

Accurate Floating Point Product

Stef Graillat

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

REC'08, Third International Workshop on Reliable Engineering Computing
Savannah, Georgia, USA, February 20-22, 2008



- Determinant of a triangle matrix

$$T = \begin{bmatrix} t_{11} & t_{12} & \cdots & t_{1n} \\ & t_{22} & & t_{2n} \\ & & \ddots & \vdots \\ & & & t_{nn} \end{bmatrix}.$$

$$\det(T) = \prod_{i=1}^n t_{ii}.$$

- Evaluation of a polynomial when represented by the root product form $p(x) = a_n \prod_{i=1}^n (x - x_i)$

What are Error-Free Transformations (EFT)?

Assume floating point arithmetic adhering IEEE 754 with **rounding to nearest** with rounding unit 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

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

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 = fl(a + b)
    y = fl((a - x) + b)
```

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

```
function [x, y] = TwoSum(a, b)
    x = fl(a + b)
    z = fl(x - a)
    y = fl((a - (x - z)) + (b - z))
```

EFT for the product (1/3)

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

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] = Split(a, b)
    factor = fl(2s + 1)           % u = 2-p, s = [p/2]
    c = fl(factor · a)
    x = fl(c - (c - a))
    y = fl(a - x)
```

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

```
function  $[x, y] = \text{TwoProduct}(a, b)$ 
```

```
   $x = \text{fl}(a \cdot b)$ 
```

```
   $[a_1, a_2] = \text{Split}(a)$ 
```

```
   $[b_1, b_2] = \text{Split}(b)$ 
```

```
   $y = \text{fl}(a_2 \cdot b_2 - (((x - a_1 \cdot b_1) - a_2 \cdot b_1) - a_1 \cdot b_2))$ 
```

EFT for the product (3/3)

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

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

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

Theorem 1

Let $a, b \in \mathbb{F}$ and let $x, y \in \mathbb{F}$ such that $[x, y] = \text{TwoSum}(a, b)$. 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 `TwoSum` requires 6 flops.

Let $a, b \in \mathbb{F}$ and let $x, y \in \mathbb{F}$ such that $[x, y] = \text{TwoProduct}(a, b)$. 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 `TwoProduct` requires 17 flops.

Classic method for computing product

The classic method for evaluating a product of n numbers

$$a = (a_1, a_2, \dots, a_n)$$

$$p = \prod_{i=1}^n a_i$$

is the following algorithm.

Algorithm 6 (Product evaluation)

```
function res = Prod(a)
    p1 = a1
    for i = 2 : n
        pi = fl(pi-1 · ai)    % rounding error πi
    end
    res = pn
```

This algorithm requires $n - 1$ flops

$$\gamma_n := \frac{n\mathbf{u}}{1 - n\mathbf{u}} \quad \text{for } n \in \mathbb{N}.$$

A forward error bound is

$$|a_1 a_2 \cdots a_n - \text{res}| \leq \gamma_{n-1} |a_1 a_2 \cdots a_n|$$

A validated error bound is

$$|a_1 a_2 \cdots a_n - \text{res}| \leq \text{fl} \left(\frac{\gamma_{n-1} |\text{res}|}{1 - 2\mathbf{u}} \right)$$

Algorithm 7 (Product evaluation with a compensated scheme)

```
function res = CompProd(a)
    p1 = a1
    e1 = 0
    for i = 2 : n
        [pi, πi] = TwoProduct(pi-1, ai)
        ei = fl(ei-1ai + πi)
    end
    res = fl(pn + en)
```

This algorithm requires $19n - 18$ flops

Algorithm 8 (Product evaluation with a compensated scheme with TwoProductFMA and FMA)

```
function res = CompProdFMA(a)
    p1 = a1
    e1 = 0
    for i = 2 : n
        [pi, πi] = TwoProductFMA(pi-1, ai)
        ei = FMA(ei-1, ai, πi)
    end
    res = fl(pn + en)
```

