

# Robust Adaptive Floating-Point Geometric Predicates

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

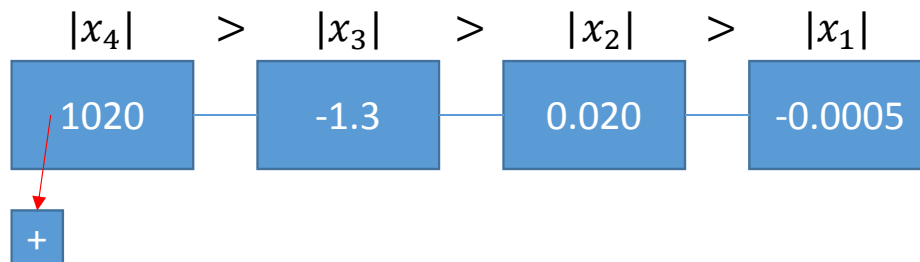
+ additional notes by Petr Felkel, CTU Prague, 2020

Version from 8.10.2020

# Expansion

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

```
1020
- 1.3
-----
1018.7
+ 0.02
-----
1018.7200
- 0.0005
-----
1018.7195
```



represents  $x = +1018.7195$   
approximated  $x \sim +1020 = x_4$

# Expansions are not unique

binary

decimal

$$1100 + (-10.1)$$

$$\dots 12 + (-2.5)$$

$$= 1100.0 - 10.1$$

$$\dots 12 - 2.5$$

$$= 1001 + 0.1$$

$$\dots 9 + 0.5$$

$$= 1000 + 1 + 0.1$$

$$\dots 8 + 1 + 0.5$$

All represent the value 1001.1 ... 9.5

# Meaning of symbols

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

- $\oplus$  addition
- $\ominus$  subtraction
- $\otimes$  multiplication

# 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

$$010 \times 011 = 100$$

$$2 \times 3 = 6$$

$$111 \times 101 = 100011$$

$$7 \times 5 = 35$$

With exact rounding to 4-bits

$$010 \otimes 011 = 100$$

$$2 \otimes 3 = 6$$

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

$$7 \otimes 5 = 36$$

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.

