



HAL
open science

Accurate Complex Multiplication in Floating-Point Arithmetic

Vincent Lefèvre, Jean-Michel Muller

► **To cite this version:**

Vincent Lefèvre, Jean-Michel Muller. Accurate Complex Multiplication in Floating-Point Arithmetic. 2019. hal-02001080v1

HAL Id: hal-02001080

<https://hal.science/hal-02001080v1>

Preprint submitted on 31 Jan 2019 (v1), last revised 4 Apr 2019 (v2)

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

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

Accurate Complex Multiplication in Floating-Point Arithmetic

Vincent Lefèvre
Univ Lyon, Inria,
CNRS, ENS de Lyon, Université Claude Bernard Lyon 1, LIP

Jean-Michel Muller
Univ Lyon, CNRS,
ENS de Lyon, Inria, Université Claude Bernard Lyon 1, LIP

January 31, 2019

Abstract

We deal with accurate complex multiplication in binary floating-point arithmetic, with an emphasis on the case where one of the operands in a “double-word” number. We provide an algorithm that returns a complex product with normwise relative error bound close to the best possible one, i.e., the rounding unit u .

Keywords. Floating-point arithmetic, Complex multiplication, Rounding error analysis.

1 Introduction and notation

This paper deals with accurate (from a normwise point of view) complex multiplication in binary floating-point (FP) arithmetic, with an emphasis on the case where one of the operands in a “double-word” number. This is of interest, for instance, when that operand is a root of one (which is the case in fast Fourier transforms), precomputed and stored in higher precision than standard FP precision. This can also be of interest for computing iterated products of several complex numbers accurately. In the following, we assume a radix-2, precision- p FP arithmetic, with correctly-rounded (to nearest) arithmetic operations. We assume that an FMA (Fused Multiply-Add) instruction is available and, to simplify our study, we also assume an unbounded exponent range. This means that the results presented in this paper apply to “real-life” IEEE 754 Floating-Point arithmetic [6] as long as underflows and overflows do not occur.

We will denote $u = 2^{-p}$ the “rounding unit”. For instance in the binary64 format of the IEEE 754 Standard (a.k.a. “double precision”), $u = 2^{-53}$. RN is the round-to-nearest function (with any choice in case of a tie). For instance,

when performing the operation $a + b$, where a and b are floating-point numbers, the obtained result is $\text{RN}(a + b)$, and it satisfies:

$$|\text{RN}(a + b) - (a + b)| \leq \frac{u}{1 + u} \cdot |a + b| < u \cdot |a + b|.$$

Assuming one wishes to compute an approximation $\hat{z}^R + i\hat{z}^I$ to $z^R + iz^I = (x^R + ix^I) \cdot (y^R + iy^I)$, where $x^R, x^I, y^R, y^I, \hat{z}^R$, and \hat{z}^I are FP numbers, one may consider the following “naive” formulas:

- if no FMA instruction is available

$$\begin{cases} \hat{z}^R &= \text{RN}(\text{RN}(x^R y^R) - \text{RN}(x^I y^I)), \\ \hat{z}^I &= \text{RN}(\text{RN}(x^R y^I) + \text{RN}(x^I y^R)). \end{cases} \quad (1)$$

- if an FMA instruction is available

$$\begin{cases} \hat{z}^R &= \text{RN}(x^R y^R - \text{RN}(x^I y^I)), \\ \hat{z}^I &= \text{RN}(x^R y^I + \text{RN}(x^I y^R)). \end{cases} \quad (2)$$

When trying to minimize the relative error, one may be interested in minimizing the *componentwise error*

$$\max \left\{ \left| \frac{\hat{z}^R - z^R}{z^R} \right|; \left| \frac{\hat{z}^I - z^I}{z^I} \right| \right\},$$

or the *normwise error*

$$\left| \frac{\hat{z} - z}{z} \right|.$$

Our goal here is to obtain small *normwise* relative errors, i.e., small values of $|(\hat{z} - z)/z|$, where z is the exact product and \hat{z} is the computed product. Note that getting small *componentwise* relative errors would be a completely different problem: formulas (1) and (2), as well as the algorithm we give in this paper (Algorithm 3), can lead to large componentwise relative errors, and to obtain small componentwise errors, one needs to use significantly different algorithms, such as an algorithm attributed to Kahan by Higham [5, p. 65], analyzed in [8], or Cornea et al’s algorithm for $ab + cd$ presented in [2]. We will give anyway a componentwise error bound for our algorithm (in Lemma 1, below), but only because this will be an intermediate step for establishing our main result (Theorem 1).

An algorithm that would always return the best possible result (i.e., a real part equal to the FP number nearest to the exact real part, and an imaginary part equal to the FP number nearest to the exact imaginary part) would have worst case relative error $u/(1 + u) \approx u$: this is therefore the best error bound achievable by an algorithm that returns the real and imaginary part of the result in floating-point. Asymptotically optimal bounds on the normwise relative error of formulas (1) and (2) are known: Brent et al. [1] show a bound $\sqrt{5} \cdot u$ for (1), and Jeannerod et al. [7] show a bound $2 \cdot u$ for (2).