This algorithm requires $3n - 2$ flops

Theorem 2

Suppose Algorithm CompProd is applied to floating point number $a_i \in \mathbb{F}$, $1 \leq i \leq n$, and set $p = \prod_{i=1}^n a_i$. Then,

$$|\text{res} - p| \leq \mathbf{u}|p| + \gamma_n \gamma_{2n}|p|$$

Condition number of the product evaluation :

$$\text{cond}(a) = \limsup_{\varepsilon \rightarrow 0} \left\{ \frac{|(a_1 + \Delta a_1) \cdots (a_n + \Delta a_n) - a_1 \cdots a_n|}{\varepsilon |a_1 a_2 \cdots a_n|} : |\Delta a_i| \leq \varepsilon |a_i| \right\}$$

A standard computation yields

$$\text{cond}(a) = n$$

Numerical experiments

Laptop with a Pentium M processor at 1.73GHz with gcc version 4.0.2.

Tab.: Measured computing times with Prod normalised to 1.0

n	Prod	CompProd
100	1.0	3.5
500	1.0	4.4
1000	1.0	5.0
10000	1.0	4.9
100000	1.0	5.5

The theoretical ratio for CompProd is 19.

Compensated algorithms are generally faster than the theoretical performances

→ due to a better instruction-level parallelism.

Lemma 1

Suppose Algorithm CompProd is applied to floating point numbers $a_i \in \mathbb{F}$, $1 \leq i \leq n$ and set $p = \prod_{i=1}^n a_i$. Then, the absolute forward error affecting the product is bounded according to

$$|\text{res} - p| \leq \text{fl} \left(\left(\mathbf{u}|\text{res}| + \frac{\gamma_n \gamma_{2n} |a_1 a_2 \cdots a_n|}{1 - (n+3)\mathbf{u}} \right) / (1 - 2\mathbf{u}) \right).$$

Faithful rounding (1/3)

Floating point predecessor and successor of a real number r satisfying $\min\{f : f \in \mathbb{R}\} < r < \max\{f : f \in \mathbb{F}\}$:

$$\text{pred}(r) := \max\{f \in \mathbb{F} : f < r\} \quad \text{and} \quad \text{succ}(r) := \min\{f \in \mathbb{F} : r < f\}.$$

Definition 1

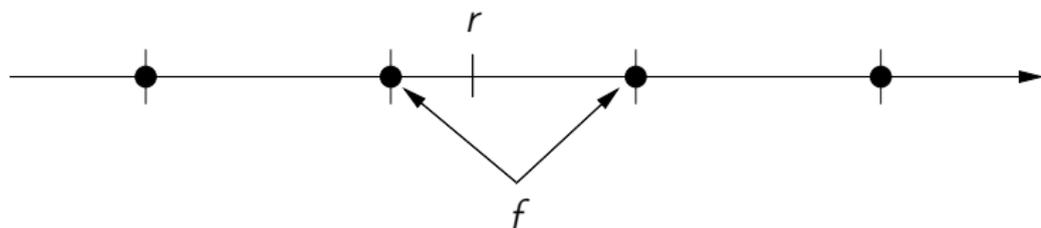
A floating point number $f \in \mathbb{F}$ is called a faithful rounding of a real number $r \in \mathbb{R}$ if

$$\text{pred}(f) < r < \text{succ}(f).$$

We denote this by $f \in \square(r)$. For $r \in \mathbb{F}$, this implies that $f = r$.

Faithful rounding means that the computed result is equal to the exact result if the latter is a floating point number and otherwise is one of the two adjacent floating point numbers of the exact result.

Faithful rounding (2/3)



Lemma 2 (Rump, Ogita and Oishi, 2005)

Let $r, \delta \in \mathbb{R}$ and $\tilde{r} := \text{fl}(r)$. Suppose that $2|\delta| < \mathbf{u}|\tilde{r}|$. Then $\tilde{r} \in \square(r + \delta)$, that means \tilde{r} is a faithful rounding of $r + \delta$.