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

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( $a, b$ )

```

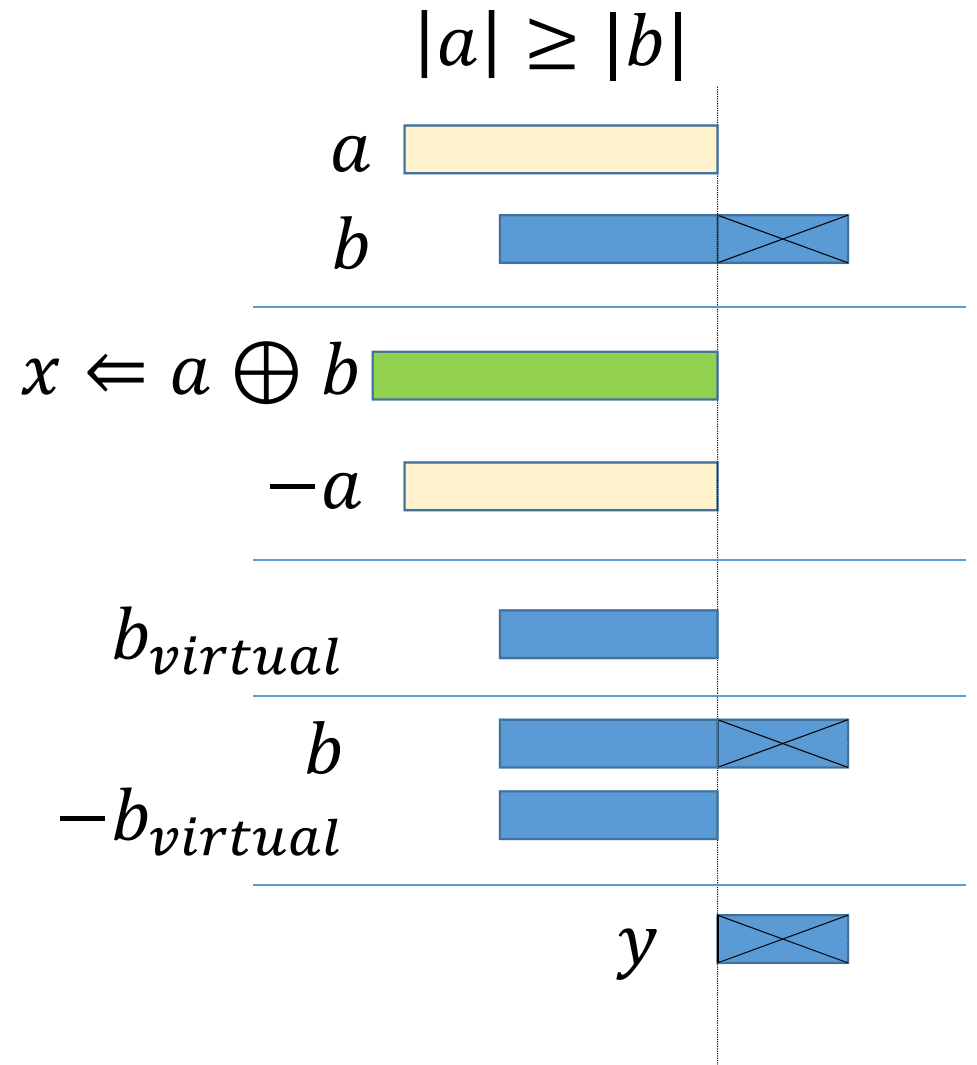
1       $x \leftarrow a \oplus b$            // Rounded sum = approximation
2       $b_{\text{virtual}} \leftarrow x \ominus a$  // What was truly added - Rounded
3       $y \leftarrow b \ominus b_{\text{virtual}}$  // round-off error
4      return ( $x, y$ )

```

# FAST-TWO-SUM( $a, b$ )

- 1  $x \leftarrow a \oplus b$
- 2  $b_{\text{virtual}} \leftarrow x \ominus a$
- 3  $y \leftarrow b \ominus b_{\text{virtual}}$
- 4 **return** ( $x, y$ )

$$\begin{aligned}
 a + b &= x + y \\
 &= a \oplus b + b \ominus b_{\text{virtual}}
 \end{aligned}$$





## Fast TwoSum with result rounded up

Correct	Rounded	Really added	Correction
$a = 5081$ $b = 93.5$ <hr/> $5174.5$	$a = 5081$ $b = 93.5$ <hr/> $x = 5175$	$x = 5175$ $-a = -5081$ <hr/> $b_{virtual} = 94$	$b = 93.5$ $-b_{virtual} = -94$ <hr/> $y = -0.5$
	$(a + b) = (x + y)$		
	$5081 + 93.5 = (5175 - 0.5)$		

## Fast TwoSum with result rounded down

Correct	Rounded	Really added	Correction
$a = 5081$ $b = 93.4$ <hr style="width: 100px; margin-left: 0;"/> $5174.4$	$a = 5081$ $b = 93.4$ <hr style="width: 100px; margin-left: 0;"/> $x = 5174$	$x = 5174$ $-a = -5081$ <hr style="width: 100px; margin-left: 0;"/> $b_{virtual} = 93$	$b = 93.4$ $-b_{virtual} = -94$ <hr style="width: 100px; margin-left: 0;"/> $y = 0.4$
$(a + b) = (x + y)$			
$5081 + 93.4 = (5174 + 0.4)$			

