

3.5 F.p. addition/subtraction with rounding.

- Rounding error is not correlated with operands' difference.

$$1.75 \times 2^{32} + 1.75 \times 2^{29} \rightarrow 4$$

$$1.75 \times 2^{32} - 1.75 \times 2^{31} \rightarrow 1$$

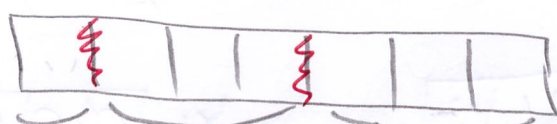
- Use a simplified floating point format
- inspired by IEEE 754

IEEE 754
- single precision
1 - sign
8 - exponent
24 - significand

- fewer bits are used:
= scaled down.
fast

- 1 bit for sign
- 3 bits for exponent field, 1 hidden bit
- 4 bits for the significant
- 3 bits for mantissa

Packed float



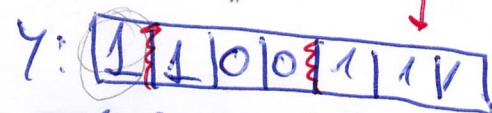
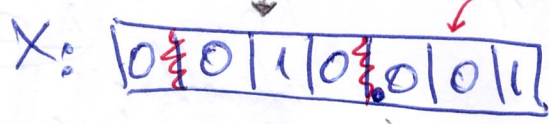
sign exponent packed part of significand

under: bias = $2^{e-1} - 1 = 2^{3-1} - 1 = 3$

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

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

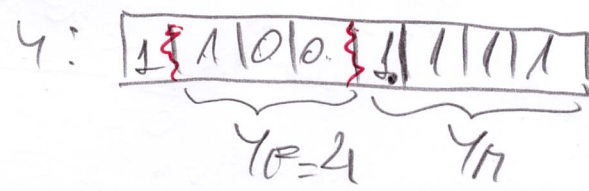
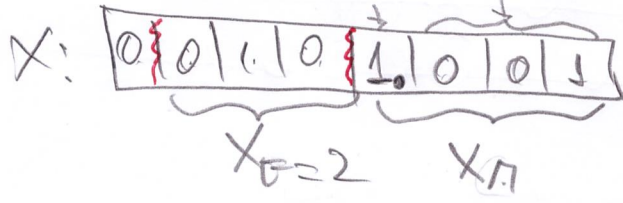
Packed operands



F.p. addition with rounding - Algorithm

Step 1 Unpack operands

- adding the hidden bit (supraunitary bit)
- check for exceptions:
- one or both operands are 0, $\pm\infty$, NaN, ...



Step 2 Compute exponents' difference, $d = X_e - Y_e$ align
 - if $d < 0 \Rightarrow |X| < |Y|$
 \rightarrow SWAP $X \Leftrightarrow Y$
 \rightarrow choose Z_e as Y_e $X \pm Y = (X_n \pm Y_n \cdot 2^{-d}) \cdot 2^{Y_e}$
if $d \geq 0$

- if $d \geq 0$
 \rightarrow choose Z_e as X_e

Save area \rightarrow reduce complexity by only allowing Y_n to be Right Shifted for alignment

$$d = X_e - Y_e = 2 - 4 = -2!$$

\Rightarrow SWAP $X \Leftrightarrow Y$, $Z_e = 4(Y_e)$

X:

1	1	0	0	1	1	1	1
---	---	---	---	---	---	---	---

Y:

0	1	0	1	0	0	1	1
---	---	---	---	---	---	---	---

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

Save area \rightarrow reduce complexity by only allowing Y_n to be Complement of 2

$\text{sign}(X) = 1 (-)$
 $\text{sign}(Y) = 0 (+)$
 \Rightarrow Complement of 2 for Y_n

Y:

0	1	0	1	0	0	1	1
---	---	---	---	---	---	---	---

$Y_n = 0.111$

$$\begin{aligned} X < 0 &\rightarrow X = -|X| \\ Y > 0 &\rightarrow Y = |Y| \\ X + Y &= -|X| + |Y| \end{aligned}$$

subtraction
 $-A \equiv$ Complement of 2 of value of A

for BCD addition:
 $-10_{(10)} \equiv +6_{(10)}$

$$10_{(10)} = \begin{pmatrix} 1 & 0 & 1 & 0 \end{pmatrix}$$

Complement of 2:
 $6_{(10)} = \begin{pmatrix} 0 & 1 & 1 & 0 \end{pmatrix}$

$$\begin{aligned} |X| - |Y| &\text{complement of 2} \\ &= -(-|X| + |Y|) \end{aligned}$$

Step 4. Align Y_n : RShift of Y_n by 1011 bits.

