# La malédiction des preuves longues et ennuyeuses
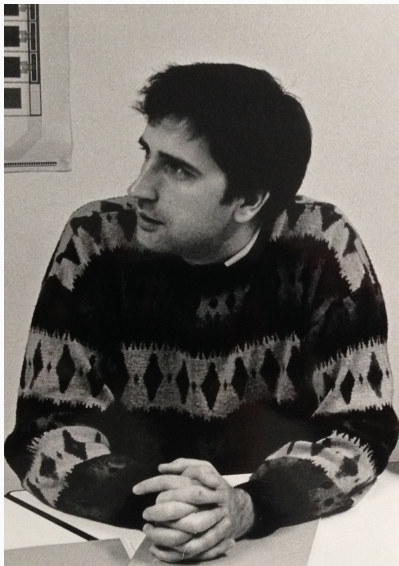
Jean-Michel Muller

avec des contributions de M. Joldes, V. Popescu, L. Rideau, B. Salvy

*Numération, Algorithmes et Cryptographie*, Paris, Février 2024

CNRS - Laboratoire LIP

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

- je suis chargé d'organiser une séance photos pour réaliser la première plaquette du LIP;
- on fait appel aux top-models locaux.

Outre JCB, on reconnait Christine Paulin et Pierre Fraigniaud.

We wish to prove error bounds of medium-size, "atomic" algorithms in FP arithmetic, but...

- error bounds... *for what purpose?*
- proofs... *for what purpose?*

A Floating-Point number (FPN) $x$ is represented by two integers:

- Floating-Point number:
$$x = \left( \frac{M}{2^{p-1}} \right) \cdot 2^e$$
  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.

- largest finite FPN $\Omega = 2^{e_{\max}+1} - 2^{e_{\max}-p+1}$;

- unit roundoff: $u = 2^{-p}$.

# Error bounds. . .

- FP system parametered by precision $p$ or unit round-off $u = 2^{-p}$;
- for a given algorithm relative error bound $\mathcal{B}(u)$;
- the (most likely unknown) worst case error is $\mathcal{W}(u)$.

The bound $\mathcal{B}$ is

- certain (for $u \leq u_0$) if $\mathcal{W}(u) \leq \mathcal{B}(u)$ for $u \leq u_0$;
- asymptotically optimal if $\mathcal{W}(u)/\mathcal{B}(u) \to 1$ as $u \to 0$;
- tight (for $u \leq u_0$) if $\mathcal{W}(u)$ is close to $\mathcal{B}(u)$ for $u \leq u_0$.

# Mais pourquoi veut-il tant calculer des bornes d'erreur ?

- choice between different algorithms:
  - an informed choice of the algorithm that has the best balance performance/accuracy requires tight bounds;
  - certainty not that important;
- guaranteeing the behavior of a possibly critical software:
  - need to prove that the error is not $\geq$ some threshold
  - $\rightarrow$ certainty important,
  - tightness not always needed;
- fully validated set of "atomic" algorithms:
  - the most common transcendental functions such as exp, ln;
  - simple algebraic functions such as $1/\sqrt{x}$, $\text{hypot}(x, y) = \sqrt{x^2 + y^2}$
  
  are "basic building blocks" of numerical computing : users expect same behavior as for $+$, $-$, $\times$, $\div$, $\sqrt{}$.
- $\rightarrow$ having bounds that are both certain and tight is desirable.

# Proofs. . . for what purpose ?

1. to check, by following the proof step by step, that the claimed property holds;

2. to have a deep and "global" understanding of what is behind the claimed property.

Rather antagonistic goals:

- goal 1 requires many details,
- goal 2 needs a focus on the "big things" (hence many "*without loss of generality. . .*" or "*the second case is similar*").

In general our "paper proofs" are in between: is this the right solution?

# The two examples considered in this talk

- "double word" arithmetic: formal proofs helped to
  - strengthen claimed results,
  - improve them,
  - find (hmmm...embarassing) bugs.
- hypotenuse function $\sqrt{x^2 + y^2}$: computer algebra helped to
  - obtain tight bounds,
  - explore several variants.

Before presenting that: additional notions on FP arithmetic (roundings, error-free transforms, double-word arithmetic).

# Correct rounding, ulp (unit in the last place)

- 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 $c = \text{RN}(a+b)$.

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.

