

# A priori Worst-Case Error Bounds for Floating-Point Computations

Walter Krämer  
Institut für Wissenschaftliches Rechnen  
und Mathematische Modellbildung (IWRMM)  
Universität Karlsruhe  
D-76128 Karlsruhe, Germany

## Abstract

*A new technique for the a priori calculation of rigorous error bounds for floating-point computations is introduced. The theorems given in the paper combined with interval arithmetic lead to the implementation of reliable software-routines, which enables the user to compute the desired error bounds automatically by a suitable computer program.*

*As a prominent example a table-lookup algorithm for calculating the function  $\exp(x) - 1$  that has been published by Tang [15] is analyzed using these new tools. The result shows the high quality of the new approach.*

## 1. Introduction

It is a well known fact, that the results of floating-point computations may be totally wrong: For the numerical solution of the following simple linear system of equations IEEE double arithmetic [4, 8] is used (64 bit data format, best possible (semimorphic) floating-point operations). The components of the true solution of the system  $Ax = b$  with

$$A = (a_{ij}) := \begin{pmatrix} 64919121 & -159018721 \\ 41869520.5 & -102558961 \end{pmatrix}, b := \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

are given by  $x_1 = a_{22}/(a_{11}a_{22} - a_{12}a_{21})$ ,  $x_2 = -a_{21}/(a_{11}a_{22} - a_{12}a_{21})$ . The corresponding floating-point computations result in the machine numbers  $\tilde{x}_1 = 102558961.0$ ,  $\tilde{x}_2 = 41869520.5$ . All computed digits are wrong! Although the data are exactly representable and only 5 floating-point operations are done, the computed "solution" has nothing to do with the exact solution<sup>1</sup>.

Such "numerical catastrophes" are extremely rare in practice, but to cite Kahan [7]:

Significant discrepancies [between the computed and the true result] are very rare, too rare to worry about all the time, yet not rare enough to ignore.

<sup>1</sup>The solution of the system is  $x_1 = 205117922$ ,  $x_2 = 83739041$

In the present paper a methodology is described, which allows to decide whether the floating-point computation of a given real valued expression delivers an accurate result. Here, the independent variables as well as the parameters of the expression may vary within given intervals (the width of which may be large). Then the conclusion drawn from the application of the new error-analyzing tools will be valid for floating-point evaluations of the expression with arbitrary point data for the independent variables and the parameters lying in their corresponding intervals.

Moreover, the tools which are discussed enable the user to analyze complete program parts including loops, iterations, and recursion. Also intermediate results in the underflow range are treated in a reliable manner.

As an example the result of a rigorous worst-case error estimation for an implementation of an accurate table-lookup algorithm for  $\exp(x)-1$  described by Tang [15] is shown. Finding a reliable worst case error bound for such an algorithm by hand is very cumbersome and error prone. A lot of different cases have to be analyzed separately. One must be sure that all different branches are considered when performing the error estimation<sup>2</sup>.

## 2. Some Notations and Conventions

$S = S(b, l, \underline{e}, \bar{e})$  denotes a floating-point system with base  $b$ , with  $l$  mantissa digits, and exponent  $e$  in the range  $\underline{e} \leq e \leq \bar{e}$ .  $S(2, 53, -1022, 1023)$  corresponds to the IEEE double data format (64 bit).  $\circ \in \{+, -, \cdot, /\}$  denotes a real operation and  $\boxtimes, \circ \in \{+, -, \cdot, /\}$  the corresponding floating-point operation (round to nearest). Notice the implicit priority rules used in the notation

<sup>2</sup>In [13] Tang proposed a table-lookup algorithm for logarithms. The user is supplied with a very small over all error bound. Nevertheless, some constants in the table shown in the cited paper are wrong and must be corrected! Apparently this fact has not been detected by the error analysis done by hand. (The correct values for the constants  $C_{\text{lead}}(j)$ ,  $j=17, 57$ , and  $73$  are 3FBF EC91 31DC 0000, 3FD7 92A5 5FDD 8000, and 3FDC E1AF 0B86 0000, respectively. Here, the same notation as in [13], pages 397 and 398 is used.)

$|a \circ b - a \boxplus b| := |(a \circ b) - (a \boxplus b)|$ .  $\text{MinReal}$  is the smallest positive normalized floating-point number,  $\text{MaxReal}$  the largest floating-point number. For the IEEE double format these numbers are  $\text{MinReal} := 2.22 \dots \cdot 10^{-308}$ ,  $\text{MaxReal} := 1.78 \dots \cdot 10^{308}$ .  $\tilde{a}$ ,  $\tilde{b}$  are quantities computed in floating-point arithmetic,  $\varepsilon^*$  means Wilkinson's epsilon  $\varepsilon^* := \frac{1}{2}b^{1-l}$ , i. e.  $\varepsilon^* = \frac{1}{2}2^{1-53} = 1.11022 \dots \cdot 10^{-16}$  for the IEEE double format.  $\mathbb{R}^+$  denotes the set of all positive real numbers,  $IR$  the set of all closed real (finite) intervals,  $IS = \{[\underline{a}, \bar{a}] | \underline{a}, \bar{a} \in S, \underline{a} \leq \bar{a}\}$  the set of all intervals bounded by floating-point numbers,  $|A| := \max_{a \in A} |a|$  the maximum absolute value,  $\langle A \rangle := \min_{a \in A} |a|$  the minimum absolute value and  $\text{diam}(A) := \sup(A) - \inf(A)$  the diameter (width) of an interval  $A \in IR$ .

