# Error-free transformations in real and complex floating point arithmetic

**Stef Graillat**

Joint work with Valérie Ménissier-Morain

LIP6/PEQUAN - Université Pierre et Marie Curie (Paris 6)

International Symposium on Nonlinear Theory and its Applications

Vancouver, Canada, September 16-19, 2007

## What are Error-Free Transformations (EFT) ?

Assume floating point arithmetic adhering IEEE 754 with rounding to nearest with rounding unit $\mathbf{u}$ (no underflow nor overflow)

Error free transformations are properties and algorithms to compute the generated elementary rounding errors,

$$a, b \text{ entries } \in \mathbb{F}, \quad a \circ b = \text{fl}(a \circ b) + e, \text{ with } e \in \mathbb{F}$$

Key tools for accurate computation

- fixed length expansions libraries : double-double (Briggs, Bailey, Hida, Li), quad-double (Bailey, Hida, Li)
- arbitrary length expansions libraries : Priest, Shewchuk
- compensated algorithms (Kahan, Priest, Ogita-Rump-Oishi, Graillat-Langlois-Louvet)

# EFT for the summation

$$x = \text{fl}(a \pm b) \quad \Rightarrow \quad a \pm b = x + y \quad \text{with } y \in \mathbb{F},$$

Algorithms of Dekker (1971) and Knuth (1974)

## Algorithm 1 (EFT of the sum of 2 floating point numbers with $|a| \geq |b|$)

function $[x, y] = $ `FastTwoSum`$(a, b)$
  $x = \text{fl}(a + b)$
  $y = \text{fl}((a - x) + b)$

## Algorithm 2 (EFT of the sum of 2 floating point numbers)

function $[x, y] = $ `TwoSum`$(a, b)$
  $x = \text{fl}(a + b)$
  $z = \text{fl}(x - a)$
  $y = \text{fl}((a - (x - z)) + (b - z))$

# EFT for the product (1/3)

$$x = \mathrm{fl}(a \cdot b) \quad \Rightarrow \quad a \cdot b = x + y \quad \text{with } y \in \mathbb{F},$$

Algorithm `TwoProduct` by Veltkamp and Dekker (1971)

$$a = x + y \quad \text{and} \quad x \text{ and } y \text{ non overlapping with } |y| \leq |x|.$$

---

### Algorithm 3 (Error-free split of a floating point number into two parts)

function $[x, y] = \mathtt{Split}(a, b)$
  factor $= \mathrm{fl}(2^s + 1)$          % $\mathbf{u} = 2^{-p}$ , $s = \lceil p/2 \rceil$
  $c = \mathrm{fl}(\mathtt{factor} \cdot a)$
  $x = \mathrm{fl}(c - (c - a))$
  $y = \mathrm{fl}(a - x)$

## Algorithm 4 (EFT of the product of 2 floating point numbers)

function $[x, y] = \mathtt{TwoProduct}(a, b)$
  $x = \mathrm{fl}(a \cdot b)$
  $[a_1, a_2] = \mathtt{Split}(a)$
  $[b_1, b_2] = \mathtt{Split}(b)$
  $y = \mathrm{fl}(a_2 \cdot b_2 - (((x - a_1 \cdot b_1) - a_2 \cdot b_1) - a_1 \cdot b_2))$

Given $a, b, c \in \mathbb{F}$,

- FMA$(a, b, c)$ is the nearest floating point number $a \cdot b + c \in \mathbb{F}$

**Algorithm 5 (EFT of the product of 2 floating point numbers)**

function $[x, y] = \texttt{TwoProductFMA}(a, b)$
$\quad x = \text{fl}(a \cdot b)$
$\quad y = \texttt{FMA}(a, b, -x)$

The `FMA` is available for example on PowerPC, Itanium, Cell processors.

