

RA

**stichting  
mathematisch  
centrum**



REKENAFDELING

RA

MR 118/70 NOVEMBER

T.J. DEKKER  
A FLOATING-POINT TECHNIQUE FOR EXTENDING  
THE AVAILABLE PRECISION

2nd (unrevised) edition

**2e boerhaavestraat 49 amsterdam**

# A FLOATING-POINT TECHNIQUE FOR EXTENDING THE AVAILABLE PRECISION

by

T.J. Dekker<sup>\*)</sup>

## Abstract

A technique is described for expressing multilength floating-point arithmetic in terms of singlelength floating-point arithmetic, i.e. the arithmetic for an available (say: single or double precision) floating-point number system. The basic algorithms are exact addition and multiplication of two singlelength floating-point numbers, delivering the result as a doublelength floating-point number. A straightforward application of the technique yields a set of algorithms for doublelength arithmetic which are given as ALGOL 60 procedures.

<sup>\*)</sup> Part of this research was done while the author was visiting Bell Telephone Laboratories, Murray Hill, New Jersey.

## Contents

1	Introduction	page	1
2	Floating-point number systems		4
3	Floating-point operations		6
4	Exact addition		8
5	Exact multiplication		16
6	Splitting into halflength numbers		21
7	Doublelength arithmetic		24
8	Appendix. ALGOL 60 procedures		32
9	Acknowledgements		36
10	References		36

## 1. Introduction

Multilength floating-point arithmetic is often described in terms of integer arithmetic. The mantissa of a multilength floating-point number is then represented by means of some integers and another integer is used for the exponent.

In contrast with this approach, we present some basic algorithms which enable us to describe multilength floating-point arithmetic in terms of singlelength floating-point arithmetic. A multilength floating-point number is then represented by means of some singlelength floating-point numbers. In particular, we represent a doublelength floating-point number as the sum of two singlelength floating-point numbers, one of them being (almost) negligible in singlelength precision. The basic algorithms are exact addition and multiplication of two singlelength floating-point numbers, delivering the result as a doublelength floating-point number.

After two introductory sections on floating-point systems and arithmetic, we deal with exact addition in section 4.

Let  $x$  and  $y$  be singlelength floating-point numbers and let

$$z = \text{fl}(x+y);$$

i.e.  $z$  is the result of a singlelength floating-point addition of  $x$  and  $y$ .

Let  $zz$  be the correction term exactly satisfying

$$z + zz = x + y.$$

It will be shown that, under various conditions,  $zz$  can be obtained by the formula

$$zz = \text{fl}((x-z)+y).$$

For example, this is true if the singlelength floating-point system is binary, addition and subtraction are optimal and  $|x| \geq |y|$ . Similar formulas have been used by Møller (1965), Kahan (1965), Babuška (1968) and Knuth (1969) for reducing the error in calculating the sum of several

terms.

Møller (1965) gives a more elaborate formula which turns out to be a double application of ours. Knuth (1969) shows that this formula yields the correction term  $zz$  if floating-point addition and subtraction are optimal.

In sections 5 and 6, we deal with exact multiplication of two single-length floating-point numbers.

First, each factor is splitted into two "halflength" numbers, and then the exact product is formed in a rather obvious way. Splitting the factors as well as forming the exact product are expressed in terms of singlelength floating-point operations which, under certain reasonable conditions, are shown to yield the correct results.

Veltkamp (1968) developed a similar technique with a slightly modified algorithm for exact multiplication.

In section 7, we give a straightforward application of the foregoing results to obtain algorithms for doublelength arithmetic. In the appendix, the algorithms are given as ALGOL 60 procedures.

An important application (in fact, the main incentive for this research), is the calculation of the doublelength scalar product of two vectors of singlelength floating-point numbers.

Application of our results to multilength arithmetic we hope to present in a future publication.

A drawback of our technique is that it works only if the singlelength floating-point system and arithmetic satisfy certain conditions (which are fulfilled, for example, on the General Electric 635 and Philips Electronica X8 computers).

Our technique is then, of course, most effective if we take as single-length system the highest precision floating-point system available.

Therefore, an essential feature of our algorithms is that, with respect to this singlelength system, no doublelength accumulator is used.

Our technique is especially attractive if the mantissa length of the highest precision floating-point system available is larger than the capacity of an available fixed-point (i.e. integer) system, because then less work is needed to obtain a certain required precision.

Moreover, since the floating-point operations used take care of the shifting and normalizing needed at various stages, a reasonable efficiency can be reached by writing our algorithms entirely in a high-level language such as ALGOL or FORTRAN.

## 2. Floating-point number systems

The starting-point for our technique is a certain system,  $R$ , of floating-point numbers, and the floating-point arithmetical operations defined on  $R$ .

Let the base,  $\beta$ , of  $R$  be an integer larger than 1.

We assume that the mantissa as well as the exponent of the elements of  $R$  are integers within a certain range.

More precisely, we assume:

- 1) a floating-point number  $x$  has the form

$$(2.1) \quad x = mx \beta^{ex},$$

where  $mx$ , or  $m(x)$ , denotes the mantissa, and  $ex$ , or  $e(x)$ , the exponent of  $x$ ;

- 2) the system  $R$  of floating-point numbers is

$$(2.2) \quad R = \{x \mid x = mx \beta^{ex}, |mx| < M, -D < ex < E\},$$

where  $M$  is a positive integer (usually  $M = \beta^t$ , where  $t$  is the number of mantissa digits) and  $D, E$  are positive integers or infinite (the latter if we want to disregard underflow and/or overflow).

Floating-point number systems of the form (2.2) were proposed by Grau (1962) and are actually used in some computers, e.g. Burrough B 5500 (octal) and Philips Electrologica X8 (binary). Apart from overflow and underflow, the floating-point systems have this form in most, if not all, computers, since the mantissa can always be interpreted as an integer by subtracting a suitable constant from the exponent.

### Remarks

- 1) Representation (2.1) is not always unique. In most machines, a certain standardization is defined in order to make the representation unique. The most common standardization is the normalization in which the magnitude of a non-zero mantissa has a lower bound  $M/\beta$ . Some elements of  $R$  having an

exponent (nearly) equal to  $-D$  do not have a normalized representation due to underflow in the exponent range. A system of the form (2.2) without requiring normalization has the advantage, that addition and subtraction never lead to underflow.

2) We shall assume that  $M$  is not congruent to 1 (mod  $\beta$ ), because this slightly simplifies some theorems and proofs. This assumption is no true restriction if we disregard overflow. Indeed, if  $M = 1 + \beta K$  for some integer  $K$ , then  $M$  may be replaced by  $\beta K$  without changing the system  $R$ , except that (for finite  $E$ ) the elements  $\pm M \beta^E$  are removed from the system. This remark also justifies the restriction to a symmetric mantissa range, since, for  $M = \beta K$ , the not uncommon asymmetric mantissa range  $-(1+M) < mx < M$  defines the same system, apart from overflow, as the symmetric mantissa range given in (2.2).



### 3. Floating-point operations

We use the following notation from Wilkinson (1963, p. 4). If  $A$  is an expression involving floating-point numbers and the arithmetical operations  $+$ ,  $-$ ,  $\times$ ,  $/$ , then  $fl(A)$  is the corresponding expression obtained by replacing the arithmetical operations by the corresponding floating-point operations.

Let  $*$  represent  $+$ ,  $-$ ,  $\times$  or  $/$ , and let  $x$  and  $y$  be elements of  $R$ .