### 3. Theoretical Results

Let  $f$  denote an expression and  $f(a)$  its evaluation at the point  $a \in A \in IR$ . Let  $\tilde{f}(\tilde{a})$  be the result of the corresponding floating-point computations. One is interested in an a priori error bound  $\Delta(f)$  for the maximal absolute error

$$|f(a) - \tilde{f}(\tilde{a})| \leq \Delta(f)$$

valid for all  $a \in A$  and all  $\tilde{a} = a + \Delta_a$  with  $\Delta_a \leq \Delta(a)$ . Here the interval  $A$  and the error bound  $\Delta(a)$  for  $\tilde{a}$  are assumed to be known.

To achieve this, the original expression is broken down into basic operations like  $+$ ,  $-$ ,  $\cdot$ ,  $/$  and  $\sqrt{x}$  for which the following theorems will give reliable worst case error bounds. These operations are used as basics because they are required by the IEEE Standard for Binary Floating-Point Arithmetic [4]. For the fundamental operations  $\circ \in \{+, -, \cdot, /\}$  this standard guarantees:

$$\bigwedge_{\substack{a, b \in S \text{ with} \\ |a \circ b| \in [\text{MinReal}, \text{MaxReal}]}} \left| \frac{a \circ b - a \boxplus b}{a \circ b} \right| \leq \varepsilon^* .$$

However, the methodology is not restricted to this set of basic operations. It can be extended, for example, to a complete set of operations (functions like  $\sin$ ,  $\cos$  etc.) as it is commonly available in a higher programming language like FORTRAN or C (see Section 6).

It should be emphasized that in contrast to conventional error estimations also higher order error terms and underflow situations are treated rigorously by the following theorems. Using an interval arithmetic as well as directed rounded operations makes it possible to build subroutines preserving these properties. Such software tools are supplied by the PASCAL-XSC module `eps.ari` (see Section 4).

The following lemma bounds the maximal absolute error for results lying in the underflow range.

**Lemma 1 (Underflow Range)** For results lying in the underflow range it holds

$$\bigwedge_{\substack{a, b \in S \\ |a \circ b| < \text{MinReal}}} |a \boxplus b - a \circ b| \leq \text{MinReal} . \quad (1)$$

**Proof:**  $\text{sign}(a \boxplus b) = \text{sign}(a \circ b)$  and  $\text{diam}([0, \text{MinReal}]) = \text{MinReal} \square$ .

What follows is a discussion of the error propagation for the basic binary operations with arguments that are afflicted by rounding errors.

Let  $\circ \in \{+, -, \cdot, /\}$  and  $a \in A \in IR$ ,  $b \in B \in IR$ ,

$$S \ni \tilde{a} = a + \Delta_a \text{ with } |\Delta_a| \leq \Delta(a),$$

$$S \ni \tilde{b} = b + \Delta_b \text{ with } |\Delta_b| \leq \Delta(b).$$

It is assumed that the quantities  $A, B, \Delta(a)$  and  $\Delta(b)$  are given. The goal is to find a bound  $\Delta(\circ) \in \mathbb{R}^+$  with

$$\bigwedge_{\substack{a \in A, b \in B \\ a \circ b, \tilde{a} \circ \tilde{b} \in [-\text{MaxReal}, \text{MaxReal}]}} \bigwedge_{\substack{|a - \tilde{a}| \leq \Delta(a) \\ |b - \tilde{b}| \leq \Delta(b)}} |a \circ b - \tilde{a} \circ \tilde{b}| \leq \Delta(\circ)$$

$$|a \circ b - \tilde{a} \circ \tilde{b}| \leq \Delta(\circ) .$$

The case of addition is treated in

**Theorem 1 (Addition)** For the addition  $\circ := +$  it holds:

$$\bigwedge_{\substack{a \in A \\ b \in B}} \bigwedge_{\substack{|a - \tilde{a}| \leq \Delta(a) \\ |b - \tilde{b}| \leq \Delta(b)}} |a + b - \tilde{a} \boxplus \tilde{b}| \leq \Delta(\text{add})$$

with

$$\Delta(\text{add}) := \text{MinReal} + \varepsilon^* |A + B| + (1 + \varepsilon^*) (\Delta(a) + \Delta(b))$$

**Proof:** Let  $a \in A, b \in B$ ,

$\tilde{a} = a + \Delta_a \in A + [-\Delta(a), \Delta(a)]$ ,  $|\Delta_a| \leq \Delta(a)$  and  $\tilde{b} = b + \Delta_b \in B + [-\Delta(b), \Delta(b)]$ ,  $|\Delta_b| \leq \Delta(b)$  arbitrary but fixed.

