

# Apporter de la confiance aux calculs en arithmétique virgule flottante



Jean-Michel Muller

CNRS - Laboratoire LIP  
Septembre 2024

<http://perso.ens-lyon.fr/jean-michel.muller/>

# Floating-Point Arithmetic

- by far the **most frequent solution** for manipulating real numbers in computers;
- comes from the “scientific notation” used for 3 centuries by the scientific community;

Sometimes a bad reputation... for bad reasons:

- **intrinsically approximate...**
    - but most data is approximate;
    - but most numerical problems we deal with have no closed-form solution;
    - and in a subtle way (correct rounding), FP arithmetic is **exact**.
  - part of the literature comes from times when it was poorly specified;
- too often, viewed as a mere set of cooking recipes.

## We wish to show that

- it is a **well specified** arithmetic, on which one can build trustable calculations;
- one can **prove useful properties** and **build efficient algorithms** on FP arithmetic;
- and yet the proofs are complex: **formal proof** is helpful.

# Desirable properties of an arithmetic system

- **Speed:** tomorrow's weather must be computed in less than 24 hours;
- **Reliability:** all numerical computing is built upon basic arithmetic. If the arithmetic collapses, everything collapses;
- **Accuracy;**
- **Range:** represent big and tiny numbers as well;
- **Size:** silicon area for hardware, memory consumption for software;
- **Power consumption;**
- **Easiness of implementation and use:** If a given arithmetic is too arcane, nobody will use it. . .

. . . of course, you can't win on all fronts.

# Much change since the 70's and 80's: i) applications

## Numerical simulation



- trillions of operations
- crash? just start again the simulation (but not too often)
- **reproducibility** may be useful.

## Finance

# Much change since the 70's and 80's: i) applications

## Numerical simulation



- trillions of operations
- crash? just start again the simulation (but not too often)
- **reproducibility** may be useful.

## Finance

### Embedded computing



- speed: yes, but no need to be faster than real time;
- crash? ahem...

→ **certified calculations.**

## Entertainment



- Supermario's pizza: no need to carefully follow the laws of physics;
- **fluidity** matters;
- **reproducibility**: each player must see the same game landscape.

## Artificial intelligence

- neural net training: huge amount of **very low precision** calculations.

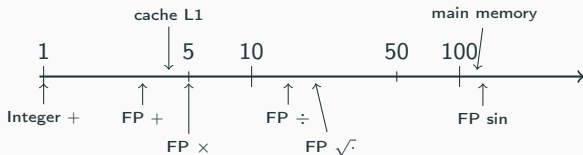
## Much change since the 70's and 80's: ii) performance

- the ratio

$$\frac{\text{time to read/write in memory}}{\text{time to perform } +, \times, \div, \sqrt{}}$$

has increased by a factor **around 140** between 1986 and 2000;

- It has continued to increase after 2000, but at a somehow slower pace;
  - the challenge is no longer to design fast arithmetic operators, but to be able to feed them with data at a very high rate;
- first consequence: many **new architectural concepts** (multiple levels of cache, pipelining, vector instructions, branch prediction);
- second consequence: incentive to use **small formats** whenever possible.



# Much change since the 70's and 80's: iii) FP formats

single precision (a.k.a. binary32)  
double precision (a.k.a. binary64)  $\Rightarrow$   $\left\{ \begin{array}{l} \text{8-bit emerging formats for IA} \\ \text{BFloat16} \\ \text{binary16} \\ \text{binary32} \\ \text{binary64} \\ \text{binary128 (quad)} \end{array} \right.$

- Combinatorial explosion of all the possible arithmetic operators of the form  $\text{Format 1} \times \text{Format 2} \rightarrow \text{Format 3}$
- Need to develop and maintain math function libraries for all these formats.

Cleverly using these formats:

Locate when  
low precision  
puts us  
at an unacceptable risk.



Locate when  
big precision  
totally destroys  
performance.

Numerical analysis, abstract interpretation, compilation, computer architecture, formal proof, ...



# A few weird arithmetic things

- Excel'2007 (first releases), compute  $65535 - 2^{-37}$ , you get 100000;
- 2020: in a competition, robotic car crash due to bad handling of floating-point exception



- if you have a Casio FX-92 pocket calculator, compute  $11^6/13$ , you will get

$$\frac{156158413}{3600}\pi$$

compute  $97027288/89521$ , you will get

$$345\pi.$$

## Base 2, precision- $p$ FP arithmetic

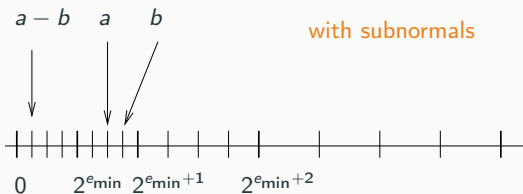
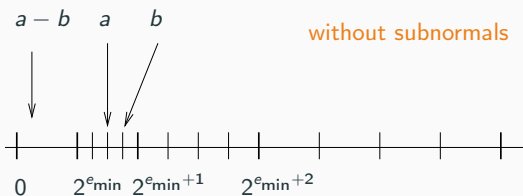
In **binary, precision- $p$  Floating-Point (FP) arithmetic**, a number  $x$  is represented by two integers  $M$  (integral significand) and  $e$  (exponent):

$$x = \left( \frac{M}{2^{p-1}} \right) \cdot 2^e = m_0.m_1m_2 \cdots m_{p-1} \cdot 2^e$$

where  $M, e \in \mathbb{Z}$ , with  $|M| \leq 2^p - 1$  and  $e_{\min} \leq e \leq e_{\max}$ . Additional requirement:  $e$  **smallest under these constraints**.

- $x$  is **normal** if  $|x| \geq 2^{e_{\min}}$  (implies  $|M| \geq 2^{p-1}$ , i.e.,  $m_0 = 1$ );
- $x$  is **subnormal** otherwise ( $m_0 = 0$ ).

Subnormal numbers complicate the implementation of FP multiplication, but...



If  $a$  and  $b$  are FPN,  $a \neq b$  equivalent to “computed  $a - b \neq 0$ ”.