(3.1) Definition. The floating-point operation corresponding to  $*$  is faithful if, for all  $x$  and  $y$ , the result  $fl(x * y)$  equals either the largest element of  $R$  smaller than or equal to  $x * y$ , or the smallest element of  $R$  larger than or equal to  $x * y$ .

Thus, when  $x * y$  lies between two successive elements of  $R$ , either one will do, but when  $x * y \in R$ , the result must be exact (and when  $x * y$  is outside the range of  $R$ , the result must be the nearest, i.e. the largest or the smallest, element of  $R$ ).

(3.2) Definition. The floating-point operation corresponding to  $*$  is optimal (or properly rounding) if, for all  $x$  and  $y$ , the result  $fl(x * y)$  is an element of  $R$  nearest to  $x * y$ .

Note that this definition uniquely determines the result, except when  $x * y$  lies halfway between two successive elements of  $R$ , in which case an optimal operation may round up or down.

Related to this definition is the following notation which will be used below.

(3.3) Notation. For any real  $\xi$ ,  $\text{round}(\xi)$  denotes an integer closest to  $\xi$ .

The result  $z = fl(x * y)$  of a properly rounding operation can be represented in the form (2.1), where mantissa  $mz$  and exponent  $ez$  satisfy

$$(3.4) \quad mz = \text{round}(x * y \beta^{-ez}),$$

provided that no overflow or underflow occurs.

For addition and subtraction, we also need the following definitions.

(3.5) Definition. Floating-point addition is properly truncating if it is commutative (i.e.  $fl(x+y) = fl(y+x)$ ) and, for all  $x$  and  $y$  satisfying  $|x| \geq |y|$ , the result  $fl(x+y)$  equals the largest element of  $R$  smaller than or equal to  $x + y$  if  $y \geq 0$  or the smallest element of  $R$  larger than or equal to  $x + y$  if  $y < 0$ .

Thus, when  $x + y$  does not belong to  $R$ , it is truncated in the direction of  $-y$ . Note that this definition uniquely determines the result.

(3.6) Definition. Floating-point subtraction is properly truncating if, for all  $x$  and  $y$ , we have  $fl(x-y) = fl(x+y')$ , where  $y' = -y$  and the floating-point addition is properly truncating.

(3.7) Definition. Floating-point addition and subtraction are super-faithful if, for each  $x$  and  $y$ , the result  $fl(x \pm y)$  is obtained by properly rounding or by properly truncating.

#### Remark

Arithmetical operations having these properties are not difficult to realize.

In fact, there are machines in which the floating-point addition and subtraction are optimal (e.g. General Electric 635, Philips Electrologica X8). In order to obtain faithful addition and subtraction, the result must only be normalized before it is truncated or rounded, cf. Kahan (1965). To obtain a super-faithful, or even optimal, addition and subtraction, it is by no means necessary to have a (nearly) doublelength accumulator. Optimal addition and subtraction can be perfectly well formulated using an accumulator having no more than two guarding digits; see Knuth (1969, p.183 Algorithms A and N and p. 194 exercise 5).

#### 4. Exact addition

Let  $x$  and  $y$  be given elements of  $R$  and let

$$(4.1) \quad z = \text{fl}(x+y).$$

We want to find the correction term  $zz$  satisfying the (exact) relation

$$(4.2) \quad z + zz = x + y.$$

We shall derive some formulas for calculating  $zz$  which use only single-length floating-point addition and subtraction.

First we consider the formula

$$(4.3) \quad w = \text{fl}(z-x), \quad zz = \text{fl}(y-w)$$

and prove some theorems stating sufficient conditions for the validity of this formula. For practical computation, formulas (4.1) and (4.3) can be written as the following sequence of ALGOL 60 statements:

$$(4.4) \quad "z := x + y; zz := y - (z-x)",$$

in which  $w$  remains anonymous.

Let  $x, y \in R$  be representable such that their exponents satisfy

$$(4.5) \quad e_x \geq e_y.$$

In particular, this holds if  $x, y \in R$  satisfy

$$(4.6) \quad |x| \geq |y|.$$

(4.7) Theorem. If  $R$  has the form (2.2), where  $\beta = 2$  or  $3$  and  $M$  is a multiple of  $\beta$ , and if, moreover, floating-point addition is optimal and subtraction faithful, then, for all  $x$  and  $y$  satisfying (4.5) and for  $z$

obtained according to (4.1), formula (4.3) yields the correction term  $zz$  defined by (4.2).

Proof. According to (4.5), we may assume

$$x = mx \beta^{ex}, y = my \beta^{ey}, ex \geq ey.$$

We prove the theorem by showing

- 1)  $w = z - x,$
- 2)  $zz = y - w.$

Proof of (1). Since floating-point subtraction is faithful, we need to show only that  $z - x$  is an element of  $R$ . Obviously,  $z$  can be represented such that  $ez \leq ex + 1$ . We first assume that no overflow occurs in forming  $z$  and distinguish between two cases.

1a)  $ez \leq ex.$

Then we have (for definition of "round" see 3.3)

$$mz = mx \beta^{ex-ez} + \text{round}(my \beta^{ey-ez}),$$

so that

$$z - x = (mz - mx \beta^{ex-ez}) \beta^{ez} = \text{round}(my \beta^{ey-ez}) \beta^{ez}.$$

Hence,  $z - x \in R$ .

1b)  $ez = ex + 1$

Let  $d = ex - ey$ . Then

$$mz = \text{round}(mx/\beta + my/\beta^{d+1}),$$

$$z - x = (\beta mz - mx) \beta^{ex} = \mu \beta^{ex} \text{ (say),}$$

where  $\mu = \beta mz - mx$  satisfies

$$|\mu| \leq |\beta mz - mx - my/\beta^d| + |my/\beta^d|$$

$$\leq \beta/2 + |my/\beta^d| < \beta/2 + M.$$

Hence, since  $\beta \leq 3$  and  $\mu$  is integral, we have

$$|\mu| \leq M,$$

so that  $z - x$  is representable as an element of  $R$  either with mantissa  $\mu$  and exponent  $ex$  if  $|\mu| < M$ , or, since  $M$  is a multiple of  $\beta$ , with mantissa  $\mu/\beta$  and exponent  $ex + 1$  if  $|\mu| = M$ .

We now assume that overflow occurs in forming  $z$ .

Then  $ez = ex$  and  $|mz| = M - 1$ , so that

$$z - x = (mz - mx) \beta^{ex}.$$

Hence,  $z - x \in R$ .

This completes the proof of (1).

Proof of (2). Since floating-point subtraction is faithful, we need to show only that  $y - w$  is an element of  $R$ .

From  $ex \geq ey$ , it follows that  $y - w$  is equal to an integer times  $\beta^{ey}$ .

Moreover,  $|y - w| \leq |y|$ , because otherwise  $x$  would be closer to  $x + y$  than  $z$  contradicting the assumption that floating-point addition is optimal.

Hence,  $y - w \in R$ .

This completes the proof of (2) and the theorem.

Theorem (4.7) does not hold for  $\beta \leq 3$  if  $M$  is not a multiple of  $\beta$  (which, according to section 2 remark 2 can only occur for  $\beta = 3$ ), nor for any  $\beta > 3$ .