# Relative error due to rounding

- if $2^{e_{\min}} \leq |x| \leq \Omega$, then

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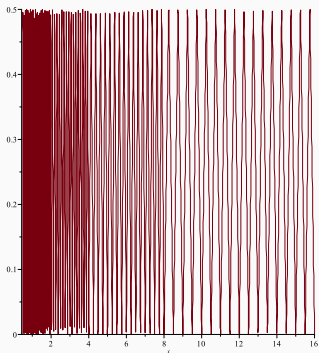
  therefore,

$$|x - \text{RN}(x)| \leq u \cdot |x|, \tag{1}$$
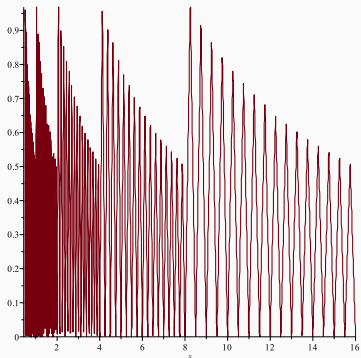
  with $u = 2^{-p}$ . Hence the relative error

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

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

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

Absolute error (in ulps) of rounding to nearest a real number $x \in [1/2, 16]$, assuming a binary FP "toy" system with $p = 5$.



Relative error (in multiples of $u = 2^{-p}$) of rounding to nearest a real number $x \in [1/2, 16]$, assuming a binary FP "toy" system with $p = 5$.

The relative error bound $u$ is tight only slightly above a power of 2.

## Error-free transforms and double-word arithmetic

2Sum($a, b$)

$\quad s \leftarrow \text{RN}(a + b)$
$\quad a' \leftarrow \text{RN}(s - b)$
$\quad b' \leftarrow \text{RN}(s - a')$
$\quad \delta_a \leftarrow \text{RN}(a - a')$
$\quad \delta_b \leftarrow \text{RN}(b - b')$
$\quad t \leftarrow \text{RN}(\delta_a + \delta_b)$
$\quad$**return** $(s, t)$

Fast2Sum($a, b$)

$\quad s \leftarrow \text{RN}(a + b)$
$\quad z \leftarrow \text{RN}(s - a)$
$\quad t \leftarrow \text{RN}(b - z)$
$\quad$**return** $(s, t)$

Barring overflow:

- the pair $(s, t)$ returned by 2Sum satisfies $s = \text{RN}(a + b)$ and $t = (a + b) - s$;

- if $|a| \geq |b|$ then the pair $(s, t)$ returned by Fast2Sum satisfies $s = \text{RN}(a + b)$ and $t = (a + b) - s$.

Such algorithms: Error-free transforms.

13

2Prod($a, b$)

$\quad \pi \leftarrow \text{RN}\,(ab)$

$\quad \rho \leftarrow \text{RN}\,(ab - \pi)$

$\quad$ **return** $(\pi, \rho)$

Barring overflow, if the exponents $e_a$ and $e_b$ of $a$ and $b$ satisfy $e_a + e_b \geq e_{\min} + p - 1$ then then the pair $(\pi, \rho)$ returned by Fast2Sum satisfies $\pi = \text{RN}\,(ab)$ and $\rho = (ab) - \pi$.

- Fast2Sum, 2Sum and 2Prod: return $x$ represented by a pair $(x_h, x_\ell)$ of FPN such that $x_h = \text{RN}\,(x)$ and $x = x_h + x_\ell$;
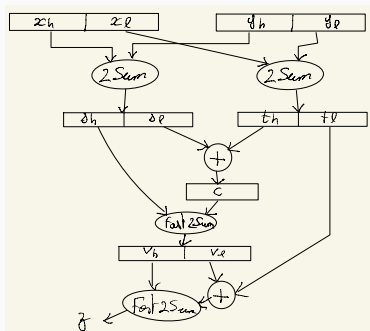
- Such pairs: double-word numbers (DW).

Algorithms for manipulating DW suggested by various authors since 1971.

# DW+DW: "accurate version"

Sum of two DW numbers. There also exists a "quick & dirty" algorithm, but its relative error is unbounded.

## DWPlusDW

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

We have (after a rather tedious proof):

**Theorem (Joldeş, Popescu, M., 2017)**

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

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

That theorem has an interesting history...

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

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

Li et al. (2000, 2002) claim that in binary64 arithmetic ($p = 53$) the relative error of Algorithm 6 is upper bounded by $2 \cdot 2^{-106}$. This bound is incor...

then the relative error of Algorithm 6 is

2.2499999...

Note that this example is somehow "gene...
$2^p - 1$, $x_\ell = -(2^p - 1) \cdot 2^{-p-1}$, $y_h = -(2^p - 5)/2$,...
that is asymptotically equivalent (as $p$ goes to infini...

Now let us try to find a relative error bound. We are g...

THEOREM 3.1. *If $p \ge 3$, then the relative error of Algorithm 6 u,*
*by*

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

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

Note that the conditions on $p$ ($p \ge 3$ for the bound (3) to hold, $p \ge 6$ 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 2Sum$(x_\ell, y_\ell)$, which is equal to $x + y$. And the... the computation is errorless Now, without loss of generality, we assume $1 \le x_h \le 2$, $x \ge |y|$ (which implies $x_h \ge |y_h|$), and $x_h + y_h$ nonzero. Notice that $1 \le x_h <$... 2 implies $1 \le x_h \le 2 - 2u$, since $x_h$ is a FP number.