### Theorem 1 (Hauser)

*If the absolute value of the sum/difference of two FP numbers is  $\leq 2^{e_{\min}+1}$  then it is a floating-point number (i.e., it is exactly representable in FP arithmetic).*

## Before 1985: a total mess...

Machine	Underflow $\lambda$	Overflow $\Lambda$
DEC PDP-11, VAX, F and D formats	$2^{-128} \approx 2.9 \times 10^{-39}$	$2^{127} \approx 1.7 \times 10^{38}$
DEC PDP-10; Honeywell 600, 6000; Univac 110x single; IBM 709X, 704X	$2^{-129} \approx 1.5 \times 10^{-39}$	$2^{127} \approx 1.7 \times 10^{38}$
Burroughs 6X00 single	$8^{-51} \approx 8.8 \times 10^{-47}$	$8^{76} \approx 4.3 \times 10^{68}$
H-P 3000	$2^{-256} \approx 8.6 \times 10^{-78}$	$2^{256} \approx 1.2 \times 10^{77}$
IBM 360, 370; Amdahl1; DG Eclipse M/600; ...	$16^{-65} \approx 5.4 \times 10^{-79}$	$16^{63} \approx 7.2 \times 10^{75}$
Most handheld calculators	$10^{-99}$	$10^{100}$
CDC 6X00, 7X00, Cyber	$2^{-976} \approx 1.5 \times 10^{-294}$	$2^{1070} \approx 1.3 \times 10^{322}$
DEC VAX G format; UNIVAC, 110X double	$2^{-1024} \approx 5.6 \times 10^{-309}$	$2^{1023} \approx 9 \times 10^{307}$

Source: Kahan, *Why do we need a Floating-Point Standard*, 1981.

## Before 1985: a total mess...

- Some **Cray computers**: overflow in FP  $\times$  detected just from the exponents of the entries, in parallel with the actual computation of the product;  
→  $1 * x$  could overflow;
- still on the Crays, only 12 bits of  $x$  were examined to detect a division by 0 when computing  $y/x$   
→ `if (x = 0) then z := 17.0 else z := y/x`  
could lead to zero divide error message... but since the multiplier too examined only 12 bits to decide if an operand is zero,  
`if (1.0 * x = 0) then z := 17.0 else z := y/x`  
was just fine.
- many systems, not enough “guard bits” for FP  $+$  → for  $x \approx 1$ , experts knew that  $(0.5 - x) + 0.5$  was much better than  $1.0 - x$ .

Writing **reliable** and **portable** numerical software was a challenge!

# IEEE-754 Standard for FP Arithmetic (1985, 2008, 2019)

- put an end to a mess (no portability, variable quality);
- leader: W. Kahan (father of the arithmetic of the HP35 and the Intel 8087);
- formats (in radices 2 and 10);
- **specification of operations** and conversions;
- exception handling ( $\max+1$ ,  $1/0$ ,  $\sqrt{-2}$ ,  $0/0$ , etc.);
- successive versions of the standard: 2008, 2019, and **2029 is already in preparation.**

# Correct rounding

- the sum, product, ... of two FP numbers is not, in general, a FP number  
→ must be rounded;
- the IEEE 754 Std for FP arithmetic specifies several rounding functions;
- the default function is **RN ties to even**.

**Correctly rounded operation:** returns what we would get by **exact operation followed by rounding**.

- correctly rounded  $+$ ,  $-$ ,  $\times$ ,  $\div$ ,  $\sqrt{\cdot}$  are required;

→ when  $c = a + b$  appears in a program, we get  $c = \text{RN}(a + b)$ .

→ somehow deterministic arithmetic (more later).

## ulp (unit in the last place), $u$ (unit round-off)

Binary, precision- $p$  FP arithmetic.

- If  $|x| \in [2^e, 2^{e+1})$ , then  $\text{ulp}(x) = 2^{\max\{e, e_{\min}\} - p + 1}$ .
  - Frequently used for expressing errors of **atomic** functions;
  - distance between consecutive FP numbers near  $x$ ;
- if  $2^{e_{\min}} \leq |x| \leq \Omega$ , then

$$|x - \text{RN}(x)| \leq \frac{1}{2} \text{ulp}(x) = 2^{\lfloor \log_2 |x| \rfloor - p},$$

therefore,

$$|x - \text{RN}(x)| \leq u \cdot |x|, \tag{1}$$

with  $u = 2^{-p}$ . Hence the **relative error**

$$\frac{|x - \text{RN}(x)|}{|x|}$$

(for  $x \neq 0$ ) is  $\leq u$ .

- $u$ , called **unit round-off** is frequently used for expressing relative errors.



# With the math functions work still needs to be done

library version	GNU libc	IML	AMD	Newlib	OpenLibm	Musl	Apple	LLVM	MSVC	FreeBSD	ArmPL	CUDA	ROCm
	2.40	2024.0.2	4.2	4.4.0	0.8.3	1.2.5	14.5	18.1.8	2022	14.1	24.04	12.2.1	5.7.0
acos	<b>0.523</b>	0.531	1.36	0.930	0.930	0.930	1.06		0.934	0.930	1.52	1.53	0.772
acosh	2.25	<b>0.509</b>	1.32	2.25	2.25	2.25	2.25		3.22	2.25	2.66	2.52	0.661
asin	<b>0.516</b>	0.531	1.06	0.981	0.981	0.981	0.709		1.05	0.981	2.69	1.99	0.710
asinh	1.92	<b>0.507</b>	1.65	1.92	1.92	1.92	1.58		2.05	1.92	2.04	2.57	0.661
atan	<b>0.523</b>	0.528	0.863	0.861	0.861	0.861	0.876		0.863	0.861	2.24	1.77	1.73
atanh	1.78	<b>0.507</b>	1.04	1.81	1.81	1.80	2.01		2.50	1.81	3.00	2.50	0.664
cbrt	3.67	0.523	1.53e22	0.670	0.668	0.668	0.729		1.86	0.668	1.79	<b>0.501</b>	<b>0.501</b>
cos	<b>0.516</b>	0.518	0.919	0.887	0.834	0.834	0.948	Inf	0.897	0.834		1.52	0.797
cosh	1.93	<b>0.516</b>	1.85	2.67	1.47	1.04	0.523		1.91	1.47	1.93	1.40	0.563
erf	1.43	<b>0.773</b>	1.00	1.02	1.02	1.02	6.41		4.62	1.02	2.29	1.50	1.12
erfc	5.19	<b>0.826</b>	4.08	4.08	3.72	3.72	10.7		8.46	4.08	1.71	4.51	4.08
exp	0.511	0.530	1.01	0.949	0.949	0.511	0.521	<b>0.500</b>	1.50	0.949	0.511	0.928	0.929

Largest errors in ulps for double-precision calculation of some math functions.  $\text{ulp}(x)$  is the distance between two FP numbers in the neighborhood of  $x$  (so the largest values should be 0.5 – which is the case with  $+$ ,  $-$ ,  $\times$ ,  $\div$ , and  $\sqrt{\cdot}$ ).