For example, if  $\beta = 3$ ,  $M = \frac{1}{2}(3^t + 1)$  (this being the most natural range in the "balanced" ternary system using the digits  $+1, 0, -1$ , see Knuth (1969, p. 173)) and  $x = y = M - 1$ , then  $z = 3^t$  and  $z - x = \frac{1}{2}(3^t + 1)$ , which is not an element of  $R$ ; if  $\beta = 10$ ,  $M = 100$  and  $x = y = 98$ , then  $z = 200$  and  $z - x = 102$ , which is not an element of  $R$ .

The theorem can be extended to any  $\beta$  and  $M$ , provided that an enlarged mantissa range (an extra digit, say) is available for  $w$ . Such an extended theorem has practical applications, because  $w$  may remain anonymous (see 4.4) and some systems have an enlarged mantissa range for the anonymous floating-point values.

Therefore, we define an enlarged system  $R^*$  for  $w$ , analogous to (2.2), as follows.

$$(4.8) \quad R^* = \{x \mid x = mx \beta^{ex}, |mx| < M^*, -D < ex < E\},$$

where  $M^* \geq M + \beta/2$ . (This is certainly true if  $R^*$  has at least one more digit for the mantissas of its elements than  $R$ .)

Then the following theorem holds without any restriction on  $\beta$  and  $M$ .

(4.9) Theorem. Let  $R$  have the form (2.2) and  $R^*$  the form (4.8). Let  $x, y \in R$  satisfy (4.5) and let  $z \in R$  be obtained according to (4.1), where the floating-point addition is optimal. Furthermore, let formula (4.3) be calculated such that the subtraction producing  $w \in R^*$  is faithful with respect to  $R^*$  and the other operation is faithful with respect to  $R$ . Then formula (4.3) yields the correction term  $zz \in R$  defined by (4.2).

Proof. The proof is analogous to that of theorem (4.7). We have to show (1) and (2). The proof of (2) is unchanged and that of (1) is modified as follows.

We need to show only that  $z - x$  is an element of  $R^*$ .

In case (1a) and in the case that overflow occurs in forming  $z$ , we obtain  $z - x \in R$ .

Hence,  $z - x \in R^*$ , since  $R \subset R^*$ .

In case (1b) we obtain  $z - x = \mu \beta^{ex}$ , where  $|\mu| < \beta/2 + M \leq M^*$ .

Hence  $z - x \in R^*$ .

This completes the proof of (1) and of the theorem.

For properly truncating addition, the following theorem holds without any restriction on  $\beta$  and  $M$  and without requiring an enlarged mantissa range for  $w$ .

(4.10) Theorem. If  $R$  has the form (2.2), floating-point addition is properly truncating and subtraction is faithful, then, for all  $x$  and  $y$  satisfying (4.5) and for  $z$  obtained according to (4.1), formula (4.3) yields the correction term  $zz$  defined by (4.2).

Proof. The proof is analogous to that of theorem (4.7). We have to show (1) and (2).

Proof of (1). We need to show only that  $z - x$  is an element of  $R$ . Let  $\text{trunc}(\xi)$  denote the largest integer  $\leq \xi$  if  $y \geq 0$  or the smallest integer  $\geq \xi$  if  $y < 0$  (i.e. the properly truncated value of  $\xi$ ; see 3.5).

Obviously,  $z$  can be represented such that  $ez \leq ex + 1$ . We first assume that no overflow occurs in forming  $z$  and distinguish between two cases.

1a)  $ez \leq ex$ .

Then we have

$$\begin{aligned} mz &= mx \beta^{ex-ez} + \text{trunc}(my \beta^{ey-ez}), \\ z - x &= (mz - mx \beta^{ex-ez}) \beta^{ez} = \text{trunc}(my \beta^{ey-ez}) \beta^{ez}. \end{aligned}$$

Hence,  $z - x \in R$ .

1b)  $ez = ex + 1$ .

Let  $d = ex - ey$ . Then

$$\begin{aligned} mz &= \text{trunc}(mx/\beta + my/\beta^{d+1}), \\ z - x &= (\beta mz - mx) \beta^{ex} = \mu \beta^{ex} \text{ (say)}. \end{aligned}$$

Since  $mz$  is obtained by properly truncating, we have

$$|\mu| \leq |my/\beta^d| < M.$$

Hence  $z - x \in R$ , because  $\mu$  is integral.

In the case that overflow occurs in forming  $z$ , we can show that  $z - x \in R$  in the same way as in the proof of (4.7).

This completes the proof of (1).

Proof of (2). We need to show only that  $y - w$  is an element of  $R$ .

From  $ex \geq ey$ , it follows that  $w$  equals an integer times  $\beta^{ey}$ . Since  $z$  is obtained by properly truncating the mantissa of  $y$  (case 1a) or  $x + y$  (case 1b), or by truncating the mantissa  $y$  to some value smaller in magnitude (in case of overflow), we have  $|w| \leq |y|$ . Moreover,  $w$  and  $y$  do

not have opposite signs, because otherwise  $x$  would be between  $z$  and  $x + y$  contradicting the assumption that floating-point addition is properly truncating.

Hence,  $|y-w| \leq |y|$ , so that  $y - w \in R$ .

This completes the proof of (2) and the theorem.

#### Remarks

1) As pointed out in section 2 remark 1, addition and subtraction never lead to underflow in a system of the form (2.2). If we restrict ourselves to normalized representations, then theorems (4.7), (4.9) and (4.10) do not remain valid, since underflow may occur in forming  $zz$ . Of course, the theorems remain valid in cases where no underflow occurs.

2) The theorems do not hold if addition is only faithful, because the correction term  $zz$  is then not always an element of  $R$ . For example, if  $\beta = 2$ ,  $M = 16$ ,  $x = 15$ ,  $y = 15/32$  and  $z = 16$ , then  $zz = -17/32$ , which is not an element of  $R$ .

For super-faithful addition (see definition 3.7), we immediately obtain by combining theorems (4.7) and (4.10):

(4.11) Corollary. Theorem (4.7) remains valid, if "optimal" is replaced by "super-faithful".

Similarly, combining theorems (4.9) and (4.10), we obtain

(4.12) Corollary. Theorem (4.9) remains valid, if "optimal" is replaced by "super-faithful".

Instead of formula (4.3), we now consider the formula

$$(4.13) \quad w = fl(x-z), \quad zz = fl(w+y).$$

For practical computation, formulas (4.1) and (4.13) can be written as the following sequence of ALGOL 60 statements:



$$(4.14) \quad "z := x + y; zz := x - z + y",$$

in which  $w$  again remains anonymous. Here, we make use of the fact that, in ALGOL 60, subsequent additions and subtractions are performed from left to right, so that the expression " $x - z + y$ " is equivalent to " $(x-z) + y$ ", see Naur (1962, 3.3.5.). For applications, we prefer this formula above (4.4), because many compilers produce a slightly faster code for (4.14) than for (4.4).

Since the floating-point number system  $R$  is symmetric (i.e.  $x \in R$  implies  $-x \in R$ , see (2.2) and section 2 remark 2), and  $R^*$  also (see 4.8), and since optimal, properly truncating or super-faithful addition is certainly faithful, we immediately obtain

(4.15) Corollary. Theorems (4.7), (4.9) and (4.10), and corollaries (4.11) and (4.12) remain valid if formula (4.3) is replaced by formula (4.13).

The formulas (4.3) and (4.13) are numerically equivalent but for the sign of  $w$ . (Since the second subtraction in (4.3) corresponds to the addition in (4.13), we talk about "other operation" in theorem (4.9).)