Define $e_1$ as the error committed at Line 3:

$$e_1 = c - (s_\ell + t_h) \quad (4)$$

and $e_2$ as the error committed at Line 5:

$$e_2 = w - (t_\ell + v_\ell). \quad (5)$$

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

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

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

$$\sigma u \le s_h \le 1 - \sigma u. \quad (6)$$

We have $|x_\ell| \le u$ and $|y_\ell| \le \frac{\sigma}{2} u$, so

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

From Equation (6), we deduce that the floating-point exponent of $s_h$ 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

$$v_h + v_\ell = s_h + c = s_h + t_h = x + y - t_\ell.$$

Equations (6) and (7) imply

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

so $|v_h| \le 1$ and $|v_\ell| \le \frac{u}{2}$. From the bounds on $|t_\ell|$ and $|v_\ell|$, we obtain:

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

and

$$|e_2| \le \frac{1}{2}\text{ulp}\left[\frac{1}{2}\text{ulp}(x_\ell + y_\ell) + \frac{1}{2}\text{ulp}\left((x + y) + \frac{1}{2}\text{ulp}(x_\ell + y_\ell)\right)\right]. \quad (9)$$

Lemma 2.1 and 3.2 $\ge \sigma u$ imply that either $s_h + t_h = 0$, or $|v_h| = |$RN$(s_h + c)| = |$RN$(s_h + t_h)| \ge \sigma u^2$. If $s_h + t_h = 0$, then $v_h = v_\ell = 0$ and the sequel of the proof is straightforward. Therefore, in the following, we assume $|v_h| \ge \sigma u^2$.

Now,

- If $|v_h| = \sigma u^2$, then, $|t_\ell + v_\ell| \le u|v_h| + u^2 = \sigma u^2 + u^2$, which implies $|w| = |$RN$(t_\ell + v_\ell)| \le \sigma u^2 = |v_h|$;
- 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 $(1 + 2u)\sigma u^2$. Hence $|v_h| \ge \sigma u^2/(1 - u)$, so $|v_\ell| = u \cdot |v_h| + u^2 \dots \sigma u^2 \le |v_\ell|$. So, $|w| = |$RN$(t_\ell + v_\ell)| \le |v_h|$.

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

$$z_h + z_\ell = v_h + w = x + y + e_2.$$

Directly 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 $|x_\ell + y_\ell| \le \frac{1}{2}\sigma u$, we have $|x_\ell + y_\ell| \le \frac{1}{2}(x_h + y_h)$. Let us now consider the two possible cases:

- If $-\frac{1}{2}(x_h + y_h) \le x_\ell + y_\ell \le -\frac{1}{2}(x_h + y_h)$, which implies $-\frac{1}{2}s_h \le t_h \le -\frac{1}{2}s_h$, 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 $v_\ell = 0$. An immediate consequence is $e_2 = 0$, so $z_h + z_\ell = v_h + w = x + y$ the computation of $x + y$ is errorless;

**ALGORITHM 6 – AccurateDWPlusDW**$(x_h, x_\ell, y_h, y_\ell)$. Calculation of $(z_h, z_\ell) = (y_h, y_\ell)$ in binary, precision-$p$, floating-point arithmetic.

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

Li et al. (2000, 2002) claim that
is upper bounded by $2 \cdot 2^{-106}$. This

==d the sequel of the proof is straightforward. Th==

then the relative error of Algorithm 6 is

$$2.2499999999999956\ldots$$

Note that this example is somehow "generic": In precision-$p$ FP arithmetic, $2^p - 1, x_\ell = -(2^p - 1) \cdot 2^{-p-5}, y_h = -(2^p - 5)/2,$ and $y_\ell = -(2^p - 1) \cdot 2^{-106}$ 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. We are going to show the following result.

**THEOREM 3.1.** If $p \ge 2$, then the relative error of Algorithm 6 (AccurateDWPlusDW) is bounded by

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

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

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

**PROOF.** First

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

and $z_2$ in the

.ame applies to the floating-poi

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

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

-point exponent of $s_h$ is at least $-p + \sigma - 1$. Therefore, the Fast2Sum algorithm, which implies

$$= x + y - t_\ell.$$

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

$$|z_\ell| \le \frac{1}{2}\mathrm{ulp}(t_\ell + v_\ell) \le \frac{1}{2}\mathrm{ulp}\left(u^2 + \frac{u}{2}\right) \le \frac{u^2}{2} \qquad (8)$$

and

$$|z_\ell| \le \frac{1}{2}\mathrm{ulp}\left[\frac{1}{2}\mathrm{ulp}(x_\ell + y_\ell) + \frac{1}{2}\mathrm{ulp}\left((x + y) + \frac{1}{2}\mathrm{ulp}(x_\ell + y_\ell)\right)\right]. \qquad (9)$$