I) Error bound for  $|\tilde{a} + \tilde{b} - \tilde{a} \boxplus \tilde{b}|$ :

a)  $\tilde{a} + \tilde{b} \in U$ , the exact sum of the disturbed arguments lies in the underflow range:

$$\tilde{a} + \tilde{b} \in U \implies \tilde{a} \boxplus \tilde{b} \in U$$

$$\stackrel{\text{Lemma 1}}{\implies} |\tilde{a} + \tilde{b} - \tilde{a} \boxplus \tilde{b}| \leq \text{MinReal}$$

b)  $\tilde{a} + \tilde{b} \notin U$ , i.e. the exact sum  $|\tilde{a} + \tilde{b}|$  is an element of  $[\text{MinReal}, \text{MaxReal}]$ :

$$\tilde{a} + \tilde{b} \notin U \implies \tilde{a} \boxplus \tilde{b} \notin U$$

$$\implies |\tilde{a} + \tilde{b} - \tilde{a} \boxplus \tilde{b}|$$

$$\leq \varepsilon^* |\tilde{a} + \tilde{b}|$$

$$\leq \varepsilon^* (|A + B| + \Delta(a) + \Delta(b)).$$

This holds, because from  $\tilde{a} \in A + [-\Delta(a), \Delta(a)]$  it follows  $|\tilde{a}| \leq |A| + \Delta(a)$ .

II) Error bound for  $|a + b - (\tilde{a} + \tilde{b})|$ :

It holds:

$$|a + b - (\tilde{a} + \tilde{b})|$$

$$\leq |a - \tilde{a}| + |b - \tilde{b}|$$

$$\leq \Delta(a) + \Delta(b)$$

Over all error bound:

$$|a + b - \tilde{a} \boxplus \tilde{b}|$$

$$\leq |a + b - (\tilde{a} + \tilde{b})| + |\tilde{a} + \tilde{b} - \tilde{a} \boxplus \tilde{b}|$$

$$\stackrel{\text{I),II)}}{\leq} \Delta(a) + \Delta(b) +$$

$$\begin{cases} \text{MinReal}, & \text{Case I) a)} \\ \varepsilon^* (|A + B| + \Delta(a) + \Delta(b)), & \text{Case I) b)} \end{cases}$$

$\leq$

$$\begin{cases} \Delta(a) + \Delta(b) + \text{MinReal}, & \text{Case I) a)} \\ \varepsilon^* |A + B| + (1 + \varepsilon^*) (\Delta(a) + \Delta(b)), & \text{Case I) b)} \end{cases}$$

$$\leq \varepsilon^* |A + B| + (1 + \varepsilon^*) (\Delta(a) + \Delta(b)) + \text{MinReal}$$

$$= \Delta(\text{add}).$$

The quantities  $a, b, \tilde{a}, \tilde{b}$  have been chosen arbitrarily. This proves Theorem 1  $\square$

**Theorem 2 (Subtraction)** The error bound  $\Delta(\text{sub})$  for the subtraction is very similar to the error bound  $\Delta(\text{add})$  for the addition. It holds:

$$\Delta(\text{sub}) := \text{MinReal} +$$

$$\varepsilon^* |A - B| + (1 + \varepsilon^*) (\Delta(a) + \Delta(b))$$

**Proof:** Use the representation  $a - b = a + (-b)$  and the

proof of Theorem 1  $\square$

**Theorem 3 (Multiplication)** For the multiplication  $\circ := \bullet$  it holds:

$$\bigwedge_{\substack{a \in A \\ b \in B}} \bigwedge_{\substack{|a - \tilde{a}| \leq \Delta(a) \\ |b - \tilde{b}| \leq \Delta(b)}} |a \cdot b - \tilde{a} \boxtimes \tilde{b}| \leq \Delta(\text{mul})$$

with

$$\Delta(\text{mul}) := \text{MinReal} + |A||B|\varepsilon^* +$$

$$(|A|\Delta(b) + |B|\Delta(a) + \Delta(a)\Delta(b)) \cdot (1 + \varepsilon^*)$$

**Proof:** Let  $a \in A, b \in B$ ,

$\tilde{a} = a + \Delta_a \in A + [-\Delta(a), \Delta(a)]$ ,  $|\Delta_a| \leq \Delta(a)$  and  $\tilde{b} = b + \Delta_b \in B + [-\Delta(b), \Delta(b)]$ ,  $|\Delta_b| \leq \Delta(b)$  arbitrary but fixed.

