

# Accurate summation, dot product and polynomial evaluation in complex floating point arithmetic

Stef Graillat\*, Valérie Ménissier-Morain<sup>1</sup>

UPMC Univ Paris 06, UMR 7606, LIP6, 4 place Jussieu, F-75252, Paris cedex 05, France

## ARTICLE INFO

### Article history:

Available online 30 March 2012

### Keywords:

Complex floating point arithmetic  
Error-free transformations  
Accurate summation  
Accurate dot product  
Accurate polynomial evaluation  
Horner's scheme  
High precision

## ABSTRACT

Several different techniques and softwares intend to improve the accuracy of results computed in a fixed finite precision. Here we focus on methods to improve the accuracy of summation, dot product and polynomial evaluation. Such algorithms exist real floating point numbers. In this paper, we provide new algorithms which deal with complex floating point numbers. We show that the computed results are as accurate as if computed in twice the working precision. The algorithms are simple since they only require addition, subtraction and multiplication of floating point numbers in the same working precision as the given data.

© 2012 Elsevier Inc. All rights reserved.

## 1. Introduction

It is well known that computing with finite precision implies some rounding errors. These errors can lead to inexact results for a computation. An important tool to try to avoid this are *error-free transformations*: to compute not only a floating point approximation but also an exact error term without overlapping. This can be viewed as a *double-double* floating point numbers [1] but without the renormalization step.

Error-free transformations have been widely used to provide some new accurate algorithms in real floating point arithmetic (see [2,3] for accurate sum and dot product and [4] for polynomial evaluation). Complex error-free transformations are then the next step for providing accurate algorithms using complex numbers.

The rest of the paper is organized as follows. In Section 2, we recall some results on real floating point arithmetic and error-free transformations. In Section 3, we present the complex floating point arithmetic and we propose some new error-free transformations for this arithmetic. In Section 4, we propose some accurate algorithms to compute summation and dot product of complex floating point vectors. These algorithms are derived by applying the real floating point algorithms to both the real and the imaginary parts. In Section 5, we study different polynomial evaluation algorithms. We first describe the Horner scheme in complex floating point arithmetic. We then present the compensated Horner scheme in complex arithmetic. We provide an error analysis for both versions of the Horner scheme and we conclude by presenting some numerical experiments confirming the accuracy of our algorithm.

This paper is an extended version of the paper [5]. The paper [5] was only dealing with accurate polynomial evaluation. Here we also consider accurate summation and dot product in complex floating point arithmetic.

\* Corresponding author.

E-mail addresses: stef.graillat@lip6.fr (S. Graillat), valerie.menissier-morain@lip6.fr (V. Ménissier-Morain).

<sup>1</sup> This work has been done while the second author was visiting University of Waterloo, 200 University Avenue West, Waterloo, Ontario, N2L 3G1, Canada.

## 2. Real floating point arithmetic

In this section, we first recall the principle of real floating point arithmetic. Then we present the well-known error-free transformations associated with the classical operations addition, subtraction, multiplication.

### 2.1. Notations and fundamental property of real floating point arithmetic

Throughout the paper, we assume to work with a floating point arithmetic adhering to IEEE 754 floating point standard [6]. We assume that no overflow nor underflow occur. The set of floating point numbers is denoted by  $\mathbb{F}$ , the relative rounding error by  $\text{eps}$ . For IEEE 754 double precision, we have  $\text{eps} = 2^{-53}$  and for single precision  $\text{eps} = 2^{-24}$ .

We denote by  $\text{fl}(\cdot)$  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 [7]

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

This implies that

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

We use standard notation for error estimations. The quantities  $\gamma_n$  are defined as usual [7] by

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

where we implicitly assume that  $n \text{eps} \leq 1$  and we will use inequality  $\text{eps} \leq \sqrt{2}\gamma_2$  in the following proofs.

### 2.2. Error-free transformations in real floating point arithmetic

One can notice that  $a \circ b \in \mathbb{R}$  and  $\text{fl}(a \circ b) \in \mathbb{F}$  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 [8]):

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

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.

We use Matlab-like notations to describe the algorithms.

#### 2.2.1. Addition

For addition, we can use the following algorithm by Knuth [9, Thm. B, p. 236].

**Algorithm 2.1** (*Error-free transformation of the sum of two floating point numbers*). (See Knuth [9].)

```
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 [8]. 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** (*Error-free transformation of the sum of two floating point numbers with  $|a| \geq |b|$* ). (See Dekker [8].)

```
function [x, y] = FastTwoSum(a, b)
    x = fl(a + b);   y = fl((a - x) + b)
```

#### 2.2.2. Multiplication

For the error-free transformation of a product, we first need to split the input argument into two parts.

**Splitting** Let  $p$  be the integer number given by  $\text{eps} = 2^{-p}$  and define  $s = \lceil p/2 \rceil$  and  $\text{factor} = \text{fl}(2^s + 1)$ . For example, if the working precision is IEEE 754 double precision, then  $p = 53$  and  $s = 27$  and  $\text{factor} = 1.\underbrace{00\dots 00}_{26}1\underbrace{00\dots 00}_{26}2^{27}$  allows by multiplying a number to split the mantissa of this number into its most and least significant halves. The quantities  $p$ ,  $s$  and  $\text{factor}$  are constants of the floating point arithmetic.

The following algorithm by Dekker [8] splits a floating point number  $a \in \mathbb{F}$  into two parts  $x$  and  $y$  such that

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

Two floating point values  $x$  and  $y$  with  $|y| \leq |x|$  are nonoverlapping if the least significant nonzero bit of  $x$  is more significant than the most significant nonzero bit of  $y$ .

**Algorithm 2.3** (Error-free split of a floating point number into two parts). (See Dekker [8].)

```
function [x, y] = Split(a, b)
    c = fl(factor * a);   x = fl(c - (c - a));   y = fl(a - x)
```

**Product** An algorithm from Veltkamp (see [8]) makes it possible to compute an error-free transformation for the product of two floating point numbers by splitting the two arguments.

This algorithm returns two floating point numbers  $x$  and  $y$  such that

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

**Algorithm 2.4** (Error-free transformation of the product of two floating point numbers). (See Veltkamp [8].)

```
function [x, y] = TwoProduct(a, b)
    x = fl(a * b)
    [a1, a2] = Split(a);   [b1, b2] = Split(b)
    y = fl(a2 * b2 - ((x - a1 * b1) - a2 * b1) - a1 * b2)
```

### 2.2.3. Properties

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

**Theorem 2.1.** (See Ogita, Rump and Oishi [2].)

**Addition** 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.

**Product** 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.

### 2.2.4. Multiplication with FMA

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 [10,11]. 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 operator satisfies

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

**Algorithm 2.5** (Error-free transformation of the product of two floating point numbers using an FMA). (See Ogita, Rump and Oishi [2].)