We aim at obtaining smaller normwise relative errors, closer to the best possible one, at the cost of more complex algorithms. We consider the product

$$\omega \times x,$$

with

$$x = x^R + i \cdot x^I,$$

where x^R and x^I are floating-point numbers, and

$$\omega = \omega^R + i \cdot \omega^I,$$

where ω^R and ω^I are either “double-word” numbers (frequently called “double-double” in the literature, because the underlying floating-point format is, in general, the binary64 format) [9, 4] or floating-point numbers. We will first consider the “double-word” case, i.e.,

$$\omega = \omega^R + i \cdot \omega^I = (\omega_h^R + \omega_\ell^R) + i \cdot (\omega_h^I + \omega_\ell^I),$$

where $\omega_h^R, \omega_\ell^R, \omega_h^I,$ and ω_ℓ^I are floating-point numbers that satisfy:

- $|\omega_\ell^R| \leq (1/2)\text{ulp}(\omega^R) \leq u \cdot |\omega^R|;$
- $|\omega_\ell^I| \leq (1/2)\text{ulp}(\omega^I) \leq u \cdot |\omega^I|.$

In Section 3.2, we deal with the case where ω^R and ω^I are floating-point numbers (i.e., $\omega_\ell^R = \omega_\ell^I = 0$).

We will need two well-known algorithms of the floating-point literature: Algorithm 2Sum (Algorithm 1 below), that takes two FP numbers a and b as input and returns two FP numbers s and t such that $s = \text{RN}(a + b)$ and $t = a + b - s$ (that is, t is the error of the floating-point addition of a and b), and Algorithm Fast2Mult (Algorithm 2 below), that requires the availability of an FMA instruction, and takes two FP numbers a and b as input and returns two FP numbers π and ρ such that $\pi = \text{RN}(ab)$ and $\rho = ab - \pi$ (that is, ρ is the error of the floating-point multiplication of a and b).

ALGORITHM 1: 2Sum(a, b). The 2Sum algorithm [12, 11].

```

s ← RN(a + b)
a' ← RN(s - b)
b' ← RN(s - a')
δa ← RN(a - a')
δb ← RN(b - b')
t ← RN(δa + δb)

```

ALGORITHM 2: Fast2Mult(a, b). The Fast2Mult algorithm (see for instance [10, 14, 13]). It requires the availability of a fused multiply-add (FMA) instruction for computing $\text{RN}(ab - \pi)$.

$\pi \leftarrow \text{RN}(ab)$
 $\rho \leftarrow \text{RN}(ab - \pi)$

2 The multiplication algorithm

For performing the complex multiplication $\omega \cdot x$, we introduce Algorithm 3 below. The real part (lines 1 to 9) and the imaginary part (lines 10 to 18) can obviously be computed in parallel, and within these parts, additional parallelism is possible. For instance lines 3 and 5 can run in parallel with line 1, and line 7 can run in parallel with line 6. This parallelism is easily exploited by recent compilers. This explains the good performance we obtain (see Section 4).

ALGORITHM 3: Accurate complex multiplication $\omega \cdot x$, where the real and imaginary parts of $\omega = (\omega_h^R + \omega_\ell^R) + i \cdot (\omega_h^I + \omega_\ell^I)$ are double-word numbers, and the real and imaginary parts of x are FP numbers.

1: $t^R \leftarrow \text{RN}(\omega_\ell^I x^I)$
2: $\pi_\ell^R \leftarrow \text{RN}(\omega_\ell^R x^R - t^R)$
3: $(P_h^R, P_\ell^R) \leftarrow \text{Fast2Mult}(\omega_h^I, x^I)$
4: $r_\ell^R \leftarrow \text{RN}(\pi_\ell^R - P_\ell^R)$
5: $(Q_h^R, Q_\ell^R) \leftarrow \text{Fast2Mult}(\omega_h^R, x^R)$
6: $s_\ell^R \leftarrow \text{RN}(Q_\ell^R + r_\ell^R)$
7: $(v_h^R, v_\ell^R) \leftarrow 2\text{Sum}(Q_h^R, -P_h^R)$
8: $\gamma_\ell^R \leftarrow \text{RN}(v_\ell^R + s_\ell^R)$
9: **return** $z^R = \text{RN}(v_h^R + \gamma_\ell^R)$ (real part)
10: $t^I \leftarrow \text{RN}(\omega_\ell^I x^R)$
11: $\pi_\ell^I \leftarrow \text{RN}(\omega_\ell^R x^I + t^I)$
12: $(P_h^I, P_\ell^I) \leftarrow \text{Fast2Mult}(\omega_h^I, x^R)$
13: $r_\ell^I \leftarrow \text{RN}(\pi_\ell^I + P_\ell^I)$
14: $(Q_h^I, Q_\ell^I) \leftarrow \text{Fast2Mult}(\omega_h^R, x^I)$
15: $s_\ell^I \leftarrow \text{RN}(Q_\ell^I + r_\ell^I)$
16: $(v_h^I, v_\ell^I) \leftarrow 2\text{Sum}(Q_h^I, P_h^I)$
17: $\gamma_\ell^I \leftarrow \text{RN}(v_\ell^I + s_\ell^I)$
18: **return** $z^I = \text{RN}(v_h^I + \gamma_\ell^I)$ (imaginary part)