I) Estimation for  $|\tilde{a} \cdot \tilde{b} - \tilde{a} \boxtimes \tilde{b}|$ :

$$\text{a) } \tilde{a} \cdot \tilde{b} \in U \implies \tilde{a} \boxtimes \tilde{b} \in U$$

$$\stackrel{\text{Lemma 1}}{\implies} |\tilde{a} \cdot \tilde{b} - \tilde{a} \boxtimes \tilde{b}| \leq \text{MinReal}$$

$$\text{b) } \tilde{a} \cdot \tilde{b} \notin U \implies |\tilde{a} \boxtimes \tilde{b}| \notin U$$

$$\implies |\tilde{a} \cdot \tilde{b} - \tilde{a} \boxtimes \tilde{b}| \leq \varepsilon^* |\tilde{a} \cdot \tilde{b}|$$

$$\leq \varepsilon^* (|A| + \Delta(a)) (|B| + \Delta(b))$$

$$= \varepsilon^* (|A||B| + |A|\Delta(b) + |B|\Delta(a) + \Delta(a)\Delta(b)).$$

II) Error bound for  $|a \cdot b - \tilde{a} \cdot \tilde{b}|$ :

$$|a \cdot b - \tilde{a} \cdot \tilde{b}|$$

$$= |a \cdot b - (a + \Delta_a)(b + \Delta_b)|$$

$$= |a \Delta_b + b \Delta_a + \Delta_a \Delta_b|$$

$$\leq |A|\Delta(b) + |B|\Delta(a) + \Delta(a)\Delta(b).$$

Over all error bound:

$$|a \cdot b - \tilde{a} \boxtimes \tilde{b}| \leq |a \cdot b - \tilde{a} \cdot \tilde{b}| + |\tilde{a} \cdot \tilde{b} - \tilde{a} \boxtimes \tilde{b}|$$

$$\stackrel{\text{I),II)}}{\leq}$$

$$\begin{cases} \text{MinReal} \\ + |A|\Delta(b) + |B|\Delta(a) + \Delta(a)\Delta(b), & \text{I) a)} \\ (|A|\Delta(b) + |B|\Delta(a) + \Delta(a)\Delta(b))(1 + \varepsilon^*) \\ + |A||B|\varepsilon^*, & \text{I) b)} \end{cases}$$

$$\leq \text{MinReal} + |A||B|\varepsilon^* + (|A|\Delta(b) + |B|\Delta(a) + \Delta(a)\Delta(b))(1 + \varepsilon^*)$$

$$= \Delta(\text{mul}) \quad \square$$

To prove Theorem 4 below, the following lemma will be useful:

**Lemma 2 (Inverse)**

$$|\eta| \leq \bigwedge_{\eta^* < 0.5} \left| \frac{1}{1+\eta} \right| \leq 1 + \varepsilon(\text{inv})$$

with  $\varepsilon(\text{inv}) := (1 + 2\eta^*) \cdot \eta^*$ .

**Proof:** With  $|\eta| \leq \eta^* < 0.5$  one gets

$$1 - \eta^* \leq \frac{1}{1+\eta} \leq 1 + (1 + 2\eta^*) \cdot \eta^* \quad \square$$

**Theorem 4 (Division)** In case of the division  $\circ := /$  it holds:

$$\bigwedge_{\substack{a \in A \\ b \in B}} \bigwedge_{\substack{|a - \tilde{a}| \leq \Delta(a) \\ |b - \tilde{b}| \leq \Delta(b)}} \left| \frac{a}{b} - \tilde{a} \upharpoonright \tilde{b} \right| \leq \Delta(\text{div}) \quad \text{with}$$

$$\Delta(\text{div}) := \text{MinReal} + \frac{1}{\langle B \rangle - \Delta(b)}$$

$$\left( \Delta(a) + (|A| + \Delta(a)) \cdot \left( \varepsilon^* + \frac{\Delta(b)}{\langle B \rangle} + 2 \left( \frac{\Delta(b)}{\langle B \rangle} \right)^2 \right) \right)$$

Here we have to make the additional assumption

$$\Delta(b) < 0.5 \cdot \langle B \rangle \quad (2)$$

Assumption (2) means, that the interval  $B$  is sufficiently far away from zero. In practice this is not a serious restriction.

**Proof:** Let  $a \in A, b \in B,$

$\tilde{a} = a + \Delta_a \in A + [-\Delta(a), \Delta(a)], \Delta_a \leq \Delta(a)$  and  $\tilde{b} = b + \Delta_b \in B + [-\Delta(b), \Delta(b)], \Delta_b \leq \Delta(b)$  arbitrary but fixed.