```
function [x, y] = TwoProductFMA(a, b)
    x = fl(a * b);   y = FMA(a, b, -x)
```

The `TwoProductFMA` algorithm requires only 2 flops.

### 3. Complex floating point arithmetic

#### 3.1. Notations and fundamental property of complex floating point arithmetic

We denote by  $\mathbb{F} + i\mathbb{F}$  the set of complex floating point numbers. As in the real case, we denote by  $\text{fl}(\cdot)$  the result of a floating point computation, where all operations inside parentheses are done in floating point working precision in the obvious way [7, p. 71]. The following properties hold [7,12] for  $x, y \in \mathbb{F} + i\mathbb{F}$ ,

$$\text{fl}(x \circ y) = (x \circ y)(1 + \varepsilon_1) = (x \circ y)/(1 + \varepsilon_2), \quad \text{for } \circ = \{+, -\} \text{ and } |\varepsilon_v| \leq \text{eps}, \quad (3.5)$$

and

$$\text{fl}(x \cdot y) = (x \cdot y)(1 + \varepsilon_1), \quad |\varepsilon_1| \leq \sqrt{2}\gamma_2. \quad (3.6)$$

This implies that

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

and

$$|x \cdot y - \text{fl}(x \cdot y)| \leq \sqrt{2}\gamma_2|x \cdot y|.$$

For the complex multiplication, we can replace the term  $\sqrt{2}\gamma_2$  by  $\sqrt{5}\text{eps}$  which is nearly optimal (see [13]). As a consequence, in the sequel, all the bounds for algorithms involving a multiplication can be improved by a small constant factor.

We will also use the notation  $\tilde{\gamma}_n$  for the quantities

$$\tilde{\gamma}_n := \frac{n\sqrt{2}\gamma_2}{1 - n\sqrt{2}\gamma_2}.$$

And we will use inequalities  $(1 + \sqrt{2}\gamma_2)(1 + \tilde{\gamma}_n) \leq (1 + \tilde{\gamma}_{n+1})$  and  $(1 + \sqrt{2}\gamma_2)\tilde{\gamma}_{n-1} \leq \tilde{\gamma}_n$ .

#### 3.2. Sum and product

The error-free transformations presented hereafter were first described in [14]. The sum requires still only one error term as for the real case but the product needs three error terms.

##### 3.2.1. Addition

**Algorithm 3.1** (Error-free transformation of the sum of two complex floating point numbers  $x = a + ib$  and  $y = c + id$ ).

```
function [s, e] = TwoSumCplx(x, y)
    [s1, e1] = TwoSum(a, c); [s2, e2] = TwoSum(b, d)
    s = s1 + is2; e = e1 + ie2
```

**Theorem 3.1.** Let  $x, y \in \mathbb{F} + i\mathbb{F}$  and let  $s, e \in \mathbb{F} + i\mathbb{F}$  such that  $[s, e] = \text{TwoSumCplx}(x, y)$  (Algorithm 3.1). Then,

$$x + y = s + e, \quad s = \text{fl}(x + y), \quad |e| \leq \text{eps}|s|, \quad |e| \leq \text{eps}|x + y|. \quad (3.7)$$

The algorithm `TwoSumCplx` requires 12 flops.

**Proof.** From Theorem 2.1 with `TwoSum`, we have  $s_1 + e_1 = a + c$  and  $s_2 + e_2 = b + d$ . It follows that  $s + e = x + y$  with  $s = \text{fl}(x + y)$ . From (3.5), we derive that  $|e| \leq \text{eps}|s|$  and  $|e| \leq \text{eps}|x + y|$ .  $\square$

##### 3.2.2. Multiplication

Algorithm 2.4 cannot be straightforward generalized to complex multiplication. We need the new following algorithm.

**Algorithm 3.2** (Error-free transformation of the product of two complex floating point numbers  $x = a + ib$  and  $y = c + id$ ).

```
function [p, e, f, g] = TwoProductCplx(x, y)
    [z1, h1] = TwoProduct(a, c); [z2, h2] = TwoProduct(b, d)
    [z3, h3] = TwoProduct(a, d); [z4, h4] = TwoProduct(b, c)
    [z5, h5] = TwoSum(z1, -z2); [z6, h6] = TwoSum(z3, z4)
    p = z5 + iz6; e = h1 + ih3; f = -h2 + ih4; g = h5 + ih6
```

**Theorem 3.2.** Let  $x, y \in \mathbb{F} + i\mathbb{F}$  and let  $p, e, f, g \in \mathbb{F} + i\mathbb{F}$  such that  $[p, e, f, g] = \text{TwoProductCplx}(x, y)$  (Algorithm 2.4). Then,

$$x \cdot y = p + e + f + g, \quad p = \text{fl}(x \cdot y), \quad |e + f + g| \leq \sqrt{2}\gamma_2|x \cdot y|. \quad (3.8)$$

The algorithm `TwoProductCplx` requires 80 flops.

**Proof.** From Theorem 2.1, it holds that  $z_1 + h_1 = a \cdot c$ ,  $z_2 + h_2 = b \cdot d$ ,  $z_3 + h_3 = a \cdot d$ ,  $z_4 + h_4 = b \cdot c$ ,  $z_5 + h_5 = z_1 - z_2$  and  $z_6 + h_6 = z_3 + z_4$ . By the definition of  $p, e, f, g$ , we conclude that  $x \cdot y = p + e + f + g$  with  $p = \text{fl}(x \cdot y)$ . From (3.6), we deduce that  $|e + f + g| = |x \cdot y - \text{fl}(x \cdot y)| \leq \sqrt{2}\gamma_2|x \cdot y|$ .  $\square$

*Optimization of the algorithm* In Algorithm 3.2, in each call to `TwoProduct`, we have to split the two arguments. Yet, we split the same numbers  $a, b, c$  and  $d$  twice. With only one split for each of these numbers, the cost is 64 flops. The previous algorithm can be expanded as follows:

**Algorithm 3.3** (Error-free transformation of the product of two complex floating point numbers  $x = a + ib$  and  $y = c + id$  with single splitting).

```
function [p, e, f, g] = TwoProductCplxSingleSplitting(x, y)
    [a1, a2] = Split(a), [b1, b2] = Split(b), [c1, c2] = Split(c), [d1, d2] = Split(d)
    z1 = fl(a · c), z2 = fl(b · d), z3 = fl(a · d), z4 = fl(b · c)
    h1 = fl(a2 · c2 - (((z1 - a1 · c1) - a2 · c1) - a1 · c2))
    h2 = fl(b2 · d2 - (((z2 - b1 · d1) - b2 · d1) - b1 · d2))
    h3 = fl(a2 · d2 - (((z3 - a1 · d1) - a2 · d1) - a1 · d2))
    h4 = fl(b2 · c2 - (((z4 - b1 · c1) - b2 · c1) - b1 · c2))
    [z5, h5] = TwoSum(z1, -z2), [z6, h6] = TwoSum(z3, z4)
    p = z5 + iz6, e = h1 + ih3, f = -h2 + ih4, g = h5 + ih6
```

### 3.2.3. Multiplication with FMA

Of course we obtain a much faster algorithm if we use `TwoProductFMA` instead of `TwoProduct`. In that case, the numbers of flops falls down to 20.

**Algorithm 3.4** (Error-free transformation of the product of two complex floating point numbers  $x = a + ib$  and  $y = c + id$  using FMA).

```
function [p, e, f, g] = TwoProductFMACplx(x, y)
    [z1, h1] = TwoProductFMA(a, c); [z2, h2] = TwoProductFMA(b, d)
    [z3, h3] = TwoProductFMA(a, d); [z4, h4] = TwoProductFMA(b, c)
    [z5, h5] = TwoSum(z1, -z2); [z6, h6] = TwoSum(z3, z4)
    p = z5 + iz6; e = h1 + ih3; f = -h2 + ih4; g = h5 + ih6
```

The 8.5:1 ratio between the cost of `TwoProduct` and `TwoProductFMA` algorithms and the 3.2:1 ratio between the cost of `TwoProductCplxSingleSplitting` and `TwoProductFMACplx` algorithms show that the availability of an FMA is crucial for fast error-free transformations in real and complex arithmetic.

## 4. Accurate summation and dot product

In this section, we first recall an accurate algorithm for the summation of real floating point numbers and we present the error bound analysis. We then show that we can apply this algorithm to both the real and imaginary part of a complex floating point number vector. We give an error bound analysis for this new algorithm. We do similar analysis for dot product.

### 4.1. Accurate summation

*Real floating point numbers case* The following algorithm makes it possible to accurately compute the sum of real floating point numbers. By accurately, we mean as if computed in twice the working precision which is sum up in the following result.

**Algorithm 4.1** (Summation in twice the working precision for real floating point numbers vectors). (See Ogita, Rump and Oishi [2].)

