

A Constructive Criticism of the C/C++ Proposal for Complex Arithmetic

FRANKY BACKELJAUW and ANNIE CUYT

Department of Mathematics and Computer Science, University of Antwerp, Campus Middelheim, Middelheim 1, B-2020 Antwerpen, Belgium, e-mail: {franky.backeljauw,annie.cuyt}@ua.ac.be

(Received: 28 June 2004; accepted: 9 January 2005)

Abstract. The IEEE 754 and 854 standards regulate the behaviour of real floating-point arithmetic, as implemented in most current hard- and software systems. Although a myriad of libraries for complex floating-point arithmetic is available and in use, there is no general consensus on their implementation. The International C Standard describes in its Annex G guidelines for the implementation of complex arithmetic, in order to achieve a similar behaviour of complex floating-point arithmetic across C-language compliant implementations. In Section 2 we summarize its recommendations and outline the problems inherent to this approach. In Section 3 we describe how the lack of reliability, when computing certain complex-valued expressions, can be overcome. Throughout the discussion the rounding mode is assumed to be round-to-nearest, as in Annex G.

1. Introduction

The IEEE standards for real floating-point arithmetic are a formal model for floating-point arithmetic in which as many properties as possible have successfully been transferred from real arithmetic. These standards realize a closed system for the basic operations, remainder and square root function, including a limited number of properties (such as commutativity) and identities (such as $1 / (1 / x) = x$). We want to analyze whether the Annex G proposal for complex arithmetic [1] is equally successful. The interested reader is also referred to [3]–[5] for introductory material on complex floating-point arithmetic.

We assume to have at our disposal an IEEE compliant implementation of floating-point arithmetic, in base β , of precision t and with the exponent of the normalized numbers ranging between L and U . The signed zeroes and denormal numbers carry exponent $L - 1$. Overflow resulting in a signed infinity or an invalid expression resulting in Not-a-Number (NaN) give rise to results carrying exponent $U + 1$. Based on this floating-point implementation we realize implementations for purely imaginary, purely real and true complex arithmetic, following the Annex G guidelines. Nevertheless it is not difficult to construct examples

of simple but problematic complex-valued expressions. Take for instance, with $2 \otimes x \geq \beta^{U+1}$,

$$(\text{rem}(x, +0) + 2i) \otimes (x + 1i) = ?$$

Clearly this expression should evaluate to a complex NaN, and indeed, using plain IEEE compliant floating-point arithmetic,

$$(\text{NaN} \ominus 2) + (\text{NaN} \oplus +\infty)i = \text{NaN} + \text{NaN}i = \text{complex NaN}.$$

But when using the recommendation for the implementation of the product of complex numbers as found in Annex G,

```
#include <math.h>
#include <complex.h>

/* Multiply z * w ... */
double complex _Cmultd(double complex z, double complex w)

{ #pragma STDC FP_CONTRACT OFF
  double a, b, c, d, ac, bd, ad, bc, x, y;
  a = creal(z); b = cimag(z);
  c = creal(w); d = cimag(w);
  ac = a * c; bd = b * d;
  ad = a * d; bc = b * c;
  x = ac - bd; y = ad + bc;
  if (isnan(x) && isnan(y)) {
    /* Recover infinities that computed as NaN+iNaN ... */
    int recalc = 0;
    if (isinf(a) || isinf(b)) { // z is infinite
      /* "Box" the infinity and change NaNs in the other factor to 0 */
      a = copysign(isinf(a) ? 1.0 : 0.0, a);
      b = copysign(isinf(b) ? 1.0 : 0.0, b);
      if (isnan(c)) c = copysign(0.0, c);
      if (isnan(d)) d = copysign(0.0, d);
      recalc = 1;
    }
    if (isinf(c) || isinf(d)) { // w is infinite
      /* "Box" the infinity and change NaNs in the other factor to 0 */
      c = copysign(isinf(c) ? 1.0 : 0.0, c);
      d = copysign(isinf(d) ? 1.0 : 0.0, d);
      if (isnan(a)) a = copysign(0.0, a);
      if (isnan(b)) b = copysign(0.0, b);
      recalc = 1;
    }
  }
}
```

```

if (!recalc && (isinf(ac) || isinf(bd) || isinf(ad) || isinf(bc)))
{
    /* Recover infinities from overflow by changing NaNs to 0 ... */
    if (isnan(a)) a = copysign(0.0, a);
    if (isnan(b)) b = copysign(0.0, b);
    if (isnan(c)) c = copysign(0.0, c);
    if (isnan(d)) d = copysign(0.0, d);
    recalc = 1;
}
if (recalc) {
    x = INFINITY * ( a * c - b * d );
    y = INFINITY * ( a * d + b * c );
}
}
return x + I * y; }

```

then the result needs to be recomputed because one of the four intermediate results, called ac , bd , ad , bc , evaluates to either $+\infty$ or $-\infty$. The recomputation delivers the final result

$$(\text{NaN} + 2i) \otimes (x + 1i) \rightarrow |\infty| \times [(+0 + 2i) \otimes (x + 1i)] = -\infty + \infty i$$

which is incorrect. Let us now analyze the origin of the problem and formulate a possible cure for it.