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

Now,

• If $|v_h| = \sigma u^2$, then $|v_\ell + t_\ell| \le u|v_h| + u^2 = \sigma u^3 + u^2$, which implies $|w| = |\mathrm{RN}(t_\ell + v_\ell)| \le \sigma u^2 = |v_h|$, 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 $\sigma(1 + 2u)u^2$. Hence $|v_h| \ge \sigma u^2/(1 - u)$, so $u \cdot |v_h| + u^2 \ge 2\sigma u^3 + |t_\ell|$. So, $|w| = |\mathrm{RN}(t_\ell + v_\ell)| \le |v_h|$.

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

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

By using Equation (10) and the bound $u^2/2$ on $|z_\ell|$ 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 $|x_h + y_\ell| \le \frac{1}{2}(x_h + y_h)$, we have $|x_\ell + y_\ell| \le \frac{1}{2}(x_h + y_h)$. Let us now consider the two possible cases:

==• $-\frac{3}{2}(x_h + y_h) \le x_\ell + y_\ell \le -\frac{1}{2}(x_h + y_h)$, which implies $-\frac{3}{2}s_h \le t_h \le -\frac{1}{2}s_h$, 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 $v_\ell = 0$. An immediate consequence is $z_\ell = 0$, so $z_h + z_\ell = v_h + w = x + y$; the computation of $x + y$ is errorless.

18

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

$$|\epsilon_2| \leq \frac{1}{2}\mathrm{ulp}\left(\frac{3}{2}\mathrm{ulp}(x+y)\right) \leq 2^{-p}\mathrm{ulp}(x+y) \leq 2 \cdot 2^{-2p} \cdot (x+y).$$

**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 we have $|x_\ell| \leq u$.

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

We have

When $\sigma = 1$, we have

$x_h \leq 2 - 2u$ implies $|y_\ell$,
$(1 + \sigma/2)u$, therefore

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

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

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

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

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

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

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

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 $|v_h| \geq u$, so Algorithm Fast2Sum introduces no error at line 6 of the algorithm, that is, $z_h + z_\ell = v_h + w$. Therefore,

$$z_h + z_\ell = x + y + \eta,$$

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

$$x + y \geq (x_h - u) + (y_h - u/2) > \begin{cases} \frac{1}{2} - \frac{3}{2}u & \text{if} \quad \sigma = 1, \\ 1 - 4u & \text{if} \quad \sigma = 2, \end{cases}$$

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 $2 - 4u \leq s_h \leq \mathrm{RN}(2x_h) = 2x_h \leq 4 - 4u$ and $|s_\ell| \leq 2u$. 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, $|s_\ell + t_h| \leq 4u$, so $|c| \leq 4u$, and $|\epsilon_1| \leq 2u^2$. 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 $|e_2| \leq \frac{1}{2}\mathrm{ulp}(t_\ell + v_\ell) \leq u^2$, or $2u \leq t_\ell + v_\ell \leq 2u + u^2$, in which case $w = \mathrm{RN}(t_\ell + v_\ell) = 2u$ 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 $v_h \geq 2 - 8u$, so $-1 \leq 2u + w \leq 2u$. Hence if $p \geq 3$, then Algorithm Fast2Sum introduces no error at line 6 of the algorithm, this gives

$$z_h + z_\ell = v_h + w = x + y + \eta,$$

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 $3u^2 + 13u^3$.   □

The bound (3) is probably not optimal. The largest relative error we have obtain through many tests is around $2.25 \times 2^{-2p} = 2.25u^2$. An example is the input values given in Equation (2), for which, with $p = 53$ (binary64 arithmetic), we obtain a relative error equal to $2.24999999999999956\cdots \times 2^{-106}$.

**Elementary calculus shows that fo**

**The bound (3) is probably**

19

So the theorem gives an error bound

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

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 started to write a formal proof in Coq;
- of course, that led to finding a (minor) flaw in our proof. . .

(I hate Coq people)

# 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

$$
\begin{aligned}
x_h &= 1 \\
x_\ell &= u - u^2 \\
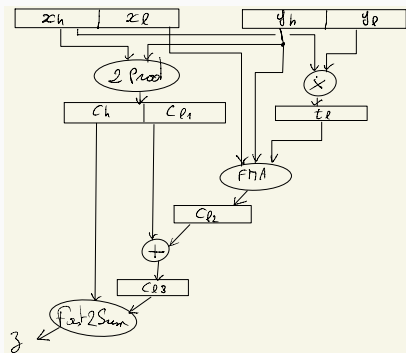y_h &= -\tfrac{1}{2} + \tfrac{u}{2} \\
y_\ell &= -\tfrac{u^2}{2} + u^3.
\end{aligned}
$$

error $\frac{3u^2 - 2u^3}{1 + 3u - 3u^2 + 2u^3}$ is attained. With $p = 53$ (binary64 arithmetic), gives error $2.99999999999999877875\cdots \times u^2$.

