

Robust Adaptive Floating-Point Geometric Predicates

Jonathan Richard Shewchuk

School of Computer Science

Carnegie Mellon University

Pittsburgh, Pennsylvania 15213

`jrs@cs.cmu.edu`

Expansion

- **Sorted** sequence of non-overlapping machine native numbers (float, double)
- Sorted by **absolute values**
- Signum of the highest FP number is the signum of the expansion
- Estimate the expansion by summing from the least significant to the most significant member
- Zero members of the expansion will be deleted.

Expansions not unique

$$1100 + -10.1$$

$$(= 1100.0 - 10.1)$$

$$= 1001 + 0.1$$

$$= 1000 + 1 + 0.1$$

Operations on expansions

IEEE 754 standard on floating point format and computing rules.

Operations on expansions require rounding of each operation to 32 / 64bit.

Fast-Two-Sum: $(a \geq b) \rightarrow (x, y), \quad a+b=x+y$

Two-Sum $(a, b) \rightarrow (x, y)$

Linear-Expansion-Sum (exp_a interleaved with exp_b) \rightarrow expansion

Split $(a) \rightarrow (a_{hi}, a_{lo}), \quad a=a_{hi}+a_{lo}$

Two-Product $(a,b) \rightarrow (x, y)$

Theorem 1 (Dekker [4]) *Let a and b be p -bit floating-point numbers such that $|a| \geq |b|$. 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 . ■*

FAST-TWO-SUM(a, b)

```

1       $x \leftarrow a \oplus b$            // Rounded sum = approximation
2       $b_{\text{virtual}} \leftarrow x \ominus a$  // What was truly added - Rounded
3       $y \leftarrow b \ominus b_{\text{virtual}}$  // round-off error
4      return ( $x, y$ )

```

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)

- 1 $\rightarrow x \leftarrow a \oplus b$ // Rounded sum = approximation
- 2 $\rightarrow b_{\text{virtual}} \leftarrow x \ominus a$ // What b was truly added - Rounded
- 3 $a_{\text{virtual}} \leftarrow x \ominus b_{\text{virtual}}$ // What a was truly added - Rounded
- 4 $\rightarrow b_{\text{roundoff}} \leftarrow b \ominus b_{\text{virtual}}$ // round-off error of b
- 5 $a_{\text{roundoff}} \leftarrow a \ominus a_{\text{virtual}}$ // round-off error of a
- 6 $y \leftarrow a_{\text{roundoff}} \oplus b_{\text{roundoff}}$
- 7 \rightarrow **return** (x, y)

Sum of two expansions (4-bit arithmetic)

Input: $1111+0.1001$ and $1100 + 0.1$

Output: $11100 + 0 + 0.0001$

Zeros slow down the computation

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

$1111+ 1100 + 0.1 + 0.1001$

Use LINEAR-EXPANSION-SUM (g)