I) Estimation for  $|\tilde{a}/\tilde{b} - \tilde{a} \upharpoonright \tilde{b}|$ :

$$\text{a) } \tilde{a}/\tilde{b} \in U \implies \tilde{a} \upharpoonright \tilde{b} \in U$$

$$\stackrel{\text{Lemma 1}}{\implies} |\tilde{a}/\tilde{b} - \tilde{a} \upharpoonright \tilde{b}| \leq \text{MinReal}$$

$$\text{b) } \tilde{a}/\tilde{b} \notin U \implies |\tilde{a} \upharpoonright \tilde{b}| \notin U$$

$$\implies |\tilde{a}/\tilde{b} - \tilde{a} \upharpoonright \tilde{b}| \leq \varepsilon^* |\tilde{a}/\tilde{b}|$$

$$\leq \varepsilon^* (|A| + \Delta(a)) \cdot \frac{1}{\langle B \rangle - \Delta(b)}$$

Here it is used that  $\tilde{b} \in B + [-\Delta(b), \Delta(b)]$  yielding

$$\frac{1}{|\tilde{b}|} \leq \frac{1}{\langle B \rangle - \Delta(b)}$$

II) Error bound for  $|a/b - \tilde{a}/\tilde{b}|$ :

$$|a/b - \tilde{a}/\tilde{b}| = \left| \frac{a}{b} - \frac{(a + \Delta_a)}{b} \cdot \frac{1}{1 + \frac{\Delta_b}{b}} \right|$$

$$= \left| \frac{a}{b} - \frac{(a + \Delta_a)}{b} \cdot (1 + \varepsilon_{\text{inv}}) \right|$$

$$= \frac{1}{|b|} |\Delta_a + (a + \Delta_a) \cdot \varepsilon_{\text{inv}}|$$

$$\stackrel{(*)}{\leq} \frac{1}{\langle B \rangle} \left( \Delta(a) + (|A| + \Delta(a)) \cdot \varepsilon(\text{inv}) \right)$$

with  $|\varepsilon_{\text{inv}}| \leq \varepsilon(\text{inv}) := (1 + 2 \cdot \frac{\Delta(b)}{\langle B \rangle}) \cdot \frac{\Delta(b)}{\langle B \rangle}$ . The last inequality (\*) results from Lemma 2 and assumption (2).

Over all error bound:

$$|a/b - \tilde{a} \upharpoonright \tilde{b}| \leq |a/b - \tilde{a}/\tilde{b}| + |\tilde{a}/\tilde{b} - \tilde{a} \upharpoonright \tilde{b}|$$

$$\stackrel{\text{I),II}}{\leq} \frac{1}{\langle B \rangle} \left( \Delta(a) + (|A| + \Delta(a)) \cdot \varepsilon(\text{inv}) \right) +$$

$$\begin{cases} \text{MinReal,} & \text{Case I) a)} \\ \varepsilon^* (|A| + \Delta(a)) \cdot \frac{1}{\langle B \rangle - \Delta(b)}, & \text{Case I) b)} \end{cases}$$

$$\leq \text{MinReal} + \frac{1}{\langle B \rangle - \Delta(b)}$$

$$\left( \Delta(a) + (|A| + \Delta(a)) \cdot (\varepsilon^* + \varepsilon(\text{inv})) \right)$$

$$= \Delta(\text{div}) \quad \square$$

In the following theorem the square root function (a unary operation) is considered.

**Theorem 5 (Square root)** Let the auxiliary function  $h(a)$  be defined by

$$h(a) = \varepsilon^* \sqrt{a} + \frac{\Delta(a)(1 + \varepsilon^*)}{2\sqrt{a - \Delta(a)}} \quad (3)$$

with  $a \in A$  and  $a > \Delta(a)$ . Then it holds:

$$\bigwedge_{\Delta(a) < a \in A} \bigwedge_{|a - \tilde{a}| \leq \Delta(a)} \left| \sqrt{a} - \sqrt{\tilde{a}} \right| \leq \Delta(\text{sqrt})$$

where

$$\Delta(\text{sqrt}) := \max \{ h(\inf(A)), h(\sup(A)) \}$$

**Proof:** Let  $a \in A, \tilde{a} = a + \Delta_a$ , with  $|\Delta_a| \leq \Delta(A)$  arbitrary but fixed. Using the Mean-Value Theorem one finds the following representation

$$\sqrt{\tilde{a}} = \sqrt{a + \Delta_a} = \sqrt{a + \Delta_a} (1 + \varepsilon_{\sqrt{\cdot}})$$

$$= \left( \sqrt{a} + \frac{\Delta_a}{2\sqrt{\xi}} \right) (1 + \varepsilon_{\sqrt{\cdot}})$$

with

$$\xi = \xi(a, \tilde{a}) \in [\min\{a, \tilde{a}\}, \max\{a, \tilde{a}\}] \subseteq A + [-\Delta(a), \Delta(a)].$$

The square root function is assumed to conform to the IEEE-754 standard, i.e.  $|\varepsilon_{\sqrt{\cdot}}| \leq \varepsilon^*$ , yielding

$$\left| \sqrt{a} - \sqrt{\tilde{a}} \right| \leq h(a).$$

Now the auxiliary function  $h$  has exactly one minimum in the range  $(0, \infty)$ , i. e.  $\max\{h(x) | x \in A\} = \max\{h(\inf(A)), h(\sup(A))\}$   $\square$

