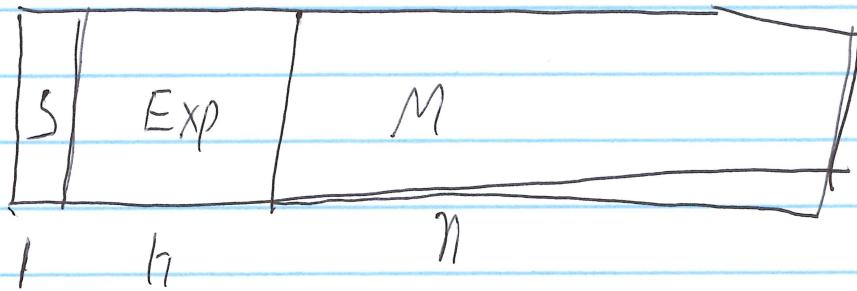


CS2318

Floating Point Operations

01/17/2016



$$(-1)^S \times (1, M) \times 2^{E - \text{Bias}}$$

Bias is $2^{h-1} - 1$

In IEEE Not in assignment

The Exp of $2^{E - \text{Bias}}$ for $\text{mx}(E)$

is considered as the indication of

exception

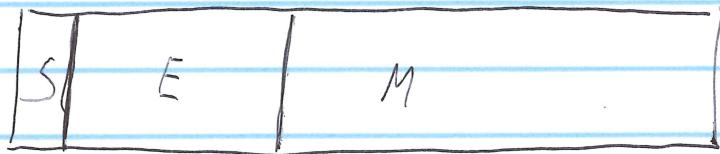
E is "treated" as unsigned

$$\text{mx}(E) = 2^{(k-1)} + 2^{-1} \text{ Bias } 2^{-1}$$

$\Rightarrow 2^{h-1}$ is the exception

2

Conversion from FP to Decimal



$$1 \cdot 5 \quad 4 \quad (-1)^5 \times (1.M) \times 2^4 \quad E = 15$$

FP to Dec

$$\underline{1.10110} \quad 1101 \Rightarrow -(1\frac{13}{16})_2^7$$

$$0.000075 \quad (1.M) = ? \quad \underline{1M}$$

Dec. to FP 5

stop at $1 + |M| + 1$ ($* 2$)

Significand

3

0.00075	0	not in numly
0.0015	0	
0.0030	0	not probably not
0.0060	0	accurate
0.012	0	
0.024	0	
0.048	0	
0.096	0	$l, M = 1,100 \times 2$ - 8
0.192	0	
0.384	0	-8 bc we moved
0.768	0	
0.156	1	the point to the right
0.062	1	
	0	
	0	?
	0	

↓
0 0 0 0 0 0 0 1 1 0 0 0 | 0]

I placed 8 ↑ |

instead of

11 might

bc to small

1 M [0]
10 [1]

4

$$A \equiv S_A E_A M_A$$

On paper

$$B \equiv S_B E_B M_B$$

5 + M rep. might

$$C \equiv S_C E_C M_C$$

be or

For computers

$$S_A M_A \xrightarrow{2^5} X$$

$$S_B M_B \xrightarrow{2^5} Y$$

$$+ 0.5 \times 10^{-3}$$

$$S_C M_C \xrightarrow{SYN} Z$$

D) Align Exp Assume $E_B < E_A$

$$E_A - E_B = h \quad SHR [l, M_B] \text{ by } h$$

Ex) $mls \quad 3.5 \times 10^5 + 0.5 \times 10^3$
 $= 0.005 \times 10^5$

5

1) Align Exp current Exp is E_A

2) produce $l_1 M_A + l_1 M_B$

Get upto $|M_A| + 2$ significant

Lets In the case of $|M_A| + L$

Normalize

Shift point for the right Left

=> shifting A to the right

=> increment Exp $E_C = E_A + 1$

In the case of $|M_A| + 1$ $E_C = E_A$

3) ~~E_C~~ ~~E_{C+1}~~

$C = \sum E_C M_C$ $\sum M_C$ obtained in
the addition

E_C $E_A + 1$

4) Check for ovf ($BCE_A + 1$ is too big)

5) Normalize again (might be needed)

Exempl

6) $A \oplus B^+ = \underline{\underline{010011110}} + \left(1\frac{7}{8}\right)_X 2^4$
 $+ \underline{\underline{0100101100}} + \left(1\frac{3}{4}\right)_X 2^3$

$$E_A = 4 \quad E_B = 3$$

$$l.M_A = 1.1110 \quad l.M_B = 1.1100$$

$$E'_C = 4$$

$$\text{After Align} \quad l.M_B \quad 0.1110$$

Add in 2's result in 5+1

$$\begin{array}{r} 1.1110 \\ 0.1110 \\ \hline 1.01100 \end{array} \Rightarrow 1.0\overset{1}{1}10$$

inc E'_C $E_C = 5$ un biased

0101000110

7

Addition / Subtraction

$$+ \begin{array}{c} S_A \\ E_A \\ M_A \end{array} \quad \begin{array}{c} S_C \\ E_C \\ M_C \end{array}$$

$$- \begin{array}{c} S_B \\ E_B \\ M_B \end{array}$$

- 1) Align EXP
 - 2) Add/Sub $l_M A \pm l_M B$
 - 3) Normalize and Round }
 - 4) check for overflow
 - 5) Normalize and round again

$$10.5 \times 10^3 \times 0.05 \times 10^4$$

8

$$\star S_A E_A M_A$$

$$1 S_B E_B M_B$$

$$S_C E_C M_C$$

1) $S_C = S_A \oplus S_B$

2) $E_C = E_A + E_B$ check for OVF

3) $I.M_C = (I.M_A) \star (I.M_B)$

Result is $2, (1+|M_A|)$

4) Normalize (select the 5 significant bits)

Add Round

5) check for OVF

6) Normalize and round

9

 010011110
 0100101100

$$S_C = 1$$

$$\begin{aligned} E_A &= x - \text{bias} & E_B &= y - \text{bias} \\ \text{Computer: } E_A + E_B &= \text{at the biased notation} \\ \text{Subtract } i \text{ bias from } E_A + E_B \end{aligned}$$

$$Vn \text{ ficsel } E_A = 4 \quad E_B = 3$$

in Li_3 $7 + 15$ 10110

110110 xxxx

$$\begin{array}{r}
 \text{M}_A 1.M_A 1.1110 \\
 \text{M}_B 1.M_B 1.1100 \\
 \hline
 \text{D} \cdot 0000 \\
 \hline
 \begin{array}{r}
 1 \quad 0 \quad 1 \quad 1 \quad 1 \\
 1 \quad 1 \quad 1 \quad 1 \\
 \hline
 1 \quad 1 \quad 1 \quad 1 \\
 \hline
 1 \quad 1.01 \quad 0 \quad 1 \quad 0
 \end{array}
 \end{array}$$

10

11.010010

$$1.\underset{E_C}{\textcircled{0}} \underset{E}{\textcircled{1}} 0$$
$$E_C = 8$$

point to the left

Enc E_C

1101111010

Normalize might be needed infinite

- 1) Convert to Binary
 - 2) on addition / multiplication
- use only n bits
for computing

when we have more than

typical n bits $n = |M_A|$

assume we get

\dots, X

0

000...01, XXXXX

Select the $n+1$ bits 1.XXXXX...

ii

Example

1.1110 need to select 4

1.111
^

1.11011xxx
^

Add for 1, n

1.11011

1.11011.0 round down cut

1.11011 Round up off 1
^ odd bit n+2 for round

1.11011 becomes 1.1110
^