- if Y_n was complemented of 2 in Step 3:
 - introduce bits of 1's in Y_n 's mob while RShifting instead of 0's
- preserve the sticky bits: g, r, s

$$Y_{n2} = \boxed{0.111}$$

$$d = -2$$

was Y_n C2 in Step 3? YES

\Rightarrow RShift Y_{n2} by $|-2| = 2$ bits

- introduce values of 1's while RShifting.

$$Y_{nreal} = \boxed{1.101} \mid \boxed{11} \boxed{0}$$

$g \quad r \quad s$

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

- if $\text{sign}(x) = \text{sign}(y)$
 - if car is generated, preserve it
- if $\text{sign}(x) \neq \text{sign}(y)$
 - if no car is generated \Rightarrow result's signficand is negative
 \rightarrow Complement of 2 the Z_n (result's signficand)
 - if car is generated \Rightarrow result's signficand is positive
 \rightarrow discard car bit

$$\begin{array}{r} X_n = 1.111 \mid \underline{g} \underline{r} \underline{s} \\ Y_{nreal} = 1.101 \mid \underline{1} \underline{1} \underline{0} \end{array} \quad +$$

$$Z_n = \cancel{1}.100 \mid 110$$

car

Z_n is positive and $\text{sign}(x) \neq \text{sign}(y) \Rightarrow$ discard car

Step 6 Renormalisation.

- according to the rules from 3.4.
- determine z_n , can update z_e
- exception checking:

- if $z_e == z_{e_{max}}$ ($2^3 - 2 = 6$) and z_n requires a 1-bit RShift \Rightarrow Overflow

- if $z_e == z_{e_{min}}$ (1) and z_n requires 1 ~~and~~ more bits LShift \Rightarrow Underflow

ignored cont from Step 1

$$z_n = 01.100 \mid \underset{9}{1} \underset{2}{1} \underset{1}{0}$$

\Rightarrow Rule 2) from 3.4 $\begin{matrix} z_{2m} & z_{2n} & z_{2m} & z_{2m} \\ z_n = & 1 & . & 1 & 0 & 0 & \end{matrix} \quad z_e = 4$

Step 7 Calculate values of R, S

- get them from the one rule of 3.4 as the one applied in Step 6

Rule 2) from 3.4.

$$\Rightarrow \underline{R} = \underline{q} = 1$$

$$\underline{S} = (\underline{r} \text{ or } \underline{s}) = \underset{\text{logic OR}}{1 + 0} = \underline{1}$$

Step 8 Rounding of z_n : determine z_n^* from z_n

- takes into account the IEEE 754 rounding mode set for the f.p. addition.

- according to the rounding rules from 3.4.

- if rounding generates a cont \Rightarrow postnormalisation.

Postnormalisation: - 1-bit RShift of z_n^*

- $z_e + 1 \Rightarrow$ possible OVERFLOW

Consider the rounding mode to be set to
round to nearest even

(3)

condition is: $R \text{ odd } (S \text{ or } zero)$

$$\left. \begin{array}{l} Z_{\text{even}} = 0 \\ R = 1 \\ S = 1 \end{array} \right\} R \cdot (S + Z_{\text{even}}) = 1(1 + 0) = 1 \quad \text{TRUE}$$

$$Z_{\text{even}} = 1.100 + \Rightarrow Z_{\text{even}} + 1$$

$$Z_n^* \quad \begin{array}{r} \text{carry} \\ 01.101 \end{array}$$

no carry out was generated
 $\Rightarrow Z_n$ is not modified.

$$Z_n^* = 1.101 \quad Z_n = 4$$

Step 9 Choose the result's sign.

- if $\text{sign}(x) == \text{sign}(y) \Rightarrow \text{sign}(z) = \text{sign}(x)$
- otherwise.

SWAP (Step 2)	Complementing of 2 (Step 5)	$\text{sign}(x)$	$\text{sign}(y)$	$\text{sign}(z)$
YES		+	-	-
YES		-	+	+
NO	YES	+	-	-
NO	YES	-	+	+
NO	NO	+	-	+
NO	NO	-	+	-

$y_e > x_e$
 $|y| > |x|$

$x_e \geq y_e$

SWAP? YES $\Rightarrow \text{sign}(z) = \text{sign}(y)$
 \hookrightarrow before the SWAP

