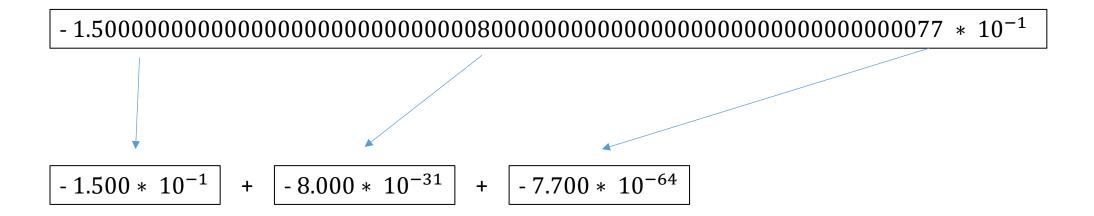
Robust Adaptive Floating-Point Geometric Predicates

Jonathan Richard Shewchuk School of Computer Science Carnegie Mellon University Pittsburgh, Pennsylvania 15213 jrs@cs.cmu.edu

+ additional notes by Petr Felkel, CTU Prague, 2020-2023

Precise floats represented as expansions

Just the idea, not using IEEE float, but 4-digit decimal numbers ...



Expansion

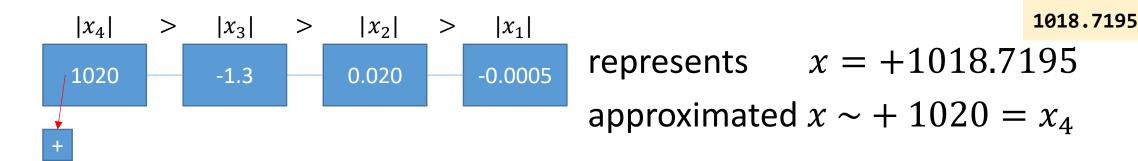
• **Sorted** sequence of **non-overlapping** machine native numbers (float, double) – each with its own exponent and significand (mantissa) 1020

1018.7

+ 0.02

1018,7200

- Sorted by absolute values
- Signum of the highest FP number is the signum of the expansion
- Zero members of the expansion will be not added.



Expansions are not unique

binary

1001.1

decimal (overlap)

... 9.5

Possibly stored as

$$1100 + (-10.1)$$

$$= 1100.0 - 10.1$$

$$= 1001 + 0.1$$

$$= 1000 + 1 + 0.1$$

$$... 12 + (-2.5)$$

$$... 12 - 2.5$$

$$...$$
 9 + 0.5

$$...$$
 8 + 1 + 0.5

Meaning of symbols

p-bit floating point operations with exact rounding (float, double):

- addition
- → subtraction
- ⊗ multiplication

Perform the operation with higher precision Round the result to the representable number

Exact rounding

Operations with exact rounding to p-bits (32 / 64) store result: exact results store exact, and non-precise results store rounded

More than 4-bits arithmetic (precise) With exact rounding to 4-bits

$$010 \times 011 = 100$$

 $2 \times 3 = 6$

$$111 \times 101 = 100011$$

 $7 \times 5 = 35$

$$010 \otimes 011 = 100$$
$$2 \otimes 3 = 6$$

$$111 \otimes 101 = 1.001 \times 2^5$$

 $7 \otimes 5 = 36$

if (possible) store exact else store rounded

Operations on expansions

IEEE 754 standard on floating point format and computing rules.

Operations on expansions require *exact rounding* of each op. to 32 / 64bit.

Theorem 1 (Dekker [4]) Let a and b be p-bit floating-point numbers such that $|a| \ge |b|$. Then the following algorithm will produce a nonoverlapping expansion x + y such that a + b = x + y, where x is an approximation to a + b and y represents the roundoff error in the calculation of x.

| 2-bits mantissa | FAS | Γ - Two - $Sum(a,b)$ | |
|-----------------------------|-----|--|-----------------------------------|
| $8 \Leftarrow 6 \oplus 1$ | 1 | $x \Leftarrow a \oplus b$ | // Rounded sum = approximation |
| $2 \Leftarrow 8 \ominus 6$ | 2 | $b_{\text{virtual}} \Leftarrow x \ominus a$ | // What was truly added - Rounded |
| $-1 \Leftarrow 1 \ominus 2$ | 3 | $y \leftarrow b \ominus b_{	ext{virtual}}$ return (x, y) | // round-off error |
| return $(8, -1)$ | 4 | return (x,y) | |