# DW × DW

- Product $z = (z_h, z_\ell)$ of two DW numbers $x = (x_h, x_\ell)$ and $y = (y_h, y_\ell)$;
- several algorithms $\rightarrow$ tradeoff speed/accuracy. We just give one of them.

### DWTimesDW

1: $(c_h, c_{\ell 1}) \leftarrow 2\text{Prod}(x_h, y_h)$
2: $t_\ell \leftarrow \text{RN}(x_h \cdot y_\ell)$
3: $c_{\ell 2} \leftarrow \text{RN}(t_\ell + x_\ell y_h)$
4: $c_{\ell 3} \leftarrow \text{RN}(c_{\ell 1} + c_{\ell 2})$
5: $(z_h, z_\ell) \leftarrow \text{Fast2Sum}(c_h, c_{\ell 3})$
6: **return** $(z_h, z_\ell)$

We have

**Theorem (M. and Rideau, 2022)**

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

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

and that theorem too has an interesting (hmmm... a bit more annoying?) history!

- in 2017, I participated to the proof of an initial relative error bound $6u^2$;
- again, Laurence tried translating the proof in Coq... and it turned out the proof was based on a wrong lemma (and this was *after* publication).

(what did I say about Coq people?)

- after a few nights of bad sleep, turn-around... that also improved the bound: $6u^2 \to 5u^2$!

- no proof of asymptotic optimality, but in binary64 arithmetic, we have examples with error $> 4.98u^2$;

- *real consolation or lame excuse?* Maybe without the flaw, we would never have found the better bound.

Full set of validated DW algorithms for the arithmetic operations and the square root (M. and Rideau, 2022; Lefèvre, Louvet, Picot, M. and Rideau, 2023).

That class of algorithms really needs formal proof:

- Proofs have too many subcases to be certain you have not forgotten one;
- they are boring: almost nobody reads them.

Alternate—or complementary—solution? try to automatically compute bounds:

- short-term goal: limit human intervention (and therefore, human error); and make simpler the exploration of many variants;
- long-term goal: bounds correct by construction.

# An example: hypotenuse function $\sqrt{x^2 + y^2}$

NaiveHypot

1: $s_x \leftarrow \text{RN}(x^2)$
2: $s_y \leftarrow \text{RN}(y^2)$
3: $\sigma \leftarrow \text{RN}(s_x + s_y)$
4: $\rho_1 = \text{RN}(\sqrt{\sigma})$

- classical relative error bound $2u + \mathcal{O}(u^2)$;
- refinement: $2u$ (Jeannerod & Rump);
- asymptotically optimal (Jeannerod, M., Plet).

Major drawback: "spurious" overflow/underflow

Examples in binary64/double precision arithmetic ($p = 53$):

- if $x = 2^{600}$ and $y = 0$, returned result $+\infty$, exact result $2^{600}$;
- if $x = 65 \times 2^{-542}$ and $y = 72 \times 2^{-542}$, returned result $96 \times 2^{-542}$, exact result $97 \times 2^{-542}$.

$\Rightarrow$ need to scale the operands.

1: **if** $|x| < |y|$ **then**
2:     swap $(x, y)$
3: **end if**
4: $r \leftarrow \mathsf{RN}\,(y/x)$
5: $t \leftarrow \mathsf{RN}\,(1 + r^2)$
6: $s \leftarrow \mathsf{RN}\,(\sqrt{t})$
7: $\rho_2 = \mathsf{RN}\,(|x| \cdot s)$

- relative error bounded by $\frac{5}{2}u + \frac{3}{8}u^2$;
- asymptotically optimal.

$\Rightarrow$ avoiding spurious overflow has a significant cost in terms of accuracy.

Improvements?

# Simple scaling with compensation (Nelson Beebe, 2017)

1: **if** $|x| < |y|$ **then**
2:     swap$(x, y)$
3: **end if**
4: $r \leftarrow$ RN$(y/x)$
5: $t \leftarrow$ RN$(1 + r^2)$
6: $s \leftarrow$ RN$(\sqrt{t})$
7: $\epsilon \leftarrow$ RN$(t - s^2)$
8: $c \leftarrow$ RN$(\epsilon/(2s))$
9: $\nu \leftarrow$ RN$(|x| \cdot c)$
10: $\rho_3 \leftarrow$ RN$(|x| \cdot s + \nu)$

- this version: requires an FMA;
- one Newton-Raphson iteration;
- relative error bound $\frac{8}{5}u + \frac{7}{5}u^2$ (Salvy & M., 2023);
- known case with error $1.5999739u$ in binary64 FP arithmetic.

