

2.3. Rounding & normalization rules:

- sticky bits: g, r, s

$$X_n + Y_{n+1} = Z_n$$

Consider Z_n : **Case**

$$Z_n = \overline{z}_n \overline{z}_{n-1} \overline{z}_{n-2} \overline{z}_{n-3} \dots \overline{z}_1 \overline{z}_0 | g \ r \ s$$

Normalization:

$$Z_n = \overline{z}_n | \overline{z}_{n-1} \overline{z}_{n-2} \overline{z}_{n-3} \dots \overline{z}_1 \overline{z}_0 | \overline{g} \ \overline{r} \ \overline{s}$$

Normalization cases:

- 1) $\overline{z}_n = 1 \rightarrow$ 1-bit RShift
- 2) $\overline{z}_n = 0, \overline{z}_{n-1} = 1$
 \overline{z}_n is already normalized
- 3) $\overline{z}_n = 0, \overline{z}_{n-1} = 0, \overline{z}_{n-2} = 1$ 1-bit LShift
- 4) $\overline{z}_n = \overline{z}_{n-1} = \overline{z}_{n-2} = 0, \overline{z}_{n-3} = 1$ 3-bit LShift

! if Z_n requires a LShift by 2 or more bits then after rounding g Z_n is completed with bits of 0's; additionally $R=S=0$

Normalization is followed by rounding: use R, S

- using $R, S \equiv$ "fractional part of 2 bits"

Rounding mode	$Z_n \geq 0$	$Z_n < 0$
to 0	(discard R, S)	(discard R, S)
towards	(discard R, S)	if $(R, S) \neq (0, 0)$ then $\overline{z}_n = 1$
towards	if $(R, S) \neq (0, 0)$ then $\overline{z}_n = 1$	if $(R, S) \neq (0, 0)$ then $\overline{z}_n = 1$
towards even	if $(R, S) \neq (0, 0)$ then $\overline{z}_n = 1$	if $(R, S) \neq (0, 0)$ then $\overline{z}_n = 1$

$$1.0001 \rightarrow 2$$

$$1 \times 2$$

! $Z_n + 1$ can generate carry \Rightarrow post-normalization

2.4. Fp. addition/subtraction with rounding

- use a IEEE 754 format scaled down:

$$\frac{a}{2^4} - \frac{b}{2^4}$$

1 bit: sign $e=3$ bits: exponent field. bias = $2^{e-1} - 1 = 3$

3 bits: significand \rightarrow 3 bits for fractional part of significand

$$X = 0.5625_{10} = 0.1001_2 \times 2^0 = 1.001_2 \times 2^{-1}$$

$$Y = -3.75_{10} = -11.11_2 \times 2^0 = -1.11_2 \times 2^1$$

Packed operands: $X_E = -1 + \text{bias} = -1 + 3 = 2$ $Y_E = 1 + \text{bias} = 1 + 3 = 4$

$$X: 0 \ 0 \ 1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 1 \ 1 \quad Y: 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1$$

Fp. addition with rounding algorithm

Step 1: unpack operands: \rightarrow adding hidden bit \rightarrow check for exceptions (0, $\pm \infty$, NaN, ...)

$$X: 0 \ 0 \ 1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 1 \ 1 \quad Y: 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1$$

$$X_E = 010110011 \quad Y_E = 11100110011$$

$$d = X_E - Y_E = 010110011 - 11100110011 = -2 < 0$$

$$\Rightarrow \text{SWAP } X \leftrightarrow Y$$

$$Z_E = Y_E = 4$$

$$X: 1 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0 \ 0 \ 1 \ 1 \quad Y: 0 \ 0 \ 1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 1 \ 1$$

Step 3: if $\text{sign}(X) \neq \text{sign}(Y) \rightarrow$ Complement Y_n of 2

$$\text{sign}(X) = 1 \quad \text{sign}(Y) = 0 \quad \text{complement } Y_n \text{ of 2}$$

$$Y: 0 \ 0 \ 1 \ 1 \ 0 \ 1 \ 0 \ 0 \ 1 \ 1$$

$$Y_{nc2} = 0.111$$

Step 4: Align of Y_n : RShift by $|d|$

\rightarrow if Y_n was complemented of 2 in Step 3, then a RShift introduces bits of 1's in Y_n , instead of 0's

\rightarrow preserve sticky bits: g, r, s

$$Y_{nc2} = 0.111 \quad d = -2$$

$$Y_{nc2d} = 1.101110$$

Step 5: Add the 2 significands: $Z_n = X_n + Y_{nc2d}$

\rightarrow if $\text{sign}(X) = \text{sign}(Y)$ if carry is generated from msb of Z_n : PRESERVE it

\rightarrow if $\text{sign}(X) \neq \text{sign}(Y)$ if carry is NOT generated: Z_n is negative. \Rightarrow Complement Z_n of 2

\rightarrow if carry is generated: Z_n is positive \Rightarrow IGNORE carry

$$X_n = 1.1111 \quad Y_{nc2d} = 1.101110$$

$$Z_n = 1.1111 + 1.101110 = 1.1001110$$

Step 6: Pre-normalization: according to the rules in section 2.3 \Rightarrow determine Z_n , update Z_E

\rightarrow exception checking:

- if $Z_E = Z_{E_{max}} (2^3 - 2 = 6)$ AND Z_n requires 1-bit RShift \Rightarrow OVERFLOW

- if $Z_E = Z_{E_{min}} (1)$ AND Z_n requires 1-bit LShift \Rightarrow UNDERFLOW

$$Z_n = 1.1001110 \quad Z_E = 4$$

Rule 2) (already normalized) $\Rightarrow Z_n = 1.100 \quad Z_E = 4$

Step 7: Calculate R, S : according to the same rules in 2.3.

$$\text{Rule 2): } R = g = 1$$

$$S = r \text{ OR } s = 1 \text{ OR } 0 = 1$$

Step 8: Round to Z_n : according to the rules in 2.3.