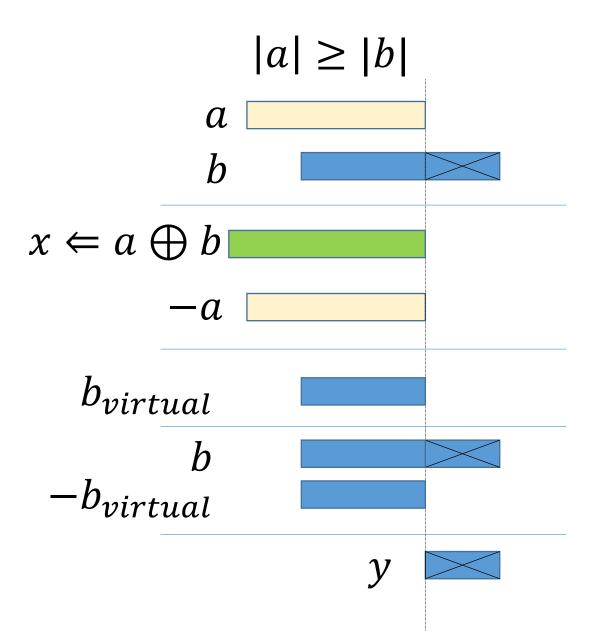
FAST-TWO-SUM(a, b)

$$\begin{array}{ccc}
1 & x \Leftarrow a \oplus b \\
b & x \rightleftharpoons x \ominus a
\end{array}$$

$$\begin{array}{ll} 2 & b_{\text{virtual}} \Leftarrow x \ominus a \\ 3 & y \Leftarrow b \ominus b_{\text{virtual}} \\ 4 & \text{return } (x,y) \end{array}$$

4 return
$$(x, y)$$

$$a + b = x + y$$
$$= a \oplus b + b \ominus b_{virtual}$$



Fast TwoSum with result rounded up (on 4-digits decimal numbers)

Correct Rounded up Really added Correction
$$a = 5081$$
 $a = 5081$ $x = 5175$ $b = 93.5$ $b = 93.5$ $a = -5081$ a

Fast TwoSum with result rounded down (on 4-digits decimal numbers)

Correct Rounded down Really added Correction
$$a = 5081$$
 $a = 5081$ $x = 5174$ $b = 93.4$ $b = 93.4$ $a = 5081$ $a = 6081$ $a = 6081$

Theorem 2 (Knuth [10]) Let a and b be p-bit floating-point numbers, where $p \geq 3$. Then the following algorithm will produce a nonoverlapping expansion x + y such that a + b = x + y.

```
TWO-SUM(a, b)
     \rightarrow x \Leftarrow a \oplus b
                                                                 // Rounded sum = approximation
2 \rightarrow b_{\text{virtual}} \Leftarrow x \ominus a
                                                                 // What b was truly added – Rounded
                                                                    for a > b
             a_{	ext{virtual}} \Leftarrow x \ominus b_{	ext{virtual}} // What a was truly added – Rounded
    \rightarrow b_{\text{roundoff}} \Leftarrow b \ominus b_{\text{virtual}} \text{ for } b > a
// round-off error of b
             a_{\text{roundoff}} \leftarrow a \ominus a_{\text{virtual}} // round-off error of a
             y \Leftarrow a_{\text{roundoff}} \oplus b_{\text{roundoff}}
     \rightarrow return (x,y)
```

Sum of two expansions (4-bit arithmetic)