# Summary

## Theorem 1

*Let $a, b \in \mathbb{F}$ and let $x, y \in \mathbb{F}$ such that $[x, y] = \texttt{TwoSum}(a, b)$. Then,*

$$a + b = x + y, \quad x = \text{fl}(a + b), \quad |y| \leq \mathbf{u}|x|, \quad |y| \leq \mathbf{u}|a + b|.$$

*The algorithm* $\texttt{TwoSum}$ *requires* 6 *flops.*

*Let $a, b \in \mathbb{F}$ and let $x, y \in \mathbb{F}$ such that $[x, y] = \texttt{TwoProduct}(a, b)$ . Then,*

$$a \cdot b = x + y, \quad x = \text{fl}(a \cdot b), \quad |y| \leq \mathbf{u}|x|, \quad |y| \leq \mathbf{u}|a \cdot b|,$$

*The algorithm* $\texttt{TwoProduct}$ *requires* 17 *flops.*

# Accurate sum and dot product

## Algorithm 6 (Ogita, Rump and Oishi 2005)

*Summation in twice the working precision*

function $\text{res} = \text{Sum2}(p)$
  $\pi_1 = p_1 \,;\, \sigma_1 = 0\,;$
  for $i = 2 : n$
    $[\pi_i, q_i] = \text{TwoSum}(\pi_{i-1}, p_i)$
    $\sigma_i = \text{fl}(\sigma_{i-1} + q_i)$
  $\text{res} = \text{fl}(\pi_n + \sigma_n)$

## Algorithm 7 (Ogita, Rump and Oishi 2005)

*Dot product in twice the working precision*

function $\text{res} = \text{Dot2}(x, y)$
  $[p, s] = \text{TwoProduct}(x_1, y_1)$
  for $i = 2 : n$
    $[h, r] = \text{TwoProduct}(x_i, y_i)$
    $[p, q] = \text{TwoSum}(p, h)$
    $s = \text{fl}(s + (q + r))$
  end
  $\text{res} = \text{fl}(p + s)$

# Accurate sum and dot product

## Proposition 1 (Ogita, Rump and Oishi 2005)

*Suppose Algorithm* Sum2 *is applied to floating point number* $p_i \in \mathbb{F}$, $1 \leq i \leq n$. *Let* $s := \sum p_i$, $S := \sum |p_i|$. *Then, we have*

$$|\mathtt{res} - s| \leq \mathbf{u}|s| + \gamma_{n-1}^2 S.$$

## Proposition 2 (Ogita, Rump and Oishi 2005)

*Let floating point numbers* $x_i, y_i \in \mathbb{F}, 1 \leq i \leq n$, *be given and denote by* $\mathtt{res} \in \mathbb{F}$ *the result computed by Algorithm* Dot2. *Then occurs,*

$$|\mathtt{res} - x^T y| \leq \mathbf{u}|x^T y| + \gamma_n^2 |x^T||y|.$$

$$\gamma_n = \frac{n\mathbf{u}}{1 - n\mathbf{u}}$$

# What about complex numbers?

Splitting between real and imaginary part

- Summation
  $s = \sum_{j=1}^{n} p_j$ with $p_j = a_j + ib_j$

$$\rightarrow \quad s = \underbrace{\sum_{j=1}^{n} a_j}_{\text{Sum2}} + i \underbrace{\sum_{j=1}^{n} b_j}_{\text{Sum2}}$$

- Dot product
  $x = (x_j)$ with $x_j = a_j + ib_j$ and $y = (y_j)$ with $y_j = c_j + id_j$ , $p = x^* y$

$$\rightarrow \quad p = \underbrace{\left[ \begin{array}{c} \text{Re}(x) \\ \text{Im}(x) \end{array} \right]^T \left[ \begin{array}{c} \text{Re}(y) \\ \text{Im}(y) \end{array} \right]}_{\text{Dot2}} + i \underbrace{\left[ \begin{array}{c} \text{Re}(x) \\ \text{Im}(x) \end{array} \right]^T \left[ \begin{array}{c} \text{Im}(y) \\ -\text{Re}(y) \end{array} \right]}_{\text{Dot2}}$$