- if rounding generates carry \Rightarrow post-normalization

Post-normalization: 1-bit RShift Z_E++ !! check exceptions (like in Step 6)

Consider round to nearest even mode

if $(R \text{ AND } (S \text{ OR } Z_{E_{min}}))$ then Z_n++

$$R = 1$$

$$S = 1$$

$$Z_{E_{min}} = 0$$

$$1 \text{ AND } (1 \text{ OR } 0) = 1 \Rightarrow$$

$$Z_n = 1.1001110 + 1$$

$$Z_n^* = 1.1011$$

no carry out \Rightarrow no post-normalization

Step 9: Calculate sign of result

- if $\text{sign}(X) = \text{sign}(Y) \Rightarrow \text{sign}(Z) = \text{sign}(X)$

- if $\text{sign}(X) \neq \text{sign}(Y) \Rightarrow$ initial signs of X, Y (before the SWAP)

SWAP (Step 2)	Complement of 2 (Step 3)	sign(X)	sign(Y)	sign(Z)
YES		+	-	-
YES		-	+	+
NO	YES	+	-	-
NO	YES	-	+	+
NO	NO	+	+	+
NO	NO	-	-	-

SWAP (Step 2) YES: $\text{sign}(Z) = \text{sign}(Y) = -$ (initial sign)

Step 10: Pack the result

$$\text{sign}(Z) = -; Z_E = 4; Z_n = 1.1011$$

$$Z: 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1$$

Z_E fractional part of Z_n

Verify: $X = 0.5625 \quad Y = -3.75$

$$Z = X + Y = -3.1875 \text{ (infinite precision)}$$

Unpack Z :

$$Z: 1 \ 1 \ 1 \ 0 \ 1 \ 1 \ 0 \ 1 \ 1 \ 1$$

$$Z = (-1)^{\text{sign}} \times 2^{Z_E - \text{bias}} \times Z_n = (-1) \times 2^{4-3} \times 1.1011 = -1.1011_2 = -3.1875_{10}$$



Chapter III Functional Analysis & Synthesis of Binary Multiplication Devices

3.1. Multiplication methods

Y = multiplicand X = multiplier

$$P = X \times Y \quad X = 11, Y = 12$$

a) paper & pencil

$$\begin{array}{r} 1100 \dots\dots\dots Y \\ 1011 = x_3x_2x_1x_0 \dots\dots X \\ \hline 1100 \dots\dots\dots x_0 \cdot 4 \cdot 2^0 \\ 1100 \dots\dots\dots x_1 \cdot 4 \cdot 2^1 \\ 0000 \dots\dots\dots x_2 \cdot 4 \cdot 2^2 \\ 1100 \dots\dots\dots x_3 \cdot 4 \cdot 2^3 \\ \hline 11000100 = P = \sum_{i=0}^3 x_i \cdot 4 \cdot 2^i = 132 \end{array}$$

HW investment:

- 2, 4-bit registers for X, Y
- multi-operand adder (CSA) \leftarrow latency !!!
- multiplicand gating

b) keep for the partial products

$$\begin{array}{r} 1100 \dots\dots\dots Y \\ 1011 = x_3x_2x_1x_0 \dots\dots X \\ \hline 00000000 \dots\dots\dots P_0 = 0 \\ 1100 \dots\dots\dots x_0 \cdot 4 \cdot 2^0 \quad + \\ \hline 00001100 \dots\dots\dots P_1 = P_0 + x_0 \cdot 4 \cdot 2^0 \\ 1100 \dots\dots\dots x_1 \cdot 4 \cdot 2^1 \quad + \\ \hline 00100100 \dots\dots\dots P_2 = P_1 + x_1 \cdot 4 \cdot 2^1 \\ 0000 \dots\dots\dots x_2 \cdot 4 \cdot 2^2 \quad + \\ \hline 00100100 \dots\dots\dots P_3 = P_2 + x_2 \cdot 4 \cdot 2^2 \\ 1100 \dots\dots\dots x_3 \cdot 4 \cdot 2^3 \quad + \\ \hline 11000100 \dots\dots\dots P_4 = P_3 + x_3 \cdot 4 \cdot 2^3 = P \end{array}$$

$$\text{iteration step: } P_{i+1} = P_i + x_i \cdot 4 \cdot 2^i, i \geq 0, P_0 = 0$$

HW investment:

- \rightarrow 2, 4-bit registers for X, Y
- \rightarrow 8-bit register for result / P_i
- \rightarrow 8-bit adder
- \rightarrow LShift for $x_i \cdot 4 \cdot 2^i$