(Extracted from Gladman, Innocente, Mather, and Zimmermann, *Accuracy of Mathematical Functions...*, Aug. 2024)

## Exception handling: the show must go on...

- when an exception occurs: the computation must continue (default behaviour);
- two infinities and two zeros, with intuitive rules:  $1/(+0) = +\infty$ ,  $5 + (-\infty) = -\infty \dots$ ;
- and yet, something a little odd:  $\sqrt{-0} = -0$ ;
- **Not a Number** (NaN): result of  $\sqrt{-5}$ ,  $(\pm 0)/(\pm 0)$ ,  $(\pm \infty)/(\pm \infty)$ ,  $(\pm 0) \times (\pm \infty)$ , NaN +3, etc.

$$f(x) = 3 + \frac{1}{x^5}$$

will give the very accurate answer 3 for huge  $x$ , even if  $x^5$  overflows.

One should be cautious: behavior of

$$\frac{x^2}{\sqrt{x^3 + 1}}$$

for large  $x$ .

# With correct rounding and standardized exception handling, arithmetic is almost deterministic

- watch the dependency graph of operations (beware of “optimizing” compilers);
- watch the format of the implicit variables (such as the  $x+y$  in  $(x+y)*(z+t)$ );
- math functions still a problem unless you use a correctly rounded library such as Zimmermann & Sibidanov’s Core Math,<sup>1</sup> or LLVM libc.<sup>2</sup>

With enough care we can **prove properties and build specific algorithms.**

---

<sup>1</sup><https://core-math.gitlabpages.inria.fr/>

<sup>2</sup><https://libc.llvm.org/>

# A useful property: Sterbenz' Theorem

## Theorem 2 (Sterbenz)

*Let  $a$  and  $b$  be positive FP numbers. If*

$$\frac{a}{2} \leq b \leq 2a$$

*then  $a - b$  is a FP number*

*( $\rightarrow$  computed exactly, whatever the rounding function).*

**Beware:** the “2”s in the formula are **not** the radix. In radices 10, 16 or 42, the same property holds, still with  $\frac{a}{2} \leq b \leq 2a$ .

## Example of use: implementation of trig. functions in precision- $p$ FP arithmetic

- cosine function: **range reduction** to small interval followed by **polynomial approximation** in that interval;
- range reduction:  $x \rightarrow y = x - k\pi$  such that  $|y|$  is small. If done naively this is a very inaccurate operation.
- assuming the largest value of  $k$  of interest fits in  $m < p$  bits, express  $\pi$  as the sum of two FP numbers  $\pi_1$  and  $\pi_2$  such that
  - $\pi_1$  is closest to  $\pi$  among the FP numbers whose significand fits in  $p - m$  bits;
  - $\pi_2 = \text{RN}(\pi - \pi_1)$ .

Program:  $y \leftarrow ((x - k*\pi_1) - k*\pi_2)$

By construction,  $\Delta = k*\pi_1$  is exact, and by **Sterbenz Lemma**,  $x - \Delta$  is exact.

(Cody-Waite range reduction. Many improvements are possible)

# The error of (RN) FP addition is a FPN

## Lemma 3

*Let  $a$  and  $b$  be two FP numbers. Let*

$$s = RN(a + b) \text{ and } r = (a + b) - s.$$

*If no overflow when computing  $s$ , then  $r$  is a FP number.*

**Beware:** does not always work with rounding functions  $\neq RN$ .

# Get $r$ : the fast2sum algorithm (Dekker)

## Theorem 4 (Fast2Sum (Dekker))

(only *radix 2*). Let  $a$  and  $b$  be FP numbers, s.t.  $|a| \geq |b|$ . Following algorithm:  $s$  and  $r$  such that

- $s + r = a + b$  exactly;
- $s$  is "the" FP number that is closest to  $a + b$ ;

## Algorithm 1 (FastTwoSum)

```
 $s \leftarrow RN(a + b)$   
 $z \leftarrow RN(s - a)$   
 $r \leftarrow RN(b - z)$ 
```

## C Program 1

```
s = a+b;  
z = s-a;  
r = b-z;
```

**Important remark:** Proving the behavior of such algorithms requires use of the correct rounding property.

# The TwoSum Algorithm (Moller-Knuth)

- no need to compare  $a$  and  $b$ ;
- 6 operations instead of 3 yet, on many architectures, very cheap in front of wrong branch prediction penalty when comparing  $a$  and  $b$ ;
- works in all bases.

## Algorithm 2 (TwoSum)

$$\begin{aligned}s &\leftarrow RN(a + b) \\ a' &\leftarrow RN(s - b) \\ b' &\leftarrow RN(s - a') \\ \delta_a &\leftarrow RN(a - a') \\ \delta_b &\leftarrow RN(b - b') \\ r &\leftarrow RN(\delta_a + \delta_b)\end{aligned}$$

**Knuth:** if no underflow nor overflow occurs then  $a + b = s + r$ , and  $s$  is nearest  $a + b$ .

**Boldo et al:** formal proof + underflow does not hinder the result (overflow does).

**TwoSum is optimal** (no way of always obtaining the same result with less than  $6 \pm$  operations).



## Example of application: computing $x_1 + x_2 + x_3 + \dots + x_n$

### Naive algorithm:

```
s ← x1
for i = 2 to n do
  s ← RN(s + xi)
end for
return s
```

### Pichat, Ogita, Rump, and Oishi:

```
s ← x1
e ← 0
for i = 2 to n do
  (s, ei) ← 2Sum(s, xi)
  e ← RN(e + ei)
end for
return RN(s + e)
```

### Error bounds:

$$(n-1) \cdot u \sum |x_i|$$

(remember:  $u = 2^{-p}$ )

$$u \left| \sum_{i=1}^n x_i \right| + \left( \frac{(n-1)u}{1 - (n-1)u} \right)^2 \sum_{i=1}^n |x_i|$$

## What about products ?

- If  $a$  and  $b$  are FP numbers, then (under mild conditions),  
 $r = ab - \text{RN}(ab)$  is a FP number;
- We use the *fused multiply-add* (fma) instruction. It computes  $\text{RN}(ab + c)$ . First appeared in IBM RS6000, Intel/HP Itanium, PowerPC... Specified since 2008.
- obtained with algorithm **TwoMultFMA**  $\begin{cases} p = \text{RN}(ab) \\ r = \text{RN}(ab - p) \end{cases}$   
→ 2 operations only, gives  $p + r = ab$ .