The pair of formulas (4.1) and (4.13), or the equivalent formulation (4.14), is our basic algorithm for exact addition of two floating-point numbers. In the applications, it will often be necessary to interchange the roles of the terms  $x$  and  $y$  when  $|x| < |y|$ , in order to ensure that (4.5) holds. In the subsequent sections, we shall use this basic algorithm to obtain algorithms for doublelength arithmetical operations.

We conclude this section with a formula for the correction term  $zz$  satisfying (4.2), in the case that floating-point addition and subtraction are optimal,  $R$  has the form (2.2) with arbitrary  $\beta$  and  $M$ , but no enlarged mantissa range is available for anonymous real values (in particular,  $w$ ). Since we cannot guarantee that  $w = z - x$  in (4.3), we apply formula (4.13) to obtain the correction term,  $z_2$ , for  $w$ . This leads to the formula

$$(4.16) \quad \left\{ \begin{array}{l} w = fl(z-x), z_1 = fl(y-w), \\ v = fl(z-w), z_2 = fl(v-x), \\ zz = fl(z_1-z_2). \end{array} \right.$$

Since we assume that  $R$  is symmetric, this formula is equivalent to that given by Møller (1965, p. 42 process A) and Knuth (1969, p. 203 formula 48), the only difference being that there  $-z_2$  instead of  $z_2$  is calculated and added to  $z_1$ .

For this formula, the following theorem holds even without requiring (4.5).

(4.17) Theorem of Møller-Knuth.

If  $R$  has the form (2.2) and floating-point addition and subtraction are optimal, then for all  $x, y \in R$  and for  $z$  obtained according to (4.1), formula (4.16) yields the correction term  $zz$  defined by (4.2).

Møller (1965) has a weaker assumption for floating-point addition and subtraction and obtains a weaker result. Knuth (1969, p. 203 Theorem B) assumes that floating-point addition and subtraction are optimal (p. 197, formula 11), but excludes overflow and underflow. In fact, overflow does not invalidate the theorem, whereas underflow does not occur in a number system  $R$  of the form (2.2) (see section 2 remark 1).

For a proof of this theorem we refer to Knuth (1969, p. 201-203). An alternative proof can be given along the lines of the proof of Theorem (4.7).

### 5. Exact multiplication

Let  $x$  and  $y$  be given elements of  $R$ . We want to calculate their exact product and to deliver it as a pair  $(z, zz)$  of elements of  $R$  satisfying the (exact) relation

$$(5.1) \quad z + zz = x \times y,$$

with some extra condition that  $zz$  be (almost) negligible within machine precision with respect to  $x \times y$ .

For simplicity, we assume that the system  $R$  of floating-point numbers has the form (2.2), with the restrictions

$$(5.2) \quad \beta = 2, \quad M = 2^t,$$

where  $t$  is the number of binary digits in the mantissa, and  $D$  and  $E$  are infinite. Thus  $R$  obtains the form

$$(5.3) \quad R = R(t) = \{x \mid x = mx \cdot 2^{ex}, |mx| < 2^t\}.$$

In other words, we restrict ourselves to binary floating-point  $t$ -digit arithmetic disregarding overflow and underflow. The following results on exact multiplication can be generalized, however, to nonbinary systems. We use the notation  $R(t)$ , because we shall also refer to binary floating-point  $k$ -digit number systems  $R(k)$  for some values  $k \neq t$ .

Moreover, we assume that floating-point addition and subtraction are optimal and multiplication is faithful (see definitions 3.1 and 3.2).

In order to form the exact product of  $x$  and  $y$ , we first split  $x$  and  $y$  each into two "half-length" numbers. Let  $t_1$  and  $t_2$  be (roughly equal) integers such that

$$(5.4) \quad t = t_1 + t_2.$$

Let  $hx$  (the "head" of  $x$ ) be an element of  $R(t_2)$  as near to  $x$  as possible; i.e. if  $x = mx \cdot 2^{ex}$ , where  $mx$  is normalized when  $x \neq 0$  (thus,  $ex$  is minimal and  $|mx| \geq 2^{t-1}$ ), then

$$(5.5) \quad hx = \text{round}(mx 2^{-t1}) 2^{ex+t1}.$$

(For the definition of round see 3.3.) Let tx (the "tail" of x) be the remaining part of x:

$$(5.6) \quad tx = x - hx.$$

Since hx is obtained by rounding, tx is an element of  $R(t1-1)$ . Let y similarly split into hy and ty.

Then we obviously have

$$\begin{aligned} hx \times hy &\in R(2t2), \\ hx \times ty, tx \times hy &\in R(t-1), \\ tx \times ty &\in R(2t1-2). \end{aligned}$$

So, in order to make sure, that these numbers fit in  $R = R(t)$ , we choose (cf. 5.4)

$$(5.7) \quad t2 = \text{entier}(t/2), t1 = t - t2.$$

Then  $2t2 \leq t$  and  $2t1 - 2 \leq t - 1$ . Since  $hx \times ty$  and  $tx \times hy$  are representable as elements of  $R(t-1)$  with the same exponent, their sum is an element of  $R(t)$ . So, the calculations

$$(5.8) \quad \begin{cases} p = \text{fl}(hx \times hy), \\ q = \text{fl}(hx \times ty + tx \times hy), \\ r = \text{fl}(tx \times ty) \end{cases}$$

are exact, because the floating-point operations involved are faithful. We find z and the correction term z1(say) by performing an exact addition of p and q:

$$(5.9) \quad z = \text{fl}(p+q), z1 = \text{fl}((p-z) + q).$$

Since certainly  $|p| \geq |q|$ , we have  $ep \geq eq$ , in other words, the relation corresponding to (4.5) holds. Hence,  $z + z1 = p + q$  according to theorem (4.7) and corollary (4.15), since we assume that R has the form (5.3) and floating-point addition and subtraction are optimal.

From these assumptions, it also follows that  $z_1$  is representable as an element of  $R(t-1)$  with the same exponent as  $r$ . Hence,  $z_1 + r \in R$ , so that  $zz$  defined by (5.1) is obtained from

$$(5.10) \quad zz = fl(z_1 + r).$$

So we have proved the following

(5.11) Theorem. If  $R$  has the form (5.3), floating-point addition and subtraction are optimal and multiplication is faithful, then for all  $x, y \in R$  splitted into head and tail according to (5.5) and (5.6), where  $t_1$  and  $t_2$  are given by (5.7), the formulas (5.8), (5.9) and (5.10) yield  $z$  and  $zz$  satisfying (5.1).

Formulas (5.8), (5.9) and (5.10) can be written as the following sequence of ALGOL 60 statements (cf. 4.14):

$$(5.12) \quad \left\{ \begin{array}{l} "p := hx \times hy; q := hx \times ty + tx \times hy; \\ z := p + q; zz := p - z + q + tx \times ty". \end{array} \right.$$

Note that, in the last assignment statement, the additions and subtractions are performed from left to right (see (4.14) and Naur (1962, 3.3.5.)). We use (5.12) in our algorithm for exact multiplication (see ALGOL 60 procedure "mul 12" in appendix).

As to the extra condition that " $zz$  be (almost) negligible within machine precision with respect to  $x \times y$ ", we shall show for  $t \geq 2$  that

$$(5.13) \quad |zz| \leq |x \times y| \frac{\tau 2^{-t}}{1 + \tau 2^{-t}},$$

where  $\tau = 2$  if  $t$  even and  $\tau = 3$  otherwise.

