

## (1) Unsigned

→ Carry always indicate overflow.

## (2) Signed

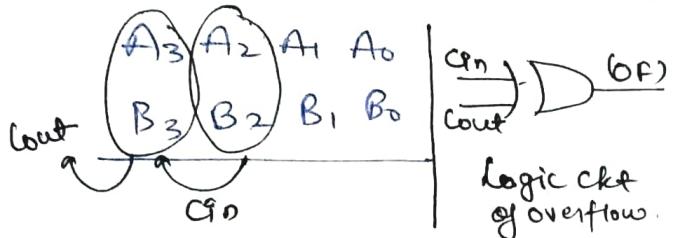
→ Two cases for overflow cond'n

- two (+ve) add to give (-ve) no.
- two (-ve) add to give (+ve) no.

Note → In case of signed, carry doesn't indicate overflow.

Trick

Assume two 4 bit numbers



If  $Cin \neq Cout$  then overflow

If  $Cin = Cout$  then  $\Rightarrow$  no overflow

→ Verify all the previously solved examples by above  $Cin$ ,  $Cout$  tricks

## Classification of binary Codes

### Binary Codes

Alphanumeric  
Codes

|                   |                    |
|-------------------|--------------------|
| 1                 | ↓                  |
| ASCII<br>(7 bits) | EBCDIC<br>(8 bits) |

Numeric  
Codes

↓  
BCD

BCD

↓

Weighted

8421

↓

Self  
complementary

Non-weighted

Gray code

Xs-3 code

Gray code called as unit distance codes  
Xs-3 code is sequential and self complementary

## Self Complementary Codes

In self complementary codes,

1's complement of patterns represents the 0's complement of digits

Suppose a weighted code

631-1

$$\begin{array}{r} 0 \ 1 \ 0 \ 1 \\ \hline \text{Pattern} \end{array} = \frac{2}{\substack{\downarrow 1's \text{comp.} \\ \downarrow 0's \text{comp.}}} \quad \downarrow$$

$$1 \ 0 \ 1 \ 0 = (2-2)=7$$

and  $1010 = 7$  in (631,-1) code

So, 631-1 is a self comp. codes

Ex

2 4 2 1

0 0 1 0 = 2

$\downarrow 1's \text{comp.}$        $\downarrow 0's \text{comp.}$

1 1 0 1 = 7

$\boxed{2 \times 1 + 4 \times 1 + 1 \times 1 = 7}$

## 8421, Xs-3, 3321 Codes

| Decimal no. | 8421 | Xs-3 | 3321  |
|-------------|------|------|-------|
| 0           | 0000 | 0011 | 00000 |
| 1           | 0001 | 0100 | 00001 |
| 2           | 0010 | 0101 | 00100 |
| 3           | 0011 | 0110 | 00111 |
| 4           | 0100 | 0111 | 01011 |

|        |
|--------|
| 3321   |
| 2421   |
| 84-2-1 |
| 631-1  |

|   |                                   |                          |      |
|---|-----------------------------------|--------------------------|------|
| 5 | 0101                              | 1000                     | 1010 |
| 6 | 0110                              | 1001                     | 1100 |
| 7 | 0111                              | 1010                     | 1101 |
| 8 | 1000                              | 1011                     | 1110 |
| 9 | 1001                              | 1100                     | 1111 |
|   | (10-15)<br>are<br>invalid<br>here | Add 3 to<br>8421<br>code |      |

→ In 3321 you might get end up with different combn for same example.

'3' can be represented as

$$\begin{array}{r} \text{Bcd} \quad \begin{array}{cccc} 3 & 3 & 2 & 1 \end{array} \\ \hline 3 \leftarrow 0 \ 0 \quad | \quad | \quad | \\ 4 \leftarrow 0 \quad 1 \quad 0 \quad 0 \\ 8 \leftarrow 1 \quad 0 \quad 0 \quad 0 \end{array} \xrightarrow{\text{3 combn}}$$

Go for min. BCD number

→ Choose the combination that gives min. BCD number upto Decimal number 4.