#### 4. Usage of the PASCAL-XSC Module eps\_ari

The theorems of the preceding section are coded in PASCAL-XSC [9] using interval operations. Of course any other programming language that provides an interval arithmetic can be used. Five functions named DeltaAdd, DeltaSub, DeltaMul, DeltaDiv, and DeltaSqrt, respectively are supplied by the module eps\_ari. Calling the function DeltaAdd in the following assignment

```
DeltaRes:=
  DeltaAdd( A, B, DeltaA, DeltaB );
```

gives as result a machine number  $\text{DeltaRes} \in S$  with

$$\bigwedge_{\substack{a \in A \\ b \in B}} \bigwedge_{\substack{|a - \tilde{a}| \leq \Delta(a) \\ |b - \tilde{b}| \leq \Delta(b)}} \left| a + b - \tilde{a} \boxplus \tilde{b} \right| \leq \Delta(\text{add})$$

$$\leq \text{DeltaRes}$$

(see Theorem 1). Here  $A \subseteq \mathbb{A} \in IS$ ,  $B \subseteq \mathbb{B} \in IS$ ,  $\Delta(a) \leq \text{DeltaA} \in S$ , and  $\Delta(b) \leq \text{DeltaB} \in S$ . The basic function DeltaAdd supplied by the module eps\_ari looks like

```
global const Epsilon = 1.110224E-16;
{ Upper bound for 2**(-53) which is }
{ Wilkinson's epsilon for          }
{ the IEEE double format           }

global const
  MinReal = 2.2250738585072013E-308;
{ Smallest positive normalized      }
{ floating-point number            }

global function DeltaAdd(
  A, B          : interval;
  DeltaA, DeltaB: real      ): real;
var u, v: real;
begin
  if { ... some special cases ... }
  ...
```

```
else begin      { See Theorem 1 }
  u:= Epsilon *> MaxAbs( A + B );
  v:= (DeltaA +> DeltaB)*>(1 +> Epsilon);
  DeltaAdd:= MinReal +> u +> v;
end;
end;
```

The symbols +> and \*> denote the machine operations addition and multiplication with rounding towards +∞. The operator + in the expression A + B means interval addition. All these operators are predefined in PASCAL-XSC.

The following simple example shows, how the basic routines DeltaAdd, DeltaSub, ... can be combined to compute a worst-case error bound for the evaluation of the polynomial  $p(x) = \sum_{i=0}^n p_i x^i$ . As usual the computation is done according to Horner's scheme. It is assumed that  $x \in X \in IS$ ,  $\tilde{x} = x + \Delta_x$  with  $|\Delta_x| \leq \Delta(x) \leq \text{DeltaX} \in S$ ,  $p_i, |\tilde{p}_i| \in p[i] \in IS$ , where X, DeltaX, and p[i] are known quantities. Of course  $p_i, \tilde{p}_i \in p[i]$  implies  $|\tilde{p}_i - p_i| \leq \text{diam}(p[i])$ , i. e.  $\Delta(p_i) \leq \text{diam}(p[i])$ .

```
H:= p[n]; DeltaH:= diam( p[n] );
for i:= n-1 downto 0 do begin
  DeltaH:=
    DeltaMul(H, X, DeltaH, DeltaX);
  H:= H * X;
  DeltaH:=
    DeltaAdd(H, p[i], DeltaH, diam(p[i]));
  H:= H + p[i];
end;
PolX := H;
AbsErr:= DeltaH;
```

Now the inequality

$$|p(x) - \tilde{p}(\tilde{x})| \leq \text{AbsErr}$$

holds for all  $x \in X$  and all  $\tilde{x} = x + \Delta_x$  with  $|\Delta_x| \leq \text{DeltaX}$ . Moreover  $p(x) \in \text{PolX}$  for  $x \in X$  and  $\text{PolX} \ni \tilde{p}(x) := (\dots((\tilde{p}_n \boxtimes x \boxplus \tilde{p}_{n-1}) \boxtimes x \boxplus \tilde{p}_{n-2}) \boxtimes x \boxplus \dots \boxplus \tilde{p}_1) \boxtimes x \boxplus \tilde{p}_0$  for any set of polynomial coefficients  $\tilde{p}_i \in p[i], i = 0(1)n$ . Examples producing numerical results are given in [5].

#### 5. A priori Error Bound for an Accurate Table Method Realizing $e^x - 1$

To show that the methods discussed in this paper are relevant for practical purposes an error estimation for an algorithm implemented in ANSI-C which computes the function  $e^x - 1$  is performed. It is assumed that the implementation of  $\exp(x) - 1$  uses operations which comply with the IEEE standard [4].

**Some remarks on the table-lookup algorithm:** Only some parts of the algorithm are described informally. This will

be sufficient to show that an error estimation by hand is very expensive (many different cases) and error prone. The detailed description as well as the complete ANSI-C source code can be found in [5]. The complete algorithm is described in [5, 15].