**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  →  $x \leftarrow a \oplus b$  // Rounded sum = approximation
2  →  $b_{\text{virtual}} \leftarrow x \ominus a$  // What  $b$  was truly added – Rounded
   // for  $a > b$ 
3  →  $a_{\text{virtual}} \leftarrow x \ominus b_{\text{virtual}}$  // What  $a$  was truly added – Rounded
   // for  $b > a$ 
4  →  $b_{\text{roundoff}} \leftarrow b \ominus b_{\text{virtual}}$  // 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  → 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 – removed afterwards

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

$1111+ 1100 + 0.1001 + 0.1$

Use LINEAR-EXPANSION-SUM ( $g$ )

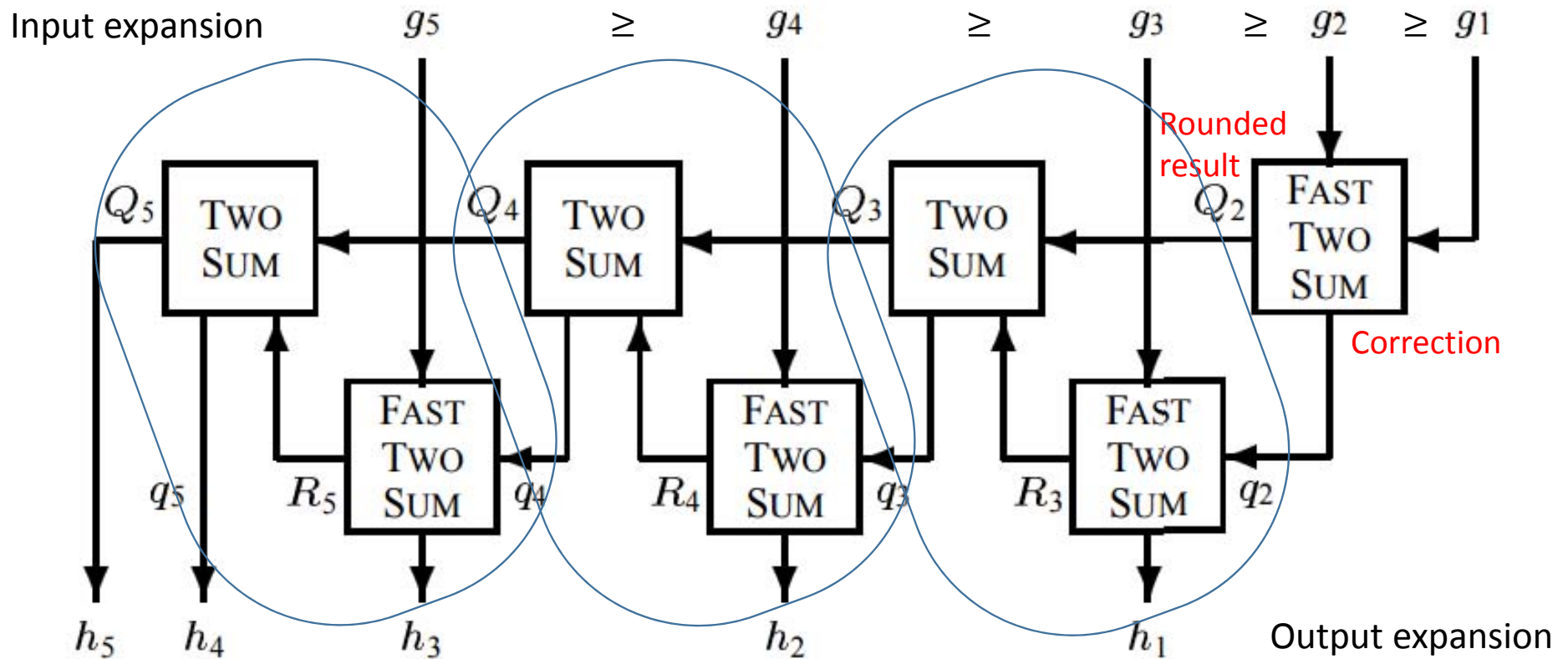


Figure 1: Operation of LINEAR-EXPANSION-SUM. The expansions  $g$  and  $h$  are illustrated with their most significant components on the left.  $Q_i + q_i$  maintains an approximate running total. The FAST-TWO-SUM operations in the bottom row exist to clip a high-order bit off each  $q_i$  term, if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