## Just an example: $ad - bc$ with fused multiply-add

Kahan's algorithm for  $x = ad - bc$ :

$$\hat{w} \leftarrow \text{RN}(bc)$$

$$e \leftarrow \text{RN}(\hat{w} - bc)$$

$$\hat{f} \leftarrow \text{RN}(ad - \hat{w})$$

$$\hat{x} \leftarrow \text{RN}(\hat{f} + e)$$

Return  $\hat{x}$

- we have proven (2011):

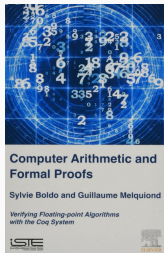
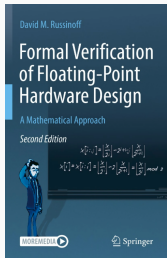
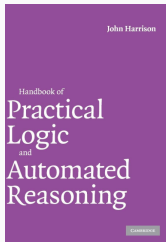
$$|\hat{x} - x| \leq 2u|x|$$

“asymptotically optimal” error bound.

- $\rightarrow$  rotations, complex arithmetic.

# Formal verification of FP algorithms

- starting point: the Pentium division bug (1994)
- **J. Harrison** formalized FP arithmetic in HOL Light, formally proved the division and sqrt algorithms of the Intel Itanium, and some elementary function algorithms (around 1999);
- **D. Russinoff**: similar things for AMD (more on the hardware side);
- **Sylvie Boldo and Guillaume Melquiond** use the Coq proof assistant (Flocq library, Gappa tool).



# Double-Word arithmetic

- Fast2Sum, 2Sum and 2MultFMA return their result as the unevaluated sum of two FP numbers.
  - **idea:** manipulate such unevaluated sums to perform more accurate calculations in critical parts of a numerical program.
- “double word” or “double-double” arithmetic. Most recent avatar: Rump and Lange’s “pair arithmetic” (2020).

## Definition 5

A **double-word** (DW) number  $x$  is the unevaluated sum  $x_h + x_\ell$  of two floating-point numbers  $x_h$  and  $x_\ell$  such that

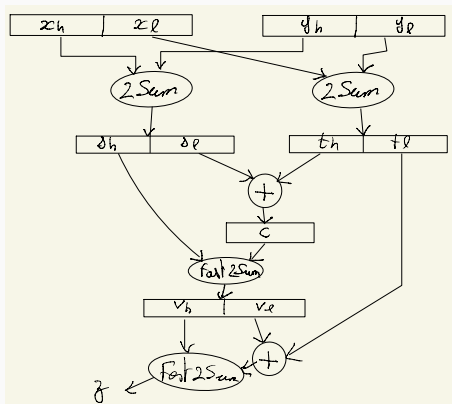
$$x_h = \text{RN}(x).$$

# DW+DW: “accurate version”

Sum of two DW numbers. There exist a “quick & dirty” algorithm, but its relative error is unbounded.

## DWPlusDW

- 1:  $(s_h, s_\ell) \leftarrow 2\text{Sum}(x_h, y_h)$
- 2:  $(t_h, t_\ell) \leftarrow 2\text{Sum}(x_\ell, y_\ell)$
- 3:  $c \leftarrow \text{RN}(s_\ell + t_h)$
- 4:  $(v_h, v_\ell) \leftarrow \text{Fast2Sum}(s_h, c)$
- 5:  $w \leftarrow \text{RN}(t_\ell + v_\ell)$
- 6:  $(z_h, z_\ell) \leftarrow \text{Fast2Sum}(v_h, w)$
- 7: **return**  $(z_h, z_\ell)$



## DW+DW: “accurate version”

We have (after a very long and tedious proof):

### Theorem 6

*If  $p \geq 3$ , the relative error of Algorithm DWPlusDW is bounded by*

$$\frac{3u^2}{1-4u} = 3u^2 + 12u^3 + 48u^4 + \dots, \quad (2)$$

That theorem has an interesting history. . .

**ALGORITHM 6:** *AccurateDWPlusDW*( $x_h, x_l, y_h, y_l$ ). Calculation of  $(x_h, x_l) + (y_h, y_l)$  in binary, precision- $p$ , floating-point arithmetic.

```

1:  $(s_h, s_l) \leftarrow 2\text{Sum}(x_h, y_h)$ 
2:  $(t_h, t_l) \leftarrow 2\text{Sum}(x_l, y_l)$ 
3:  $c \leftarrow \text{RN}(t_l + t_h)$ 
4:  $(v_h, v_l) \leftarrow \text{Fast2Sum}(s_h, c)$ 
5:  $w \leftarrow \text{RN}(t_l + v_l)$ 
6:  $(u_h, u_l) \leftarrow \text{Fast2Sum}(v_h, w)$ 
7: return  $(u_h, u_l)$ 
    
```

Li et al. (2006, 2002) claim that in binary64 arithmetic ( $p = 53$ ) the relative error of Algorithm 6 is upper bounded by  $2 \cdot 2^{-166}$ . This bound is inco-

```

 $x_h = 900719$ 
 $x_l = -9007$ 
 $y_h = -9007$ 
 $y_l = -9007$ 
    
```

in which one of the opera-  
the returned result is  $2\text{Sum}(x_l, y_l)$ ,

then the relative error of Algorithm 6 is

```
2.249999999
```

low, without loss of generality, we

Note that this example is somehow "gen-

$2^p - 1, x_l = -(2^p - 1), 2^{p-1}, y_h = -(2^p - 5)/2, y_l$

that is asymptotically equivalent (as  $p$  goes to infinity).