~~For no. greater than 4 take its 9's complement and then take 3's complement of pattern denoting obtained no.~~

Ex - for 5 do  $(9-5)=4$

4 represent by 0101 then take its 3's complement

$$\text{so } 5 = \underline{1010}$$

Similarly obtain 3321 code for 6, 7, 8, 9.

→ for a self complementary code sum of weights should be 9, for weighted code only.

Ex 3321 code is self Comp. So  $3+3+2+1 = 9$ .

It is necessary condition not sufficient condition.

Suppose 'a' weighted code

$$\begin{array}{cccc} a & b & c & d \\ b \leftarrow 0 \quad 1 \quad 0 \quad 0 & \rightarrow \text{say } x \text{ in BCD} \\ 1's \leftarrow 1 \quad 0 \quad 1 \quad 1 & \rightarrow (9-x) \\ \downarrow & \downarrow & \downarrow & \downarrow \\ a+c+d \end{array}$$

$$\Rightarrow (9-x+x) = a+b+c+d$$

$$\Rightarrow \underline{a+b+c+d = 9}$$

→ 8421 is sequential code because we get another number by just adding 1. Similarly X5-3 is also sequential.

→ 3321 is not sequential code

### Examples on Code

Which of the following is invalid

8421 code.

i) 0010    ii) 1100    iii) 1001

iv) 0111

i) → 2 (valid) ii) 1100 → (12) (invalid)

bcz 8421 is BCD ranges in (0-9).

Represent  $(863)_{10}$  in a) BCD  
② Binary    c) 2241    d) 3321

(e)  $\times 5 - 3$ .

$$\text{Sol}^n \quad \textcircled{1} \quad (863)$$

|      |      |      |
|------|------|------|
| 1000 | 0110 | 0011 |
|------|------|------|

$$\text{BCD} \rightarrow (1000\ 0110\ 0011)$$


---

② Binary

$$\begin{array}{r|rr} 16 & 863 \\ \hline 16 & 53 \\ \hline 16 & 3 \\ \hline & 3 \end{array} \quad \begin{array}{r|c} 15 & \rightarrow F \\ \hline 5 \\ \hline 3 \end{array}$$

$$(35F)_{16} \rightarrow (0011\ 0101\ 1111)_2$$

$$\underline{(1101011111)_2}$$

③ 2 2 4 1

~~This~~ is self complementary code.

$$\begin{array}{ccc} 8 & 6 & 3 \\ \downarrow \text{9's comp} & \downarrow \text{a's comp} & \downarrow \\ 1 & 3 & \underline{(0101)} \\ (0001) & (0101) & \\ \text{l's} \left\{ & \downarrow \text{1's comp} & \\ \underline{1110 \rightarrow 8} & (1010) \rightarrow 6 & \end{array}$$

$$(863)_{10} \rightarrow (1110\ 1010\ 0101)_{2241}$$

Since 8, 6 are greater than 4

so we go for a's comp. and then l's comp.

④ 3321

$$(1110\ 1100\ 0011)$$

Approach same as 2241  
for 3321

⑤  $\times 5 - 3$

~~$$\begin{array}{r|rr} 10 & 8 & 6 & 3 \\ \hline 10 & +3 & +3 & +3 \\ \hline 11 & 9 & 6 & \\ \hline (1011) & (1001) & (0110) \end{array}$$~~

$$(863)_{10} = (1011\ 1001\ 0110)_{\times 5 - 3}$$

BCD Addition

$$\begin{array}{r} 3 \rightarrow 0011 \\ + 6 \rightarrow 0110 \\ \hline 9 \quad \underline{1001} \quad \text{Ans} \end{array}$$

But if sum  $> 9$  then we have to do some modification bcoz (10-15) form invalid binary pattern.

$$\begin{array}{r} \text{Ex} \quad 3 \rightarrow 0011 \\ 7 \rightarrow \underline{0111} \\ \hline 10 \quad \underline{1010} \end{array}$$