Proof. If  $x = 0$  or  $y = 0$ , then obviously  $zz = z = 0$  and (5.13) holds. So we may assume that  $x$  and  $y$  are nonzero, and thus also  $hx$ ,  $hy$  and  $p+q$ . Let

$$\epsilon = tx/x, \eta = ty/y, \delta = z_1/(p+q).$$

Then a simple calculation yields

$$zz = (x \times y) ((1 - \epsilon\eta) \delta + \epsilon\eta).$$

Since floating-point addition is optimal, we have

$$|\delta| \leq \frac{2^{-t}}{1+2^{-t}}$$

(cf. Wilkinson (1963, p. 17-19) who gives the bound  $2^{-t}$ ). Similarly, it follows from (5.5) and (5.6) that

$$|\epsilon|, |\eta| \leq \frac{2^{-t2}}{1+2^{-t2}}.$$

Hence,

$$\begin{aligned} |zz| &\leq |x \times y| \left[ \left( 1 + \left( \frac{2^{-t2}}{1+2^{-t2}} \right)^2 \right) \frac{2^{-t}}{1+2^{-t}} + \left( \frac{2^{-t2}}{1+2^{-t2}} \right)^2 \right] \\ &= |x \times y| \frac{2^{-t}}{1+2^{-t}} \left[ 1 + \frac{2^{t-2t2}(1+2^{1-t})}{(1+2^{-t2})^2} \right]. \end{aligned}$$

Since, according to (5.7),

$$2^{t-2t2} = \tau - 1$$

and, for  $t \geq 2$ ,

$$1 + 2^{-t2} \geq 1 + \tau 2^{-t} \geq 1 + 2^{1-t},$$

we obtain

$$|zz| \leq |x \times y| \frac{2^{-t}}{1+2^{-t}} \left[ 1 + \frac{\tau - 1}{1 + \tau 2^{-t}} \right]$$

from which (5.13) immediately follows.

A slightly smaller bound for  $zz$  is obtained if the algorithm is modified as follows. After calculating  $z$  and  $zz$  as above (see 5.12), an exact addition of  $z$  and  $zz$  is performed:

$$(5.14) \quad "u := z; z := u + zz; zz := u - z + zz".$$

Then  $z$  and  $zz$  still satisfy (5.1) and , since addition is optimal, we now have

$$(5.15) \quad |zz| \leq |x \times y| \frac{2^{-t}}{1+2^{-t}}.$$

This bound is only about 2 or 3 times smaller than the bound given by (5.13). Therefore we did not include (5.14), requiring three extra additions or subtractions, in our algorithm "mul 12".

The following algorithm for exact multiplication, due to Veltkamp (1968), also works under the conditions stated in theorem (5.11).

$$z = fl(x \times y),$$

$$zz = fl((((hx \times hy - z) + hx \times ty) + tx \times hy) + tx \times ty).$$

If multiplication is optimal, then  $zz$  satisfies (5.15), otherwise the bound for  $zz$  is about twice as large, because multiplication is assumed to be faithful. Veltkamp's algorithm requires one more multiplication and one less addition than ours.

## 6. Splitting into halflength numbers

Let  $x$  be a given element of  $R$ . We want to calculate  $hx$  and  $tx$  as defined by (5.5) and (5.6) by means of some simple arithmetical operations. We consider the formula

$$(6.1) \quad p = fl(x \times c), \quad q = fl(x-p), \quad hx = fl(q+p),$$

where  $c$  is some constant, and we try to find a value for  $c$  such that the formula yields  $hx$  defined by (5.5). The most obvious choices  $c = \pm 2^{t_1}$  do not always work.

For example, if  $t$  is even,  $c = 2^{t_1}$  and  $|mx| = 2^{t-1} + 2^{t_1-1}$ , then  $fl(q+p)$  is not an element of  $R(t_2)$ ; if  $c = -2^{t_1}$  and  $|mx| = 2^{t-1}$ , then  $fl(q+p)$  is not properly rounded to  $t_2$  bits.

On the other hand, defining

$$(6.2) \quad c = 2^{t_1} + 1,$$

the following theorem holds.

(6.3) Theorem. If  $R$  has the form (5.3), floating-point addition and subtraction are optimal, multiplication is faithful, and  $t_1$ ,  $t_2$  and  $c$  are defined by (5.7) and (6.2), then, for all  $x \in R$ , formula (6.1) yields  $hx$  defined by (5.5).

Proof Let, for any real  $\xi$ ,  $\text{trunc}(\xi)$  denote the largest integer  $\leq \xi$  or the smallest integer  $\geq \xi$ . If  $x = 0$ , then the theorem is obvious. If  $x \neq 0$ , then we may and shall assume (for definiteness) that the floating-point numbers involved are normalized, i.e. the absolute value of the mantissas are  $\geq 2^{t-1}$ . Since  $hx$  defined by (5.5) is an element of  $R(t_2)$ , thus, certainly of  $R = R(t)$ , and addition is optimal, we need to show only that  $q + p$  obtained from (6.1) equals  $hx$ . Obviously  $ep$  equals either  $ex + t_1$  or  $ex + t_1 + 1$  which cases we now consider separately.

a)  $ep = ex + t_1$ . Since multiplication is faithful, we have

$$mp = \text{trunc}(mx + mx 2^{-t_1}) = mx + \text{trunc}(mx 2^{-t_1}).$$



Obviously,  $eq = ep$  or  $eq = ep-1$ . We show that the latter is impossible. Indeed, if  $eq = ep-1$ , then, since subtraction is optimal,

$$mq = \text{round} (mx 2^{1-t_1}) - 2mp = \text{round} (mx 2^{1-t_1}) - 2mx - 2 \text{trunc} (mx 2^{-t_1}).$$

Hence,

$$|mq| > 2|mx| - \frac{1}{2} - 2.$$

If  $|mx| > 2^{t-1}$ , then this yields  $|mq| \geq 2^t$ , since  $mq$  is an integer, and if  $|mx| = 2^{t-1}$ , then round and trunc are exact, so that  $|mq| = 2^t$ . So,  $mq$  is not within the mantissa range, which shows that  $eq = ep-1$  is impossible. Hence, we have  $eq = ep = ex + t_1$  and

$$mq = \text{round} (mx 2^{-t_1}) - mp,$$

so that

$$q+p = \text{round} (mx 2^{-t_1}) 2^{ep},$$

which equals  $hx$  defined by (5.5).

b)  $ep = ex + t_1 + 1$ . Then

$$mp = \text{trunc} (mx/2 + mx 2^{-t_1-1}).$$

Obviously,  $eq = ep-1$  or  $eq = ep$ .

If  $eq = ep-1 = ex + t_1$ , then

$$mq = \text{round} (mx 2^{-t_1}) - 2mp,$$

so that

$$q+p = \text{round} (mx 2^{-t_1}) 2^{eq},$$

which equals  $hx$  defined by (5.5).

If  $eq = ep = ex + t1 + 1$ , then

$$mq = \text{round} (mx 2^{-t1-1}) - mp = \text{round} (mx 2^{-t1-1}) - \text{trunc} (mx/2 + mx 2^{-t1-1}).$$

Hence,

$$|mq| < |mx|/2 + \frac{1}{2} + 1.$$

This is within the normalized mantissa range  $2^{t-1} \leq |mq| < 2^t$  only if  $|mx| = 2^t - \epsilon$ , where  $\epsilon = 1$  or  $2$ , but then we obtain

$$q+p = \text{round} (mx 2^{-t1-1}) 2^{ep} = \text{round} (mx 2^{-t1}) 2^{ep-1},$$

which equals  $hx$  defined by (5.5).