# The various bounds obtained

| Algorithm | reference | error bound | condition | status |
|---|---|---|---|---|
| Naive | folklore | $2u - \frac{8}{5}(9 - 4\sqrt{6})u^2$ | $p \geq 2$ | asympt. optimal |
| Simple scaling | folklore | $\frac{5}{2}u + \frac{3}{8}u^2$ | $p \geq 2$ | asympt. optimal |
| Scaling w. compensation | N. Beebe (2017) | $\frac{8}{5}u + \frac{7}{5}u^2$ | $p \geq 4$ | sharp |
| Borges "fused" | C. Borges (2020) | $u + 14u^2$ | $p \geq 5$ | asympt. optimal |
| Kahan | W. Kahan (1987) | $1.5355u + \mathcal{O}(u^2)$ ? | TBD | a bit loose |

# Goal: tight and certain relative error bounds

- Programs that at step $k$ have an instruction of the form

$$\texttt{x\_k = x\_i op x\_j} \quad \text{or} \quad \texttt{x\_k = sqrt(x\_i)}$$

  where op is +, -, * or /, and x_i and x_j are either precomputed values or input values ($i, j < k$);

- Computed values:

$$x_k = \text{RN}\left(x_i \ \text{op} \ x_j\right) \quad \text{or} \quad x_k = \text{RN}\left(\sqrt{x_i}\right);$$

- basic relations:

$$\begin{aligned} x_k &= x_i \ \text{op} \ x_j \pm \tfrac{1}{2} \, \text{ulp}\left(x_i \ \text{op} \ x_j\right), \\ x_k &= (x_i \ \text{op} \ x_j)(1 + \epsilon), \quad \text{with } |\epsilon| \leq \tfrac{u}{1+u} < u. \end{aligned} \qquad (3)$$

  (or the same with $\sqrt{x_i}$)

Optimisation problem: find the maximum and the minimum of the quantity $\rho / \sqrt{x^2 + y^2} - 1$ in the region defined by the equalities and inequalities obtained from analyzing the program (e.g., (3)) $\rightarrow$ Algebraic bound.

# Goal: tight and certain relative error bounds

- Computed values

$$x_k = \mathrm{RN}\left(x_i \ \mathrm{op} \ x_j\right) \quad \text{or} \quad x_k = \mathrm{RN}\left(\sqrt{x_i}\right);$$

- we compare the computed values $x_k$ with the exact values:

$$x_k^* = x_i^* \ \mathrm{op} \ x_j^* \quad \text{or} \quad x_k^* = \sqrt{x_i^*};$$

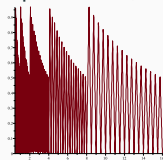(initial values: $x_i = x_i^*$ for $i \leq 0$).

- The analysis consists in iteratively computing relative error bounds $\epsilon_k^\ell(u)$ and $\epsilon_k^r(u)$ such that (here, for positive $x_k$ and $x_k^*$)

$$x_k^* \left(1 - \epsilon_k^\ell(u)\right) \leq x_k \leq x_k^* \left(1 + \epsilon_k^r(u)\right); \tag{4}$$

# Goal: tight and certain relative error bounds

- with care, iteratively computing bounds of the form (4), using at each step the "basic relations" (3) is not so difficult;

- ending up with a tight bound is difficult. Two reasons:
  - requires existence of input values for which the individual rounding errors attain their maximum (with the right sign) at each operation.

    $\rightarrow$ Not always possible: Correlations. $3 \cdot (x \cdot y)$, one cannot have both $(x \cdot y)$ and $3 \cdot (x \cdot y)$ very slightly above a power of 2;

    

    (and, indeed, $3 \cdot (x \cdot y)$ more accurate than $(3 \cdot x) \cdot y$)
  - the "basic relations" (3) are not the last word: additional properties (e.g., Sterbenz Lemma) specific to FP arithmetic.

1: **if** $|x| < |y|$ **then**
2:      swap$(x, y)$
3: **end if**
4: $r \leftarrow$ RN $(y/x)$
5: $t \leftarrow$ RN $(1 + r^2)$
6: $s \leftarrow$ RN $(\sqrt{t})$
7: $\epsilon \leftarrow$ RN $(t - s^2)$
8: $c \leftarrow$ RN $(\epsilon/(2s))$
9: $\nu \leftarrow$ RN $(|x| \cdot c)$
10: $\rho_3 \leftarrow$ RN $(|x| \cdot s + \nu)$