(1010) is invalid in BCD

To represent ~~10~~ 10 decode 10 as  $\underline{\underline{0001}}\ \underline{\underline{0000}}$  but

it not valid in decimal

So, (10-15), 6 combinations are invalid. So for

Sum  $> 9$  add 6 to it and make it valid, as,

$$\begin{array}{r}
 3 \quad 0011 \\
 7 \quad 0111 \\
 \hline
 10 \rightarrow 1010 \\
 + 6 \quad + 0110 \\
 \hline
 0001 \quad 0000 \\
 \hline
 . \quad 1 \quad 0
 \end{array}$$

9

9

18 → Max. output we can get bcoz max. no. allows is 9 but we can't go beyond 15

$$\begin{array}{r}
 1001 \\
 1001 \\
 \hline
 10010 \rightarrow 18 \\
 \downarrow \\
 0001 \quad 0010 \\
 \hline
 1 \quad 2
 \end{array} \quad (12) \text{ BCD.}$$

Add -6 and get answer

$$\begin{array}{r}
 0001 \quad 0010 \\
 0000 \quad 0110 \\
 \hline
 00011000 \\
 \hline
 1 \quad 8
 \end{array} \quad (-18) \text{ BCD.}$$

If sum > 9 then add -6 to the sum to get answers

$$\begin{array}{r}
 1 \quad 2 \quad 3 \\
 7 \quad 8 \quad 9 \\
 \hline
 9 \quad 1 \quad 2
 \end{array}$$

In BCD

|   |   |   |                      |
|---|---|---|----------------------|
| $  \begin{array}{r}  0001 \\  0111 \\  \hline  1001  \end{array}  $ | $  \begin{array}{r}  0010 \\  1000 \\  \hline  1011  \end{array}  $ | $  \begin{array}{r}  0011 \\  1001 \\  \hline  1100  \end{array}  $ | → 12 out<br>of range |
| <u>0001</u>   | <u>0010</u>   | <u>1100</u>   |                      |

0000      0001      0010

9      1      2

(912) BCD

→ Correction used in BCD addition is +6.

### XS-3 Addition

If carry occurs add (+3) as correction.

If carry doesn't occur add (-3) as correction.

Input and output should have same no. of digits.

Ex ① 7 → 0111 + 011 → 1010

$$\begin{array}{r}
 2 \rightarrow 0010 + 011 \rightarrow 010 \\
 \hline
 9
 \end{array}$$

↓  
No carry

→ Subtract 3

$$\begin{array}{r}
 1111 \\
 0011 \\
 \hline
 1100 \rightarrow \text{represent } 9 \text{ in XS-3}
 \end{array}$$

② 07 → 0011 1010

$$\begin{array}{r}
 03 \rightarrow 0011 0110 \\
 \hline
 10 \quad \text{---} \quad \text{---} \rightarrow \text{Cont}
 \end{array}$$

Because of two digits in o/p we make to add '0' to make i/p of two digits

|          |        |
|----------|--------|
| 0011     | 1010   |
| 0011     | 0110   |
| 0111     | 0000   |
| (→) 0011 | 0011   |
| 0100     | 0011   |
| (1)      | 0      |
| xs-3     | (xs-3) |

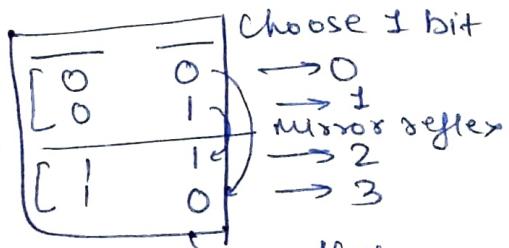
Gray Code (Reflexive, Unit distance, code  
cyclic code, non-weighted)

for 1 bit — either 0 or 1 can fill it.