This completes the proof.

After calculating  $hx$  according to (6.1), we obtain  $tx$  defined by (5.6) from

$$(6.4) \quad tx = fl(x-hx)$$

which is equivalent to (5.6), because  $tx \in R$  and subtraction is optimal. Formulas (6.1) and (6.4) can be written as the following sequence of ALGOL 60 statements, where  $c$  is assumed to have the value given by (6.2):

$$(6.5) \quad "p := x \times c; hx := x - p + p; tx := x - hx".$$

Note that, in ALGOL 60, " $x - p + p$ " is equivalent to " $(x-p) + p$ " (see (4.14) and Naur (1962, 3.3.5.)). We use formula (6.5) in our algorithm for exact multiplication (see ALGOL 60 procedure "mul 12" in appendix).

## 7. Doublelength arithmetic

We now give a straightforward application of the results of the previous sections to obtain algorithms for doublelength addition, subtraction, multiplication, division and square root. In the appendix, the algorithms are given as ALGOL 60 procedures. For simplicity, we assume that the conditions of section 5 are satisfied, viz.  $R$  has the form (5.3), floating-point addition and subtraction are optimal, and multiplication is faithful. The algorithms for addition and subtraction, however, would also work correctly under the weaker assumptions stated in theorems (4.7), (4.9) and (4.10) and corollaries (4.11), (4.12) and (4.15), provided that no overflow occurs.

In several computers, more than one floating-point system is available, e.g. a "single" and a "double" precision system. Our technique is, of course, most effective if we start from the highest precision system available (provided that it satisfies the requirements). Thus, starting from a "double" precision system and arithmetic, our technique yields a "quadruple" precision system and arithmetic.

We shall, however, call the system  $R$  used as starting-point for our technique "singlelength floating-point number system", and the arithmetical operations defined on  $R$  "singlelength floating-point arithmetical operations".

(7.1) Definition. A doublelength floating-point number is a pair  $(r,s)$  of singlelength floating-point numbers (i.e.  $r,s \in R$ ) satisfying

$$(7.2) \quad |s| \leq |r+s| \frac{2^{-t}}{1+2^{-t}}.$$

The value of the doublelength number  $(r,s)$  is, by definition, equal to  $r + s$ . We call  $r$  the head and  $s$  the tail of  $(r,s)$ .

In particular, any pair  $(r,0)$  is a doublelength floating-point number, and, since addition is optimal, also any pair  $(z,zz)$  obtained by performing an exact addition (see section 4).

Sometimes we replace (7.2) by the weaker condition

$$(7.3) \quad |s| \leq |r+s| C 2^{-t},$$

where  $C$  is some constant not much larger than 1; we call a pair  $(r,s)$  satisfying (7.3) a nearly doublelength floating-point number.

In particular, a pair  $(z,zz)$  obtained by exact multiplication (5.12) is a nearly doublelength floating-point number, because, according to (5.13) and Theorem (5.11), we can take  $C = \tau/(1 + \tau 2^{-t})$ .

On the other hand, the magnitude of the tail of a doublelength number may be much smaller than the bound given by (7.2). Thus, the mantissa of a doublelength floating-point number cannot always be represented by means of a multilength integer of fixed maximum length. This is in contrast with the usual approach in which the mantissa is represented by means of some integers.

The doublelength sum of two (nearly) doublelength floating-point numbers  $(x,xx)$  and  $(y,yy)$  is calculated as follows (see ALGOL 60 procedure "add 2" in appendix).

First, the heads  $x$  and  $y$  are added exactly (4.14). Here the roles of  $x$  and  $y$  are interchanged when  $|x| < |y|$ , in order to ensure that (4.5) holds. Thus, we obtain a doublelength number  $(r,rr)$  such that  $r + rr = x + y$ . Subsequently, the tails are added to  $rr$ :

$$(7.4) \quad s = fl((rr+yy) + xx),$$

so that  $r + s$  approximately equals the sum of  $(x,xx)$  and  $(y,yy)$ . Here again the roles of  $xx$  and  $yy$  are interchanged when  $|x| < |y|$ , in order to reduce the maximum error in (7.4) and to ensure commutativity. Finally, an exact addition of  $r$  and  $s$  is performed. Although not always  $|r| \geq |s|$  (cancellation in forming  $r$  may cause  $|r|$  to be slightly smaller than  $|s|$ ), the relation corresponding to (4.5) certainly holds, so that, for this final exact addition, we never need interchange the roles of  $r$  and  $s$ . Since singlelength addition is optimal, the final exact addition transforms the approximate sum into a doublelength floating-point number having the same value.

Doublelength subtraction is performed in a completely analogous fashion (see ALGOL 60 procedure "sub 2" in appendix).

The calculation of doublelength product, quotient and square root is rather obvious and can be sketched as follows (for details see error analysis below and ALGOL 60 procedures "mul 2", "div 2" and "sqrt 2" in appendix). First a nearly doublelength approximation,  $(c, cc)$ , of the required result is calculated. Here, besides some singlelength operations, exact multiplication is used. The pair  $(c, cc)$  satisfies the relation corresponding to (7.3), but not always (7.2). Therefore, an exact addition is performed, which transforms the result obtained into a doublelength floating-point number having the same value.

#### Error analysis

The only error in the doublelength addition is committed in forming  $s$  (7.4). Assuming that  $(x, xx)$  and  $(y, yy)$  are nearly doublelength numbers satisfying

$$(7.5) \quad |xx| \leq |x+xx| C_1 2^{-t}, \quad |yy| \leq |y+yy| C_2 2^{-t},$$

where  $C_1$  and  $C_2$  are constants not much larger than 1, we shall show that, for sufficiently large  $t$ , the error,  $E^+$ , of the doublelength addition satisfies

$$(7.6) \quad |E^+| \leq \{|x+xx|(1+C_1) + |y+yy|(1+C_2)\} 2^{1-2t}.$$

Proof From (7.4), it follows that

$$s = (rr+yy) (1+\epsilon) (1+\epsilon') + xx(1+\epsilon'),$$

where  $|\epsilon|, |\epsilon'| \leq 2^{-t}/(1+2^{-t}) < 2^{-t}$ , because addition is optimal, cf. Wilkinson (1963, p.9). Similarly,

$$|rr| \leq |x+y| 2^{-t}.$$

From these relations and (7.5), we obtain

$$\begin{aligned}
|E^+| &\leq |rr+yy|(2^{1-t}+2^{-2t}) + |xx|2^{-t} \\
&\leq |x+xx|\{(2^{1-2t}+2^{-3t})(1+C_1 2^{-t}) + C_1 2^{-2t}\} \\
&\quad + |y+yy|(2^{1-2t}+2^{-3t})(1+C_2 2^{-t}+C_2).
\end{aligned}$$

Since  $|y| \leq |x|$  (otherwise the roles of  $x$  and  $y$  would be interchanged), the sum of the terms of order of magnitude  $2^{-3t}$  and lower is smaller than  $|x+xx| C_1 2^{-2t}$  for sufficiently large  $t$ . This establishes (7.6).