2. Complex Arithmetic as Detailed in Annex G

As can be expected, the problematic situations arise from the appearance of signed zeroes, signed infinities and NaNs in expressions, especially when at least 4 representations of zero or underflow, namely $(+/- 0) + (+/- 0)i$, and 4 representations of the complex Riemann infinity, namely $(+/- \infty) + (+/- \infty)i$, exist. The fact that combinations of these special values in the real and imaginary part of a complex number are possible, gives rise to particular problems. For instance, the result of the multiplication

$$(+\infty + \infty i) \otimes (-\infty - \infty i) = \text{NaN} - \infty i \tag{2.1}$$

is interpreted in Annex G as infinity. That this is correct, can be shown as follows. Let us introduce the notation $r \exp^{[\theta_1, \theta_2]i}$ to denote a complex number with modulus r and argument θ satisfying $\theta_1 \leq \theta \leq \theta_2$. Using this notation, it is easy to see that (2.1) signifies

$$\begin{aligned} (+\infty + \infty i) \otimes (-\infty - \infty i) &= (|\infty| \exp^{[0, \pi/2]i}) \times (|\infty| \exp^{[-\pi, -\pi/2]i}) \\ &= |\infty| \exp^{[-\pi, 0]i} \end{aligned} \tag{2.2}$$

which is a complex number of infinitely large modulus lying in the lower half-plane.

As a consequence of this type of situations, Annex G proposes to also consider the 4 values $(+/-\infty) + \text{NaN}i$ and $\text{NaN} + (+/-\infty)i$ as representations of the Riemann infinity, in addition to the 4 representations already listed above. These representations actually stand for complex numbers with infinitely large modulus lying in the 4 possible halfplanes:

$$\begin{aligned} (+\infty) + \text{NaN}i &= |\infty| \exp^{[-\pi/2, \pi/2]i}, \\ (-\infty) + \text{NaN}i &= |\infty| \exp^{[\pi/2, 3\pi/2]i}, \\ \text{NaN} + (+\infty)i &= |\infty| \exp^{[0, \pi]i}, \\ \text{NaN} + (-\infty)i &= |\infty| \exp^{[-\pi, 0]i}. \end{aligned} \tag{2.3}$$

A truly invalid result is represented by $\text{NaN} + \text{NaN}i$.

Of course this situation gives rise to some new problems, such as in the expression

$$(+0) \oslash (+0) \oplus (y + \infty i) = \text{NaN} + \infty i$$

which is incorrectly being interpreted as a representation of complex infinity instead of invalid, because of $(+0) \oslash (+0)$.

The fact that now a different number of representations for zero and infinity exists, namely 4 versus 8, leads to violations of the floating-point identity $1/(1/x) = x$. The latter holds for all nonzero x and also for x being either a real signed zero or a real signed infinity. However, it is easy to see that the identity does not hold anymore for the representations (2.3). With $x = \text{NaN} - \infty i$, which represents $|\infty| \exp^{[-\pi, 0]i}$:

$$\begin{aligned} 1 \oslash (\text{NaN} - \infty i) &= (+0) + (+0)i, \\ 1 \oslash (+0 + 0i) &= +\infty + \text{NaN}i \neq \text{NaN} - \infty i. \end{aligned}$$

Moreover, the 4 additional representations for large complex numbers in one of the 4 possible halfplanes do not cover all the cases. Take for instance the expression

$$(\text{NaN} + \infty i) \otimes (\text{NaN} - \infty i)$$

which should be interpreted as

$$(|\infty| \exp^{[0, \pi]i}) \times (|\infty| \exp^{[-\pi, 0]i}) = |\infty| \exp^{[-\pi, \pi]i}. \tag{2.4}$$

Its result is a complex number that can lie in any of the four quadrants. For this result Annex G does not provide a representation. After recomputation, as recommended in the cited guideline for the multiplication, the product returned by Annex G is erroneously $-\infty + \infty i$.