Now let us try to find a relative error bound. We are by  
**THEOREM 3.1.** If  $p \geq 3$ , then the relative error of Algorithm 6 is

by

$$\frac{3u^2}{1-4u} = 3u^2 + 12u^3 + 48u^4 + \dots,$$

which is less than  $3u^2 + 13u^3$  as soon as  $p \geq 6$ .

Note that the conditions on  $p$  ( $p \geq 3$  for the bound (3) to hold,  $p \geq 6$  for the simplified bound  $3u^2 + 13u^3$ ) are satisfied in all practical cases.

**PROOF.** First, we exclude the straightforward case in which one of the operands is zero. We can also quickly proceed with the case  $x_h + y_h = 0$ . The returned result is  $2\text{Sum}(x_l, y_l)$ , which is equal to  $x + y$ , that is, the computation is error-free. Now, without loss of generality, we assume  $1 \leq x_h < 2, x \geq |y|$  (which implies  $x_h \geq |y_h|$ ), and  $x_h + y_h$  nonzero. Notice that  $1 \leq x_h < 2$  implies  $1 \leq x_h \leq 2 - 2u$ , since  $x_h$  is a FP number.

Define  $e_1$  as the error committed at Line 3 of the algorithm:  
$$e_1 = c - (t_l + t_h) \tag{4}$$

and  $e_2$  as the error committed at Line 5:  
$$e_2 = w - (t_l + v_l). \tag{5}$$

**1.** If  $-x_h < y_h \leq -x_h/2$ . Sterbenz Lemma, applied to the first line of the algorithm, implies  $s_h = x_h + y_h, s_l = 0$ , and  $c = \text{RN}(t_h) = 0$ .

Define  
$$\sigma = \begin{cases} 2 & \text{if } y_h \leq -1, \\ 1 & \text{if } -1 < y_h \leq -x_h/2. \end{cases}$$

ACM Transactions on Mathematical Software, Vol. 44, No. 2, Article 15res. Publication date: October 2017.

We have  $-x_h < y_h \leq (1 - \sigma) + \frac{\sigma}{2}(x - 2)$ , so  $0 \leq x_h + y_h \leq 1 + \sigma \cdot (\frac{x}{2} - 1) \leq 1 - \sigma u$ . Also, since  $x_h$  is a multiple of  $2u$  and  $y_h$  is a multiple of  $\sigma u, x_h = x_h + y_h$  is a multiple of  $\sigma u$ . Since  $x_h$  is nonzero, we finally obtain

$$\sigma u \leq x_h \leq 1 - \sigma u. \tag{6}$$

We have  $|x_l| \leq u$  and  $|y_l| \leq \frac{u}{2}$ , so

$$|t_h| \leq \left(1 + \frac{\sigma}{2}\right)u \quad \text{and} \quad |t_l| \leq u^2. \tag{7}$$

From Equation (6), we deduce that the floating-point exponent of  $x_h$  is at least  $-p + \sigma - 1$ . From Equation (7), the floating-point exponent of  $c = t_h$  is at most  $-p + \sigma - 1$ . Therefore, the Fast2Sum algorithm introduces no error at line 4 of the algorithm, which implies

$$v_h + v_l = s_h + c = x_h + t_h = x + y - t_l.$$

Equations (6) and (7) imply

$$|s_h + t_h| \leq 1 + \left(1 - \frac{\sigma}{2}\right)u \leq 1 + \frac{u}{2},$$

so  $|x_h| \leq 1$  and  $|v_l| \leq \frac{u}{2}$ . From the bounds on  $|t_l|$  and  $|v_l|$ , we obtain:

$$|e_1| \leq \frac{1}{2} \text{ulp}(t_l + v_l) \leq \frac{1}{2} \text{ulp}\left(u^2 + \frac{u}{2}\right) = \frac{u^2}{2} \tag{8}$$

and

$$|e_2| \leq \frac{1}{2} \text{ulp}\left(\frac{1}{2} \text{ulp}(x_l + y_l) + \frac{1}{2} \text{ulp}\left((x + y) + \frac{1}{2} \text{ulp}(x_l + y_l)\right)\right). \tag{9}$$

**Lemma 2.1** and  $|s_h| \geq \sigma u$  imply that either  $x_h + t_h = 0$ , or  $|v_h| = |\text{RN}(x_h + c)| = |\text{RN}(x_h + t_h)| \geq \sigma u^2$ . If  $x_h + t_h = 0$ , then  $v_h = v_l = 0$  and the sequel of the proof is straightforward. Therefore, in the following, we assume  $|v_h| \geq \sigma u^2$ .

Now,

- If  $|v_h| = \sigma u^2$ , then  $|v_l + t_l \leq u|v_h| + u^2 = \sigma u^3 + u^3$ , which implies  $|w| = |\text{RN}(t_l + v_l)| \leq \sigma u^2 + |v_h|$ ;
- If  $|v_h| > \sigma u^2$ , then, since  $v_h$  is a FP number,  $|v_h|$  is larger than or equal to the FP number immediately above  $\sigma u^2$ , which is  $\sigma(1 + 2u)\sigma u^2$ . Hence  $|v_h| \geq \sigma u^2/(1 - u)$ , so  $|v_h| \geq u \cdot |v_h| + \sigma u^2 \geq |v_l| + |t_l|$ . So,  $|w| = |\text{RN}(t_l + v_l)| \leq |v_h|$ .

Therefore, in all cases, Fast2Sum introduces no error at line 6 of the algorithm, and we have

$$z_h + z_l = v_h + w = x + y + e_2. \tag{10}$$

Directly using Equation (10) and the bound  $u^2/2$  on  $|e_1|$  to get a relative error bound would result in a large bound, because  $x + y$  may be small. However, when  $x + y$  is very small, some simplification occurs thanks to Sterbenz Lemma. First,  $x_h + y_h$  is a nonzero multiple of  $\sigma u$ . Hence, since  $|x_l + y_l| \leq (1 + \frac{\sigma}{2})u$ , we have  $|x_l + y_l| \leq \frac{1}{2}(x_h + y_h)$ . Let us now consider the two possible cases:

- If  $-\frac{1}{2}(x_h + y_h) \leq x_l + y_l \leq -\frac{1}{2}(x_h + y_h)$ , which implies  $-\frac{1}{2}x_h \leq t_h \leq -\frac{1}{2}y_h$ , then Sterbenz Lemma applies to the floating-point addition of  $s_h$  and  $c = t_h$ . Therefore line 4 of the algorithm results in  $v_h = s_h$  and  $v_l = 0$ . An immediate consequence is  $e_2 = 0$ , so  $x_h - z_h = v_h + w = x + y$ : the computation of  $x + y$  is errorless;

ACM Transactions on Mathematical Software, Vol. 44, No. 2, Article 15res. Publication date: October 2017.