Input: 1111+0.1001 and 1100+0.1

Output: 11100 + 0 + 0.0001

Zeroes slow down the computation – removed afterwards

1. Merge both input expansions into a single sequence *g* respecting the order of magnitudes

numbers in the sequence overlap

2. Use LINEAR-EXPANSION-SUM (g) to create a correct expansion

$$g5 + g4 + g3 + g2 + g1 \rightarrow h5 + h4 + h3 + h2 + h1$$

overlapping input \rightarrow non-overlapping output

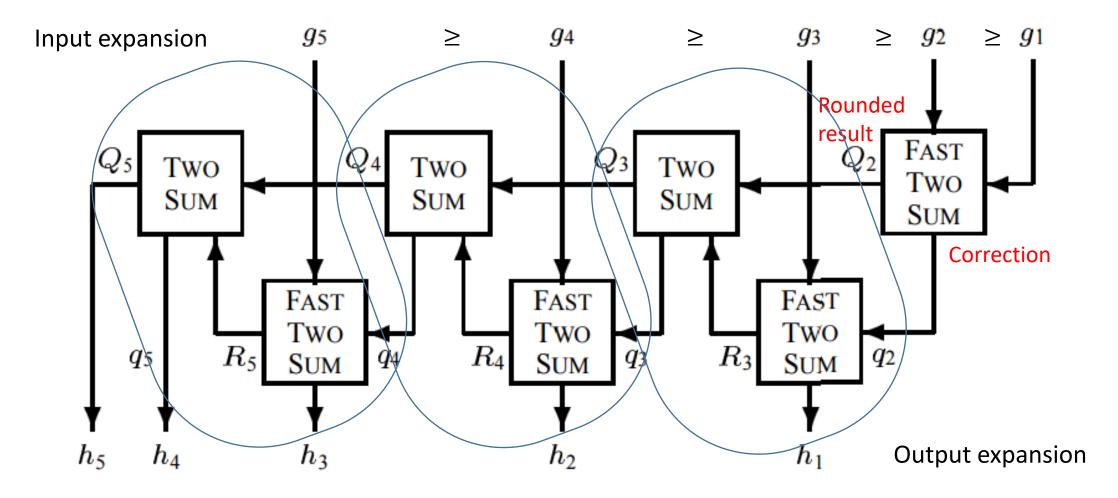


Figure 1: Operation of Linear-Expansion-Sum. The expansions g and h are illustrated with their most significant components on the left. $Q_i + q_i$ maintains an approximate running total. The Fast-Two-Sum operations in the bottom row exist to clip a high-order bit off each q_i term, if necessary, before outputting it.

LINEAR-EXPANSION-SUM

| 1111 1100 0 1001 0 1 | | | |
|---------------------------|--------------------|--------|------------|
| 1111+ 1100 + 0.1001 + 0.1 | 1111 | 1100 | 0.1001 |
| 1111 | 1101 | 1 | 0.1 |
| 1100 | | | |
| 0.1001 | 11100 | 1101+0 | 1 + 0.0001 |
| 0.1 | | | |
| | 11100 + 0 + 0.0001 | | |
| 11100.0001 | 11100 + 0.0001 | | |
| | | | |

Multiplication

Multiplies two p-bit values a and b

- 1. Split both p-bit values into two halves (with $^{\sim}$ p/2 bits)
- 2. perform four exact multiplications on these fragments.

$$a_{hi} \times b_{hi}$$
 $a_{hi} \times b_{lo}$ $a_{lo} \times b_{hi}$ $a_{lo} \times b_{lo}$

The trick is to find a way to split a floating-point value into two.

SPLIT(a) operation

- Splits p bits into two non-overlapping halves $\left(\left|\frac{p}{2}\right| \text{ bits a}_{\text{hi}} \text{ and } \left|\frac{p}{2}\right| 1 \text{ bits a}_{lo}\right)$
- ullet Missing bit is hidden in the signum of a_{lo}
- Example

7bit number splits to two 3 bit significands 1001001 splits to 1010000 (101×2^4) and -111 73 = 80 - 7 **Theorem 4 (Dekker [4])** Let a be a p-bit floating-point number, where $p \ge 3$. The following algorithm will produce $a \lfloor \frac{p}{2} \rfloor$ -bit value a_{hi} and a nonoverlapping $(\lceil \frac{p}{2} \rceil - 1)$ -bit value a_{lo} such that $|a_{\text{hi}}| \ge |a_{\text{lo}}|$ and $a = a_{\text{hi}} + a_{\text{lo}}$.

$$1 \qquad c \Leftarrow (2^{\lceil p/2 \rceil} + 1) \otimes a$$

$$a_{\mathsf{big}} \Leftarrow c \ominus a$$

$$a_{hi} \Leftarrow c \ominus a_{big}$$

$$a_{lo} \Leftarrow a \ominus a_{hi}$$

return
$$(a_{hi}, a_{lo})$$

Theorem 5 (Veltkamp) Let a and b be p-bit floating-point numbers, where $p \ge 4$. The following algorithm will produce a nonoverlapping expansion x + y such that ab = x + y.

```
TWO-PRODUCT(a, b)
          x \Leftarrow a \otimes b
           (a_{hi}, a_{lo}) = SPLIT(a)
           (b_{hi}, b_{lo}) = SPLIT(b)
          err_1 \Leftarrow x \ominus (a_{hi} \otimes b_{hi})
          err_2 \Leftarrow err_1 \ominus (a_{lo} \otimes b_{hi})
          err_3 \Leftarrow err_2 \ominus (a_{hi} \otimes b_{lo})
          y \Leftarrow (a_{10} \otimes b_{10}) \ominus err_3
           return (x, y)
```