```

function res = Sum2(p)
  pi_1 = p_1; sigma_1 = 0;
  for i = 2 : n
    [pi_i, qi_i] = TwoSum(pi_{i-1}, p_i)
    sigma_i = fl(sigma_{i-1} + qi_i)
  end
  res = fl(pi_n + sigma_n)

```

**Proposition 4.1.** (See Ogita, Rump and Oishi [2].) Suppose Algorithm Sum2 is applied to floating point number  $p_i \in \mathbb{F}$ ,  $1 \leq i \leq n$ . Let  $s := \sum p_i$ ,  $S := \sum |p_i|$ . Then, we have

$$|\text{res} - s| \leq \text{eps}|s| + \gamma_{n-1}^2 S.$$

*Complex floating point numbers case* If the inputs are now complex floating point numbers  $p_j = a_j + ib_j$ , we want to compute  $s = \sum_{j=1}^n p_j$ . This can easily be done by compute the sum of the real part and the imaginary part with Sum2.

**Algorithm 4.2** (Summation in twice the working precision for complex floating point numbers vectors).

```

function res = Sum2cplx(p)
  Let a and b be the vectors representing the real and imaginary parts of p
  res_r = Sum2(a)
  res_i = Sum2(b)
  res = res_r + i res_i

```

**Proposition 4.2.** Suppose Algorithm Sum2cplx is applied to floating point number  $p_j = a_j + ib_j \in \mathbb{F} + i\mathbb{F}$ ,  $1 \leq j \leq n$ . Let  $s := \sum p_j$ ,  $S := \sum |p_j|$ . Then, we have

$$|\text{res} - s| \leq \sqrt{2}\text{eps}|s| + 2\gamma_{n-1}^2 S.$$

**Proof.** We have  $|\text{res} - s|^2 = |\text{res}_r - \sum_{i=1}^n a_i|^2 + |\text{res}_i - \sum_{i=1}^n b_i|^2$ . We know from Proposition 4.1 that

$$\left| \text{res}_r - \sum_{i=1}^n a_i \right| \leq \text{eps} \left| \sum_{i=1}^n a_i \right| + \gamma_{n-1}^2 \sum_{i=1}^n |a_i| \quad \text{and} \quad \left| \text{res}_i - \sum_{i=1}^n b_i \right| \leq \text{eps} \left| \sum_{i=1}^n b_i \right| + \gamma_{n-1}^2 \sum_{i=1}^n |b_i|.$$

As a consequence,

$$|\text{res} - s|^2 \leq \left( \text{eps} \left| \sum_{i=1}^n a_i \right| + \gamma_{n-1}^2 \sum_{i=1}^n |a_i| \right)^2 + \left( \text{eps} \left| \sum_{i=1}^n b_i \right| + \gamma_{n-1}^2 \sum_{i=1}^n |b_i| \right)^2.$$

Since for all numbers  $x$  and  $y$ , we have  $(x + y)^2 \leq 2(x^2 + y^2)$ , it follows

$$|\text{res} - s|^2 \leq 2\text{eps}^2 \left( \left| \sum_{i=1}^n a_i \right|^2 + \left| \sum_{i=1}^n b_i \right|^2 \right) + 2\gamma_{n-1}^4 \left( \left( \sum_{i=1}^n |a_i| \right)^2 + \left( \sum_{i=1}^n |b_i| \right)^2 \right).$$

Since  $|\sum_{i=1}^n a_i|^2 + |\sum_{i=1}^n b_i|^2 = |\sum_{i=1}^n (a_i + ib_i)|^2 = |s|^2$ , we have

$$|\text{res} - s|^2 \leq 2\text{eps}^2 |s|^2 + 2\gamma_{n-1}^4 \left( \left( \sum_{i=1}^n |a_i| \right)^2 + \left( \sum_{i=1}^n |b_i| \right)^2 \right).$$

Furthermore, since for all numbers  $x, y \geq 0$ , we have  $x^2 + y^2 \leq (x + y)^2$ , we get

$$|\text{res} - s|^2 \leq 2\text{eps}^2 |s|^2 + 2\gamma_{n-1}^4 \left( \sum_{i=1}^n (|a_i| + |b_i|) \right)^2.$$

As for all numbers  $x$  and  $y$ ,  $x + y \leq \sqrt{2}\sqrt{x^2 + y^2}$ , we have  $|a_i| + |b_i| \leq \sqrt{2}|p_i|$  and hence,

$$|\text{res} - s|^2 \leq 2\text{eps}^2 |s|^2 + 4\gamma_{n-1}^4 \left( \sum_{i=1}^n |p_i| \right)^2.$$

Since for all  $x, y \geq 0$ ,  $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$ , it follows

$$|\text{res} - s| \leq \sqrt{2}\text{eps}|s| + 2\gamma_{n-1}^2 \sum_{i=1}^n |p_i|.$$

This concludes the proof.  $\square$

#### 4.2. Accurate dot product

*Real floating point numbers case* Let  $x = (x_j)$  and  $y = (y_j)$ , we compute  $p = x^T y = \sum x_j y_j$

**Algorithm 4.3** (*Dot product in twice the working precision for real floating point numbers vectors*). (See Ogita, Rump and Oishi [2].)

```
function res = Dot2(x, y)
[p, s] = TwoProduct(x1, y1)
for i = 2 : n
    [h, r] = TwoProduct(xi, yi)
    [p, q] = TwoSum(p, h)
    s = fl(s + (q + r))
end
res = fl(p + s)
```

**Proposition 4.3.** (See Ogita, Rump and Oishi [2].) *Let floating point numbers  $x_i, y_i \in \mathbb{F}$ ,  $1 \leq i \leq n$ , be given and denote by  $\text{res} \in \mathbb{F}$  the result computed by Algorithm `Dot2`. Then occurs,*

$$|\text{res} - x^T y| \leq \text{eps}|x^T y| + \gamma_n^2 |x^T| |y|.$$

*Complex floating point numbers case* Let  $x = (x_j)$  with  $x_j = a_j + ib_j$  and  $y = (y_j)$  with  $y_j = c_j + id_j$ , we compute  $p = x^* y = \sum x_j \bar{y}_j = \sum (a_j c_j + b_j d_j) + i \sum (b_j c_j - a_j d_j)$ . These two sums will be each one computed by a dot product of real floating point numbers vectors with double length: let  $X$  be the vector with the  $a_j$  as first elements followed by the  $b_j$ ,  $Y$  be the vector equivalent for  $y$  with the  $c_j$  as first elements followed by the  $d_j$  and  $Y'$  be the vector equivalent for  $-iy$  with the  $d_j$  as first elements followed by the  $-c_j$ . With a block vector notation, we have

$$X = \begin{bmatrix} \text{Re}(x) \\ \text{Im}(x) \end{bmatrix}, \quad Y = \begin{bmatrix} \text{Re}(y) \\ \text{Im}(y) \end{bmatrix} \quad \text{and} \quad Y' = \begin{bmatrix} \text{Im}(y) \\ -\text{Re}(y) \end{bmatrix}$$

and we have  $p = X^T Y + i X^T Y'$ .

**Algorithm 4.4** (*Dot product in twice the working precision for complex floating point numbers vectors*).