**ALGORITHM 6:** AccurateDWPPlusDW( $x_h, x_\ell, y_h, y_\ell$ ). Calculation of  $(x_h, x_\ell) + (y_h, y_\ell)$  in binary, precision- $p$ , floating-point arithmetic.

```

1:  $(s_h, t_h) \leftarrow \text{Fast2Sum}(x_h, y_h)$ 
2:  $(s_\ell, t_\ell) \leftarrow \text{Fast2Sum}(x_\ell, y_\ell)$ 
3:  $c \leftarrow \text{RN}(t_\ell + t_h)$ 
4:  $(v_h, v_\ell) \leftarrow \text{Fast2Sum}(s_h, c)$ 
5:  $w \leftarrow \text{RN}(t_\ell + v_\ell)$ 
6:  $(e_2, e_1) \leftarrow \text{Fast2Sum}(v_h, w)$ 
7: return  $(e_2, e_1)$ 
    
```

Li et al. (2006, 2002) claim that it is upper bounded by  $2 \cdot 2^{-16p}$ . The

either  $s_h + t_h = 0$ , or  $|v_h| = |\text{RN}(s_h + c)| = |\text{RN}(s_h + t_h + t_\ell)|$ . In the sequel of the proof is straightforward. The

then the relative error of Algorithm 6 is

$$\frac{|e_2|}{|x|} = \frac{|v_h|}{|x|} = \frac{|s_h + t_h + t_\ell|}{|x|} \leq 2.24999999999999956 \dots$$

Note that this example is somehow "generic". In precision- $p$  FP arithmetic,  $2^p - 1, x_\ell = -(2^p - 1) \cdot 2^{p-1}, y_\ell = -(2^p - 5)/2$ , and  $y_\ell = -(2^p - 1) \cdot 2^{p-1}$  leads to a relative error that is asymptotically equivalent (as  $p$  goes to infinity) to  $2.25u^2$ .

Now let us try to find a relative error bound. We are going to show the following result. **THEOREM 3.1.** If  $p \geq 3$ , then the relative error of Algorithm 6 (AccurateDWPPlusDW) is bounded by

$$\frac{3u^2}{1 - 4u} = 3u^2 + 12u^3 + 48u^4 + \dots, \quad (3)$$

which is less than  $3u^2 + 13u^3$  as soon as  $p \geq 6$ . Note that the conditions on  $p$  ( $p \geq 3$  for the bound (3) to hold,  $p \geq 6$  for the simplified bound  $3u^2 + 13u^3$ ) are satisfied in all practical cases.

**PROOF.** First we have  $|x_\ell + y_\ell| = 2 \cdot |c|$ .

• If  $-\frac{3}{2}(x_h + y_h) \leq x_\ell + y_\ell \leq -\frac{1}{2}(x_h + y_h)$

and  $e_2$  as the error. This applies to the floating-point

1. If  $-x_h < y_h \leq -x_h/2$ , Sterbenz Lemma, applied to the first line of the algorithm, implies  $s_h = x_h + y_h, t_h = 0$ , and  $c = \text{RN}(t_h) = 0$ . Define

$$\sigma = \begin{cases} 2 & \text{if } y_h \leq -1, \\ 1 & \text{if } -1 < y_h \leq -x_h/2. \end{cases}$$

We have  $-x_h < y_h \leq (1 - \sigma)(x - 2)$ , so  $0 \leq x_h + y_h \leq 1 + \sigma \cdot (\frac{x}{2} - 1) \leq 1 - \sigma u$ . Also, since  $x_h$  is a multiple of  $2u$  and  $y_h$  is a multiple of  $\sigma u, x_h - x_h + y_h$  is a multiple of  $\sigma u$ . Since  $x_h$  is nonzero, we finally obtain

$$\sigma u \leq x_h \leq 1 - \sigma u. \quad (6)$$

Since  $|x_\ell| \leq u$  and  $|y_\ell| \leq \frac{u}{2}$ , so

$$|c| \leq \left(1 + \frac{\sigma}{2}\right)u \quad \text{and} \quad |t_\ell| \leq u^2. \quad (7)$$

point exponent of  $t_h$  is at least  $-p + \sigma - 1$ . From is at most  $-p + \sigma - 1$ . Therefore, the Fast2Sum Num, which implies

$$|v_h| \leq x + y - t_\ell.$$

$$\left(1 - \frac{\sigma}{2}\right)u \leq 1 + \frac{u}{2}.$$

bounds on  $|t_\ell|$  and  $|v_\ell|$ , we obtain:

$$|e_2| \leq \frac{1}{2}u^2p(t_\ell + v_\ell) \leq \frac{1}{2}u^2p \left(u^2 + \frac{u}{2}\right) = \frac{u^2}{2} \quad (8)$$

and

$$|e_1| \leq \frac{1}{2}u^2p \left[ \frac{1}{2}u^2p(x_\ell + y_\ell) + \frac{1}{2}u^2p \left( (x + y) + \frac{1}{2}u^2p(x_\ell + y_\ell) \right) \right]. \quad (9)$$

Lemma 2.1 and  $|s_h| \geq \sigma u$  imply that either  $x_h = 0$  or  $|v_h| = |\text{RN}(s_h + c)| = |\text{RN}(s_h + t_h)| \geq \sigma u^2$ . If  $x_h + t_h = 0$ , then  $v_h = v_\ell = 0$  and the sequel of the proof is straightforward. Therefore, in the following, we assume  $|v_h| \geq \sigma u^2$ .