Our main result is

Theorem 1. *The normwise relative error η of Algorithm 3 satisfies*

$$\eta \leq u + 33u^2 + \mathcal{O}(u^3).$$

The higher order terms hidden in the “ $\mathcal{O}(u^3)$ ” term of Theorem 1 can be made totally explicit: Eq. (18) below gives a bound on η^2 with all higher order

terms, and one easily shows from (18) that as soon as $u \leq 1/16$ (i.e., the precision p is larger than or equal to 4, which always holds in practice), $\eta \leq u + 34u^2$. For instance, in binary64 arithmetic, (18) implies an error bound $1.0000000000000037u$, and in binary32 arithmetic, the obtained error bound is $1.000002u$.

Note that Algorithm 3 can easily be transformed into an algorithm that returns the real and imaginary parts of the product as double-word numbers: it suffices to replace the floating-point additions of lines 9 and 18 by a call to 2Sum (or to the somehow simpler Fast2Sum algorithm, see for instance [13]). We deal with this solution in Section 3.1.

Let us now prove Theorem 1.

Proof. In the following, to simplify the notation, we will write

$$N^R = |\omega^R x^R| + |\omega^I x^I|,$$

and

$$n^R = |\omega^R x^R - \omega^I x^I|.$$

Later on, we will need similar notation for the imaginary part, i.e.,

$$N^I = |\omega^R x^I| + |\omega^I x^R|,$$

and

$$n^I = |\omega^R x^I + \omega^I x^R|.$$

We will focus on the calculation of the real part of the complex product (i.e., lines 1 to 9 of the algorithm), since the results for the real part hold for the imaginary part through a simple symmetry argument (namely, the imaginary part of $(a + ib) \cdot (c + id)$ is equal to the real part of $(b - ia) \cdot (c + id)$).

2.1 Lines 1-2 of Algorithm 3: computation of an approximation π_ℓ^R to $(\omega_\ell^R x^R - \omega_\ell^I x^I)$.

We have

$$|t^R - (\omega_\ell^I x^I)| \leq u^2 \cdot |\omega^I x^I|, \quad (3)$$

and $|t^R| \leq (u + u^2) \cdot |\omega^I x^I|$, so that

$$\begin{aligned} |\omega_\ell^R x^R - t^R| &\leq u \cdot |\omega^R x^R| + (u + u^2) \cdot |\omega^I x^I| \\ &\leq (u + u^2) \cdot (|\omega^R x^R| + |\omega^I x^I|) \\ &= (u + u^2) \cdot N^R, \end{aligned}$$

a consequence of which is

$$|\pi_\ell^R - (\omega_\ell^R x^R - t^R)| \leq (u^2 + u^3) \cdot N^R, \quad (4)$$

and

$$|\pi_\ell^R| \leq (u + 2u^2 + u^3) \cdot N^R.$$

From (3) and (4), we obtain

$$|\pi_\ell^R - (\omega_\ell^R x^R - \omega_\ell^I x^I)| \leq \epsilon_1, \quad (5)$$

with

$$\epsilon_1 = (2u^2 + u^3) \cdot N^R.$$

2.2 Line 3.

We have

$$P_h^R + P_\ell^R = \omega_h^I x^I,$$

with

$$|P_\ell^R| \leq u \cdot |\omega_h^I x^I| \leq u(1+u) \cdot |\omega^I x^I|,$$

and

$$|P_h^R| \leq (1+u)^2 \cdot |\omega^I x^I|.$$

2.3 Line 4.

From

$$\begin{aligned} |\pi_\ell^R - P_\ell^R| &\leq (u+u^2) \cdot |\omega^I x^I| + (u+2u^2+u^3) \cdot N^R \\ &\leq (2u+3u^2+u^3) \cdot N^R, \end{aligned}$$

we obtain

$$|r_\ell^R - (\pi_\ell^R - P_\ell^R)| \leq \epsilon_2, \quad (6)$$

with

$$\epsilon_2 = (2u^2 + 3u^3 + u^4) \cdot N^R,$$

and

$$|r_\ell^R| \leq (2u + 5u^2 + 4u^3 + u^4) \cdot N^R, \quad (7)$$

and, using (5) and (6),

$$\begin{aligned} |r_\ell^R - (\omega_\ell^R x^R - \omega_\ell^I x^I - P_\ell^R)| \\ \leq (4u^2 + 4u^3 + u^4) \cdot N^R. \end{aligned} \quad (8)$$