If  $(x,xx)$  and  $(y,yy)$  are doublelength numbers, then  $C_1, C_2 < 1$  so that (7.6) reduces to

$$(7.7) \quad |E^+| \leq (|x+xx| + |y+yy|) 2^{2-2t}$$

If, however, doublelength addition is used, in combination with our exact multiplication algorithm, for calculating the doublelength scalar product of two vectors of singlelength floating-point numbers, then one of the constants  $C_1, C_2$  is smaller than  $\tau$  (i.e. 2 if  $t$  even and 3 otherwise (5.13)) and the other is smaller than 1. Repeated application of (7.6) then yields an upper bound for the error of the doublelength scalar product, cf. Wilkinson (1963, p.18).

Formula (7.6) means that the error is small in doublelength precision with respect to the sum of the magnitudes of the terms, the loss being at most a few bits. The relative error of the doublelength sum, however, is not always small, because severe cancellation may take place in forming  $r$ . Since this can happen only if  $rr = 0$ , the algorithm for doublelength addition can easily be modified such that a small relative error is ensured (Veltkamp, 1968). We did not include this modification in our algorithm, because it is of only limited value if the terms are not exact.

The doublelength product, quotient and square root produced by our algorithms all have a small relative error. In other words, if the operands are (nearly) doublelength numbers, then the corresponding absolute errors,  $E^x$ ,  $E^/$  and  $E^{\sqrt{\phantom{x}}}$ , satisfy

$$(7.8) \quad \begin{cases} |E^x| \leq |x+xx| \times |y+yy| C^x 2^{-2t}, \\ |E'| \leq (|x+xx| / |y+yy|) C' 2^{-2t}, \\ |E^\sqrt{}| \leq \sqrt{(x+xx)} C^\sqrt{} 2^{-2t}, \end{cases}$$

where  $C^x$ ,  $C'$  and  $C^\sqrt{}$  are constants not much larger than 1. (Here we assume, of course, that the denominator of the quotient is nonzero and the argument of the square root is nonnegative.)

If the operands are doublelength floating-point numbers (7.1) and single-length division and square root are faithful, then, for  $t \geq 10$  (say), the constants are bounded by

$$(7.9) \quad C^x \leq 9 + \tau, \quad C' \leq 21.1, \quad C^\sqrt{} \leq 12.7.$$

where  $\tau = 2$  if  $t$  even and  $\tau = 3$  otherwise (5.13). If, moreover, single-length multiplication and division are optimal (which holds for the Philips Electrologica X8 computer), then

$$(7.10) \quad C^x \leq 7 + \tau, \quad C' \leq 12.1, \quad C^\sqrt{} \leq 10.2.$$

Proof Let

$$P = (x+xx) \times (y+yy), \quad P' = x \times y + x \times yy + xx \times y.$$

Since the operands are doublelength numbers (7.1), we have

$$(7.11) \quad |P' - P| = |xx \times yy| \leq |P| 2^{-2t}.$$

In "mul 2" the following approximation,  $P''$ , is calculated:

$$P'' = c + \text{fl}((x \times yy + xx \times y) + cc),$$

where  $c$  and  $cc$ , obtained by exactly multiplying  $x$  and  $y$ , satisfy (cf. section 5)

$$c + cc = x \times y, \quad |cc| \leq |x \times y| \frac{\tau 2^{-t}}{1 + \tau 2^{-t}}.$$

Hence, we obtain (cf. Wilkinson (1963, p. 7-11))

$$P'' = c + \left[ \{(x \times yy)(1+\epsilon) + (xx \times y)(1+\epsilon')\}(1+\epsilon_1) + c\bar{c} \right] (1+\epsilon_2),$$

where

$$|\epsilon_1|, |\epsilon_2| \leq 2^{-t}/(1+2^{-t}),$$

because addition is optimal, and, for some constant  $\mu$ ,

$$|\epsilon|, |\epsilon'| \leq \mu 2^{-t}/(1+2^{-t}).$$

From the assumption that multiplication is faithful, it follows that  $\mu \leq 2$ ; if multiplication is optimal, then  $\mu = 1$ .

From the formulas for  $P'$  and  $P''$ , we easily derive

$$(7.12) \quad |P''-P'| \leq |P| (2\mu+4+\tau) 2^{-2t}.$$

So, combining (7.11) and (7.12), we obtain

$$(7.13) \quad |E^x| = |P''-P| \leq |P| (2\mu+5+\tau) 2^{-2t}.$$

This establishes the relations given above for  $C^x$ .

Let

$$Q = (x+xx)/(y+yy), \quad Q' = (x+xx - x \times yy/y)/y.$$

Since the operands are doublelength numbers (7.1), we have

$$(7.14) \quad |Q'-Q| \leq |Q| 2^{1-2t}.$$

In "div 2" the following approximation,  $Q''$ , is calculated:

$$Q'' = c + fl((((x-u) - uu) + xx) - c \times yy)/y),$$

where

$$c = fl(x/y) = (x/y)(1+\epsilon) \text{ (say)}, \quad u + uu = c \times y,$$



$u$  and  $uu$  being obtained by exactly multiplying  $c$  and  $y$ . Since cancellation occurs in calculating  $\text{fl}(x-u)$ , this is exact. Moreover, since division is faithful, we have  $x - c \times y \in R$ , so that

$$\text{fl}((x-u) - uu) = x - u - uu = -x\varepsilon.$$

Hence (cf. Wilkinson (1963, p. 7-11)),

$$\begin{aligned} Q'' &= (x/y)(1+\varepsilon) + \left[ \{(x/y)(-\varepsilon) + xx/y\} (1+\varepsilon_1) \right. \\ &\quad \left. - (x/y)(yy/y)(1+\varepsilon)(1+\varepsilon'') \right] (1+\varepsilon_2)(1+\varepsilon'), \end{aligned}$$

where

$$|\varepsilon_1|, |\varepsilon_2| \leq 2^{-t}/(1+2^{-t}),$$

because addition and subtraction are optimal, and, for some constant  $\mu$ ,

$$|\varepsilon|, |\varepsilon'|, |\varepsilon''| \leq \mu 2^{-t}.$$

Since we assume that multiplication and division are faithful, we obviously have  $\mu \leq 2$ ; if these operations are optimal, then  $\mu < 1$ . Denoting the sum of the terms of order  $2^{-3t}$  and lower by  $\theta 2^{-2t}$ , we obtain

$$(7.15) \quad |Q''-Q'| \leq |Q| (\mu^2+6\mu+3+\theta)2^{-2t}.$$

Further analysis shows that, for  $t \geq 10$ , we can take  $\theta = 0.1$ . So, from this and (7.14), we obtain

$$(7.16) \quad |E'| = |Q''-Q| \leq |Q| (\mu^2+6\mu+5.1)2^{-2t}.$$

This establishes the relations given above for  $C'$ .