Simplification: $x \geq y > 0$

1: $r \leftarrow \text{RN}(y/x)$
2: $t \leftarrow \text{RN}(1 + r^2)$
3: $s \leftarrow \text{RN}(\sqrt{t})$
4: $\epsilon \leftarrow \text{RN}(t - s^2)$
5: $c \leftarrow \text{RN}(\epsilon/(2s))$
6: $\nu \leftarrow \text{RN}(x \cdot c)$
7: $\rho_3 \leftarrow \text{RN}(x \cdot s + \nu)$

Main idea: Newton-Raphson iteration

$$\frac{\epsilon}{2s} + s = \frac{t - s^2}{2s} + s = \sqrt{t} + \frac{(s - \sqrt{t})^2}{2s},$$

so that

$$\left(\frac{\epsilon}{2s} + s\right) - \sqrt{t} = \frac{(s - \sqrt{t})^2}{2s}.$$

1: $r \leftarrow \text{RN}(y/x)$
2: $t \leftarrow \text{RN}(1 + r^2)$
3: $s \leftarrow \text{RN}(\sqrt{t})$
4: $\epsilon \leftarrow \text{RN}(t - s^2)$
5: $c \leftarrow \text{RN}(\epsilon/(2s))$
6: $\nu \leftarrow \text{RN}(x \cdot c)$
7: $\rho_3 \leftarrow \text{RN}(x \cdot s + \nu)$

- define $\alpha$ by $y = \alpha x$, so that $r = \text{RN}(\alpha)$;

- $r = \alpha + u\epsilon_r$, with

$$|\epsilon_r| \leq \begin{cases} \frac{1}{4}, & \text{if } \alpha \leq 1/2, \\ \frac{1}{2}, & \text{if } \alpha > 1/2. \end{cases}$$

- $t = 1 + r^2 + u\epsilon_t$, with $|\epsilon_t| \leq 1$ (comes from $1 + r^2 \leq 2$);

- $s = \sqrt{t} + u\epsilon_s$, with $|\epsilon_s| \leq 1$ (comes from $t < 2$);

- $\epsilon = t - s^2$ (comes from Sterbenz Lemma).

35

$$\left|\frac{\epsilon}{2s}\right| = \left|\frac{t-s^2}{2s}\right|$$

$$= \left|\frac{(s-u\epsilon_s)^2-s^2}{2s}\right| \tag{5}$$

$$= \left|-u\epsilon_s + \frac{u^2\epsilon_s^2}{2s}\right| \le u + \frac{u^2}{2}.$$

1: $r \leftarrow \text{RN}(y/x)$
2: $t \leftarrow \text{RN}(1+r^2)$
3: $s \leftarrow \text{RN}(\sqrt{t})$
4: $\epsilon \leftarrow \text{RN}(t-s^2)$
5: $c \leftarrow \text{RN}(\epsilon/(2s))$
6: $\nu \leftarrow \text{RN}(x \cdot c)$
7: $\rho_3 \leftarrow \text{RN}(x \cdot s + \nu)$

- If $|\epsilon/(2s)| \le u$ then the error committed by rounding $\frac{\epsilon}{2s}$ to nearest is $\le u^2/2$;

- If $|\epsilon/(2s)| > u$, then since the FPN above $u$ is $u+2u^2$, (5) implies $\text{RN}(\epsilon/(2s)) = \pm u$ $\Rightarrow$ again the rounding error is $\le u^2/2$.

Hence in all cases, $|c| \le u$ and

$$c = \frac{\epsilon}{2s} + \epsilon_c \frac{u^2}{2},$$

with $|\epsilon_c| \le 1$.

1: $r \leftarrow \text{RN}(y/x)$
2: $t \leftarrow \text{RN}(1 + r^2)$
3: $s \leftarrow \text{RN}(\sqrt{t})$
4: $\epsilon \leftarrow \text{RN}(t - s^2)$
5: $c \leftarrow \text{RN}(\epsilon/(2s))$
6: $\nu \leftarrow \text{RN}(x \cdot c)$
7: $\rho_3 \leftarrow \text{RN}(x \cdot s + \nu)$

- $\nu = xc(1 + u\epsilon_\nu)$ with $|\epsilon_\nu| \leq 1/(1 + u)$;
- $\rho_3 = (\nu + xs)(1 + u\epsilon_\rho)$ with $|\epsilon_\rho| \leq 1/(1 + u)$;

Putting all this together:

$$
\begin{aligned}
\rho_3 &= (\nu + xs)(1 + u\epsilon_\rho), \\
&= x\left((-u\epsilon_s + \tfrac{u^2}{2}(\epsilon_c + \epsilon_s^2/s))(1 + u\epsilon_\nu) + \sqrt{t} + u\epsilon_s\right)(1 + u\epsilon_\rho), \\
&= x\left(\sqrt{t} + \tfrac{u^2}{2}((\epsilon_c + \epsilon_s^2/s)(1 + u\epsilon_\nu) - 2\epsilon_s\epsilon_\nu)\right)(1 + u\epsilon_\rho) \\
&= x\sqrt{1+r^2}\sqrt{1 + \tfrac{u\epsilon_t}{1+r^2}}\left(1 + \tfrac{u^2}{2\sqrt{t}}((\epsilon_c + \epsilon_s^2/s)(1 + u\epsilon_\nu) - 2\epsilon_s\epsilon_\nu)\right)(1 + u\epsilon_\rho),
\end{aligned}
$$