Now,

- If  $|v_h| = \sigma u^2$ , then  $|v_\ell + t_\ell| \leq u|v_h| + u^2 = \sigma u^2 + u^2$ , which implies  $|w| = |\text{RN}(t_\ell + v_\ell)| \leq \sigma u^2 = |v_h|$ .
- If  $|v_h| > \sigma u^2$ , then, since  $v_h$  is a FP number,  $|v_h|$  is larger than or equal to the FP number immediately above  $\sigma u^2$ , which is  $\sigma(1 + 2u)\sigma^2$ . Hence  $|v_h| \geq \sigma u^2/(1 - u)$ , so  $|v_h| \geq u \cdot |v_h| + u^2 \geq |v_\ell| + |t_\ell|$ . So,  $|w| = |\text{RN}(t_\ell + v_\ell)| \leq |v_h|$ .

All cases, Fast2Sum introduces no error at line 6 of the algorithm, and we have

$$z_h + z_\ell = v_h + w = x + y + e_2. \quad (10)$$

Using Equation (10) and the bound  $u^2/2$  on  $|e_2|$  to get a relative error bound would result in a large bound, because  $x + y$  may be small. However, when  $x + y$  is very small, some simplification occurs thanks to Sterbenz Lemma. First,  $x_h + y_h$  is a nonzero multiple of  $\sigma u$ . Hence, since  $|x_\ell + y_\ell| \leq (1 + \frac{\sigma}{2})u$ , we have  $|x_\ell + y_\ell| \leq \frac{1}{2}(x_h + y_h)$ . Let us now consider the two possible cases:

- If  $-\frac{3}{2}(x_h + y_h) \leq x_\ell + y_\ell \leq -\frac{1}{2}(x_h + y_h)$ , which implies  $-\frac{3}{2}x_h \leq t_h \leq -\frac{1}{2}x_h$ , then Sterbenz Lemma implies to the floating-point addition of  $s_h$  and  $c = t_h$ . Therefore line 4 of the algorithm results in  $v_h = s_h$  and  $v_\ell = 0$ . An immediate consequence is  $e_2 = 0$ , so  $x_h + z_\ell = v_h + w = x + y$ : the computation of  $x + y$  is errorless;

- If  $-\frac{1}{2}(x_h + y_h) < x_\ell + y_\ell \leq \frac{3}{2}(x_h + y_h)$ , then  $\frac{2}{3}(x_\ell + y_\ell) \leq \frac{2}{3}(x_h + y_h + x_\ell + y_\ell) = \frac{2}{3}(x + y)$ , and  $-\frac{1}{2}(x + y) < \frac{2}{3}(x_\ell + y_\ell)$ . Hence,  $|x_\ell + y_\ell| < |x + y|$ , so  $\text{ulp}(x_\ell + y_\ell) \leq \text{ulp}(x + y)$ . Combined with Equation (9), this gives

$$|e_1| \leq \frac{1}{2} \text{ulp} \left( \frac{3}{2} \text{ulp}(x + y) \right) \leq 2^{-p} \text{ulp}(x + y) \leq 2 \cdot 2^{-p} \cdot (x + y).$$

## 2. If $-x_h/2 < y_h \leq x_h$

Notice that we have  $x_h/2 < x_h + y_h \leq 2x_h$ , so  $x_h/2 \leq s_h \leq 2x_h$ . Also notice that  $-\dots$  have  $|x_\ell| \leq u$ .

- If  $\frac{1}{2} < x_h + y_h \leq 2 - 4u$ . Define

We have

When  $\sigma = 1$ , we i.  
 $x_h \leq 2 - 2u$  implies  $|y_\ell|$   
 $(1 + \sigma/2)u$ , therefore

Elementary calculus shows that fo

The bound (3) is probabl

Now,  $|x_\ell + t_h| \leq (1 + \sigma)u$ , so

$$|c| \leq (1 + \sigma)u \quad \text{and} \quad |e_1| \leq \sigma u^2. \quad (13)$$

Since  $s_h \geq 1/2$  and  $|c| \leq 3u$ , if  $p \geq 3$ , then Algorithm Fast2Sum introduces no error at line 4 of the algorithm, that is,

$$v_h + v_\ell = s_h + c.$$

Therefore  $|v_h + v_\ell| = |s_h + c| \leq \sigma(1 - 2u) + (1 + \sigma)u \leq \sigma$ . This implies

$$|v_h| \leq \sigma \quad \text{and} \quad |v_\ell| \leq \frac{\sigma}{2}u. \quad (14)$$

Thus  $|t_\ell + v_\ell| \leq u^2 + \frac{\sigma}{2}u$ , so

$$|w| \leq \frac{\sigma}{2}u + u^2 \quad \text{and} \quad |e_2| \leq \frac{\sigma}{2}u^2. \quad (15)$$

From Equations (11) and (13), we deduce  $s_h + c \geq \frac{\sigma}{2} - u(2\sigma + 1)$ , so  $|v_h| \geq \frac{\sigma}{2} - u(2\sigma + 1)$ . If  $p \geq 3$ , then  $|v_h| \geq |w|$ , so Algorithm Fast2Sum introduces no error at line 6 of the algorithm, that is,  $x_h + x_\ell = v_h + w$ .

Therefore,

$$x_h + x_\ell = x + y + \eta,$$

with  $|\eta| = |e_1 + e_2| \leq \frac{3\sigma}{2}u^2$ . Since

$$x + y \geq (x_h - u) + (y_h - u/2) > \begin{cases} \frac{1}{2} - \frac{1}{2}u & \text{if } \sigma = 1, \\ 1 - 4u & \text{if } \sigma = 2, \end{cases}$$

the relative error  $|\eta|/(x + y)$  is upper bounded by

$$\frac{3u^2}{1 - 4u}.$$

- If  $2 - 4u < x_h + y_h \leq 2x_h$ , then  $2 - 4u \leq s_h \leq \text{RN}(2x_h) = 2x_h \leq 4 - 4u$  and  $|x_\ell| \leq 2u$ . We have

$$t_h + t_\ell = x_\ell + y_\ell,$$

with  $|x_\ell + y_\ell| \leq 2u$ , hence  $|t_h| \leq 2u$ , and  $|t_\ell| \leq u^2$ . Now,  $|x_\ell + t_h| \leq 4u$ , so  $|c| \leq 4u$ , and  $|e_1| \leq 2u^2$ . Since  $s_h \geq 2 - 4u$  and  $|c| \leq 4u$ , if  $p \geq 3$ , then Algorithm Fast2Sum introduces no error at line 4 of the algorithm. Therefore,