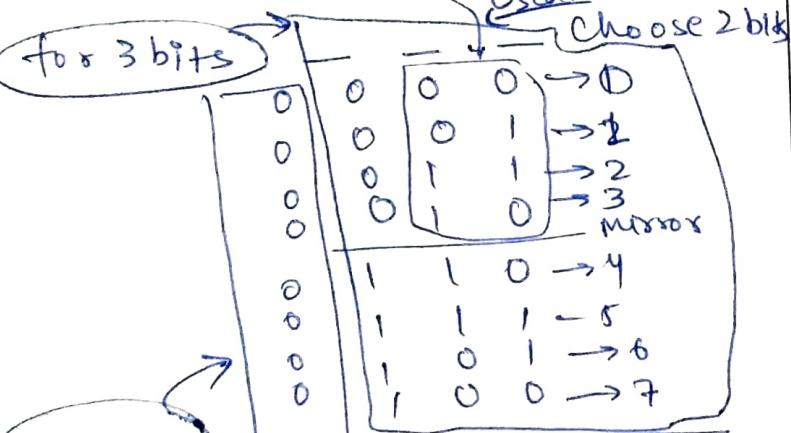
for 2 bit

Put 0 in first half

Put 1 in 2nd half.



for 3 bits



for 4 bits  
Choose bits

|   |          |                   |
|---|----------|-------------------|
| 1 | 0 0 → 8  | Mirror for 4 bits |
| 1 | 0 1 → 9  |                   |
| 1 | 1 1 → 10 |                   |
| 1 | 1 0 → 11 |                   |
| 1 | 0 1 → 12 |                   |
| 1 | 0 1 → 13 |                   |
| 1 | 0 0 → 14 |                   |
| 1 | 0 0 → 15 |                   |

that is why it is called as reflexive

Ex → Represent 12 in

① BCD

(0001 0010)

② 3321

(0001, 0010)

③ Gray (1010)

④ Binary (1100)

→ Unit distance code because on going from one no. to another just one bit is changed only.

→ Cyclic bcoz distance between last number and first number is also one.

Ex  $0 \rightarrow 0000$  only 1 bit change  
 $15 \rightarrow 1000$

Binary to Gray and Vice Versa.

12 → Binary (1100)

Binary to Gray

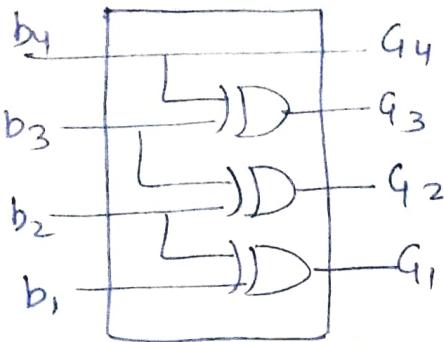
Binary      1    1    0    0  
↓      ⊕    ↓      ⊕    ↓      ↓  
Gray      1    0    1    0

(1010) is 12 in Gray.

$b_4 \ b_3 \ b_2 \ b_1$   
↓      ⊕    ↓      ⊕    ↓  
G<sub>4</sub>    G<sub>3</sub>    G<sub>2</sub>    G<sub>1</sub>

$G_4 = b_4, G_3 = b_4 \oplus b_3$

$G_2 = b_2 \oplus b_3, G_1 = b_1 \oplus b_2$



Binary to Gray Converter.

$$\text{Binary } 15 \rightarrow (1111)$$

$$\text{Gray } 15 \rightarrow (1000)$$

Binary 200 to gray 200.

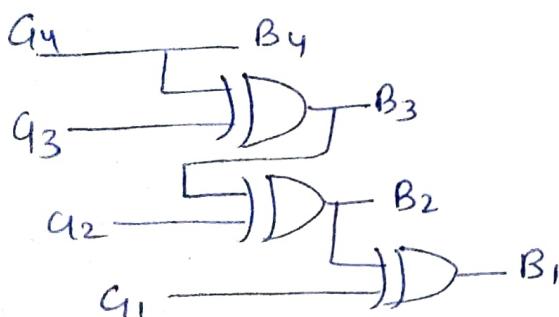
$$\begin{array}{cccccc} 1 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow & \downarrow \\ 1 & 0 & 1 & 0 & 1 & 1 & 0 & 0 \end{array} \xrightarrow{\text{Binary (200)}} \xrightarrow{\text{Gray (200)}}$$