# 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.

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

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

# SPLIT(a) operation

- Splits  $p$  bits into two non-overlapping halves  
( $\lfloor \frac{p}{2} \rfloor$  bits  $a_{hi}$  and  $\lfloor \frac{p}{2} \rfloor - 1$  bits  $a_{lo}$ )
- Missing bit is hidden in the signum of  $a_{lo}$
- Example

7bit number splits to two 3 bit significands

1001001 splits to 1010000 ( $101 \times 2^4$ ) and -111

$$73 = 80 - 7$$

**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  $a_{\text{hi}}$  and a nonoverlapping  $(\lceil \frac{p}{2} \rceil - 1)$ -bit value  $a_{\text{lo}}$  such that  $|a_{\text{hi}}| \geq |a_{\text{lo}}|$  and  $a = a_{\text{hi}} + a_{\text{lo}}$ . ■*

**SPLIT**( $a$ )

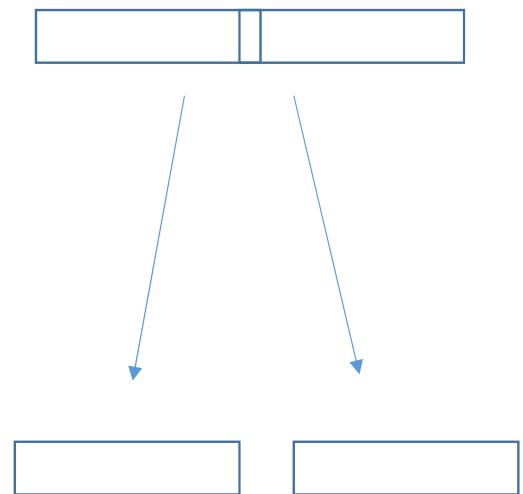
1  $c \leftarrow (2^{\lceil p/2 \rceil} + 1) \otimes a$

