

SCAN 2008 EI Paso

Staggered Correction Computations with Enhanced Accuracy and Extremely Wide Exponent Range

University of Wuppertal - Germany
Frithjof Blomquist

Summary

1. Classic staggered correction arithmetic
2. Disadvantages of this arithmetic
 - Restricted exponent range
 - Low accuracy near the underflow range
3. New extended staggered interval arithmetic
4. Numerical examples

There are two possibilities for implementing an interval arithmetic in higher precision:

1. **Multiprecision-arithmetic** with base $B = 2^{32}$ and integer-valued mantissas. \leadsto high memory requirements and an enlarged runtime, caused by the simulation of the basic operations by software.
2. **Classic Interval staggered arithmetic**

$$\mathbf{x} := \sum_{i=1}^{p-1} x_i + [x_p, x_{p+1}] = [\text{Inf}(\mathbf{x}), \text{Sup}(\mathbf{x})],$$

$$\text{Inf}(\mathbf{x}) := \sum_{i=1}^{p-1} x_i + x_p, \quad \text{Sup}(\mathbf{x}) := \sum_{i=1}^{p-1} x_i + x_{p+1};$$

x_i : IEEE double numbers; $p \geq 1$: precision.

The basic arithmetic operations are based on the data type `dotprecision` and the x_i are generated by reading out the accumulator. So the x_i are in general non-overlapping IEEE numbers.

$$\mathbf{x} := \sum_{i=1}^{p-1} x_i + [x_p, x_{p+1}] = [\text{Inf}(\mathbf{x}), \text{Sup}(\mathbf{x})],$$

$$\text{Inf}(\mathbf{x}) := \sum_{i=1}^{p-1} x_i + x_p, \quad \text{Sup}(\mathbf{x}) := \sum_{i=1}^{p-1} x_i + x_{p+1};$$

x_i : IEEE double numbers; $p \geq 1$: precision.

$\text{Inf}(\mathbf{x}) \sim 10^k$:

$$\boxed{10^k} \quad \boxed{10^{k-16}} \quad \boxed{10^{k-32}} \quad \dots \quad \boxed{10^{k-16(p-1)}}$$

$$x_1 \quad x_2 \quad x_3 \quad \dots \quad x_p$$

Precision: $(16 \cdot p)$ decimal digits;

In case of $x_p = x_{p+1}$ (\mathbf{x} is a point interval) the maximum accuracy: $(16 \cdot p)$ decimal digits.

Thus, **precision and accuracy seem to be unbounded.**

IEEE-double number 10^k : $-308 \leq k \leq +308$.

$$10^{k-16(p-1)} \geq 10^{-308} \quad \iff \quad p \leq 20 + k/16.$$

Restriction of the exponent: $-308 \leq k \leq +308$.

Restriction of p : $1 \leq p \leq 20 + k/16 =: p_{max}$

Some examples for $p_{max} = p_{max}(k)$:

k	$1 \leq p \leq p_{max}$	$p_{max} \cdot 16$ dec. digits
+308	$p_{max} = 39$	624
+150	$p_{max} = 29$	464
0	$p_{max} = 20$	320
-308	$p_{max} = 1$	16

Thus, using the classic staggered interval arithmetic intermediate results near the underflow range **must** be avoided!

The next example will demonstrate how to manage the described difficulties near the underflow range:

$$\mathbf{x} \sim 10^{-150}, \mathbf{y} \sim 10^{-150},$$

$$\text{Product: } \mathbf{x} \cdot \mathbf{y} \sim 10^{-300}, \text{ (only 20 dec. digits)}$$

$$\tilde{\mathbf{x}} := 10^{300} \cdot \mathbf{x} \sim 10^{150}, \tilde{\mathbf{y}} := 10^{300} \cdot \mathbf{y} \sim 10^{150};$$

$$\tilde{\mathbf{x}}, \tilde{\mathbf{y}} \sim 10^{150} \rightsquigarrow \text{max. accuracy: 464 decimal digits.}$$

$$\text{Product: } \tilde{\mathbf{x}} \cdot \tilde{\mathbf{y}} \sim 10^{300}, \text{ without overflow!}$$

$$\mathbf{x} = 10^{-300} \cdot \tilde{\mathbf{x}}, \quad \mathbf{y} = 10^{-300} \cdot \tilde{\mathbf{y}}.$$

$$\mathbf{x} \cdot \mathbf{y} = (10^{-300} \cdot \tilde{\mathbf{x}}) \cdot (10^{-300} \cdot \tilde{\mathbf{y}}) = 10^{-600} \cdot (\tilde{\mathbf{x}} \cdot \tilde{\mathbf{y}}).$$

Now the scaled intervals $\tilde{\mathbf{x}}, \tilde{\mathbf{y}}$ can be multiplied in high accuracy and the above factor 10^{-600} or better -600 should be stored in some way.