Demonstration of SPLIT splitting a five-bit number into two two-bit numbers

Demonstration of TWO-PRODUCT in six-bit arithmetic

The resulting expansion is $110110 \times 2^6 + 11001$ $54 * 2^6 + 25$ $56^2 = 3481 \rightarrow (3456 + 25)$

Adaptive arithmetic

- Expensive avoid when possible
- Some applications need results with absolute error below a threshold
- Set of procedures with different precision (& speed) + error bounds
- For each input compute the error bounds and choose the procedure
 But
- Sometimes hard to determine error before computation
- Especially when relative error needed like sign of expression compar.
 - Result can be much larger than error bound rounded arithmetic will suffice
 - Result can be near zero must be evaluated exactly

Shewchuk predicates

- Compute a sequence of increasingly accurate results
- Testing each for accuracy
- Not using separate procedures BUT
- Using intermediate results as steps to more accurate results (work already done is not discarded, but refined)
- Idea: presented routines can be split to two parts
 - Line 1 gives an approximate result run each time
 - Remaining lines compute the roundoff error delayed until needed, if ever ...

Principle of adaptive computation

Distance of two points
$$(b_x - a_x)^2 + (b_y - a_y)^2$$

Store
$$b_x - a_x$$
 as $x_1 + y_1$

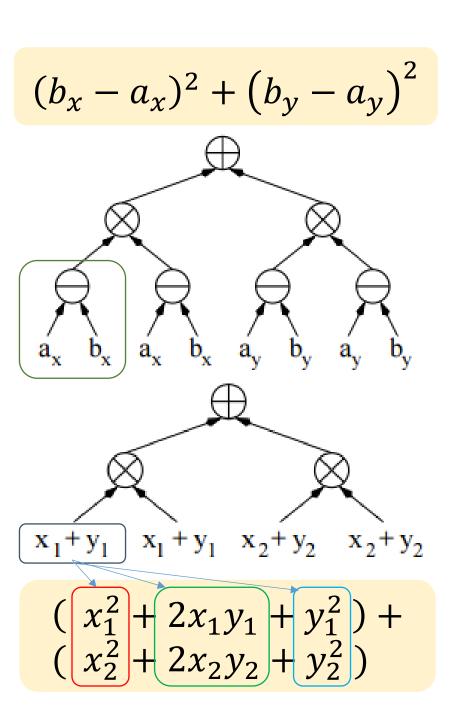
and
$$b_y - a_y$$
 as $x_2 + y_2$

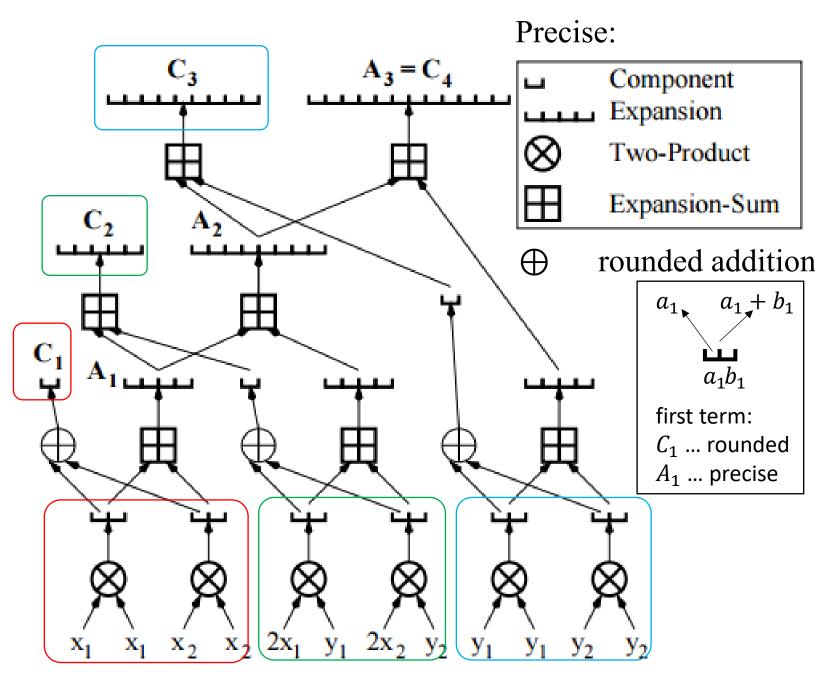
$$(x_1^2 + 2x_1y_1 + y_1^2) + (x_2^2 + 2x_2y_2 + y_2^2)$$