# Error bounds

## Proposition 3

*Suppose Algorithm* Sum2cplx *is applied to floating point number*
$p_j = a_j + i b_j \in \mathbb{F} + i\mathbb{F}$, $1 \leq j \leq n$. *Let* $s := \sum p_j$, $S := \sum |p_j|$. *Then, we have*

$$|\text{res} - s| \leq \sqrt{2}\mathbf{u}|s| + 2\gamma_{n-1}^2 S.$$

## Proposition 4

*Let floating point numbers* $x = (x_j)$ *with* $x_j = a_j + i b_j$ *and* $y = (y_j)$ *with* $y_j = c_j + i d_j$ *be given and denote by* $\text{res} \in \mathbb{F} + i\mathbb{F}$ *the result computed by Algorithm* Dot2cplx. *Then occurs,*

$$|\text{res} - x^* y| \leq \sqrt{2}\mathbf{u}|x^* y| + 2\gamma_{2n}^2 |x|^T |y|.$$

$$p(z) = \sum_{j=0}^{n} a_j z^j, \qquad a_j \in \mathbb{C}, z = x + iy \in \mathbb{C}$$

$\rightarrow$ Write $p(z) = p_r(x, y) + iq_i(x, y)$ with $p_r$ and $q_r$ with real coefficients and evaluate $p_r$ and $q_r$ with Horner scheme

Problem : need formal manipulations

$\Rightarrow$ need new EFT for complex floating point arithmetic

Given $x, y \in \mathbb{F} + i\mathbb{F}$,

$$\mathrm{fl}(x \circ y) = (x \circ y)(1 + \varepsilon_1), \text{ for } \circ \in \{+, -\} \text{ and } |\varepsilon_\nu| \leq \mathbf{u},$$

and

$$\mathrm{fl}(x \cdot y) = (x \cdot y)(1 + \varepsilon_1), |\varepsilon_1| \leq \sqrt{2}\gamma_2.$$

**Algorithm 8 (EFT of the sum of 2 complex floating point numbers $x = a + ib$ and $y = c + id$ )**

function $[s, e] = \texttt{TwoSumCplx}(x, y)$
  $[s_1, e_1] = \texttt{TwoSum}(a, c)$
  $[s_2, e_2] = \texttt{TwoSum}(b, d)$
  $s = s_1 + is_2$
  $e = e_1 + ie_2$

# Complex EFT (2/2)

## Algorithm 9 (EFT of the product of two complex floating point numbers $x = a + ib$ and $y = c + id$)

function $[p, e, f, g] = $ TwoProductCplx$(x, y)$
  $[z_1, h_1] = $ TwoProduct$(a, c)$
  $[z_2, h_2] = $ TwoProduct$(b, d)$
  $[z_3, h_3] = $ TwoProduct$(a, d)$
  $[z_4, h_4] = $ TwoProduct$(b, c)$
  $[z_5, h_5] = $ TwoSum$(z_1, -z_2)$
  $[z_6, h_6] = $ TwoSum$(z_3, z_4)$
  $p = z_5 + iz_6$
  $e = h_1 + ih_3$
  $f = -h_2 + ih_4$
  $g = h_5 + ih_6$

## Summary

### Theorem 2

Let $x, y \in \mathbb{F} + i\mathbb{F}$ and let $s, e \in \mathbb{F} + i\mathbb{F}$ such that
$[s, e] = \texttt{TwoSumCplx}(x, y)$. Then,

$$x + y = s + e, \quad s = \text{fl}(x + y), \quad |e| \leq \mathbf{u}|s|, \quad |e| \leq \mathbf{u}|x + y|.$$

The algorithm $\texttt{TwoSumCplx}$ requires 12 flops.

### Theorem 3

Let $x, y \in \mathbb{F} + i\mathbb{F}$ and let $p, e, f, g \in \mathbb{F} + i\mathbb{F}$ such that
$[p, e, f, g] = \texttt{TwoProductCplx}(x, y)$ . Then,

$$x \cdot y = p + e + f + g \quad p = \text{fl}(x \cdot y), \quad |e + f + g| \leq \sqrt{2}\gamma_2|x \cdot y|,$$

The algorithm $\texttt{TwoProductCplx}$ requires 80 flops.