Faithful rounding (3/3)

Let $\text{res} = \text{CompProd}(p)$

Lemma 3

If $n < \frac{\sqrt{1-u}}{\sqrt{2}\sqrt{2+u}+2\sqrt{(1-u)u}} u^{-1/2}$ then res is a faithful rounding of p .

If $n < \alpha u^{-1/2}$ where $\alpha \approx 1/2$ then the result is faithfully rounded

In double precision where $u = 2^{-53}$, if $n < 2^{25} \approx 5 \cdot 10^7$, we get a faithfully rounded result

If

$$\text{fl} \left(2 \frac{\gamma_n \gamma_{2n} |a_1 a_2 \cdots a_n|}{1 - (n+3)\mathbf{u}} \right) < \text{fl}(\mathbf{u}|\text{res}|)$$

then we got a faithfully rounded result. This makes it possible to check a *posteriori* if the result is faithfully rounded.

Algorithm 9 (Power evaluation with a compensated scheme)

```
function res = CompLinPower(x, n)
    p1 = x
    e1 = 0
    for i = 2 : n
        [pi, πi] = TwoProduct(pi-1, x)
        ei = fl(ei-1x + πi)
    end
    res = fl(pn + en)
```

Complexity : $\mathcal{O}(n)$

In double precision where $\mathbf{u} = 2^{-53}$, if $n < 2^{25} \approx 5 \cdot 10^7$, we get a faithfully rounded result

Algorithm 10 (Multiplication of two double-double numbers)

```
function  $[r_h, r_l] = \text{prod\_dd\_dd}(a_h, a_l, b_h, b_l)$   
     $[t_1, t_2] = \text{TwoProduct}(a_h, b_h)$   
     $t_3 = \text{fl}(((a_h \cdot b_l) + (a_l \cdot b_h)) + t_2)$   
     $[r_h, r_l] = \text{TwoProduct}(t_1, t_3)$ 
```

Algorithm 11 (Multiplication of double-double number by a double number)

```
function  $[r_h, r_l] = \text{prod\_dd\_d}(a, b_h, b_l)$   
     $[t_1, t_2] = \text{TwoProduct}(a, b_h)$   
     $t_3 = \text{fl}((a \cdot b_l) + t_2)$   
     $[r_h, r_l] = \text{TwoProduct}(t_1, t_3)$ 
```

Theorem 3 (Lauter 2005)

Let be $a_h + a_l$ and $b_h + b_l$ the double-double arguments of Algorithm `prod_dd_dd`. Then the returned values r_h and r_l satisfy

$$r_h + r_l = ((a_h + a_l) \cdot (b_h + b_l))(1 + \varepsilon)$$

where ε is bounded as follows : $|\varepsilon| \leq 16\mathbf{u}^2$. Furthermore, we have $|r_l| \leq \mathbf{u}|r_h|$.

A logarithmic algorithm

Algorithm 12 (Power evaluation with a compensated scheme)

```
function res = CompLogPower(x, n)           % n = (n_t n_{t-1} \cdots n_1 n_0)_2
    [h, l] = [1, 0]
    for i = t : -1 : 0
        [h, l] = prod_dd_dd(h, l, h, l)
        if n_i = 1
            [h, l] = prod_dd_d(x, h, l)
        end
    end
    res = fl(h + l)
```

Complexity : $\mathcal{O}(\log n)$

Theorem 4

The two values h and l returned by Algorithm `CompLogPower` satisfy

$$h + l = x^n(1 + \varepsilon)$$

with

$$(1 - 16\mathbf{u}^2)^{n-1} \leq 1 + \varepsilon \leq (1 + 16\mathbf{u}^2)^{n-1}.$$

For example, in double precision where $\mathbf{u} = 2^{-53}$, if $n < 2^{49} \approx 5 \cdot 10^{14}$, then we get a faithfully rounded result.

Thank you for your attention