```
function res = Dot2cplx(x, y)
    build X, Y, Y'
    res = Dot2(X, Y) + i Dot2(X, Y')
```

**Proposition 4.4.** *Let floating point numbers  $x = (x_j)$  with  $x_j = a_j + ib_j$  and  $y = (y_j)$  with  $y_j = c_j + id_j$  be given and denote by  $\text{res} \in \mathbb{F} + i\mathbb{F}$  the result computed by Algorithm `Dot2cplx`. Then occurs,*

$$|\text{res} - x^* y| \leq \sqrt{2}\text{eps}|x^* y| + 2\gamma_{2n}^2 |x^T| |y|.$$

**Proof.** From Proposition 4.3, it follows that

$$|X^T Y - \text{Dot2}(X, Y)| \leq \text{eps}|X^T Y| + \gamma_{2n}^2 |X^T| |Y|,$$

and

$$|X^T Y' - \text{Dot2}(X, Y')| \leq \text{eps}|X^T Y'| + \gamma_{2n}^2 |X^T| |Y'|.$$

As a consequence,

$$|\text{res} - x^*y|^2 = |X^T Y - \text{Dot2}(X, Y)|^2 + |X^T Y' - \text{Dot2}(X, Y')|^2.$$

Since for all numbers  $x$  and  $y$ , we have  $(x + y)^2 \leq 2(x^2 + y^2)$ , it follows

$$|\text{res} - x^*y|^2 \leq 2(\text{eps}^2 |X^T Y|^2 + \gamma_{2n}^4 (|X|^T |Y|)^2) + 2(\text{eps}^2 |X^T Y'|^2 + \gamma_{2n}^4 (|X|^T |Y'|)^2).$$

We rearrange this inequalities and we obtain

$$|\text{res} - x^*y|^2 \leq 2\text{eps}^2 (|X^T Y|^2 + |X^T Y'|^2) + 2\gamma_{2n}^4 ((|X|^T |Y|)^2 + (|X|^T |Y'|)^2).$$

For the first part of the right member we use the definition  $|X^T Y|^2 + |X^T Y'|^2 = |x^*y|^2$ . We now have to find an upper bound for the second part  $(|X|^T |Y|)^2 + (|X|^T |Y'|)^2$ . Since still for all  $x, y \geq 0$ , we have  $x^2 + y^2 \leq (x + y)^2$ , it follows that  $(|X|^T |Y|)^2 + (|X|^T |Y'|)^2 \leq (|X|^T |Y| + |X|^T |Y'|)^2$ . As a consequence, we just have to find an upper bound for  $|X|^T |Y| + |X|^T |Y'|$ .

We have

$$|X|^T |Y| + |X|^T |Y'| = \sum_{i=1}^n (|a_i c_i| + |b_i d_i|) + \sum_{i=1}^n (|a_i d_i| + |b_i c_i|) = \sum_{i=1}^n (|a_i c_i| + |b_i d_i| + |a_i d_i| + |b_i c_i|)$$

and we factorize this internal sum

$$|X|^T |Y| + |X|^T |Y'| = \sum_{i=1}^n (|a_i| + |b_i|)(|c_i| + |d_i|).$$

By definition  $|x|^T |y| = \sum_{i=1}^n (\sqrt{a_i^2 + b_i^2} \sqrt{c_i^2 + d_i^2})$ . Since for all numbers  $x$  and  $y$ , we have  $x + y \leq \sqrt{2} \sqrt{x^2 + y^2}$ , it follows that

$$|a_i| + |b_i| \leq \sqrt{2} \sqrt{a_i^2 + b_i^2} \quad \text{and} \quad |c_i| + |d_i| \leq \sqrt{2} \sqrt{c_i^2 + d_i^2}$$

and by multiplying these two inequalities, we get

$$\sum_{i=1}^n (|a_i| + |b_i|)(|c_i| + |d_i|) \leq \sum_{i=1}^n (\sqrt{2} \sqrt{a_i^2 + b_i^2}) (\sqrt{2} \sqrt{c_i^2 + d_i^2}) = 2 \sum_{i=1}^n \sqrt{a_i^2 + b_i^2} \sqrt{c_i^2 + d_i^2},$$

and so

$$|X|^T |Y| + |X|^T |Y'| \leq 2|x|^T |y|.$$

Consequently, we obtain

$$|\text{res} - x^*y|^2 \leq 2\text{eps}^2 |x^*y|^2 + 4\gamma_{2n}^4 (|x|^T |y|)^2. \quad (4.9)$$

It is clear that for all numbers  $x, y \geq 0$ , we have  $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$ . Applying this to Eq. (4.9), we obtain

$$|\text{res} - x^*y| \leq \sqrt{2\text{eps}^2 |x^*y|^2 + 4\gamma_{2n}^4 (|x|^T |y|)^2} \leq \sqrt{2}\text{eps} |x^*y| + 2\gamma_{2n}^2 |x|^T |y|.$$

This concludes the proof.  $\square$

It is difficult to use such a scheme for polynomial evaluation. Indeed, let

$$p(z) = \sum_{j=0}^n a_j z^j, \quad a_j \in \mathbb{C}, \quad z = x + iy \in \mathbb{C}$$

be a polynomial. By separating real and imaginary parts, we can write it as  $p(z) = p_r(x, y) + iq_r(x, y)$  with  $p_r$  and  $q_r$  with real coefficients and evaluate  $p_r$  and  $q_r$  with Horner scheme. To achieve this, we need formal manipulations to compute  $p_r$  and  $q_r$  that are costly. In that case, it is easier to use the new error-free transformations for complex floating point arithmetic.

## 5. Accurate polynomial evaluation

First of all we describe the classical Horner scheme to evaluate polynomial  $p$  with complex floating point coefficients on  $x$  a complex floating point value. The computed value  $\text{res}$  is generally not the mathematical value  $p(x)$  rounded to the working precision. We want then to reduce the gap between these values so we modify this algorithm to compute  $\text{res}$  and additionally four polynomial error terms that we will have to evaluate on  $x$  to deduce a complex floating point correction term  $c$  that we have to add to  $\text{res}$ . Afterwards we will study mathematically and experimentally the improvement of the accuracy consisting in replacing  $\text{res}$  by  $\text{fl}(\text{res} + c)$ .



5.1. Classical Horner scheme for complex floating point arithmetic

The classical method for evaluating a polynomial

$$p(x) = \sum_{i=0}^n a_i x^i, \quad a_i, x \in \mathbb{F} + i\mathbb{F}$$

is the Horner scheme which consists on the following algorithm:

**Algorithm 5.1** (Polynomial evaluation with Horner's scheme).

```
function res = Horner(p, x)
    sn = an
    for i = n - 1 : -1 : 0
        si = si+1 · x + ai
    end
    res = s0
```

**Proposition 5.1.** A forward error bound is

$$|p(x) - \text{Horner}(p, x)| \leq \tilde{\gamma}_{2n} \sum_{i=0}^n |a_i| |x|^i = \tilde{\gamma}_{2n} \tilde{p}(|x|) \tag{5.10}$$

where  $\tilde{p}(x) = \sum_{i=0}^n |a_i| x^i$ .

**Proof.** This is a straightforward adaptation of the proof found in [7, p. 95] using (3.5) and (3.6) for complex floating point arithmetic. □

The classical condition number that describes the evaluation of  $p(x) = \sum_{i=0}^n a_i x^i$  at  $x$  is

$$\text{cond}(p, x) = \frac{\sum_{i=0}^n |a_i| |x|^i}{|\sum_{i=0}^n a_i x^i|} = \frac{\tilde{p}(|x|)}{|p(x)|}. \tag{5.11}$$

Thus if  $p(x) \neq 0$ , Eqs. (5.10) and (5.11) can be combined so that

$$\frac{|p(x) - \text{Horner}(p, x)|}{|p(x)|} \leq \tilde{\gamma}_{2n} \text{cond}(p, x). \tag{5.12}$$

5.2. Compensated Horner scheme

We now propose an error-free transformation for polynomial evaluation with the Horner scheme. We produce four polynomial error terms monomial-by-monomial: a monomial for each polynomial at each iteration.

**Algorithm 5.2** (Error-free transformation for the Horner scheme).