## Lemma

*The relative error of the algorithm is*

$$
\begin{aligned}
R &= \sqrt{1 + \tfrac{r^2 - \alpha^2}{1+\alpha^2}}\sqrt{1 + \tfrac{u\epsilon_t}{1+r^2}} \\
&\quad \times \left(1 + \tfrac{u^2}{2\sqrt{t}}((\epsilon_c + \epsilon_s^2/s)(1 + u\epsilon_\nu) - 2\epsilon_s\epsilon_\nu)\right)(1 + u\epsilon_\rho) - 1,
\end{aligned}
$$

*Moreover, $|\epsilon_s|, |\epsilon_t|, |\epsilon_c|$ are bounded by 1 and $|\epsilon_\nu|$ and $|\epsilon_\rho|$ by $1/(1+u)$.*

- linear term

$$\left( \frac{2\alpha\epsilon_r + \epsilon_t}{2(1 + \alpha^2)} + \epsilon_\rho \right) \cdot u$$

  - increasing function of $\epsilon_r, \epsilon_t$ and $\epsilon_\rho$,
  - $\epsilon_r \leq 1/4$ if $\alpha \leq 1/2$, $\epsilon_r \leq 1/2$ otherwise,
  - $\epsilon_t, \epsilon_\rho \leq 1$

  $\rightarrow$ max. value 8/5;

- show that for $u \in [0, 1/2]$,

$$\frac{\partial R}{\partial \epsilon_\rho} \geq 0, \quad \frac{\partial R}{\partial \epsilon_t} \geq 0, \quad \frac{\partial R}{\partial \epsilon_r} \geq 0, \quad \frac{\partial R}{\partial \epsilon_c} \geq 0.$$

$\rightarrow$ it suffices to consider the extremum values of $\epsilon_\rho$, $\epsilon_t$, $\epsilon_r$, and $\epsilon_c$;

- process the cases $\alpha < 1/2$ and $1/2 \leq \alpha \leq 1$ separately;

- in each case, lower and upper bound on $R$. . .

## Theorem

*Assuming $u \leq 1/16$ (i.e., $p \geq 4$), the relative error of Beebe's algorithm is bounded by*

$$\chi_4(u) = (1 + 2u) \sqrt{\frac{1 + u/5}{1 + u}} - 1 + u^2 \frac{(1 + 2u)^2}{(1 + u)^2} \left( \frac{\sqrt{5}}{5} + \frac{1}{\frac{5 \sqrt{(1+u)\left(1+\frac{u}{5}\right)}}{2} - u} + \frac{2\sqrt{5}}{5(1 + 2u)} \right),$$

$$\leq \frac{8}{5} u + \frac{7}{5} u^2.$$

How do we publish a proof? Have a Maple worksheet publicly available and just get a rough sketch (similar to these slides) in a paper?

## And the other algorithms?

- Another algorithm due to Borges: really painful. . . but we managed to obtain the result;
- Kahan's algorithm. . . the first result was:



We are succeeding (paper to come soon)
It seems we are approaching a limit. . .

. . . and again, as for DW arithmetic, if we fully "expand" the proofs they are terrible (probably unpublishable).

## But, really, what were we trying to do?

- obtain the best "algebraic bound": the best one could deduce from the individual bounds on the rounding errors of the operations and a few properties such as Sterbenz Lemma;

- but when the algorithms become complex, does that bound remain tight?

    - we have seen: correlations;
    - even without correlations: tightness requires that for each operation the maximum error is almost reached, with the right signs;
    - in general: probability of this decreases exponentially with number of operations;

$\rightarrow$ Rule of thumb: when the number of operations is no longer small in front of $p$, little hope of having a worst-case error close to the algebraic bound.

# Conclusion

- formal proof and computer algebra:
  - add confidence to the computed bounds;
  - allow us to get to grips with (slightly) bigger algorithms;
  - make it possible to explore many variants of an algorithm (just "replay" the calculation with small modifications);
- long-term goal: use both techniques together (have the computer algebra tool generate a certificate);
- seems we are approaching the limit (in terms of algorithm size) of what can be done "exactly";
- consolation: for larger algorithms, little hope of having a worst-case error close to the algebraic bound;
- what is a publishable proof? A human-readable rough sketch along with a Coq file and/or a Maple (or whatever tool) worksheet? What we currently do is just a stylistic exercise...