Extended Staggered Arithmetic in C-XSC

Actual implementation in C-XSC: `lx_interval`

$$X = 2^m \cdot x, \quad m : \text{double}, \quad x : \text{lx_interval};$$
$$-9007199254740991 \leq m \leq +9007199254740991;$$

Multiplication:

$$X \cdot Y = (2^{m_x} \cdot x) \cdot (2^{m_y} \cdot y)$$
$$= (2^{m_x - s_x} \{2^{s_x} \cdot x\}) \cdot (2^{m_y - s_y} \{2^{s_y} \cdot y\})$$

Condition to get maximum accuracy:

$$u := 2^{s_x} \cdot x \sim 10^{150}, \quad v := 2^{s_y} \cdot y \sim 10^{150},$$
$$\rightsquigarrow u \cdot v \sim 10^{300}, \text{ no overflow,}$$

and a maximum accuracy of about 620 decimal digits.

$$X \cdot Y = (2^{m_u} \cdot u) \cdot (2^{m_v} \cdot v) = 2^{m_u + m_v} \{u \cdot v\};$$

Division:

$$\begin{aligned} X/Y &= (2^{m_x} \cdot x) / (2^{m_y} \cdot y) \\ &= (2^{m_x - s_x} \{2^{s_x} \cdot x\}) / (2^{m_y - s_y} \{2^{s_y} \cdot y\}) \end{aligned}$$

Condition to get maximum accuracy:

$$\begin{aligned} u &:= 2^{s_x} \cdot x \sim 10^{300}, & v &:= 2^{s_y} \cdot y \sim 10^{150}, \\ &\rightsquigarrow u/v \sim 10^{150}, & &\text{no overflow,} \end{aligned}$$

and a maximum accuracy of about 470 decimal digits.

$$X/Y = (2^{m_u} \cdot u) / (2^{m_v} \cdot v) = 2^{m_u - m_v} \{u/v\};$$

Addition:

$$\begin{aligned} X + Y &= (2^{m_x} \cdot x) + (2^{m_y} \cdot y) \\ &= (2^{m_x - s_x} \{2^{s_x} \cdot x\}) + (2^{m_y - s_y} \{2^{s_y} \cdot y\}) \end{aligned}$$

1. condition: $u := 2^{s_x} \cdot x \sim 10^{+300}$;
2. condition: $v := 2^{s_y} \cdot y$; $m_y - s_y \stackrel{!}{=} m_x - s_x =: n$;

$$X + Y = (2^n \cdot u) + (2^n \cdot v) = 2^n \cdot (u + v);$$

1. Example

Two complex numbers:

$$z = a + i \cdot b, \quad w = c + i \cdot d, \quad i := \sqrt{-1};$$

$$\Im(z/w) = \frac{b \cdot c - a \cdot d}{c^2 + d^2},$$

$$a = b = 3; \quad c = 10^{300}; \quad d = 10^{300} - 1;$$

$$\begin{aligned} \Im(z/w) &= \frac{3}{2 \cdot 10^{600} - 2 \cdot 10^{300} + 1} \\ &\approx 1.5 \cdot 10^{-600}; \quad \mathbf{476} \text{ dec. digits.} \end{aligned}$$

2. Example

$$z = a + i \cdot b, \quad w = c + i \cdot d, \quad i := \sqrt{-1};$$

$$\Im(z/w) = \frac{b \cdot c - a \cdot d}{c^2 + d^2},$$

$$a = 2 \cdot 10^{5000}; \quad b = 2 \cdot a; \quad c = d = 10^{-5000};$$

$$\begin{aligned} \Im(z/w) &= \frac{2}{2 \cdot 10^{-10000}} \\ &= 1.0 \cdot 10^{+10000}; \quad \mathbf{474} \text{ decimal digits.} \end{aligned}$$

Elementary Interval Functions in C-XSC

$ x $	$\log_{10}(x)$
x^2	$\sin(x)$
\sqrt{x}	$\cos(x)$
$\sqrt[n]{x}$	$\tan(x)$
$\sqrt{1+x} - 1$	$\cot(x)$
$\sqrt{1+x^2}$	$\arcsin(x)$
$\sqrt{1-x^2}$	$\arccos(x)$
$\sqrt{x^2-1}$	$\arctan(x)$
$x^n, p \in \mathbb{Z}$	$\operatorname{arccot}(x)$
$x^p, p: \text{lx_interval}$	$\sinh(x)$
e^x	$\cosh(x)$
$e^x - 1$	$\tanh(x)$
2^x	$\operatorname{coth}(x)$
10^x	$\operatorname{arsinh}(x)$
$\ln(x)$	$\operatorname{arcosh}(x)$
$\ln(1+x)$	$\operatorname{artanh}(x)$
$\log_2(x)$	$\operatorname{arcoth}(x)$

Example Exponential function:

IEEE: $x = 709.5 \mapsto e^x = 1.35 \dots \cdot 10^{+308}$.