```
function [res, pπ, pμ, pν, pσ] = EFTHorner(p, x)
    sn = an
    for i = n - 1 : -1 : 0
        [pi, πi, μi, νi] = TwoProductCplx(si+1, x)
        [si, σi] = TwoSumCplx(pi, ai)
        Set πi, μi, νi, σi respectively as the coefficient of degree i in pπ, pμ, pν, pσ
    end
    res = s0
```

The next theorems and proofs are very similar to the ones of [4]. It is just necessary to change real error-free transformations into complex error-free transformations and to change  $\epsilon$  into  $\sqrt{2}\gamma_2$ . This leads to change the  $\gamma_n$  into  $\tilde{\gamma}_n$ .

**Theorem 5.2** (Equality). Let  $p(x) = \sum_{i=0}^n a_i x^i$  be a polynomial of degree  $n$  with complex floating point coefficients, and let  $x$  be a complex floating point value. Then Algorithm 5.2 computes both

- i) the floating point evaluation  $\text{res} = \text{Horner}(p, x)$  and
- ii) four polynomials  $p_\pi, p_\mu, p_\nu$  and  $p_\sigma$  of degree  $n - 1$  with complex floating point coefficients.

Then,

$$p(x) = \text{res} + (p_\pi + p_\sigma + p_\mu + p_\nu)(x). \quad (5.13)$$

**Proof.** Thanks to the error-free transformations, we have  $p_i + \pi_i + \mu_i + \nu_i = s_{i+1} \cdot x$  and  $s_i + \sigma_i = p_i + a_i$ . By induction, it is easy to show that

$$\sum_{i=0}^n a_i x^i = s_0 + \sum_{i=0}^{n-1} \pi_i x^i + \sum_{i=0}^{n-1} \mu_i x^i + \sum_{i=0}^{n-1} \nu_i x^i + \sum_{i=0}^{n-1} \sigma_i x^i,$$

which is exactly (5.13).  $\square$

**Proposition 5.3** (Bound on the error). Given  $p(x) = \sum_{i=0}^n a_i x^i$  a polynomial of degree  $n$  with complex floating point coefficients, and  $x$  a complex floating point value. Let  $\text{res}$  be the floating point value,  $p_\pi, p_\mu, p_\nu$  and  $p_\sigma$  be the four polynomials of degree  $n - 1$ , with complex floating point coefficients, such that  $[\text{res}, p_\pi, p_\mu, p_\nu, p_\sigma] = \text{EFTHorner}(p, x)$ . Then,

$$((p_\pi + \widetilde{p_\mu} + p_\nu) + \widetilde{p_\sigma})(|x|) \leq \widetilde{\gamma}_{2n} \widetilde{p}(|x|).$$

**Proof.** The proof is organized as follows: we prove a bound on  $|p_{n-i}| |x|^{n-i}$  and  $|s_{n-i}| |x|^{n-i}$  from which we deduce a bound on  $|\pi_i + \mu_i + \nu_i|$  and  $|\sigma_i|$  and we use these bounds on each coefficient of the polynomial error terms to obtain finally the expected bound on these polynomials.

- By definition, for  $i = 1, \dots, n$ ,  $p_{n-i} = s_{n-i+1} \cdot x$  and  $s_{n-i} = p_{n-i} + a_{n-i}$ . From Eq. (3.6), we deduce  $\text{fl}(s_{n-i+1} \cdot x) = (1 + \varepsilon_1) s_{n-i+1} \cdot x$  with  $|\varepsilon_1| \leq \sqrt{2} \gamma_2$ . From Eq. (3.5), we deduce  $\text{fl}(p_{n-i} + a_{n-i}) = (1 + \varepsilon_2)(p_{n-i} + a_{n-i})$  with  $|\varepsilon_2| \leq \text{eps} \leq \sqrt{2} \gamma_2$ . Consequently

$$|p_{n-i}| \leq (1 + \sqrt{2} \gamma_2) |s_{n-i+1}| |x| \quad \text{and} \quad |s_{n-i}| \leq (1 + \sqrt{2} \gamma_2) (|p_{n-i}| + |a_{n-i}|). \quad (5.14)$$

- These two bounds will be used in the basic case and the inductive case of the following double property: for  $i = 1, \dots, n$ ,

$$|p_{n-i}| \leq (1 + \widetilde{\gamma}_{2i-1}) \sum_{j=1}^i |a_{n-i+j}| |x^j| \quad \text{and} \quad |s_{n-i}| \leq (1 + \widetilde{\gamma}_{2i}) \sum_{j=0}^i |a_{n-i+j}| |x^j|. \quad (5.15)$$

For  $i = 1$ :

Since  $s_n = a_n$  the bound of Eq. (5.14) can be rewritten as  $|p_{n-1}| \leq (1 + \sqrt{2} \gamma_2) |a_n| |x| \leq (1 + \widetilde{\gamma}_1) |a_n| |x|$ . We combine this bound on  $|p_{n-1}|$  to (5.14) to obtain  $|s_{n-1}| \leq (1 + \sqrt{2} \gamma_2) ((1 + \widetilde{\gamma}_1) |a_n| |x| + |a_{n-1}|) \leq (1 + \widetilde{\gamma}_2) (|a_n| |x| + |a_{n-1}|)$ . Thus (5.15) is satisfied for  $i = 1$ .

Let us now suppose that (5.15) is true for some integer  $i$  such that  $1 \leq i < n$ . According to (5.14), we have  $|p_{n-(i+1)}| \leq (1 + \sqrt{2} \gamma_2) |s_{n-i}| |x|$ . Thanks to the induction hypothesis, we derive,

$$|p_{n-(i+1)}| \leq (1 + \sqrt{2} \gamma_2) (1 + \widetilde{\gamma}_{2i}) \sum_{j=0}^i |a_{n-i+j}| |x^{j+1}| \leq (1 + \widetilde{\gamma}_{2(i+1)-1}) \sum_{j=1}^{i+1} |a_{n-(i+1)+j}| |x^j|.$$

Let us combine (5.14) with this inequality, we have,

$$\begin{aligned} |s_{n-(i+1)}| &\leq (1 + \sqrt{2} \gamma_2) (|p_{n-(i+1)}| + |a_{n-(i+1)}|) \\ &\leq (1 + \sqrt{2} \gamma_2) (1 + \widetilde{\gamma}_{2(i+1)-1}) \left[ \sum_{j=1}^{i+1} |a_{n-(i+1)+j}| |x^j| + |a_{n-(i+1)}| \right] \\ &\leq (1 + \widetilde{\gamma}_{2(i+1)}) \sum_{j=0}^{i+1} |a_{n-(i+1)+j}| |x^j|. \end{aligned}$$

So (5.15) is proved by induction. We bound each of these sums by  $p(|x|)/|x|^{n-i}$  and obtain for  $i = 1, \dots, n$ ,

$$|p_{n-i}| |x|^{n-i} \leq (1 + \widetilde{\gamma}_{2i-1}) \widetilde{p}(|x|) \quad \text{and} \quad |s_{n-i}| |x|^{n-i} \leq (1 + \widetilde{\gamma}_{2i}) \widetilde{p}(|x|). \quad (5.16)$$

- From Theorem 3.1 and Theorem 3.2, for  $i = 0, \dots, n - 1$ , we have  $|\pi_i + \mu_i + \nu_i| \leq \sqrt{2}\gamma_2|p_i|$  and  $|\sigma_i| \leq \text{eps}|s_i| \leq \sqrt{2}\gamma_2|s_i|$ . Therefore,

$$((p_\pi + \widetilde{p}_\mu + p_\nu) + \widetilde{p}_\sigma)(|x|) = \sum_{i=0}^{n-1} (|\pi_i + \mu_i + \nu_i| + |\sigma_i|)|x^i| \leq \sum_{i=0}^{n-1} (\sqrt{2}\gamma_2|p_i|)|x^i| + \sum_{i=0}^{n-1} (\sqrt{2}\gamma_2|s_i|)|x^i|.$$

We now transform the summation into

$$((p_\pi + \widetilde{p}_\mu + p_\nu) + \widetilde{p}_\sigma)(|x|) \leq \sqrt{2}\gamma_2 \sum_{i=1}^n (|p_{n-i}| |x^{n-i}| + |s_{n-i}| |x^{n-i}|)$$

and use the preceding equation (5.16) and the growth of the sequence  $\tilde{\gamma}_k$  so that

$$\begin{aligned} ((p_\pi + \widetilde{p}_\mu + p_\nu) + \widetilde{p}_\sigma)(|x|) &\leq \sqrt{2}\gamma_2 \sum_{i=1}^n ((1 + \tilde{\gamma}_{2i-1})\tilde{p}(|x|) + (1 + \tilde{\gamma}_{2i})\tilde{p}(|x|)) \\ &\leq \sqrt{2}\gamma_2 \sum_{i=1}^n 2(1 + \tilde{\gamma}_{2n})\tilde{p}(|x|) = 2n\sqrt{2}\gamma_2(1 + \tilde{\gamma}_{2n})\tilde{p}(|x|). \end{aligned}$$

Since  $2n\sqrt{2}\gamma_2(1 + \tilde{\gamma}_{2n}) = \tilde{\gamma}_{2n}$ , we finally obtain  $((p_\pi + \widetilde{p}_\mu + p_\nu) + \widetilde{p}_\sigma)(|x|) \leq \tilde{\gamma}_{2n}\tilde{p}(|x|)$ .  $\square$

From Theorem 5.2 the forward error affecting the evaluation of  $p$  at  $x$  according to the Horner scheme is

$$e(x) = p(x) - \text{Horner}(p, x) = (p_\pi + p_\mu + p_\nu + p_\sigma)(x).$$

The coefficients of these polynomials are exactly computed by Algorithm 5.2, together with  $\text{Horner}(p, x)$ .

If we try to compute a complete error-free transformation for the evaluation of a polynomial of degree  $n$ , we will have to perform recursively the same computation for four polynomials of degree  $n - 1$  and so on. This will produce at the end of the computation  $\sum_{i=0}^n 4^i = \frac{4^{n+1}-1}{4-1}$  error terms (for example for a polynomial of degree 10 we would obtain more than one million error terms), almost all of which are null with underflow and the other ones do not have the essential nonoverlapping property. It will takes a very long time to compute this result (even more probably than with exact symbolic computation) and we will have to make a drastic selection on the huge amount of data to keep only a few meaningful terms as a usable result. We only consider here intentionally the first-order error term to obtain a really satisfactory improvement of the result of the evaluation with a reasonable running time.

Consequently we compute here a single complex floating point number as the first-order error term, the most significant correction term. The key is then to compute an approximate of the error  $e(x)$  in working precision, and then to compute a corrected result  $\text{res}' = \text{fl}(\text{Horner}(p, x) + e(x))$ .

Our aim is now to compute the correction term  $c = \text{fl}(e(x)) = \text{fl}((p_\pi + p_\sigma + p_\mu + p_\nu)(x))$ . For that we evaluate the polynomial  $P$  whose coefficients are those of  $p_\pi + p_\sigma + p_\mu + p_\nu$  faithfully rounded<sup>2</sup> since the sums of the coefficients  $p_i + q_i + r_i + s_i$  are not necessarily floating point numbers. We compute the coefficients of polynomial  $P$  thanks to `Accsum` algorithm [3]. This can also be done via other accurate summation algorithms (see [15] for example). We could not use `Sum2` because even with only 4 numbers to sum up, this algorithm could not guarantee a good accuracy of the result. We modify the classical Horner scheme applied to  $P$ , to compute  $P$  at the same time.

**Algorithm 5.3** (Evaluation of the sum of four polynomials with degree  $n$ ).