2.4 Lines 5 and 6.

We have

$$Q_h^R + Q_\ell^R = \omega_h^R x^R,$$

and

$$|Q_\ell^R| \leq u(1+u) \cdot |\omega^R x^R|.$$

From this bound on $|Q_\ell^R|$ and (7), we obtain

$$|Q_\ell^R + r_\ell^R| \leq (3u + 6u^2 + 4u^3 + u^4) \cdot N^R,$$

from which we deduce

$$|s_\ell^R - (Q_\ell^R + r_\ell^R)| \leq \epsilon_3, \quad (9)$$

with

$$\epsilon_3 = (3u^2 + 6u^3 + 4u^4 + u^5) \cdot N^R,$$

and

$$|s_\ell^R| \leq (3u + 9u^2 + 10u^3 + 5u^4 + u^5) \cdot N^R. \quad (10)$$

All this gives

$$\begin{aligned} (Q_h^R - P_h^R + s_\ell^R) &= (Q_h^R - P_h^R + Q_\ell^R + r_\ell^R) \\ &\quad + \xi_3 \\ &= (Q_h^R + Q_\ell^R - P_h^R + (\pi_\ell^R - P_\ell^R)) \\ &\quad + \xi_3 + \xi_2 \\ &= (\omega_h^R x^R - \omega_h^I x^I + \omega_\ell^R x^R - \omega_\ell^I x^I) \\ &\quad + \xi_3 + \xi_2 + \xi_1, \end{aligned}$$

with

$$\begin{aligned} |\xi_3| &\leq \epsilon_3 = (3u^2 + 6u^3 + 4u^4 + u^5) \cdot N^R, \\ |\xi_2| &\leq \epsilon_2 = (2u^2 + 3u^3 + u^4) \cdot N^R, \\ |\xi_1| &\leq \epsilon_1 = (2u^2 + u^3) \cdot N^R. \end{aligned}$$

This implies

$$\begin{aligned} |(Q_h^R - P_h^R + s_\ell^R) - (\omega_h^R x^R - \omega_h^I x^I)| \\ \leq (7u^2 + 10u^3 + 5u^4 + u^5) \cdot N^R. \end{aligned} \quad (11)$$

2.5 Lines 7 to 9.

We have

$$\begin{aligned} |Q_h^R - P_h^R| &= |\omega_h^R x^R - \omega_h^I x^I - Q_\ell^R + P_\ell^R| \\ &\leq n^R + (2u + u^2) \cdot N^R, \end{aligned}$$

hence,

$$|v_\ell^R| \leq u \cdot n^R + (2u^2 + u^3) \cdot N^R,$$

and

$$|v_h^R| \leq (1 + u) \cdot n^R + (2u + 3u^2 + u^3) \cdot N^R.$$

Therefore, using (10),

$$|v_\ell^R + s_\ell^R| \leq u \cdot n^R + (3u + 11u^2 + 11u^3 + 5u^4 + u^5) \cdot N^R,$$

so that

$$\begin{aligned} |\gamma_\ell^R - (v_\ell^R + s_\ell^R)| \\ \leq u^2 \cdot n^R \\ + (3u^2 + 11u^3 + 11u^4 + 5u^5 + u^6) \cdot N^R, \end{aligned} \quad (12)$$

and

$$\begin{aligned} |\gamma_\ell^R| &\leq (u + u^2) \cdot n^R \\ + (3u + 14u^2 + 22u^3 + 16u^4 + 6u^5 + u^6) \cdot N^R. \end{aligned}$$

This gives

$$\begin{aligned} |v_h^R + \gamma_\ell^R| \\ \leq (1 + 2u + u^2) \cdot n^R \\ + (5u + 17u^2 + 23u^3 + 16u^4 + 6u^5 + u^6) \cdot N^R, \end{aligned}$$

so that the error of the final addition $v_h^R + \gamma_\ell^R$ is bounded by

$$(u + 2u^2 + u^3) \cdot n^R + (5u^2 + 17u^3 + 23u^4 + 16u^5 + 6u^6 + u^7) \cdot N^R. \quad (13)$$

The bound on the total error (i.e., $|z^R - \Re(wx)|$) is obtained by adding the bounds (11), (12), and (13), i.e.,

$$(u + 3u^2 + u^3) \cdot n^R + (15u^2 + 38u^3 + 39u^4 + 22u^5 + 7u^6 + u^7) \cdot N^R. \quad (14)$$

We therefore deduce,

Lemma 1 (Componentwise absolute error of Algorithm 3). *We have*

$$\begin{aligned} |z^R - \Re(wx)| &\leq \alpha n^R + \beta N^R, \\ |z^I - \Im(wx)| &\leq \alpha n^I + \beta N^I, \end{aligned} \quad (15)$$

with

$$\alpha = u + 3u^2 + u^3$$

and

$$\beta = 15u^2 + 38u^3 + 39u^4 + 22u^5 + 7u^6 + u^7.$$

Lemma 1 gives a bound on the componentwise *absolute* error of Algorithm 3. However, it does not allow one to infer a useful bound on the componentwise *relative* error, since $|N^R/n^R|$ and $|N^I/n^I|$ can be arbitrarily large. However, the normwise relative error is always small, as we are going to see.

