk. Adaptive precision floating-point arithmetic and fast robust geometric predicates. *Discrete Comput. Geom.*, 18(3):305–363, 1997.
- [20] Pat H. Sterbenz. *Floating-point computation*. Prentice-Hall Inc., Englewood Cliffs, N.J., 1974. Prentice-Hall Series in Automatic Computation.