Reorder terms according to their size

$$(x_1^2 + x_2^2) + (2x_1y_1 + 2x_2y_2) + (y_1^2 + y_2^2)$$

Compute them only if needed





Orientation predicate - definition

orientation
$$(p,q,r) = \text{sign} \left(\det \begin{bmatrix} 1 & p_x & p_y \\ 1 & q_x & q_y \\ 1 & r_x & r_y \end{bmatrix} \right) =$$

$$= sign \left((p_x - r_x) (q_y - r_y) - (p_y - r_y) (q_x - r_x) \right),$$

where point $p = (p_x, p_y), \dots$

= third coordinate of = $(\vec{u} \times \vec{v})$,

Three points

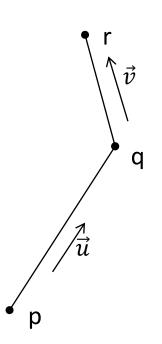
- lie on common line
- form a left turn
- form a right turn

orientation(p, q, r) =

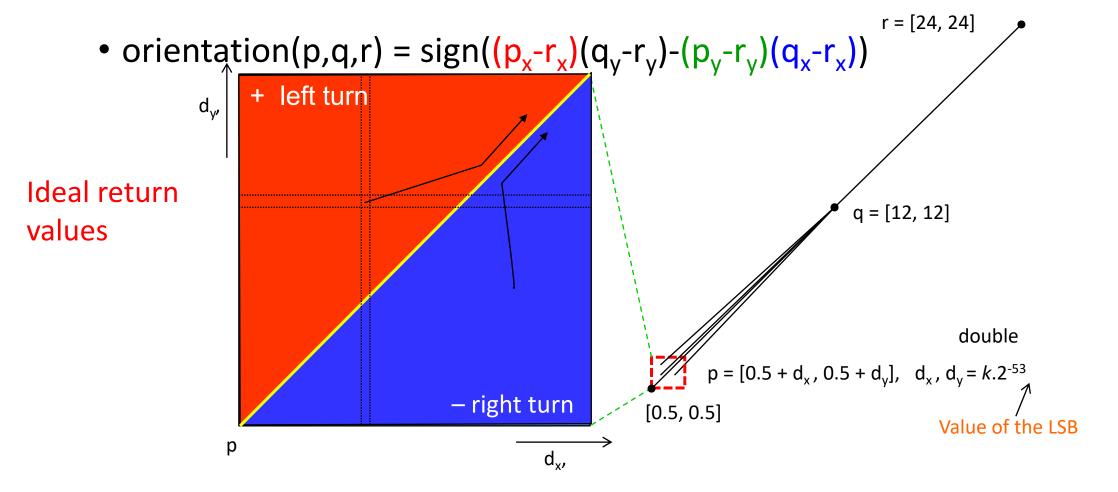
= C

= +1 (positive)

=-1 (negative)



Experiment with orientation predicate



Standard floats

| | single | double | extended double (Intel) | | |
|------------------------------------|----------------------------|----------------|-------------------------|--|--|
| celkový počet bitů | 32 | 64 | 80 | | |
| počet bitů mantisy | 24 | 53 | 64 | | |
| počet bitů exponentu | 8 | 11 | 15 | | |
| rozsah exponentu | [-126; +127] | [-1022; +1023] | [-16382; +16383] | | |
| platná desetinná | | | | | |
| místa (přibližně) | 7 | 16 | 19 | | |
| místa na lib. Čísla v o | dané přesnosti 9 | 17 | | | |
| $2^{24} = 16777216$ různých mantis | | | | | |
| $\log 2^{24} \doteq 7.2$ | | | | | |
| 1,11111 11111 11111 11 | $1111111_2 \times 2^{-12}$ | 26 | | | |

 $2,3509885615147285834557659820715330266457179855179808553659 \\ 26236850006129930346077117064851336181163787841796875 \times 10^{-38}$