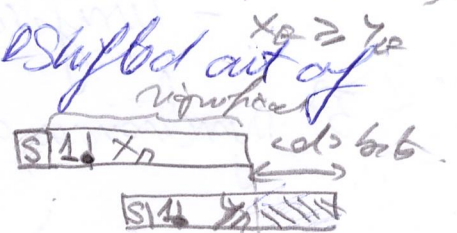


AC CIO

2.3. F.p. addition without rounding

- while aligning, keep all bits that were shifted out of Y_n 's storage capacity

$$X \pm Y = (X_n \pm Y_n) \cdot 2^{E_1 - E_2}$$



adder/subtractor

declare registers $A[(m+1):0]$, $M[(m+1):0]$, $E_1[(e-1):0]$, $E_2[(e-1):0]$, $E[(e-1):0]$, A_COUT , $ERROR$;

declare bus $INBUS[(m+e+1):0]$, $OUTBUS[(m+e+1):0]$;

BEGIN: $A_COUT := 0$, $ERROR := 0$;

INPUT: $A := INBUS(X_n)$, $E_1 := INBUS(X_e)$;

$M := INBUS(Y_n)$, $E_2 := INBUS(Y_e)$;

COMPARE: $E := E_1 - E_2$

ALIGN: if $E < 0$ then $A := RShift(A, E := E + 1, \text{go to ALIGN})$;
if $E > 0$ then $M := RShift(M, E := E - 1, \text{go to ALIGN})$;

ADD/SUBTRACT: $A := A \pm M$, $E := \max(E_1, E_2)$;

OVERFLOW: if $A_COUT = 1$ then begin

if $E == E_{max} (= 2^e - 2)$ then go to ERROR;

$A := RShift(A, 1)$, $E := E + 1$, go to END;

end

ZERO: if $A == 0$ then $E := 0$, go to END;

NORMALIZE: if $isNormalized(A) == 1$ then go to END;

UNDERFLOW: if $E > E_{min} (= -1)$ then

$A := LShift(A)$, $E := E - 1$, go to NORMALIZE;

ERROR: $ERROR := 1$;

END :

Pseudo-language:

1) declare registers \leftarrow name & width

- concatenation: $A_COUT.A$

2) declare bus \rightarrow operands
width

$$2m + 2e + 2 \text{ bits}$$

- unified bus: IDBUS $m + e + 2 \text{ bits}$

3) synchronous execution \rightarrow non-conflicting: concurrent execution,
operations separated by ;
 \rightarrow conflicting: sequential execution,
operations separated by ;

4) $:=$ assignment

- hardwired operators: RShift, isNormalised, max

5) flow control \rightarrow unconditional: go to ERROR
jump

conditional: if $B \neq \text{ERROR}$ then go to ERROR
jump

6) simultaneous read/write from/to register/bus:

$A[77] := M[77] \text{ over } [00], [10] := 0;$

Comments for algorithm:

A) operands $m + e + 2 \text{ bits}$ \rightarrow L: sign,
e: exponent
 $m + 1$: mantissa / significand
 \rightarrow s: hidden, supplant
 \rightarrow m: fractional part

B) registers A and M:

- RShift capability

- store Sign + Mantissa / Significand

- A \rightarrow LShift capabilities

extended with flag A-OUT

C) register E:

- ++/-- capabilities

D) register ERROR:

- indicate exceptions: OVERFLOW, UNDERFLOW

2.4. Rounding and normalization rules for f.p. addition (2)

$$X_M = \left| 1. \overset{2^{-1}}{x_{m-2}} \overset{2^{-2}}{x_{m-3}} \dots \overset{2^{-m+2}}{x_1} \overset{2^{-m+1}}{x_0} \right| \overset{2^{-m}}{x_{-1}} \overset{2^{-m-1}}{x_{-2}} \overset{2^{-m-2}}{x_{-3}} \dots$$

$$Y_M = \left| 1. y_{m-2} y_{m-3} \dots y_1 y_0 \right| \text{ and } x_0 \geq y_0$$

align Y_M by RShift with d bits, where $d = x_0 - y_0$

Y's alignment

- with infinite precision: keep all bits with weight $< 2^{-m+1}$
- with finite precision: keep only 3 bits with weight $< 2^{-m+1}$

the 3 bits == sticky bits

g : guard bit, weight 2^{-m} , the last bit RShifted out of Y_M
 - guards against loss of precision.
 r : round bit, weight 2^{-m+1} , the last but one bit. RShift out of Y_M
 s : sticky bit, weight 2^{-m+2} , obtained by logic OR-ing all the other bits RShifted out of Y_M (excluding g, r)

$$X_M + Y_M = Z_M$$

Consider $Z_M = \left(\overset{\text{carry out}}{z_m} \mid \overset{2^{-1}}{z_{m-1}} \overset{2^{-2}}{z_{m-2}} \dots \overset{2^{-m}}{z_1} \overset{2^{-m-1}}{z_0} \mid g \ r \ s \right)$

Normalize $Z_M \Rightarrow Z_{M_n}$

$$Z_{M_n} = \left| 1. \overset{2^{-1}}{z_{m-2_n}} \overset{2^{-2}}{z_{m-3_n}} \dots \overset{2^{-m}}{z_n} \overset{2^{-m-1}}{z_{n-1}} \right| \text{ (R S)}$$

$$z_n = (z_m, z_{m-1}, z_{m-2}, z_{m-3}, \dots, z_2, z_1, g, r, s)$$

$$z_m = 1, z_{m-1}, z_{m-2}, \dots, z_2, z_1, g, r, s$$

Normalization cases:

1) $z_m = 1$; 1-bit RShift

2) $z_m = 0, z_{m-1} = 1$

z_m is normalized.