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  $(x, y)$ 
```

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

$$\begin{array}{rcl}
 a & = & \phantom{000}11101 \\
 2^3 a & = & \phantom{00}11101 \times 2^3 \\
 c & = (2^3 + 1) \otimes a & = \begin{array}{r}
 \phantom{00}11101 \\
 \hline
 10000 \\
 \hline
 \phantom{00}11101 \\
 \hline
 11100
 \end{array} \times 2^4 \\
 a & = & \phantom{000}11101 \\
 a_{\text{big}} & = c \ominus a & = \phantom{00}11100 \times 2^3 \\
 a_{\text{hi}} & = c \ominus a_{\text{big}} & = \phantom{000}10000 \times 2^1 \\
 a_{\text{lo}} & = a \ominus a_{\text{hi}} & = \phantom{0000}11
 \end{array}$$

# Demonstration of TWO-PRODUCT in six-bit arithmetic

$$\begin{array}{rcl}
 a & = & \phantom{111011} \\
 b & = & \phantom{111011} \\
 x & = & a \otimes b = \begin{array}{r} 110110 \\ \hline 110001 \end{array} \times 2^6 \\
 err_1 & = & x \ominus (a_{hi} \otimes b_{hi}) = \begin{array}{r} 110001 \\ \hline 101000 \end{array} \times 2^3 \\
 err_2 & = & err_1 \ominus (a_{lo} \otimes b_{hi}) = \begin{array}{r} 101000 \\ \hline 100110 \end{array} \times 2^2 \\
 err_3 & = & err_2 \ominus (a_{hi} \otimes b_{lo}) = \begin{array}{r} 100110 \\ \hline -10000 \end{array} \\
