

# CH 3.5 1.3. Representation of fixed-point decimal nos.

## 1.3.1 Binary Coded Decimal (BCD): represent a digit a tetrad

Ex:  $297 = 0010 \ 1001 \ 0111$  BCD  
 tetrad = nibble = 4 bits.  
 fixed weight  $\Rightarrow$  POSITIONAL  
 weight  $10^i \cdot 2^j$   $\leftarrow i = i^{\text{th}} \text{ digit}$   
 $\leftarrow j = \text{rank in tetrad}$

## 1.3.2. Excess of Three (E3) excess == bias - represent digit on a tetrad. = amount that is added/subtr from a code. - E3 adds 3 units to the corresponding BCD digit $\Rightarrow$ E3 is not POSITIONAL. ! - facilitates addition.

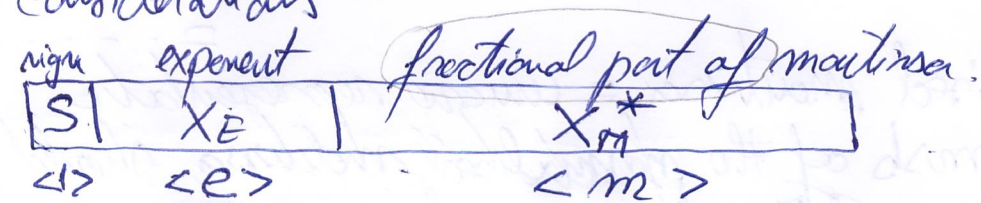
Ex:  $297 = 0101 \ 1100 \ 1010$

## 1.3.3. Two-out-of-five - represents a digit on 5 bits $\leftarrow$ 2 bits are 1s 3 are 0s $C_5^2 = \frac{5!}{2!3!} = \frac{5 \cdot 4}{2} = 10$ Ex: $297 = 00101 \ 10100 \ 10001$

! - error detection using parity codes

## 1.4. Representing floating-point (f.p.) numbers $X = X_M \cdot B^{X_E}$ $X_M$ = mantissa = fixed-point, fractional $2^{\text{SM}}$ $X_E$ = exponent = fixed-point, integer $2^{\text{SM}}$ $B = \text{radix} = 2, \text{ power of } 2$

### 1.4.1 General considerations



$X_M = S \cdot X_M^*$  - sign appended  $X_M^*$   
 $\hookrightarrow$  mantissa.



Constraints for representing f.p. nos.

- (A) Representing  $X_E$
- $X_E$  for 0 should be the smallest possible value
  - fast comparison against 0
  - normalizing errors

fixed-point integer

$X_E$  for 0 - smallest value

SR:  $-2^{e-1} + 1$

CR:  $-2^{e-1}$

$X = X_f \times 2^{X_E}$

$\langle e \rangle$

1 1 ... 1

1 0 ... 0

$X_f$  for 0 - all-0s representation

$X_f$  for 0: 00 ... 00

Putting all fields together.

SR: 0 1 1 ... 1 00 ... 00

CR: 0 1 1 ... 0 00 ... 00

$\langle e \rangle$

$\langle m \rangle$

fast 0 comparison  $\rightarrow$  require the  $X_E$  to be all-0s

Q: how to obtain all-0s config. for  $X_E$  of 0.

A: represent  $X_E$  in a bias code. (an excess code)

bias code

SR:  $+2^{e-1} - 1$

CR:  $+2^{e-1}$

$\Rightarrow 0$ :

SR: 0 00 ... 00 00 ... 00

CR: 0 00 ... 00 00 ... 00

$\langle e \rangle$

$\langle m \rangle$

(B) Represent mantissas:

- redundant! in decimal  $0.237 = 0.0237 \times 10^2$

$= 0.000237 \times 10^3$

normalised mantissa = unique representation

- msb of the fractional part of normalised mantissa is

SR: 1

CR: 5



Ex:  $0.\underline{111}sm = \frac{7}{2^3}$   
 $msb = 1$

$1.\underline{111}sm = \frac{-7}{2^3}$  (2)  
 $msb = 1$

$0.\underline{111}c2 = \frac{7}{2^3}$   
 $S \quad msb = \overline{5}$

$\underline{1}.\underline{001}c2$   
 $S \quad msb = \overline{5}$

Range of values for  $X_M$

$\frac{1}{2} \leq |X_M| < 1$

$X_M^* = 0.\underline{1011} * 2^{X_E}$   
 normalise  $X_M^*$ :  
 $X_M^* = 0.\underline{1011} * 2^{X_E-1}$