The square of the normwise relative error η is

$$\eta^2 = \frac{(z^R - \Re(\omega x))^2 + (z^I - \Im(\omega x))^2}{(\Re(\omega x))^2 + (\Im(\omega x))^2}.$$

From (15), we have

$$\begin{aligned} &(z^R - \Re(\omega x))^2 + (z^I - \Im(\omega x))^2 \\ &\leq \alpha^2 \left((n^R)^2 + (n^I)^2 \right) + 2\alpha\beta (n^R N^R + n^I N^I) \\ &\quad + \beta^2 \left((N^R)^2 + (N^I)^2 \right). \end{aligned}$$

We also have

$$\begin{aligned} &(\Re(\omega x))^2 + (\Im(\omega x))^2 = (n^R)^2 + (n^I)^2 \\ &= (\omega^R x^R)^2 + (\omega^I x^I)^2 + (\omega^R x^I)^2 + (\omega^I x^R)^2. \end{aligned}$$

Hence,

$$\begin{aligned} \eta^2 &\leq \alpha^2 + 2\alpha\beta \frac{n^R N^R + n^I N^I}{(n^R)^2 + (n^I)^2} \\ &\quad + \beta^2 \frac{(N^R)^2 + (N^I)^2}{(n^R)^2 + (n^I)^2}. \end{aligned} \quad (16)$$

We can notice that

$$\frac{n^R N^R + n^I N^I}{(n^R)^2 + (n^I)^2} \leq \frac{(N^R)^2 + (N^I)^2}{(n^R)^2 + (n^I)^2},$$

and

$$\begin{aligned} & \frac{(N^R)^2 + (N^I)^2}{(n^R)^2 + (n^I)^2} \\ &= 1 + \frac{4 \cdot |\omega^R x^R \omega^I x^I|}{(\omega^R x^R)^2 + (\omega^I x^I)^2 + (\omega^R x^I)^2 + (\omega^I x^R)^2}. \end{aligned} \tag{17}$$

In the denominator of the right-hand part of (17), the sum of the two terms $(\omega^R x^R)^2$ and $(\omega^I x^I)^2$ is larger than or equal to $2 \cdot |\omega^R x^R \omega^I x^I|$, since

$$\begin{aligned} & (\omega^R x^R)^2 + (\omega^I x^I)^2 - 2 \cdot |\omega^R x^R \omega^I x^I| \\ &= (|\omega^R x^R| - |\omega^I x^I|)^2 \geq 0. \end{aligned}$$

The same holds for the sum of the two terms $(\omega^R x^I)^2$ and $(\omega^I x^R)^2$. This immediately gives

$$1 + \frac{4 \cdot |\omega^R x^R \omega^I x^I|}{(\omega^R x^R)^2 + (\omega^I x^I)^2 + (\omega^R x^I)^2 + (\omega^I x^R)^2} \leq 2.$$

Combined with (16), this gives

$$\eta^2 \leq \alpha^2 + 4\alpha\beta + 2\beta^2.$$

From which we obtain

$$\begin{aligned} \eta^2 \leq & u^2 + 66u^3 + 793u^4 + 2958u^5 + 5937u^6 \\ & + 7696u^7 + 6982u^8 + 4596u^9 + 2216u^{10} \\ & + 772u^{11} + 186u^{12} + 28u^{13} + 2u^{14} \end{aligned} \tag{18}$$

Theorem 1 immediately follows from (18). \square

3 Two special cases

3.1 Obtaining the product as a double-word number

As explained above, one may wish to obtain the real and imaginary parts of the product $\omega \cdot x$ as double-word numbers, by replacing the floating-point addition $z^R = \text{RN}(v_h^R + \gamma_\ell^R)$ of line 9 of Algorithm 3 by a call to $2\text{Sum}(v_h^R, \gamma_\ell^R)$, and by replacing the floating-point addition $z^I = \text{RN}(v_h^I + \gamma_\ell^I)$ of line 18 by a call to $2\text{Sum}(v_h^I, \gamma_\ell^I)$. The resulting componentwise error bound is obtained by adding (11) and (12), so that the values α and β of Lemma 15 must be replaced by

$$\alpha' = u^2$$

and

$$\beta' = 10u^2 + 21u^3 + 16u^4 + 6u^5 + u^6.$$

With these new values, the square η'^2 of the normwise error now satisfies

$$\begin{aligned} \eta'^2 \leq & 241 u^4 + 924 u^5 + 1586 u^6 \\ & + 1608 u^7 + 1060 u^8 + 468 u^9 + 136 u^{10} \\ & + 24 u^{11} + 2 u^{12}, \end{aligned} \quad (19)$$

so that the normwise relative error becomes less than $\sqrt{241} \cdot u^2 + \mathcal{O}(u^3) \approx 15.53u^2 + \mathcal{O}(u^3)$.

That variant is of interest if one wishes to accurately evaluate the product

$$z_1 \times z_2 \times \cdots \times z_n$$

of n complex numbers. One can evaluate that product iteratively, keeping the real and imaginary parts of the partial product of all already considered terms as double-word numbers, and just using the unmodified Algorithm 3 (i.e., a simple FP addition and not a 2Sum at lines 9 and 18) for the last multiplication. Still assuming that overflow or underflow do not occur (which may become unlikely if n is very large and the z_i are arbitrary numbers), the total relative error is bounded by

$$(1 + \eta')^{n-2} \cdot (1 + \eta) - 1.$$

For instance, assuming binary64 arithmetic ($u = 2^{-53}$), one can multiply 1000 numbers and still have a normwise relative error bounded by $1.000000000001724u$.