 -y & = & err_3 \ominus (a_{lo} \otimes b_{lo}) = \begin{array}{r} -10000 \\ \hline -11001 \end{array}
 \end{array}$$

The resulting expansion is  $110110 \times 2^6 + 11001$

# 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 comp.
  - Result can be much larger than error bound – exact 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  $x_2 + y_2$

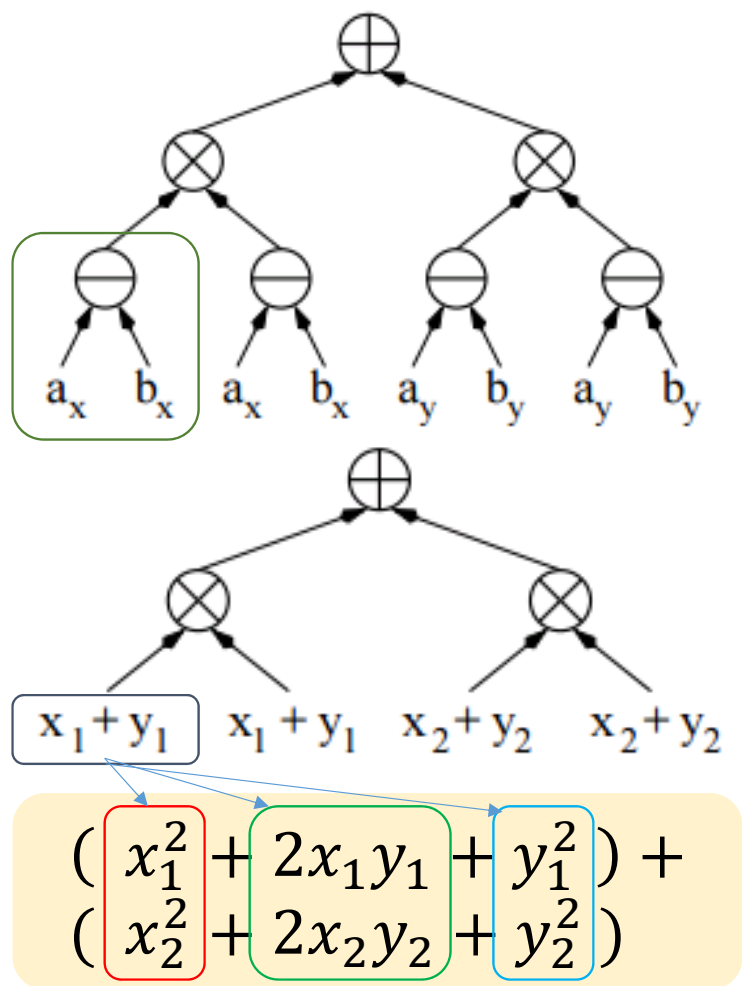
$$(x_1^2 + 2x_1y_1 + y_1^2) + (x_2^2 + 2x_2y_2 + y_2^2)$$

Reorder terms according to their size

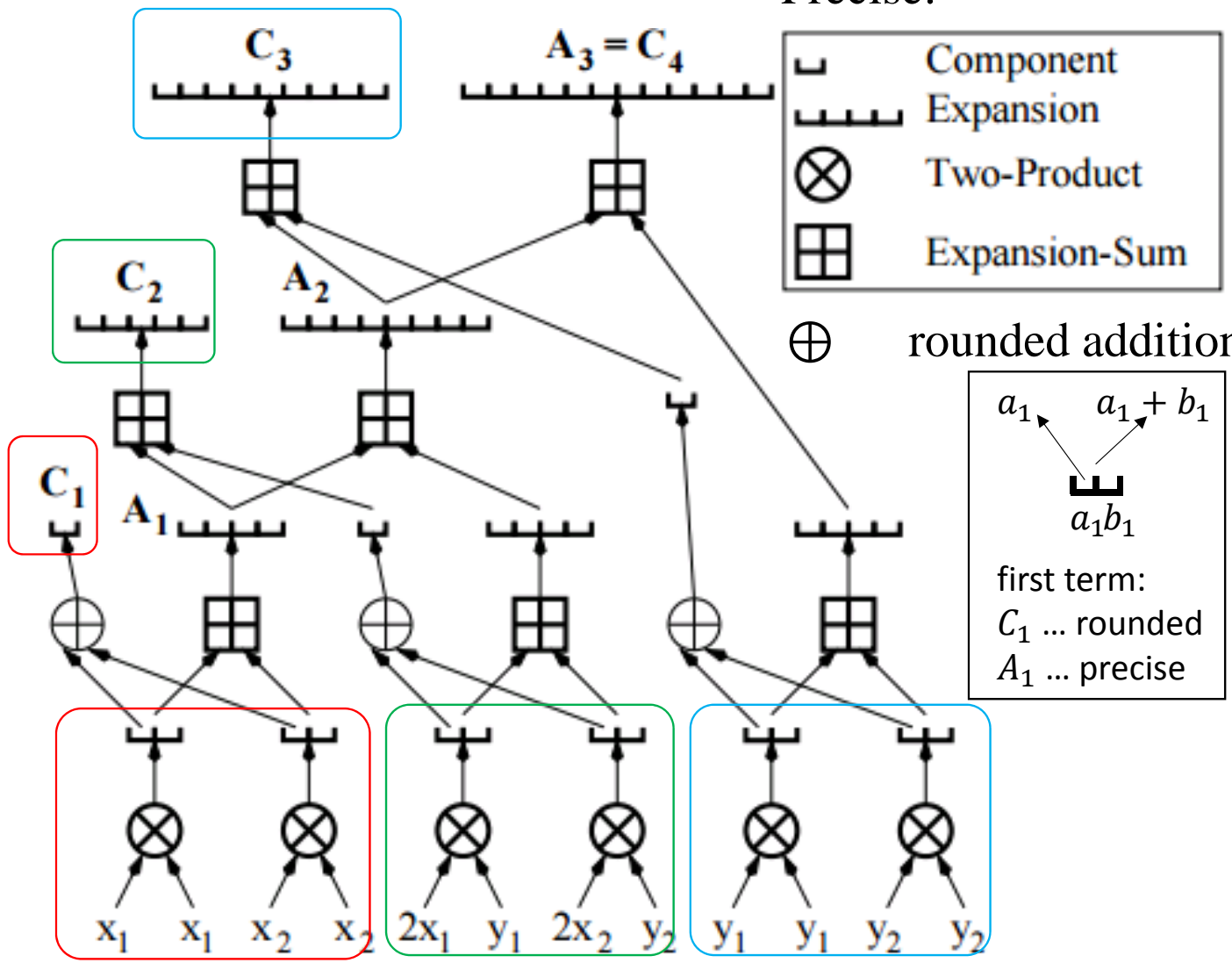
$$(x_1^2 + x_2^2) + (2x_1y_1 + 2x_2y_2) + (y_1^2 + y_2^2)$$

Compute them only if needed

$$(b_x - a_x)^2 + (b_y - a_y)^2$$



Precise:



# Orientation predicate - definition

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

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

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

## Three points

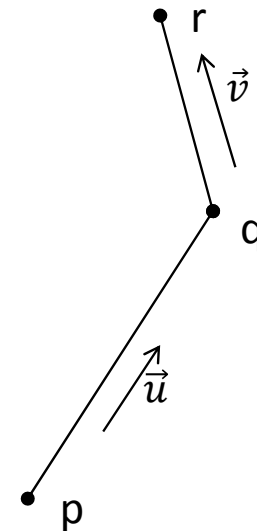
- lie on common line
- form a left turn
- form a right turn

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

$$= 0$$

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

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





pivot  $r$

# 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