→ Gray to Binary.

$$\begin{array}{c} \text{Gray } 12 \rightarrow 1010 \\ \quad \quad \quad \downarrow \quad \downarrow \quad \downarrow \quad \downarrow \\ \text{(12) Binary} \leftarrow (1100) \end{array}$$

$$G_4 = B_4, \quad B_3 = B_4 \oplus G_3$$

$$B_2 = B_3 \oplus G_2 \quad B_1 = G_1 \oplus B_2$$



$$B_3 = G_4 \oplus G_3$$

$$B_2 = G_4 \oplus G_3 \oplus G_2$$

$$B_1 = G_4 \oplus G_3 \oplus G_2 \oplus G_1$$

### Error Detection

We know that BCD have range (0-9). So if error occurred at,

$$\begin{array}{c} 0000 \\ \downarrow \\ 0001 \end{array} ] \quad \text{1 bit gets corrupted}$$

but here both 0 and 1 lies in range 0-9. It means if both the pattern before and after error are valid pattern then it is impossible to detect the error.

$$\begin{array}{c} 0000 \text{ (0)} \\ \downarrow \\ 1100 \text{ (10)} \end{array} ] \quad \begin{array}{l} \text{Here we can detect} \\ \text{error - } 1100 \text{ doesn't} \\ \text{lies in range (0-9).} \end{array}$$

So, (1100) is invalid codes

→ To detect 1 bit error the distance between 2 patterns should be atleast 2.

$$\begin{array}{c} 0000 \\ 0101 \end{array} ] \quad \text{Error can be detected}$$

$$\begin{array}{c} 0110 \\ 0010 \end{array} ] \quad \text{Error can't be detected}$$

~~error in n-bits~~  
→ To detect ~~n~~ bit error the distance should be atleast  $(n+1)$ , also called hamming distance.

### Error Correction

→ It includes detection and correction

→ for single bit error correction we need all valid combination with distance 3 and combination other than distance 3 is invalid.

→ To correct 2 bit error correction distance between two valid patterns should be atleast 5.

→ for 3 bit error correction, distance between two valid pattern should be atleast 7.

~~for t-bit error correction,~~ the min. Hamming distance should be  $(2t+1)$ .

→ Bcoz of error correction, distance between two valid patterns can be very high and it reduces the no. of valid combinations. So, to get all the valid combination we have to increase the number of bits in pattern.

→ Min. distance b/w any two pattern called Hamming distance.

Note → Just remember the imp. formula for Hamming distance in error correction and detection. All the Ques will get solved but if very much needed only then watch the video again in downloaded videos.

Hamming Code :- (1-bit error correction code)

Suppose message have  $m$ -bits and we add ' $p$ ' parity bits. So message finally contains

$m+p$  bits

Assume there is 1 bit error or no error in  $(m+p)$  bits. So, for  $(m+p)$  bits total no. of combination of errors equals to

$$\frac{(m+p)!}{p!} \downarrow \text{No error}$$

for any bit of  $(m+p)$  can get corrupted

So, condition is,

$$2^p \geq (m+p+1)$$

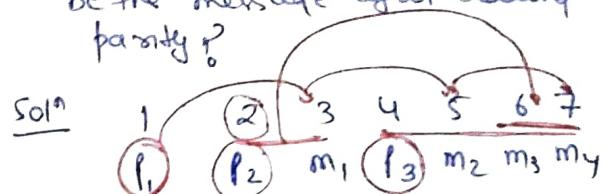
for adding no. of parity bits to check errors.

Suppose  $m=4$

$$\Rightarrow 2^p \geq (p+5)$$

So,  $p=3$  satisfy above cond<sup>n</sup> means we send total  $(m+p)=7$  bits.

Ex: If  $m=4$  and  $p=3$ , message ( $m$ ) is 0101. What should be the message after adding parity?



Add parity bits at power of 2.