3. A Proposal for Reliable Complex Arithmetic

Let us take another look at example (2.1). The result of the correct interpretation (2.2) could also be stored as (we detail our implementation further down)

$$|\infty| \exp^{[-\pi, 0]i} = (\mp\infty) + (-\infty)i$$

instead of using $\text{NaN} - \infty i$ as representation for the lower halfplane. By introducing three possible “signs,” namely positive, negative and insecure, denoted respectively by $+$, $-$, and \mp , one obtains 9 representations for complex infinity,

$$(*\infty) + (*\infty)i \quad *, * \in \{+, -, \mp\} \quad (3.1)$$

and 9 representations for complex zero,

$$(*0) + (*0)i \quad *, * \in \{+, -, \mp\}. \quad (3.2)$$

A simple verification shows that in this way also the identity $1 / (1 / x) = x$ can be restored.

Moreover, complex numbers $x + yi$ of which real and/or imaginary part equal NaN, can now be reserved for truly invalid results. In this way, the ambiguity introduced in Annex G about NaN parts in complex numbers is eliminated.

When we return to expression (2.4), it is clear that with the introduction of (3.1), we can now store

$$(|\infty| \exp^{[0, \pi]i}) \times (|\infty| \exp^{[-\pi, 0]i}) = |\infty| \exp^{[-\pi, \pi]i} = (\mp\infty) + (\mp\infty)i.$$

Tables 1 and 2 contain some more examples of expressions that cannot be dealt with correctly by Annex G and necessitate the new approach. For operands representing a halfplane, we have used the Annex G notation.

For the correct implementation of the basic operations on complex operands, making use of the indispensable additional special values, the sign ($+$, $-$, or \mp) accompanying intermediate NaN results in the real and/or imaginary part, can be useful to determine the quadrant(s) containing the result. For instance:

$$\begin{aligned} (2.0 + 0.0i) \otimes (3.0 + \infty i) &= (6.0 \ominus +\text{NaN}) + (+\infty \oplus 0.0)i \\ &= \mp\text{NaN} + \infty i \rightarrow \mp\infty + \infty i, \\ (0.0 + 3.0i) \otimes (2.0 + \infty i) &= (0.0 \ominus +\infty) + (6.0 \oplus +\text{NaN})i \\ &= -\infty + \text{NaN}i \rightarrow -\infty + \infty i. \end{aligned}$$

Moreover, while Annex G says that for the implementation of $\exp(-\infty + \infty i)$, the signs returned for the real and imaginary zero parts of the result are at the implementor’s choice, we can correctly return

$$\exp(-\infty + \infty i) = (\mp 0) + (\mp 0)i.$$

In the same way

$$\log(-0 + 0i) = -\infty + (\mp\infty)i$$

while Annex G suggests $\log(-0 + 0i) = -\infty + \pi i$.

Table 1. Some exceptional cases for the basic operations. ζ is used to indicate a mathematically incorrect result, because of its interpretation by Annex G.

operation	C++ Annex G	New Style
$(+\infty + \infty i) \otimes (-\infty - \infty i)$	NaN $-\infty i$	$(\mp\infty) - \infty i$
$(-\infty + \infty i) \otimes (3.0 + 2.0 i)$	$-\infty + \text{NaN } i$	$-\infty + (\mp\infty) i$
$(-\infty + \infty i) \oslash (3.0 + 2.0 i)$	NaN $+\infty i$	$(\mp\infty) + \infty i$
$(+\infty + \text{NaN } i) \otimes (-\infty + \infty i)$	$-\infty + \infty i \zeta$	$(\mp\infty) + (\mp\infty) i$
$(\text{NaN} - \infty i) \oslash (3.0 + 2.0 i)$	$-\infty - \infty i \zeta$	$(\mp\infty) + (\mp\infty) i$
$(+\infty + \infty i) \oplus (-\infty + \infty i)$	NaN $+\infty i$	$(\mp\infty) + \infty i$
$(\text{NaN} + \infty i) \oplus (-\infty + \infty i)$	NaN $+\infty i \zeta$	NaN $+\text{NaN } i$
$(+\infty - \infty i) \otimes (\text{NaN} + 2.0 i)$	$+\infty + \infty i \zeta$	NaN $+\text{NaN } i$
$(\text{NaN} + 2.0 i) \otimes (x + 1.0 i)^1$	$-\infty + \infty i \zeta$	NaN $+\text{NaN } i$
NaN $\oplus (5.0 + \infty i)$	NaN $+\infty i \zeta$	NaN $+\infty i$
$1.0 \oslash (+\infty - \infty i)$	$+0.0 + 0.0 i$	$+0.0 + 0.0 i$
$1.0 \oslash (\text{NaN} - \infty i)$	$+0.0 + 0.0 i \zeta$	$(\mp 0.0) + 0.0 i$
$1.0 \oslash (+0.0 + 0.0 i)$	$+\infty + \text{NaN } i \zeta$	$+\infty - \infty i$
$(+0.0 + 0.0 i) \oplus (-0.0 + 0.0 i)$	$+0.0 + 0.0 i \zeta$	$(\mp 0.0) + 0.0 i$
$(+0.0 + 0.0 i) \oplus (-0.0 - 0.0 i)$	$+0.0 + 0.0 i \zeta$	$(\mp 0.0) + (\mp 0.0) i$
$((\mp 0.0) + 0.0 i) \oplus (-0.0 + 0.0 i)$	–	$(\mp 0.0) + 0.0 i$
$(2.0 + \infty i) \otimes (0.0 + 3.0 i)$	$-\infty + \text{NaN } i \zeta$	$-\infty + \infty i$
$(7.0 + 0.0 i) \otimes (4.0 - 0.0 i)$	$+28.0 + 0.0 i \zeta$	$28.0 + (\mp 0.0) i$
$(-6.0 + 0.0 i) \oplus (-9.0 - 0.0 i)$	$-15.0 + 0.0 i \zeta$	$-15.0 + (\mp 0.0) i$