```
function c = HornerSumAcc(p, q, r, s, x)
    v_n = Accsum(p_n + q_n + r_n + s_n)
    for i = n - 1 : -1 : 0
        v_i = fl(v_{i+1} * x + Accsum(p_i + q_i + r_i + s_i))
    end
    c = v_0
```

**Lemma 5.4.** Let us consider the floating point evaluation of  $(p + q + r + s)(x)$  computed with `HornerSumAcc(p, q, r, s, x)`. Then, the computed result satisfies the following forward error bound,

$$|\text{HornerSumAcc}(p, q, r, s, x) - (p + q + r + s)(x)| \leq \tilde{\gamma}_{2n+1}((p + q + r) + \tilde{s})(|x|).$$

<sup>2</sup> Faithful rounding means that the computed result is equal to the exact result if the latter is a floating point number and otherwise is one of the two adjacent floating point numbers of the exact result.

**Proof.** We will use as in [7, p. 68] the notation  $\langle k \rangle$  to denote the product of  $k$  terms of the form  $1 + \varepsilon_i$  for some  $\varepsilon_i$  such that  $|\varepsilon_i| \leq \sqrt{2}\gamma_2$ . A product of  $j$  such terms multiplied by the product of  $k$  such terms is a product of  $j+k$  such terms and consequently we have  $\langle j \rangle \langle k \rangle = \langle j+k \rangle$ .

Considering [Algorithm 5.3](#), we have  $v_n = \text{Accsum}(p_n + q_n + r_n + s_n)$  so according to the property of the `Accsum` algorithm we have  $v_n = (p_n + q_n + r_n + s_n)\langle 1 \rangle$ .

For  $i = n-1, \dots, 0$ , the computation of  $v_i$  from  $v_{i+1}$  leads to an error term for the product and for the `Accsum` algorithm and then another on the sum and we have

$$v_i = \text{fl}(v_{i+1}x + \text{Accsum}(p_i + q_i + r_i + s_i)) = v_{i+1}x\langle 2 \rangle + (p_i + q_i + r_i + s_i)\langle 2 \rangle.$$

Therefore we can prove by induction on  $i$  that

$$v_{n-i} = (p_n + q_n + r_n + s_n)x^i\langle 2i+1 \rangle + \sum_{k=0}^{i-1} (p_{n-i+k} + q_{n-i+k} + r_{n-i+k} + s_{n-i+k})x^k\langle 2(k+1) \rangle$$

and then for  $i = n$  we obtain

$$c = v_0 = (p_n + q_n + r_n + s_n)x^n\langle 2n+1 \rangle + \sum_{k=0}^{n-1} (p_k + q_k + r_k + s_k)x^k\langle 2(k+1) \rangle.$$

Consequently we have

$$c - \sum_{i=0}^n (p_i + q_i + r_i + s_i)x^i = (p_n + q_n + r_n + s_n)x^n((2n+1) - 1) + \sum_{k=0}^{n-1} (p_k + q_k + r_k + s_k)x^k((2(k+1)) - 1).$$

Since for any  $\varepsilon$  implied in  $\langle k \rangle$  notation, we have  $|\varepsilon| \leq \sqrt{2}\gamma_2$ , we have

$$|\langle k \rangle - 1| \leq (1 + \sqrt{2}\gamma_2)^k - 1 \leq \frac{1}{1 - k\sqrt{2}\gamma_2} - 1 = \frac{k\sqrt{2}\gamma_2}{1 - k\sqrt{2}\gamma_2} = \tilde{\gamma}_k$$

and the  $\tilde{\gamma}_k$  sequence is growing, thus  $|\langle k \rangle - 1| \leq \tilde{\gamma}_k \leq \widetilde{\gamma_{2n+1}}$  pour tout  $k \leq 2n+1$ . We finally obtain

$$\left| c - \sum_{i=0}^n (p_i + q_i + r_i + s_i)x^i \right| \leq \widetilde{\gamma_{2n+1}} \sum_{i=0}^n (|p_i + q_i + r_i| + |s_i|)|x^i| \leq \widetilde{\gamma_{2n+1}}((\widetilde{p+q+r}) + \widetilde{s})(|x|). \quad \square$$

We combine now the error-free transformation for the Horner scheme that produces four polynomials and the algorithm for the evaluation of the sum of four polynomials to obtain a compensated Horner scheme algorithm that improves the numerical accuracy of the classical Horner scheme on complex numbers.

**Algorithm 5.4** (Compensated Horner scheme).