Step 10 Push the result $\Rightarrow \text{sign}(z) = 1(-)$

1 1 0 0 . 1 0 1

$z: \text{sign} = 1$
exponent = 4
significand = 1.101

Verification:

$$X = 0.5625$$

$$Y = -3.75$$

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

$$Z = \boxed{111101011011}$$

unpoising:

$$n_{\text{sig}} = 1$$

$$Z_2 = 100 = 4$$

$$Z_n = \underline{1}.101$$

hidden bit

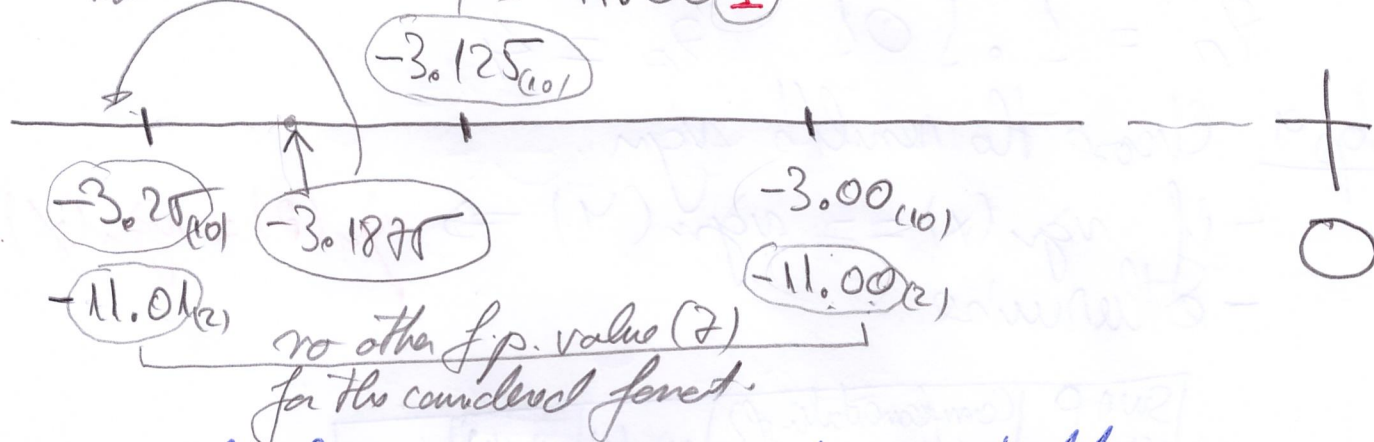
$$Z = (-1)^{n_{\text{sig}}} \times 2^{Z_2 - \text{bias}} \times Z_n =$$

$$= (-1) \times 2^{4-3} \times 1.101 = -1 \times 2 \times 1.101_2 =$$

$$= -11.01_2 = -3.25$$

Rounded:

$$\rightarrow -11.00 \text{ (1)}$$



Design of the personalization shifter.

- cases

Step 6 and Step 7

→ according to the rules from 3.4

- purely combinational design.

Result of Step 5

parallel calc.

$$Z_n = Z_n \ Z_3 \ Z_2 \ Z_1 \ Z_0 \ | \ g \ R \ S$$

Output of Step 6 & Step 7

$$Z_n = 1 \cdot Z_n \ Z_n \ Z_n \ | \ R \ S$$

Normalization cases -

- A) z_n is already normalized (l_{k0})
- B) z_n needs a 1-bit LShift (l_1)
- C) z_n needs a 2-bit LShift (l_2)
- D) z_n needs a 3-bit LShift (l_3)
- E) z_n needs a 1-bit RShift (r_1)

$$z_{2n} = \begin{array}{c|ccc|cc} 1 & z_n & z_n & z_n & R & S \end{array} \quad \text{S} \text{ (circled)} \\
\begin{array}{c|ccc|cc} 1 & z_2 & z_1 & z_0 & g & (z_{2n}) \\ 1 & z_1 & z_0 & g & r & \downarrow \\ 1 & z_0 & g & 0 & 1 & 0 & 0 \\ 1 & g & 0 & 0 & 1 & 0 & 0 \\ 1 & z_3 & z_2 & z_1 & r_0 & (g_{2n}) \end{array}$$

Associate each of the 5 cases to a boolean variable

$$\begin{aligned} z_{2n} &= z_2 \cdot l_{k0} + z_1 \cdot l_1 + z_0 \cdot l_2 + g \cdot l_3 + z_3 \cdot r_1 \\ z_{1n} &= z_1 \cdot l_{k0} + z_0 \cdot l_1 + g \cdot l_2 + z_2 \cdot r_1 \\ z_{0n} &= z_0 \cdot l_{k0} + g \cdot l_1 + z_1 \cdot r_1 \\ R &= g \cdot l_{k0} + r \cdot l_1 + z_0 \cdot r_1 \\ S &= (z_{2n}) \cdot l_{k0} + \downarrow \cdot l_1 + (g_{2n}) \cdot r_1 \end{aligned}$$