¹ $2 \otimes x \geq \beta^{U+1}$

Table 2. Some exceptional cases for the square root.

operand ²	C++ Annex G	operand ³	New Style
$\pm 0.0 + 0.0 i$	$+0.0 + 0.0 i$	$(\neq 0.0) + 0.0 i$	$+0.0 + 0.0 i$
$a + \infty i$	$+\infty + \infty i$	$a + \infty i$	$+\infty + \infty i$
NaN $+\infty i$	$+\infty + \infty i$	$(\mp\infty) + \infty i$	$+\infty + \infty i$
$-\infty + b i$	$+0.0 + \infty i$	$-\infty + b i$	$+0.0 + \infty i$
$+\infty + b i$	$+\infty + 0.0 i$	$+\infty + b i$	$+\infty + 0.0 i$
$-\infty + \text{NaN } i$	NaN $\pm \infty i \zeta$	$+\infty + (\mp\infty) i$	$(\mp\infty) + (\mp\infty) i$
$+\infty + \text{NaN } i$	$+\infty + \text{NaN } i$	$+\infty + (\mp\infty) i$	$+\infty + (\mp\infty) i$
–	–	$(\mp\infty) + (\mp\infty) i$	$(\mp\infty) + (\mp\infty) i$

² using the Annex G representations³ using the new representations

4. Concluding Remarks

The signed zeroes have been the subject of a lot of debate in the scientific community. In complex floating-point arithmetic the IEEE signed infinities raise a similar problem. While the introduction of a zero with uncertain sign, namely ∓ 0 , is not new (see [2]), its importance grows when considering complex arithmetic.

In IEEE floating-point arithmetic, it is the correct result for

$$(+0) + (-0) = \mp 0.$$

In complex arithmetic, it can be used for underflowing quantities of the type

$$\lim_{n \rightarrow \infty} (-1 / (2 + i))^n = (\mp 0) + (\mp 0)i.$$

In addition, one, two or all four halfplanes of the complex plane can be represented by supporting $\mp \infty$. Even the 3 quadrant result

$$\frac{+\infty + (\mp \infty)i}{+1 - 2i} = |\infty| \exp^{[-\pi/2, \pi]i} \subseteq |\infty| \exp^{[-\pi, \pi]i}$$

can be contained in $(\mp \infty) + (\mp \infty)i$.

While the Annex G guidelines were probably the best that could be achieved when solely making use of the IEEE real floating-point representations, it is clear that a reliable implementation of complex floating-point arithmetic needs additional special values.

References

1. ANSI/ISO/IEC 9899:1999: *Annex G: IEC 60559-Compatible Complex Arithmetic (Informative)*, ANSI, 1999, pp. 465–478.
2. Juric, Z.: *TIGCC—Routines for Floating Point Arithmetic*, Texas Instruments, <http://tigcc.ticalc.org/doc/timath.html>.
3. Kahan, W.: Branch Cuts for Complex Elementary Functions, or Much Ado about Nothing's Sign Bit, in: Iserles, A. and Powell, M. J. D. (eds), *The State of the Art in Numerical Analysis*, Oxford University Press, Oxford, 1987, pp. 165–211.
4. *Language Independent Arithmetic—Part 3: Complex Integer and Floating Point Arithmetic and Complex Elementary Numerical Functions*, Technical Report, WD 10967–3, ISO/IEC, 2001, <http://anubis.dkuug.dk/JTC1/SC22/WG11/docs/n476.pdf>.
5. Ziv, A.: Sharp Ulp Rounding Error Bound for the Hypotenuse Function, *Math. Comp.* **68** (227) (1999), pp. 1143–1148.