P<sub>1</sub>    P<sub>2</sub>    0    P<sub>3</sub>    1    0    1

(P<sub>3</sub>)    (P<sub>2</sub>)    (P<sub>1</sub>)

| C <sub>1</sub> | C <sub>2</sub> | C <sub>3</sub> |    |
|----------------|----------------|----------------|----|
| 0              | 0              | 0              | -0 |
| 0              | 0              | 1              | -1 |
| 0              | 1              | 0              | -2 |
| 0              | 1              | 1              | -3 |
| 1              | 0              | 0              | -4 |
| 1              | 0              | 1              | -5 |
| 1              | 1              | 0              | -6 |
| 1              | 1              | 1              | -7 |

P<sub>1</sub> → form pair of 1, P<sub>2</sub> form pair of 2, P<sub>3</sub> form pair of 4 and check at distance equals to no. of elements in pairs

$E_3$  is 1 at (1, 3, 5, 7)  
 $C_2$  is 1 at (2, 3, 6, 7)  
 $C_1$  is 1 at (4, 5, 6, 7)

so, every check bit will take care of respective position to correct errors.

→ Generally, by default we go for even parity.

→ So,  $P_1 = 0$  bcoz 1, 5, 7 already have 1's. so 2 (1's) means even parity.

Similarly  $P_2 = 1$  so make even no. of pattern.

$P_3 = 0$ . so message should be

0 1 0 0 1 0 1  
 ↓      ↓      ↓  
 $P_1 \quad P_2 \quad P_3$

Suppose corrupted message is

0 1 0 0 0 1  
 $\frac{C_1}{1} \quad \frac{C_2}{0} \quad \frac{C_3}{1}$  [check bits]  
 (5).

Each check bit check according to respective position for even parity

→ Location 5 have error, so we can change it to

→ Check bit is 1 means error.

→ We can have error in any (m+p) combn or 1 for no correction. So total combn possible is (m+p+1)

## Floating point Conversion

### Binary to decimal

Ex ① 110.101

$\begin{array}{ccccccc} & \uparrow & \uparrow & \uparrow & \uparrow & \uparrow & \uparrow \\ & 2^2 & 2^1 & 2^0 & 2^{-1} & 2^{-2} & 2^{-3} \end{array}$

$$= 6 + (0.5 + 0.25 + 0.125)$$

      1      0      1

$$= 6 + (0.5 + 0.125)$$

$$= \underline{\underline{6.625}}$$

② 1101.01

$$= 13 + (0.25)$$

$$= \underline{\underline{13.25}}$$

### Decimal to binary

① 0.625 to binary.

$$\begin{array}{r} 0.625 \\ + 0.625 \\ \hline 1 \leftarrow 1.250 \\ 0.250 \\ 0 \leftarrow 0.500 \\ 0.500 \\ \downarrow 1 \leftarrow 1.000 \end{array} \quad (0.\underline{101})_2$$

→ Stop when you get all zeroes.

② 0.25 to binary

$$\begin{array}{r} 0.25 \\ 0.25 \\ \hline 0 \leftarrow 0.50 \\ 0.50 \\ \downarrow 1 \leftarrow 1.00 \end{array} \quad (0.\underline{01})_2$$

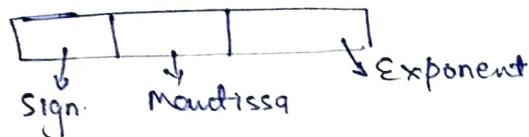
$$\begin{array}{r} 0.625 \\ 0.625 \\ \hline 1 \leftarrow 1.250 \\ 0.250 \\ 0 \leftarrow 0.500 \\ 0.500 \\ \hline 0 \leftarrow 0.904 \\ 0.904 \\ \hline 0 \leftarrow 0.808 \\ 0.808 \\ \hline \end{array} \quad \begin{array}{r} 1 \rightarrow 1.616 \\ 0.616 \\ \hline 1 \leftarrow 1.232 \\ \vdots \end{array}$$

So, for some no. it is not possible to give binary representation of fixed length.

→ In number can't represent completely in binary then we represent only some part of it.