$$v_h + v_\ell = s_h + c \leq 4 - 4u + 4u = 4,$$

so  $v_h \leq 4$  and  $|v_\ell| \leq 2u$ . Thus,  $|t_\ell + v_\ell| \leq 2u + u^2$ . Hence, either  $|t_\ell + v_\ell| < 2u$  and  $|e_2| \leq \frac{1}{2} \text{ulp}(t_\ell + v_\ell) \leq u^2$ , or  $2u \leq t_\ell + v_\ell \leq 2u + u^2$ , in which case  $w = \text{RN}(t_\ell + v_\ell) = 2u$  and  $|e_2| \leq u^2$ . In all cases,  $|e_2| \leq u^2$ . Also,  $s_h \geq 2 - 4u$  and  $|c| \leq 4u$  imply  $v_h \geq 2 - 8u$ , and  $|c| \leq 2u + u^2$  implies  $|w| \leq 2u$ . Hence if  $p \geq 3$ , then Algorithm Fast2Sum introduces no error at line 6 of the algorithm.

$$x_h + x_\ell = v_h + w = x + y + \eta,$$

with  $|\eta| = |e_1 + e_2| \leq 3u^2$ .

Since  $x + y \geq (x_h - u) + (y_h - u) > 2 - 6u$ , the relative error  $|\eta|/(x + y)$  is upper bounded by

$$\frac{3u^2}{2 - 6u}.$$

The largest bound obtained in the various cases we have analyzed is

$$\frac{3u^2}{1 - 4u}.$$

Elementary calculus shows that for  $u \in (0, 1/64]$  (i.e.,  $p \geq 6$ ) this is always less than  $3u^2 + 13u^3$ .  $\square$

The bound (3) is *provably* not optimal. The largest relative error we have obtain through many tests is around  $2.25 \times 2^{-29} = 2.25u^2$ . An example is the input values given in Equation (2), for which, with  $p = 53$  (binary64 arithmetic), we obtain a relative error equal to  $2.24999999999999956 \dots \times 2^{-116}$ .

## DW+DW: “accurate version”

So the theorem gives an error bound

$$\frac{3u^2}{1-4u} \simeq 3u^2 \dots$$

As said before, that theorem has an interesting history:

- the authors of the first paper where a bound was given (in 2000) claimed (without published proof) that the relative error was always  $\leq 2u^2$  (in binary64 arithmetic);
- when trying (without success) to prove their bound, we found an example with error  $\approx 2.25u^2$ ;
- we finally proved the theorem, and Laurence Rideau (Inria Nice) started to write a **formal proof in Coq**;
- of course, that led to finding a (minor) **flaw** in our proof. . .

## DW+DW: “accurate version”

- fortunately the flaw was quickly corrected (before final publication of the paper... Phew)!
- still, the gap between worst case found ( $\approx 2.25u^2$ ) and the bound ( $\approx 3u^2$ ) was frustrating, so I spent **months** trying to improve the theorem...
- and of course **this could not be done**: it was the worst case that needed spending time!
- we finally found that with

$$x_h = 1$$

$$x_\ell = u - u^2$$

$$y_h = -\frac{1}{2} + \frac{u}{2}$$

$$y_\ell = -\frac{u^2}{2} + u^3.$$

error  $\frac{3u^2 - 2u^3}{1 + 3u - 3u^2 + 2u^3}$  is attained. With  $p = 53$  (binary64 arithmetic), gives error  $2.999999999999999877875 \dots \times u^2$ .

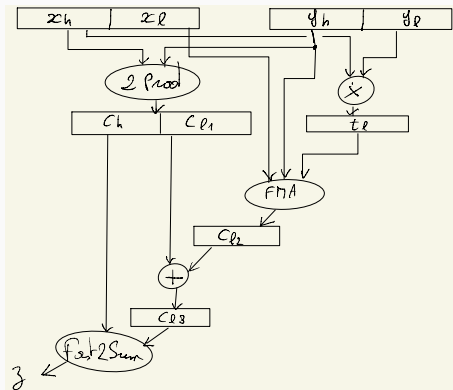
## DW+DW: “accurate version”

- We suspect the initial authors hinted their claimed bound by performing zillions of random tests
- in this domain, the worst cases are extremely unlikely: you must **build** them. Almost impossible to find them by chance.

- Product  $z = (z_h, z_l)$  of two DW numbers  $x = (x_h, x_l)$  and  $y = (y_h, y_l)$ ;
- several algorithms → tradeoff speed/accuracy. We just give one of them.

## DWTimesDW

- 1:  $(c_h, c_{l1}) \leftarrow 2\text{Prod}(x_h, y_h)$
- 2:  $t_l \leftarrow \text{RN}(x_h \cdot y_l)$
- 3:  $c_{l2} \leftarrow \text{RN}(t_l + x_l y_h)$
- 4:  $c_{l3} \leftarrow \text{RN}(c_{l1} + c_{l2})$
- 5:  $(z_h, z_l) \leftarrow \text{Fast2Sum}(c_h, c_{l3})$
- 6: **return**  $(z_h, z_l)$



We have

### Theorem 7 (Error bound for Algorithm DWTimesDW)

*If  $p \geq 5$ , the relative error of Algorithm DWTimesDW2 is less than or equal to*

$$\frac{5u^2}{(1+u)^2} < 5u^2.$$

and that theorem too has an interesting history!

- initial bound  $6u^2$ ;
- again, we tried formal proof... and it turned out the proof was based on a **wrong lemma**.

- after a few nights of very bad sleep, we found a turn-around. . . that also improved the bound !
- no proof of asymptotic optimality, but in binary64 arithmetic, we have examples with error  $> 4.98u^2$ ;
- (real consolation or lame excuse?) without the flaw, we would never have found the better bound;
- without the formal proof effort, the error would probably have remained unnoticed (in this case, without serious consequence since the property was true anyway, but. . .).



# Conclusion

- (almost) fully specified arithmetic: one can prove properties of (small enough) programs, and build algorithms;
- ongoing effort for also standardizing a **kernel of math functions** (at least exp, sin, cos, log);
- all of numerical computing is built from **computer arithmetic**: it must be reliable;
- for some algorithms (e.g., DW arithmetic, FP division algorithms) the “paper proofs” are terrible: use of **formal proof** and **computer algebra**.