```
function res' = CompHorner(p, x)
    [res, pπ, pμ, pν, pσ] = EFTHorner(p, x)
    c = HornerSumAcc(pπ, pμ, pν, pσ, x)
    res' = fl(res + c)
```

We prove hereafter that the result of a polynomial evaluation computed with the compensated Horner scheme [Algorithm 5.4](#) is as accurate as if computed by the classic Horner scheme using twice the working precision and then rounded to the working precision.

**Theorem 5.5.** Given a polynomial  $p = \sum_{i=0}^n p_i x^i$  of degree  $n$  with floating point coefficients, and  $x$  a floating point value. We consider the result `CompHorner(p, x)` computed by [Algorithm 5.4](#). Then,

$$|\text{CompHorner}(p, x) - p(x)| \leq \text{eps}|p(x)| + \widetilde{\gamma_{2n}^2} \widetilde{p}(|x|). \quad (5.17)$$

**Proof.** As  $\text{res}' = \text{fl}(\text{res} + c)$  so, according to [Theorem 3.1](#),  $\text{res}' = (1 + \varepsilon)(\text{res} + c)$  with  $|\varepsilon| \leq \text{eps} \leq \sqrt{2}\gamma_2$ . Thus we have  $|\text{res}' - p(x)| = |\text{fl}(\text{res} + c) - p(x)| = |(1 + \varepsilon)(\text{res} + c - p(x)) + \varepsilon p(x)|$ . Since  $p(x) = \text{res} + e(x)$ , we have  $|\text{res}' - p| = |(1 + \varepsilon)(c - e(x)) + \varepsilon p(x)| \leq \text{eps}|p(x)| + (1 + \text{eps})|e(x) - c|$ . By [Lemma 5.4](#) applied to four polynomials of degree  $n-1$ , we have

$$|e(x) - c| \leq \widetilde{\gamma_{2n-1}}((\widetilde{p\pi} + \widetilde{p\mu} + \widetilde{p\nu}) + \widetilde{p\sigma})(|x|).$$

By Proposition 5.3 we have also  $((p_\pi + \widetilde{p}_\mu + p_\nu) + \widetilde{p}_\sigma)(|x|) \leq \widetilde{\gamma}_{2n} \widetilde{p}(|x|)$ . We combine these two bounds and obtain  $|e(x) - c| \leq \widetilde{\gamma}_{2n-1} \widetilde{\gamma}_{2n} \widetilde{p}(|x|)$ . As a consequence,  $|res' - p(x)| \leq \text{eps}|p(x)| + (1 + \sqrt{2}\gamma_2)\widetilde{\gamma}_{2n-1}\widetilde{\gamma}_{2n}\widetilde{p}(|x|)$ . Since  $(1 + \sqrt{2}\gamma_2)\widetilde{\gamma}_{2n-1} \leq \widetilde{\gamma}_{2n}$ , it follows that  $|res' - p(x)| \leq \text{eps}|p(x)| + \widetilde{\gamma}_{2n}^2 \widetilde{p}(x)$ .  $\square$

### 5.3. Numerical experiments

In this section, we will compare our compensated algorithm to other algorithms both in term of accuracy and computing time.

#### 5.3.1. Accuracy comparisons

Eq. (5.17) can be written

$$\frac{|\text{CompHorner}(p, x) - p(x)|}{|p(x)|} \leq \text{eps} + \widetilde{\gamma}_{2n}^2 \text{cond}(p, x). \tag{5.18}$$

The comparison with the bound (5.12) for the classical Horner scheme shows that the coefficient of the condition number vanish from  $\widetilde{\gamma}_{2n}$  to  $\widetilde{\gamma}_{2n}^2$ .

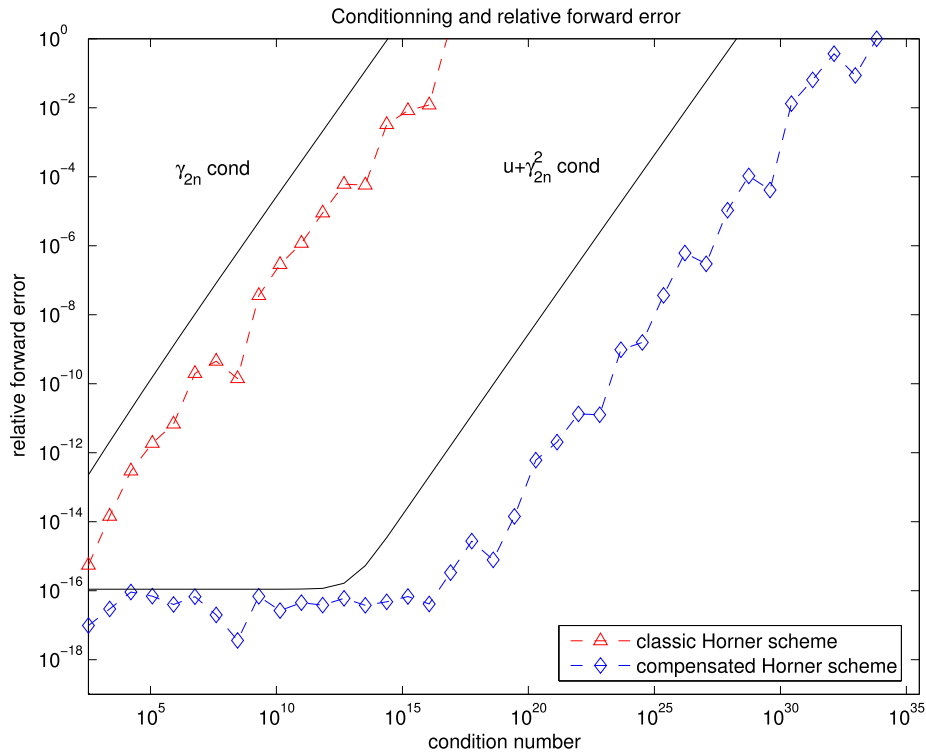
We present here comparison curves for the classical and the compensated Horner scheme.

All our experiments are performed using the IEEE 754 double precision with MATLAB 7. When needed, we use the Symbolic Math Toolbox to accurately compute the polynomial evaluation (in order to compute the relative forward error).

We test the compensated Horner scheme on the expanded form of the polynomial  $p_n(x) = (x - (1+i))^n$  at  $x = \text{fl}(1.333 + 1.333i)$  for  $n = 3 : 42$ . The condition number  $\text{cond}(p_n, x)$  varies from  $10^3$  to  $10^{33}$ .

The following figure shows the relative accuracy  $|res - p_n(x)|/|p_n(x)|$  where  $res$  is the computed value by the two Algorithms 5.1 and 5.4. We also plot the *a priori* error estimations (5.12) and (5.18).

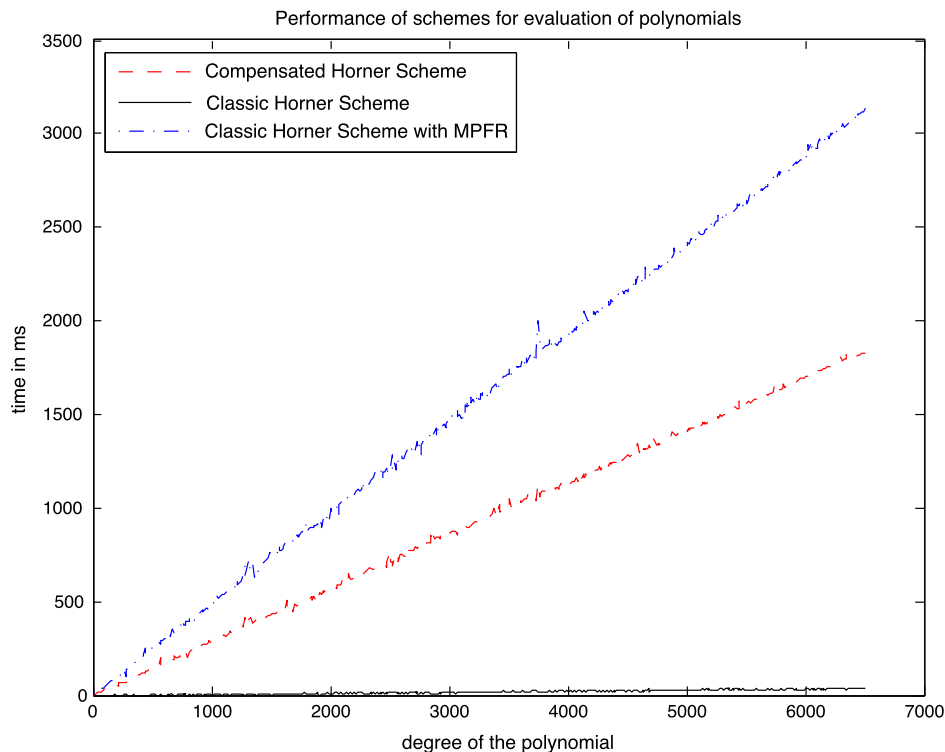
As we can see below, the compensated Horner scheme exhibits the expected behavior, that is to say, the compensated rule of thumb (5.18). As long as the condition number is less than  $\text{eps}^{-1} \approx 10^{16}$ , the compensated Horner scheme produces results with full precision (forward relative error of the order of  $\text{eps} \approx 10^{-16}$ ). For condition numbers greater than  $\text{eps}^{-1} \approx 10^{16}$ , the accuracy decreases until no accuracy at all when the condition number is greater than  $\text{eps}^{-2} \approx 10^{32}$ .



#### 5.3.2. Performance comparisons

We have compared in term of computing time three algorithms: the classic Horner scheme, the compensated Horner scheme, and the classic Horner scheme using MPFR [16] with 106 bits of mantissa. The computations were performed on quad-core Core i7 M620 at 2.67 GHz with 4 GB of RAM and 4 MB of cache memory. We used gcc-4.5.2 compiler and for multiprecision mpfr-3.0.1/gmp-5.0.2.