In the range  $x \leq \ln(1 - \frac{1}{4})$  or  $\ln(1 + \frac{1}{4}) \leq x$  one computes the value  $e^x - 1$  in three steps. First the argument is reduced to the interval  $[-\ln(2)/64, \ln(2)/64]$  using the formula

$$x = (32m + j) \cdot \frac{\ln(2)}{32} + r$$

with  $m \in \mathbf{Z}$ ,  $j \in \{0, 1, \dots, 31\}$  and  $|r| \leq \ln(2)/64$ . Here  $m$ ,  $j$  are determined and then  $r$  is computed simulating a higher precision arithmetic (staggered arithmetic). Second, for the reduced argument  $r$  an approximation to  $e^r - 1$  is computed by the polynomial approximation  $e^r - 1 \approx p(r) := r + r^2 \cdot (a_0 + a_1r + a_2r^2 + a_3r^3 + a_4r^4)$ . Third, the reconstruction of the result is done in accordance to the formula

$$e^x - 1 = 2^m(2^{j/32} + 2^{j/32}(e^r - 1)) - 1.$$

The powers  $2^{j/32}$ ,  $j \in \{0, 1, 2, \dots, 31\}$  are precomputed values and stored in a table (table method) using two double numbers for each entry.

**Automatic error analysis:** Applying the routines of the module `eps_ari` gives a worst case error bound for the algorithm described above. To avoid overestimations in interval calculations, i.e. to get sharp error bounds, the argument range has to be subdivided into thousands of small subintervals (which in fact can be done using an additional simple loop statement). Again, the complete program for the automatic error estimation can be found in [5].

With the abbreviation  $f(x) := \exp(x) - 1$  one finds the following over all upper error bound  $\varepsilon(f)$  for the relative error of the accurate table method<sup>3</sup>:

$$\bigwedge_{\substack{0 \neq x \in S \\ x \text{ causes no overflow}}} \left| \frac{(\exp(x) - 1) - \widetilde{f(x)}}{\exp(x) - 1} \right| \leq \varepsilon(f) := 1.302 \cdot 10^{-16} \leq 1.18 \cdot \varepsilon^*. \quad (4)$$

Notice, that this bound holds simultaneously for **all** floating-point arguments which do not cause an overflow error. It must be emphasized that, in general, already **one** simple floating-point operation leads to the bound  $1.0 \cdot \varepsilon^*$ . This demonstrates the high quality of the factor 1.18 in (4). Indeed, as already stated by Tang [15], the algorithm leads to an implementation with near-perfect accuracy. In comparison with the result (4) Tang's original analysis gives a slightly better error bound in ulps. But see Footnote 2.

<sup>3</sup>The approximation errors of the different polynomial approximations have been estimated rigorously using an interval tool described in [11].

## 6. Conclusions

The new technique can be used to get tight and reliable worst-case error bounds for floating-point computations (almost) automatically. In doing a priori error estimations a lot of error prone hand calculations can be omitted. The described implementation of the "absolute error arithmetic" using interval operations leads to bounds that are rigorous, realistic, and of high quality.

Some improvements and extensions are planned:

- Theoretical foundation and implementation of functions like `DeltaSin`, `DeltaArcSin`, `DeltaCos`, `DeltaArcCos`, `DeltaExp`, `DeltaLog`, ... corresponding to the functions `sin`, `arcsin`, `cos`, `arccos`, `exp`, `log` ... commonly supplied by higher programming languages<sup>4</sup>.
- Theoretical foundation and realization of a worst-case error arithmetic for *relative* error bounds.
- A C++ implementation using an object oriented approach will simplify the usage of the described error arithmetic tools even more. Such an approach can also be done in PASCAL-XSC. Introducing a new data type, e.g.

```
eps_ari_type =
  record
    Enclosure : interval;
    ErrorBound: real;
  end;
```

allows the overloading of operators and functions for this type. If  $A$  denotes a variable of type `eps_ari_type` then the record component `A.Enclosure` represents an enclosure (floating-point interval) of  $A \ni a$  and the floating-point number `A.ErrorBound` is an upper bound for  $|a - \tilde{a}|$ , i. e. for  $\Delta(a)$ .

For quantities of the new data type the routines `DeltaAdd`, `DeltaSub`, ... can be replaced by operator calls `+`, `-`, ..., respectively. This preserves the mathematical notation of expressions when using the error bound arithmetic. Horner's Scheme example from above could be simplified to (H and X have to be of type `eps_ari_type`):

```
{ Initialization: }
H.Enclosure := p[n];
```

<sup>4</sup>The numerical quality of the functions as they are supplied by the compiler or run-time library must be known (for example in the form: their accuracy is better than  $k$  units in the last place)

```

H.ErrorBound:= diam( p[n] );
{ Horner's scheme:      }
for i:= n-1 downto 0 do begin
    H:= H * X + p[i];
end;
{ Range of the polynomial }
{ over X.Enclosure:      }
PolX:= H.Enclosure;
{ Worst-case error bound: }
AbsErr:= H.ErrorBound;

```

The presented error arithmetic can also be used to control the numerical error of multi-precision computations [10]. If, for example, error bounds are known for the steps of an iterative process done in the field of the real numbers, then the additional numerical errors introduced by the multi-precision computations can be estimated rigorously in advance. Statements like "An arithmetic using  $(k + n)$ -digits with respect to base  $b$  is sufficient to get a relative accuracy better than  $b^{-k}$ " are possible. Such statements are important for reliable multi-precision calculations for example in the field of computer algebra.