3) $z_m = 0, z_{m-1} = 0, z_{m-2} = 1$

1-bit LShift

4) $z_m = 0, z_{m-1} = 0, z_{m-2} = 0, z_{m-3} = 1$

2-bit LShift

1	z_{m-1}	z_{m-2}	z_{m-3}	z_2	z_1	g	r	s
1	z_{m-1}	z_{m-2}	z_{m-3}	z_2	z_1	g	r	s
1	z_{m-2}	z_{m-3}	z_{m-4}	z_2	z_1	g	r	s
1	z_{m-3}	z_{m-4}	z_{m-5}	z_2	z_1	g	r	s
1	z_{m-4}	z_{m-5}	z_{m-6}	z_2	z_1	g	r	s

if a LShift is required by 2 or more bits of z_m
then after appending g, z_m is completed with 0s

R and S are used for rounding of z_m
- only 2 bits are needed out of those to be eliminated
for implementing the 4 rounding modes of IEEE 754

Rounding mode	$z_m \geq 0$	$z_m < 0$
to 0	(discard R, S)	(discard R, S)
towards $-\infty$	—	if (R or S) then $z_m - 1$
towards $+\infty$	if (R or S) then $z_m + 1$	—
to nearest even	if (R and (S or z_{m-1})) then $z_m + 1$	if (R and (S or z_{m-1})) then $z_m - 1$

$$-1.000005 \rightarrow -\infty \rightarrow -2$$

$$0.5 \Delta_{(2)} = 0.5 \Delta_{(10)}$$

$$0.5 \Delta_{(2)} = 0.5 \Delta_{(10)}$$

$$? (+1) \leftarrow \text{prob } p_2 > \frac{1}{2} (R=1, S=1)$$

$$\text{prob } p_2 = \frac{1}{2} (R=1, S=0)$$

$$\text{Si prob } p_2 \text{ is free input } (z_m = 1)$$

$$\text{prob } p_2 \text{ is free input } (z_m = 1)$$

$$1 \text{ } z_m \text{ } (R, S)$$

$$R \cdot S + R \cdot \overline{S} \cdot z_{m-1} = R \cdot (S + z_{m-1})$$

2.5. F.p. addition/subtraction with rounding (3)

- Rounding error is not correlated with exponents' difference.

$$1.75 \times 2^{32} + 1.25 \times 2^4$$

$$\approx 1.75 \times 2^{32} + 1.25 \times 2^{31}$$

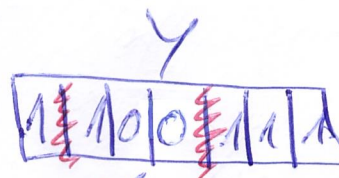
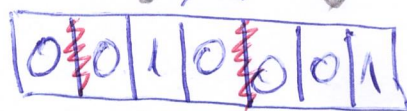
- Use IEEE 754 format **scaled down**:

1 - sign.
3 - exponent \rightarrow bias = $2^{3-1} - 1 = 3$
4 - significand (3 - fractional part of significand)

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

$$Y = -3.75_{(10)} = 11.11 \times 2^0 = 1.111 \times 2^{+1}$$

Packed operands

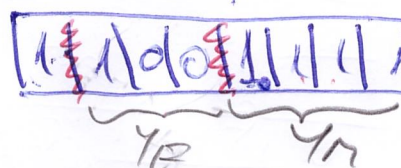
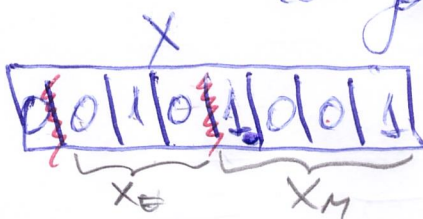


F.p. addition with rounding Algorithm

Step 1 \rightarrow unpack operands

- add the hidden bit to significand.

- check for special cases (zero, $\pm\infty$, NaN)



Step 2 Compute exponents' difference $d = X_E - Y_E$

- if $d < 0 \Rightarrow |X| < |Y| \Rightarrow$ **SWAP** $X \leftrightarrow Y$

choose Y_E as result's exponent

- if $d \geq 0 \Rightarrow$ choose X_E as result's exponent

! Only Y_M has RShift capabilities (in order to save area)

$$d = x_0 - y_0 = 0100_2 - 100_2 = 2_{(10)} - 4_{(10)} = -2 < 0$$

\Rightarrow SWAS $X \leftrightarrow Y$.

$$\textcircled{20} = y_0 \rightarrow z_0 = 4$$

X

111101011111

Y

01011011011

11110111

11010101

11110111

11010101