The results are presented in the figure below. For that, we randomly generated some polynomials (coefficients chosen in  $[-5; 5]$ ) with different degrees varying from 1 to 6500. Then we evaluated those polynomials on a randomly chosen point (sill in  $[-5; 5]$ ) and we measured the computing time. As one can see, our algorithm is, of course, less efficient than the classic Horner scheme but provides much more accuracy. Nevertheless, compared the classic Horner scheme with 106 bits en precision (via MPFR), we are faster while sharing the same accuracy.



## 6. Conclusion and future work

In this article, we derived some new error-free transformations for complex floating point arithmetic. This makes it possible to provide a complex version of the compensated Horner scheme.

Nevertheless, the error bound provided in this article is a theoretical one since it contains the quantity  $|p(x)|$ . It would be very interesting to derive a validated error bound  $\alpha \in \mathbb{F}$  that can be computed in floating point arithmetic satisfying  $|\text{CompHorner}(p, x) - p(x)| \leq \alpha$ . This can be done via a kind of running error analysis [17].

We have also provided some new algorithms to compute sum and dot product of complex floating point numbers.

## References

- [1] X.S. Li, J.W. Demmel, D.H. Bailey, G. Henry, Y. Hida, J. Iskandar, W. Kahan, S.Y. Kang, A. Kapur, M.C. Martin, B.J. Thompson, T. Tung, D.J. Yoo, Design, implementation and testing of extended and mixed precision BLAS, *ACM Trans. Math. Software* 28 (2) (2002) 152–205.
- [2] T. Ogita, S.M. Rump, S. Oishi, Accurate sum and dot product, *SIAM J. Sci. Comput.* 26 (6) (2005) 1955–1988.
- [3] S.M. Rump, T. Ogita, S. Oishi, Accurate floating-point summation, Tech. Rep. 05.12, Faculty for Information and Communication Sciences, Hamburg University of Technology, November 2005.
- [4] S. Graillat, N. Louvet, P. Langlois, Compensated Horner scheme, Research Report 04, Équipe de recherche DALI, Laboratoire LP2A, Université de Perpignan Via Domitia, France, July 2005.
- [5] S. Graillat, V. Ménéssier-Morain, Compensated Horner scheme in complex floating point arithmetic, in: *Proceedings of the 8th Conference on Real Numbers and Computers*, Santiago de Compostela, Spain, July 7–9, 2008, pp. 133–146.
- [6] IEEE Standard for Binary Floating-Point Arithmetic, ANSI/IEEE Standard 754-1985, Institute of Electrical and Electronics Engineers, New York, 1985, reprinted in *ACM SIGPLAN Notices* 22 (2) (1987) 9–25.
- [7] N.J. Higham, *Accuracy and Stability of Numerical Algorithms*, 2nd edition, Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002.
- [8] T.J. Dekker, A floating-point technique for extending the available precision, *Numer. Math.* 18 (1971) 224–242.
- [9] D.E. Knuth, *The Art of Computer Programming*, vol. 2. *Seminumerical Algorithms*, 3rd edition, Addison-Wesley, Reading, MA, 1998.
- [10] Y. Nievergelt, Scalar fused multiply-add instructions produce floating-point matrix arithmetic provably accurate to the penultimate digit, *ACM Trans. Math. Software* 29 (1) (2003) 27–48.
- [11] S. Boldo, J.-M. Muller, Some functions computable with a Fused-mac, in: *Proceedings of the 17th Symposium on Computer Arithmetic*, Cape Cod, USA, 2005.
- [12] S.M. Rump, Verification of positive definiteness, *BIT* 46 (2) (2006) 433–452.

- [13] R. Brent, C. Percival, P. Zimmermann, Error bounds on complex floating-point multiplication, *Math. Comp.* 76 (259) (2007) 1469–1481 (electronic).
- [14] S. Graillat, V. Ménissier-Morain, Error-free transformations in real and complex floating point arithmetic, in: *Proceedings of the International Symposium on Nonlinear Theory and Its Applications*, Vancouver, Canada, September 16–19, 2007, pp. 341–344.
- [15] J.W. Demmel, Y. Hida, Accurate and efficient floating point summation, *SIAM J. Sci. Comput.* 25 (4) (2003) 1214–1248, (electronic).
- [16] L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier, P. Zimmermann, MPFR: A multiple-precision binary floating-point library with correct rounding, *ACM Trans. Math. Software* 33 (2) (2007) 13:1–13:15, <http://www.mpfr.org>.
- [17] J.H. Wilkinson, *Rounding Errors in Algebraic Processes*, Prentice Hall Inc., Englewood Cliffs, NJ, 1963.