3.2 If ω^I and ω^R are floating-point numbers

If ω^I and ω^R are floating-point numbers (i.e., if $\omega_\ell^I = \omega_\ell^R = 0$), Algorithm 3 becomes simpler, and we obtain Algorithm 4 below. Each separate part (computation of the real part, lines 1 to 6, or computation of the imaginary part, lines 7 to 12) in Algorithm 4 is similar to Cornea et al's algorithm for $ab + cd$ presented in [2] (with an addition replaced here by a call to 2Sum), and to Algorithm 5.3 in [15] (with an inversion in the order of summation of P_ℓ^R , Q_ℓ^R , and v_ℓ^R for the real part, and of P_ℓ^I , Q_ℓ^I , and v_ℓ^I for the imaginary part).

Of course the error bounds given by (18) and Theorem 1 still apply. However, one can redo the calculations taking into account the zero terms, and obtain new error bounds with smaller higher-order terms, more precisely, we find

$$\begin{aligned} \eta^2 \leq & u^2 + 38 u^3 + 299 u^4 + 782 u^5 \\ & + 1025 u^6 + 768 u^7 + 336 u^8 + 80 u^9 + 8 u^{10} \end{aligned} \quad (20)$$

which gives

$$\eta \leq u + 19u^2 + \mathcal{O}(u^3),$$

and one easily shows from (20) that as soon as $u \leq 1/16$ (i.e., the precision p is larger than or equal to 4, which always holds in practice), $\eta \leq u + 20u^2$.

ALGORITHM 4: Accurate complex multiplication $\omega \cdot x$, where the real and imaginary parts of ω and the real and imaginary parts of x are FP numbers, derived from Algorithm 3.

```

1:  $(P_h^R, P_\ell^R) \leftarrow \text{Fast2Mult}(\omega^I, x^I)$ 
2:  $(Q_h^R, Q_\ell^R) \leftarrow \text{Fast2Mult}(\omega^R, x^R)$ 
3:  $s_\ell^R \leftarrow \text{RN}(Q_\ell^R - P_\ell^R)$ 
4:  $(v_h^R, v_\ell^R) \leftarrow \text{2Sum}(Q_h^R, -P_h^R)$ 
5:  $\gamma_\ell^R \leftarrow \text{RN}(v_\ell^R + s_\ell^R)$ 
6: return  $z^R = \text{RN}(v_h^R + \gamma_\ell^R)$  (real part)
7:  $(P_h^I, P_\ell^I) \leftarrow \text{Fast2Mult}(\omega^I, x^R)$ 
8:  $(Q_h^I, Q_\ell^I) \leftarrow \text{Fast2Mult}(\omega^R, x^I)$ 
9:  $s_\ell^I \leftarrow \text{RN}(Q_\ell^I + P_\ell^I)$ 
10:  $(v_h^I, v_\ell^I) \leftarrow \text{2Sum}(Q_h^I, P_h^I)$ 
11:  $\gamma_\ell^I \leftarrow \text{RN}(v_\ell^I + s_\ell^I)$ 
12: return  $z^I = \text{RN}(v_h^I + \gamma_\ell^I)$  (imaginary part)

```

4 Implementation and experiments

Algorithm 3, implemented in binary64 arithmetic (i.e., $p = 53$ and $u = 2^{-53}$), was compared with other solutions, using a loop over N random inputs, itself inside another loop doing K iterations. The goal of the external loop is to get precise timings without having to choose a large value of N , with input data that would not fit in the cache: we do not want to include memory transfers in the timings. For each test, we chose $(N, K) = (1024, 65536)$, $(2048, 32768)$ and $(4096, 16384)$.

The other considered solutions were: use of the naive formula (1) in binary64 arithmetic; use of (1) in binary128 (a.k.a. “quad precision”) arithmetic; use of GNU MPFR [3] with precision ranging from 53 to 106 bits either with fused multiplications/subtractions `fmma/fmms` (thus implementing the formulas, correctly rounded) or with separate additions, subtractions and multiplications.

The tests were run on two machines with a hardware FMA:

- an `x86_64` machine with Intel Xeon E5-2609 v3 CPUs, under Linux (Debian/unstable), with GCC 8.2.0 and a Clang 8 preversion, using `-march=native`;
- a `ppc64le` machine with POWER9 CPUs, under Linux (CentOS 7), with GCC 8.2.1, using `-mcpu=power9`.

The following optimization options were used: `-O3` and `-O2`. With GCC, `-O3 -fno-tree-slp-vectorize` was also used in order to avoid a loss of performance with some vectorized codes. In all the cases, `-static` was used to avoid the overhead due to function calls to dynamic libraries.

The tests were run on several random data sets, giving a range of timings and a range of ratios. The smallest 10 % values and largest 10 % values have been excluded to take into account inaccuracies in the timings.