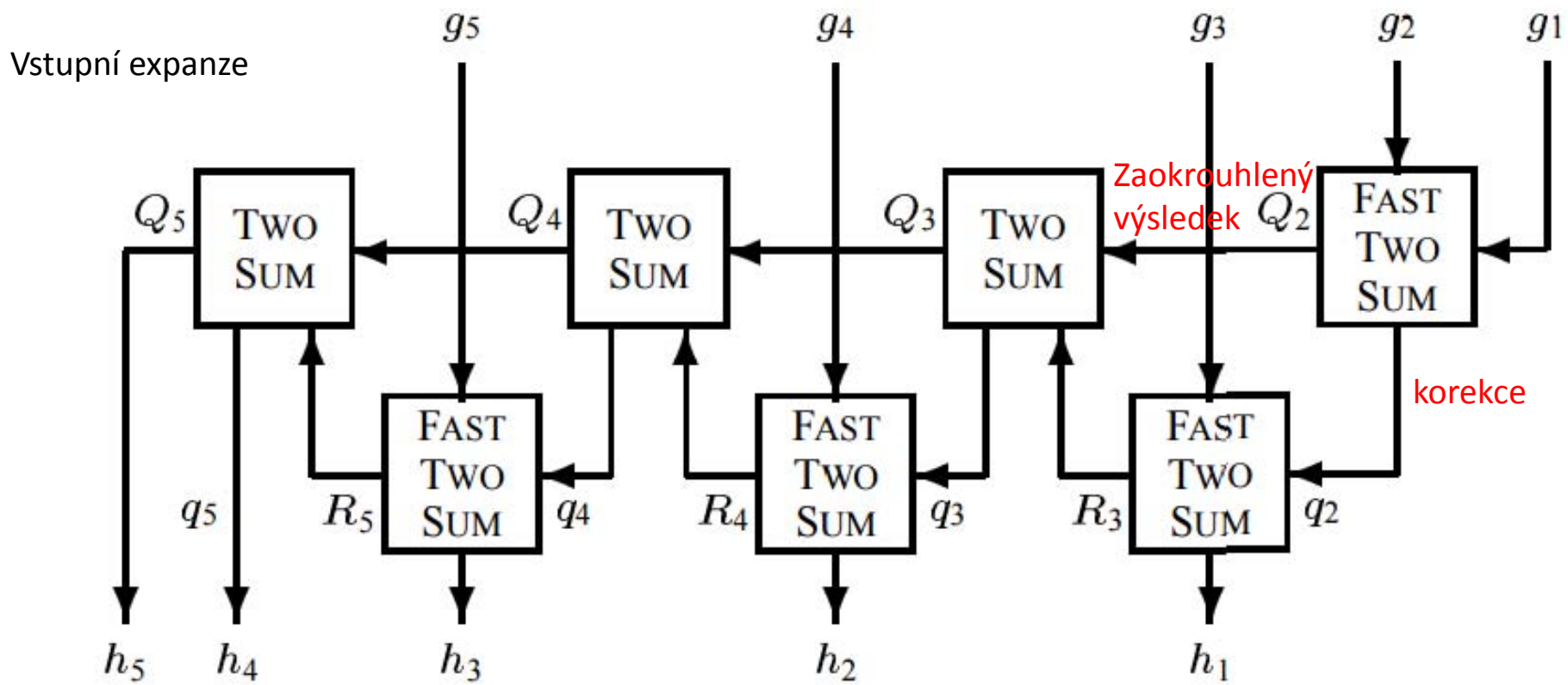


Figure 1: Operation of LINEAR-EXPANSION-SUM. The expansions g and h are illustrated with their most significant components on the left. $Q_i + q_i$ maintains an approximate running total. The FAST-TWO-SUM operations in the bottom row exist to clip a high-order bit off each q_i term, if necessary, before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

Theorem 4 (Dekker [4]) *Let a be a p -bit floating-point number, where $p \geq 3$. 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 $|a_{\text{hi}}| \geq |a_{\text{lo}}|$ and $a = a_{\text{hi}} + a_{\text{lo}}$. ■*

SPLIT(a)

```

1    $c \leftarrow (2^{\lceil p/2 \rceil} + 1) \otimes a$    
2    $a_{\text{big}} \leftarrow c \ominus a$ 
3    $a_{\text{hi}} \leftarrow c \ominus a_{\text{big}}$ 
4    $a_{\text{lo}} \leftarrow a \ominus a_{\text{hi}}$ 
5   return ( $a_{\text{hi}}, a_{\text{lo}}$ )            
```

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

TWO-PRODUCT(a, b)

```
1    $x \leftarrow a \otimes b$ 
2    $(a_{\text{hi}}, a_{\text{lo}}) = \text{SPLIT}(a)$ 
3    $(b_{\text{hi}}, b_{\text{lo}}) = \text{SPLIT}(b)$ 
4    $err_1 \leftarrow x \ominus (a_{\text{hi}} \otimes b_{\text{hi}})$ 
5    $err_2 \leftarrow err_1 \ominus (a_{\text{lo}} \otimes b_{\text{hi}})$ 
6    $err_3 \leftarrow err_2 \ominus (a_{\text{hi}} \otimes b_{\text{lo}})$ 
7    $y \leftarrow (a_{\text{lo}} \otimes b_{\text{lo}}) \ominus err_3$ 
8   return  $(x, y)$ 
```

Orientation predicate - definition

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

where point $p = (p_x, p_y), \dots$
= third coordinate of $(\vec{u} \times \vec{v})$,

Three points

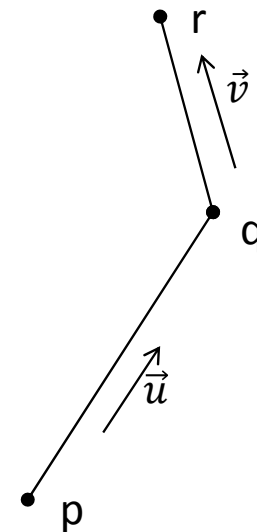
- lie on common line
- form a left turn
- form a right turn

$$\text{orientation}(p, q, r) =$$

$$= 0$$

$$= +1 \text{ (positive)}$$

$$= -1 \text{ (negative)}$$



Experiment with orientation predicate

- $\text{orientation}(p,q,r) = \text{sign}((p_x - r_x)(q_y - r_y) - (p_y - r_y)(q_x - r_x))$

Ideal return values