The new tools are also of great advantage when constructing algorithms according to Ziv [16].

Future work will show whether the presented concept of a reliable error arithmetic will give tight results in situations as discussed in [2, 3]<sup>5</sup>.

For the implementation of the new tools (machine) interval operations are of great advantage. At least directed rounded operations must be available. It is a great pity that such features of IEEE 754 lack support in programming languages and compilers (see the discussion in [8]). In the so called XSC-languages [9] interval operations are available via their usual mathematical operator notations. These languages also support directed rounded operations.

## References

- [1] Adams, E., Kulisch, U.: *Scientific Computing with Automatic Result Verification*, Mathematics in Science and Engineering, Vol. 189, Academic Press, 1993.
- [2] W. E. Ferguson, T. Brightman: *Accurate and Monotone Approximations of Some Transcendental Functions*, Proceedings to the 10th Symposium on Computer Arithmetic, IEEE Computer Society Press, pp. 237–244, 1991.
- [3] W. E. Ferguson: *Exact Computation of a Sum or Difference with Applications to Argument Reduction*, Proceedings to the 12th Symposium on Computer Arithmetic, IEEE Computer Society Press, pp. 216–221, 1995.
- [4] American National Standards Institute / Institute of Electrical and Electronics Engineers: *A Standard for Binary Floating-Point Arithmetic*. ANSI/IEEE Std. 754-1985, New York, 1985 (reprinted in SIGPLAN 22, 2, pp 9–25, 1987).
- [5] Hofschuster, W., Krämer, W.: *Ein rechnergestützter Fehlerkalkül mit Anwendung auf ein genaues Tabellenverfahren*, Preprint 96/5 des Instituts für Wissenschaftliches Rechnen und Mathematische Modellbildung, 1996:  
ftp://iamk4515.mathematik.uni-karlsruhe.de,  
Directory: /pub/iwrmm/preprints
- [6] Hofschuster, W. and Krämer, W.: *A Computer Oriented Approach to Get Sharp Reliable Error Bounds*, Reliable Computing, Issue 3, Volume 3, 1997.
- [7] Kahan, W. M.: *The Regrettable Failure of Automated Error Analysis*. A Mini-Course prepared for the conference at MIT on Computers and Mathematics, 1989.
- [8] Kahan, W. M.: *Lecture Notes on the Status of IEEE Standard 754 for Binary Floating-Point Arithmetic*, June 5, 1995.
- [9] Klatt, R., Kulisch, U., Neage, M., Ratz, D., Ullrich, Ch.: *PASCAL-XSC, Language Reference with Examples*, Springer, 1992.
- [10] Krämer, W.: *Multiple-Precision Computations with Result Verification*, in: *Scientific Computing with Automatic Result Verification* ed. by E. Adams and U. Kulisch, Academic Press, 1992.
- [11] Krämer, W.: *Sichere und genaue Abschätzung des Approximationsfehlers bei rationalen Approximationen*, Bericht des Instituts für Angewandte Mathematik, Universität Karlsruhe, 1996:  
ftp://iamk4515.mathematik.uni-karlsruhe.de,  
Directory: /pub/documents/reports
- [12] Tang, P.T.P.: *Table-Driven Implementation of the Exponential Function in IEEE Floating-Point Arithmetic*. ACM Trans. on Math. Software, Vol 15, No 2, 144-157, 1989.
- [13] Tang, P.T.P.: *Table-Driven Implementation of the Logarithm Function in IEEE Floating-Point Arithmetic*. ACM Trans. on Math. Software, Vol 16, No 4, 378-400, 1990.

<sup>5</sup>Due to pessimism on the part of interval arithmetic, which is caused by insufficient use of the correlation of different parts of the expression, the computed bounds may be not tight enough. In such cases more intelligence (e. g. monotonicity considerations) has to be added to the interval expression evaluation and/or a suitable subdivision strategy must be applied.

- [14] Tang, P.T.P.: *Table-Lookup Algorithms for Elementary Functions and Their Error Analysis*, Proceedings to the 10th Symposium on Computer Arithmetic, IEEE Computer Society Press, pp. 232–236, 1991.
- [15] Tang, P.T.P.: *Table-Driven Implementation of the  $\text{Exp}m1$  Function in IEEE Floating-Point Arithmetic*. ACM Trans. on Math. Software, Vol 18, No 2, 211-222, 1992.
- [16] Ziv, A.: *Fast Evaluation of Elementary Mathematical Functions with Correctly Rounded Last Bit*. ACM Trans. on Math. Software, Vol. 17, NO. 3, pp 410-423, 1991.