Table 1: Summary of the timings on an x86_64 machine (in seconds, for $NK = 2^{26}$ operations). “a3” stands for “Algorithm 3” (in binary64 arithmetic), “sw” corresponds to the naive formula (1) in binary64 arithmetic, “dw” corresponds to (1) in binary128 arithmetic, “cr” is GNU MPFR with fused multiplications/subtractions fmma/fmms, and “na” is GNU MPFR with separate additions, subtractions and multiplications.

$N \rightarrow$		minimums			maximums		
		1024	2048	4096	1024	2048	4096
gcc -03	a3	0.94	0.97	0.97	0.96	1.02	1.02
	sw	0.61	0.62	0.62	0.61	0.62	0.62
	dw	21.02	21.17	21.20	21.18	21.25	21.28
	cr	15.76	15.99	16.08	21.48	21.63	21.66
	na	12.46	12.88	12.99	23.16	23.23	23.22
gcc -03 -f...	a3	0.92	0.97	0.97	0.95	1.02	1.02
	sw	0.61	0.61	0.62	0.61	0.62	0.62
	dw	21.32	21.44	21.46	21.43	21.53	21.54
	cr	15.87	16.11	16.16	21.54	21.73	21.78
	na	12.59	13.01	13.12	22.72	22.85	22.80
gcc -02	a3	0.91	0.97	0.97	0.95	1.02	1.02
	sw	0.61	0.62	0.62	0.61	0.62	0.62
	dw	20.90	21.03	21.08	21.01	21.10	21.13
	cr	15.93	16.17	16.26	21.57	21.70	21.75
	na	12.31	12.74	12.85	23.11	23.20	23.18
clang -03	a3	0.86	1.09	1.10	0.96	1.15	1.15
	sw	0.39	0.61	0.63	0.47	0.65	0.66
	dw	21.65	21.77	21.81	21.74	21.87	21.88
	cr	16.00	16.24	16.32	21.46	21.69	21.71
	na	12.24	12.63	12.72	22.91	22.94	22.97
clang -02	a3	0.88	1.08	1.10	0.96	1.14	1.15
	sw	0.40	0.61	0.63	0.48	0.65	0.66
	dw	21.33	21.45	21.50	21.49	21.57	21.59
	cr	15.38	15.62	15.70	21.62	21.79	21.87
	na	12.15	12.54	12.65	23.15	23.21	23.21

Table 2: Summary of the timings on a POWER9 machine (in seconds, for $NK = 2^{26}$ operations). “a3”, “sw”, “dw”, “cr”, and “na” have the same meaning as in Table 1.

$N \rightarrow$		minimums			maximums		
		1024	2048	4096	1024	2048	4096
gcc -03	a3	0.97	0.97	0.97	0.99	0.99	1.00
	sw	0.47	0.47	0.51	0.48	0.48	0.52
	dw	2.22	2.22	2.22	2.24	2.24	2.24
	cr	19.44	19.56	19.62	23.94	24.07	24.06
	na	16.41	16.60	16.66	30.07	30.34	30.63
gcc -03 -f...	a3	0.97	0.97	0.97	0.98	0.99	1.00
	sw	0.47	0.47	0.51	0.48	0.48	0.52
	dw	2.22	2.22	2.22	2.24	2.24	2.24
	cr	19.45	19.59	19.61	24.11	24.08	24.07
	na	16.42	16.59	16.66	30.06	30.39	30.44
gcc -02	a3	0.98	0.98	0.98	0.99	1.01	1.01
	sw	0.47	0.47	0.51	0.47	0.47	0.51
	dw	2.22	2.22	2.22	2.24	2.24	2.24
	cr	19.50	19.66	19.68	24.14	24.11	24.05
	na	16.36	16.58	16.63	30.29	30.29	30.49

We checked that the various timings were globally consistent, in particular between the three chosen parameters for (N, K) , and rejected some anomalies manually: for Algorithm 3, $(N, K) = (2048, 32768)$ and $(4096, 16384)$ with Clang gave running times larger than those obtained with GCC, affecting the comparison with GNU MPFR. The timings are given in Table 1 (x86_64) and Table 2 (Power9). Note that reading the inputs is included in the timings (thus the ratios will be closer to 1 than one could expect), but these inputs are already in the right format for each implementation.

As a summary from the tables:

- **Implementation based on the naive formula (1) in binary64** (inlined code): It is about two times as fast as our implementation of Algorithm 3, but it is significantly less accurate.
- **Implementation based on the naive formula in binary128**, using the `__float128` C type (inlined code): On the x86_64 platform, it is from 19 to 25 times as slow as our implementation of Algorithm 3, the reason being that this format is implemented in software. The case of the POWER9 platform is particularly interesting as it has binary128 support in hardware. Here, the implementation is about 2.3 times as slow. This shows that even though one has hardware support for binary128, there is still interest in algorithms using a mix of binary64 and double-binary64.
- **Implementation based on GNU MPFR**, using precisions from 53 (corresponding to binary64) to 106 (roughly corresponding to double-binary64). Both codes based on `fmma/fmms` (thus implementing the formulas, correctly rounded) and based on separate additions, subtraction and multiplication operations were tested. This is from 11 to 26 times as slow as our implementation of Algorithm 3 on x86_64, and from 17 to 31 times as slow on POWER9.