### 1.4.2 IEEE 754

portable formats  $\begin{cases} 32 \text{ bits (single precision)} \\ 64 \text{ bits (double precision)} \\ 128 \text{ bits (quaduple precision)} \end{cases}$

interchange formats  $\rightarrow$  16 bit (half-precision)  $\rightarrow$  NN

Simple precision format

$\rightarrow$  fractional part of significant



$X_S = \underline{1} \cdot X_S^*$

hidden bit,  
integer bit

fractional part of significant

Range of values for  $X_S$ :

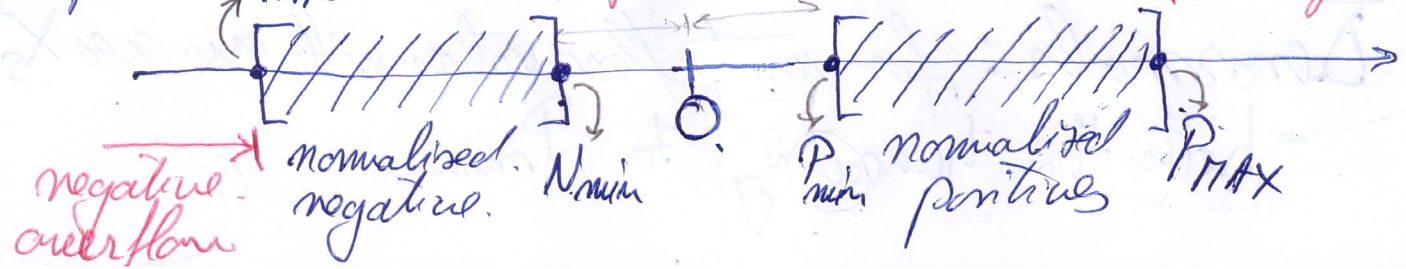
$X_E = \text{bias of } 127 (= 2^{8-1} - 1)$

$1 \leq X_S < 2$

Value of  $X = (-1)^S * 2^{X_E - \text{bias}} * (1 \cdot X_S^*)$

Examples:

negative underflow  $\rightarrow$  positive underflow  $\rightarrow$  positive overflow





$$P_{MAX}: X_E = X_{E_{max}} - 1 = 254$$

$$X_S^* = \underset{\substack{\uparrow \\ \text{binary} \\ \text{point}}}{.} \underset{\substack{\leftarrow 23 \rightarrow \\ \text{bits}}}{11} \underset{\substack{\leftarrow 23 \rightarrow \\ \text{hidden} \\ \text{bits}}}{1} = 1 - 2^{-23}$$

$$P_{MAX} = (-1)^0 * 2^{254-127} * (1.11 \dots 1) = 2^{127} * (1 + 1 - 2^{-23}) = 2^{127} * (2 - 2^{-23}) \approx 3.4 \times 10^{38}$$

$$P_{min}: X_E = X_{E_{min}} + 1 = 1$$

$$X_S^* = \underset{\substack{\leftarrow 23 \rightarrow}}{.00 \dots 0} = 0$$

$$P_{min} = (-1)^0 * 2^{1-127} * (1.00 \dots 00) = 2^{-126} \approx 1.18 \times 10^{-38}$$

$X_{mod} < 4$

Special quantities in IEEE 754

(A) Not a Number (NaN)

- result of:  $\infty \pm \infty$ ,  $0 * \infty$ ,  $\frac{\infty}{\infty}$ ,  $\frac{0}{0}$ ,  $\sqrt{\text{negative}}$ ,  $X \bmod 0$

Represent:  $X_E = X_{E_{max}} (255)$   
 $X_S^* \neq 0$  - indication of the operation that generated NaN

(B) Infinity: for overflow

- result of:  $X \div 0$ ,  $X \pm \infty$ ,  $\sqrt{+\infty}$   
 - better than truncating to  $P_{MAX}$

Represent:  $X_E = X_{E_{max}} (255)$   
 $X_S^* = 0$

(C) Zero: Represent:  $X_E = X_{E_{min}} (0)$   
 $X_S^* = 0$

(D) Denormals: for underflow values with non-zero  $X_S^*$   
 - better than truncating to  $P_{min}$

Representation  $X_E = X_{E_{min}} (10)$   
 $X_S^* \neq 0$

(3)

Value of the denormal no:

$$X_D = (-1)^S * 2^{1-bias} * (\underbrace{0.X_S^*}_{\text{hidden bit for denormals}})$$

denormal

Examples:

$$612.8046875_{10} = 1001100100.1100111_2$$

612 -	0.8046875	x2
512 $2^9$	1.6093750	x2
100 -	1.2187500	x2
64 $2^6$	0.4375000	x2
36 -	0.8750000	x2
32 $2^5$	1.7500000	x2
4 $2^2$	1.5000000	x2
	0.0000000	

$$612.8046875 = 1.0011001001100111 * 2^9$$