Let

$$R = \sqrt{x + xx}, \quad R' = \sqrt{x + xx/(2\sqrt{x})}.$$

Since the argument is a doublelength number, we have

$$(7.17) \quad |R'-R| \leq \frac{1}{8}|R| 2^{-2t}.$$

In "sqrt 2" the following approximation,  $R''$ , is calculated:

$$R'' = c + \text{fl}(((x-u) - uu) + xx) \times 0.5/c,$$

where

$$c = \text{sqrt}(x) = \sqrt{x}(1+\epsilon) \text{ (say), } u + uu = c^2,$$

$u$  and  $uu$  being obtained by exactly multiplying  $c$  and  $c$ . Since cancellation occurs in calculating  $\text{fl}(x-u)$ , this is exact. Hence,

$$\begin{aligned} \text{fl}((x-u) - uu) &= (x-u-uu)(1+\epsilon_1) = -x(2\epsilon+\epsilon^2)(1+\epsilon_1), \\ R'' &= \sqrt{x}(1+\epsilon) + \left[ \{-x(2\epsilon+\epsilon^2)(1+\epsilon_1) + xx\} \times 0.5/\sqrt{x} \right] \frac{(1+\epsilon_2)(1+\epsilon')}{1+\epsilon}, \end{aligned}$$

where, for some constant  $\mu \leq 2$ ,

$$|\epsilon_1|, |\epsilon_2| \leq 2^{-t}/(1+2^{-t}), |\epsilon'| \leq \mu 2^{-t}, |\epsilon| \leq 2^{1-t}.$$

Here again  $\mu < 1$  if division is optimal. From this we obtain

$$(7.18) \quad ||R''-R'| \leq |R| (2.5\mu+7.5+\theta)2^{-2t},$$

where  $\theta 2^{-2t}$  again denotes the sum of the terms of order  $2^{-3t}$  and lower. For  $t \geq 10$ , we can take  $\theta = 0.075$ . So, from this and (7.17), we obtain

$$|E^V| = |R''-R| \leq |R| (2.5\mu+7.7)2^{-2t}.$$

This establishes the relations given above for  $C^V$  and completes the proof.

The doublelength operations defined by our algorithms "add 2", "sub 2", "mul 2", "div 2" and "sqrt 2" are not faithful. To obtain faithful doublelength operations, one would have to calculate the result in triplelength (or maybe even quadruplelength) precision and then to round or truncate it to doublelength. This would, of course, require considerably more work and be preferable only in exceptional situations.

## 8. Appendix. ALGOL 60 procedures

In this appendix, we give a set of ALGOL 60 procedures for doublelength arithmetic and exact multiplication.

The procedures work correctly if the singlelength floating-point system is binary, singlelength floating-point addition and subtraction are optimal (3.2), multiplication is faithful (3.1) and no overflow or underflow occurs (cf. section 5).

The procedures for doublelength addition and subtraction also work correctly in a nonbinary system, provided that a guarding digit is available for the mantissas of anonymous quantities (i.e. quantities which are not assigned to a variable; for details see section 4).

The procedures have been tested on the Philips Electrologica X8 computer at Mathematical Centre, Amsterdam. This is a binary machine having  $t = 40$  binary digits in the mantissa and 12 bits for the exponent (i.e.  $D = E = 2048$  in (2.2)).

In particular, the algorithms described in "add 2" and "mul 12" have been used extensively for calculating doublelength scalar products of vectors of singlelength floating-point numbers.

Procedure "add 2" has been tested also on the GE635 computer at Bell Telephone Laboratories, Murray Hill. This is a binary machine having two floating-point systems, viz. single precision  $t = 27$  and double precision  $t = 63$ . Procedure "add 2" worked correctly in both systems.

In the comments,  $(x,xx)$ ,  $(y,yy)$  and  $(z,zz)$  denote (nearly) doublelength numbers (see section 7).

comment add 2 calculates the doublelength sum of (x,xx) and (y,yy),  
the result being (z,zz);

procedure add 2 (x,xx,y,yy,z,zz);

value x,xx, y,yy; real x,xx, y,yy, z,zz;

begin real r, s;

    r := x + y;

    s := if abs (x) > abs (y) then

        x - r + y + yy + xx else y - r + x + xx + yy;

    z := r + s;

    zz := r - z + s

end add 2;

comment sub 2 calculates the doublelength difference of (x,xx) and  
(y,yy), the result being (z,zz);

procedure sub 2 (x,xx, y,yy, z,zz);

value x,xx, y,yy; real x,xx, y,yy, z,zz;

begin real r, s;

    r := x - y;

    s := if abs (x) > abs (y) then

        x - r - y - yy + xx else - y - r + x + xx - yy;

    z := r + s;

    zz := r - z + s

end sub 2;

comment mul 12 calculates the exact product of x and y, the result being the nearly doublelength number (z,zz). The constant should be chosen equal to  $2 \uparrow (t - t \div 2) + 1$ , where t is the number of binary digits in the mantissa;

```
procedure mul 12 (x, y, z,zz);
value x, y; real x, y, z,zz;
begin real hx, tx, hy, ty, p, q;
    p := x × constant;
    hx := x - p + p; tx := x - hx;
    p := y × constant;
    hy := y - p + p; ty := y - hy;
    p := hx × hy;
    q := hx × ty + tx × hy;
    z := p + q;
    zz := p - z + q + tx × ty
end mul 12;
```

comment mul 2 calculates the doublelength product of (x,xx) and (y,yy), the result being (z,zz);

```
procedure mul 2 (x,xx, y,yy, z,zz);
value x,xx,y,yy; real x,xx, y,yy, z,zz;
begin real c,cc;
    mul 12 (x, y, c,cc);
    cc := x × yy + xx × y + cc;
    z := c + cc;
    zz := c - z + cc
end mul 2;
```

comment div 2 calculates the doublelength quotient of (x,xx) and (y,yy), the result being (z,zz). If y = 0, the effect of this procedure is undefined;

```
procedure div 2 (x,xx, y,yy, z,zz);
value x,xx, y,yy; real x,xx, y,yy, z,zz;
begin real c,cc, u,uu;
      c := x/y;
      mul 12 (c, y, u,uu);
      cc := (x-u-uu+xx-c*yy)/y;
      z := c + cc;
      zz := c - z + cc
end div 2;
```

comment sqrt 2 calculates the doublelength square root of (x,xx), the result being (y,yy). If (x,xx) is not positive, then the result equals (0,0);

```
procedure sqrt 2 (x,xx, y,yy);
value x,xx; real x,xx, y,yy;
begin real c,cc, u,uu;
      if x > 0 then
        begin c := sqrt (x);
          mul 12 (c,c, u,uu);
          cc := (x-u-uu+xx) × 0.5/c;
          y := c + cc;
          yy := c - y + cc
        end
      else y := yy := 0
end sqrt 2.
```

## 9. Acknowledgements

The author expresses his gratitude to Professor Dr. F.E.J. Kruseman Aretz, formerly at Mathematical Centre, Amsterdam, to Professor Dr. G.W. Veltkamp at Technological University, Eindhoven, and to Dr. W.S. Brown and Dr. P.L. Richman at Bell Telephone Laboratories, Murray Hill, for many valuable suggestions and inspiring discussions. The author is also grateful to Mr. H.J.W. ten Hagen, Mr. H.N. Glorie and Mr. D.F. Winter at Mathematical Centre, Amsterdam, for their assistance in testing the algorithms.

## 10. References

- I. Babuška (1968), Numerical stability in mathematical analysis, IFIP congress 68, Invited papers, p. 1-13.
- A.A. Grau (1962), On a floating-point number representation for use with algorithmic languages, Comm. ACM 5, p. 160-161.
- W. Kahan (1965), Further remarks on reducing truncation errors, Comm. ACM 8, p. 40.
- D.E. Knuth (1969), The art of computer programming, vol. 2 (Addison Wesley).
- O. Møller (1965), Quasi double-precision in floating-point addition, BIT 5, p. 37-50.
- P. Naur, ed. (1962), Revised report on the algorithmic language ALGOL 60.
- G.W. Veltkamp (1968), Private communications (see also RC Informatie nr. 21 & 22, Technological University, Eindhoven).
- J.H. Wilkinson (1963), Rounding errors in algebraic processes (Her Majesty's Stationary Office).

T.J. Dekker  
Mathematical Centre  
Amsterdam  
Netherlands.