$\texttt{TwoProductCplx}$ requires 20 flops if one uses $\texttt{TwoProductFMA}$.

# The Horner scheme

## Algorithm 10 (Horner scheme)

function $\mathtt{res} = \mathtt{Horner}(p, x)$

$\quad s_n = a_n$

$\quad$ for $i = n - 1 : -1 : 0$

$\quad\quad p_i = \mathrm{fl}(s_{i+1} \cdot x)$ $\qquad\qquad$ % rounding error

$\quad\quad s_i = \mathrm{fl}(p_i + a_i)$ $\qquad\qquad$ % rounding error

$\quad$ end

$\quad \mathtt{res} = s_0$

# EFT for the polynomial evaluation

We now propose an EFT for the polynomial evaluation with the Horner scheme.

## Algorithm 11 (EFT for the Horner scheme)

$\text{function } [h, p_\pi, p_\mu, p_\nu, p_\sigma] = \texttt{EFTHornerCplx}(p, x)$

   $s_n = a_n$

   $\text{for } i = n - 1 : -1 : 0$

      $[p_i, \pi_i, \mu_i, \nu_i] = \texttt{TwoProductCplx}(s_{i+1}, x)$

      $[s_i, \sigma_i] = \texttt{TwoSumCplx}(p_i, a_i)$

      $\text{Let } \pi_i \text{ be the coefficient of degree } i \text{ in } p_\pi$

      $\text{Let } \mu_i \text{ be the coefficient of degree } i \text{ in } p_\mu$

      $\text{Let } \nu_i \text{ be the coefficient of degree } i \text{ in } p_\nu$

      $\text{Let } \sigma_i \text{ be the coefficient of degree } i \text{ in } p_\sigma$

   $\text{end}$

   $h = s_0$
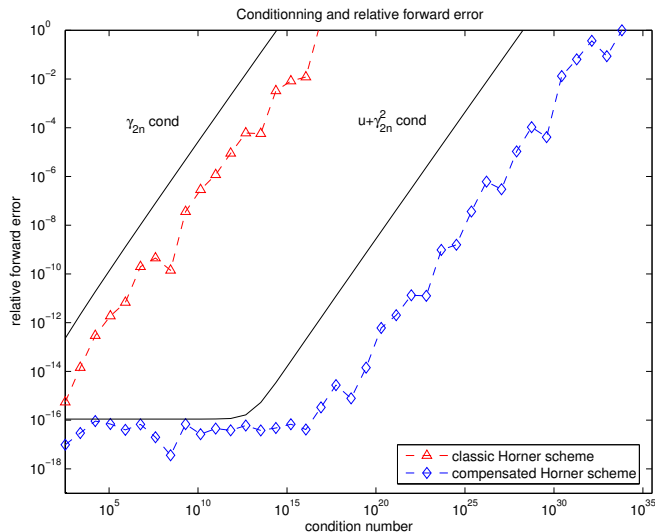
# Complex compensated Horner scheme

$$p(x) = h + (p_\pi + p_\sigma + p_\mu + p_\nu)(x)$$

## Algorithm 12 (Complex compensated Horner scheme)

*function* res $=$ CompHornerCplx$(p, x)$
  $[h, p_\pi, p_\mu, p_\nu, p_\sigma] = $ EFTHornerCplx$(p, x)$
  $c = $ HornerSumAcc$(p_\pi, p_\mu, p_\nu, p_\sigma, x)$
  res $= $ fl$(h + c)$

# Numerical experiment

$p(x) = (x - (1 + i))^n$ evaluated at $x = \mathrm{fl}(1.333 + 1.333i)$ and $n = 3 : 42$

# Conclusion

- Compensated algorithms in complex floating point arithmetic :
  - use of real EFT when possible
  - use of complex EFT otherwise
- Future work
  - complex version of the Compensated Horner Scheme
  - validation in complex floating point arithmetic

# Thank you for your attention