Algorithm 3 has also been tested on random inputs to search for large norm-wise relative errors. For binary32, the input values (in ISO C99 / IEEE 754-2008 hexadecimal format) and the corresponding largest error found until now are

$$\begin{aligned}
 \omega^R &= 0x1.3ec518p-37 + 0x1.371ep-64 \\
 \omega^I &= 0x1.38c03ap-1 + -0x1.265e7ep-26 \\
 x^R &= 0x1.a31812p-1 \\
 x^I &= 0x1.59cdcap-17 \\
 \eta &\simeq 0.9999991420274016184 u,
 \end{aligned}$$

and for binary64, we have obtained

$$\begin{aligned}
 \omega^R &= 0x1.d1ef9ea4aa013p-1 + 0x1.ae88ba2a277ep-56 \\
 \omega^I &= 0x1.f5c28321df365p-81 + 0x1.c4c3e7b506d06p-135 \\
 x^R &= 0x1.194f298b4d152p-1 \\
 x^I &= 0x1.5c1fdca444f7cp-14 \\
 \eta &\simeq 0.99999900913907117123 u.
 \end{aligned}$$

This corroborates the bound given by Theorem 1.

Conclusion

We have given algorithms for complex multiplication in floating-point arithmetic, that either return the real and imaginary parts of the product as floating-point numbers with a normwise relative error bound close to the best one that one can guarantee, namely $u/(1+u)$, or as double-word numbers. Our implementation is only twice as slow as a significantly less accurate naive implementation. It is much faster than an implementation based on binary128 or multiple-precision software.

References

- [1] R. P. Brent, C. Percival, and P. Zimmermann. Error bounds on complex floating-point multiplication. *Mathematics of Computation*, 76:1469–1481, 2007.
- [2] M. Cornea, J. Harrison, and P. T. P. Tang. *Scientific Computing on Itanium[®]-based Systems*. Intel Press, Hillsboro, OR, 2002.
- [3] L. Fousse, G. Hanrot, V. Lefèvre, P. Pélissier, and P. Zimmermann. MPFR: A multiple-precision binary floating-point library with correct rounding. *ACM Transactions on Mathematical Software*, 33(2), 2007. 15 pages. Available at <http://www.mpfr.org/>.
- [4] Y. Hida, X. S. Li, and D. H. Bailey. C++/fortran-90 double-double and quad-double package, release 2.3.17, March 2012. Accessible electronically at <http://crd-legacy.lbl.gov/~dhbailey/mpdist/>.
- [5] N. J. Higham. *Accuracy and Stability of Numerical Algorithms*. SIAM, Philadelphia, PA, 2nd edition, 2002.
- [6] IEEE Computer Society. *IEEE Standard for Floating-Point Arithmetic*. IEEE Standard 754-2008, August 2008. Available at <https://doi.org/10.1109/IEEESTD.2008.4610935>.
- [7] Claude-Pierre Jeannerod, Peter Kornerup, Nicolas Louvet, and Jean-Michel Muller. Error bounds on complex floating-point multiplication with an FMA. *Mathematics of Computation*, 86(304):881–898, 2017.
- [8] Claude-Pierre Jeannerod, Nicolas Louvet, and Jean-Michel Muller. Further analysis of Kahan’s algorithm for the accurate computation of 2×2 determinants. *Mathematics of Computation*, 82(284):2245–2264, October 2013.
- [9] Mioara Joldeş, Jean-Michel Muller, and Valentina Popescu. Tight and rigorous error bounds for basic building blocks of double-word arithmetic. *ACM Transactions on Mathematical Software*, 44(2), October 2017.

- [10] W. Kahan. Lecture notes on the status of IEEE-754. Available at <http://www.cs.berkeley.edu/~wkahan/ieee754status/IEEE754.PDF>, 1997.
- [11] D. E. Knuth. *The Art of Computer Programming*, volume 2. Addison-Wesley, Reading, MA, 3rd edition, 1998.
- [12] O. Møller. Quasi double-precision in floating-point addition. *BIT*, 5:37–50, 1965.
- [13] Jean-Michel Muller, Nicolas Brunie, Florent de Dinechin, Claude-Pierre Jeannerod, Mioara Joldes, Vincent Lefèvre, Guillaume Melquiond, Nathalie Revol, and Serge Torres. *Handbook of Floating-Point Arithmetic*. Birkhäuser Boston, 2018. ACM G.1.0; G.1.2; G.4; B.2.0; B.2.4; F.2.1., ISBN 978-3-319-76525-9.
- [14] Y. Nievergelt. Scalar fused multiply-add instructions produce floating-point matrix arithmetic provably accurate to the penultimate digit. *ACM Transactions on Mathematical Software*, 29(1):27–48, 2003.
- [15] T. Ogita, S. M. Rump, and S. Oishi. Accurate sum and dot product. *SIAM Journal on Scientific Computing*, 26(6):1955–1988, 2005.