Inclusion with an **extended** staggered interval of type `lx_interval`

$$x = \overbrace{6243314768166065}^{16 \text{ decimal digits}}$$

$$e^x \in 10^{\overbrace{+2711437152599601}^{16 \text{ decimal digits}}} \cdot \underbrace{9.90 \dots 663}_{451 \text{ dec. digits}} \overset{83\dots}{05\dots}$$

Mathematica:

$$\text{Exp}[\overbrace{1488521748}^{10 \text{ dec. digits}}] \rightsquigarrow \text{Overflow message!}$$

Next application **Taylor Arithmetic**

In C-XSC a complete Taylor arithmetic, based on the class `lx_interval`, is implemented.

$$f(x) = x^4 \cdot \sin(x^4);$$

Task:

Find a guaranteed inclusion of the fiftieth derivative for $x = 10^5$,

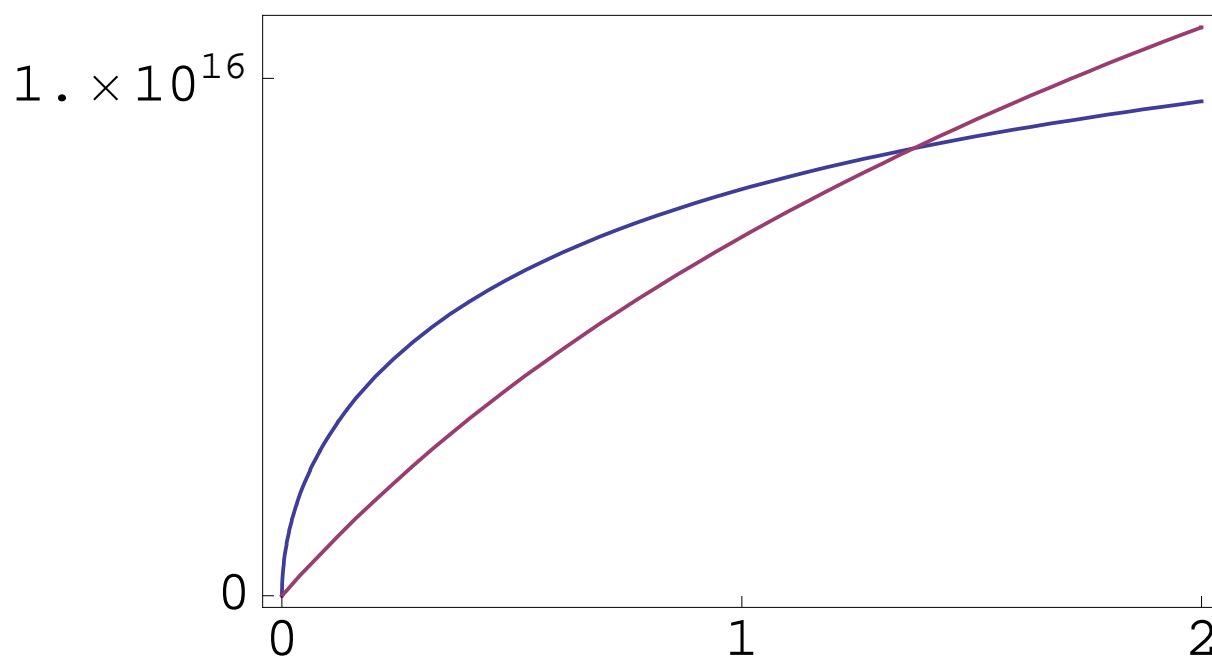
$$f^{(50)}(10^5) \in ?$$

$$f^{(50)}(10^5) \in 10^{799} \cdot \underbrace{8.179531 \dots 70843}_{461 \text{ dec. digits}} \frac{44\dots}{20\dots}$$

The inclusion interval lies **deep** in the IEEE overflow range and so the inclusion can only be calculated with the **extended** interval staggered arithmetic!

Next application: **Inclusion of (all) zeros**

$$g(x) := 10^{16} \cdot \arctan(\sqrt{x}) - 10^{16} \cdot \ln(1+x);$$



$$f(x) := 10^{16} \cdot \arctan(\sqrt{x}) - 10^0 \cdot \ln(1+x).$$

$$f(x_0) = 0, \quad x_0 \in 10^{44707} \cdot \underbrace{7.66 \dots 564}_{471 \text{ dec. digits}} \frac{96 \dots}{80 \dots}$$

Complex-valued Interval Functions in C-XSC

$z = x + i \cdot y, \quad z : \text{lx_cinterval}$	
$ z $	$\arctan(z)$
z^2	$\operatorname{arccot}(z)$
\sqrt{z}	$\sinh(z)$
$\sqrt[n]{z}$	$\cosh(z)$
$\arg(z)$	$\tanh(z)$
$\ln(z)$	$\operatorname{coth}(z)$
$z^n, n \in \mathbb{Z}$	$\operatorname{arsinh}(z)$
$z^p, p : \text{lx_interval}$	$\operatorname{arcosh}(z)$
$z^w, w : \text{lx_cinterval}$	$\operatorname{artanh}(z)$
e^z	$\operatorname{arcoth}(z)$
$\sin(z)$	$\sqrt{1+z} - 1$
$\cos(z)$	$\sqrt{1+z^2}$
$\tan(z)$	$\sqrt{1-z^2}$
$\cot(z)$	$\sqrt{z^2-1}$
$\arcsin(z)$	$e^z - 1$
$\arccos(z)$	$\ln(1+z)$