### Floating point representation - I

$$(4.5)_{10} = (100.1)_2$$



$(100.1)_2$  can be written as  
 $(0.1001 \times 2^3)$

Evaluate the value as

$$0.1001 \times 2^3$$

↑      ↑  
2^{-1}    2^{-4}

$$\Rightarrow (2^4 + 2^1) \times 2^3$$

$$\Rightarrow 2^2 + 2^1$$

$$\Rightarrow 4 + 0.5 = (4.5)$$

So,  $(0.1001 \times 2^3)$  is a correct representation, as

$0 | 1001 | 11$

So here is a problem that  $(0.01001 \times 2^4)$  also represent  $(4.5)$  so we will get more than one combination for same number.

→ So, we will use mantissa, exponent method with normalization to avoid more combinations problem.

→ Normalization says use only one combination for a number.

→ Normalization says that move the decimal upto more significant '1'.

~~elsewhere~~ → followed by anything after dot

1 0 0 . 1

So,  $(0.1001 \times 2^3)$  is only one valid combination

$0 | 1001 | 11$

Also called as 'explicit normalization.'

→ Mantissa and exponent can replace their positions

Suppose mantissa can hold 5 bits and exponent needs 4 bits then  $(0.1001 \times 2^3)$  can be written as

$0 | 10010 | 1001$

↑      ↑  
Add zero    Add '0's  
after      before

Ex:  $0.00101$

Above no. be represented

$$\underline{0.101 \times 2^{-2}}$$

We can represent  $(-2)$  (Eve) number using 1's or 2's complement but we use 2's complement

$0 | 10100 | 1110$   
(-2)

$$2 = 0010$$

$$1\text{'s comp} = 1101$$

$$2\text{'s comp} = 1110$$

→  $(-1)^S 0.M \times 2^E$  use to get the original no. from representation.

Biasing → Converts every (+ve) number to (-ve).

→ It is difficult to check (+ve) exponent. So we generally prefer to convert it into (-ve) value and then write its

$$Ex \rightarrow 0.101 \times 2^{-2}$$

-2 will represent as  
 $-2 + 16 = 14$ .

| M(4)     | E(5)  |
|----------|-------|
| 0   1010 | 01110 |

↳ biased exponent

Suppose you have exponent of n bits then in 2's comp. you can represent the range  $-2^{n-1} + (2^n - 1)$  then, bias used is  $\underline{(2^{n-1})}$  means max. (-ve) no. .

$$\rightarrow \text{on 5 bit } 2^{n-1} = 16$$

in case of biased exponent we can get the number as

$$(-1)^S \times 0.M \times 2^{E-\text{bias}}$$
$$\Rightarrow (-1)^S \times 0.1010 \times 2^{14-16}$$
$$= \underline{0.101 \times 2^{-2}} \text{ Ans.}$$

Ex:  $4.875$   
↓  
 $(100.111)$

→ Explicit normalization

$$0.100111 \times 2^3$$

|  |  |      |
|--|--|------|
|  |  | M(5) |
|--|--|------|

To save 6 bit mantissa in 5 bit field we have to neglect a bit.

$$10011(1) \text{ - 4 bits}$$

| E(4)     | M(5)  |
|----------|-------|
| 0   1011 | 10011 |

$$\text{for 4 bit } -2^3 = -8 \text{ (max. (+ve) no.)}$$

for exponent 3 in  $2^3$  we get the final exponent after adding 8 (biasing)

$$\text{so, } E = 8 + 3 = (1011)11$$

$$(-1)^S \times 0.M \times 2^{E-\text{Bias}}$$

$$\Rightarrow (-1)^0 \times 0.10011 \times 2^3$$
$$= 100.11$$
$$= (11.75) \neq 4.875$$

So, we have two methods to avoid this,

- ① Increase no. of bits
- ② Implicit Normalization.

100.111 can be written in implicit normalization as

$$1.00111 \times 2^2.$$

It means move left upto 1 and don't cross the one, comes last to

| S        | E(4)  | M(5) |
|----------|-------|------|
| 0   1010 | 00111 | 1    |

$$\frac{2+8=10}{\text{Biased exponent}}$$

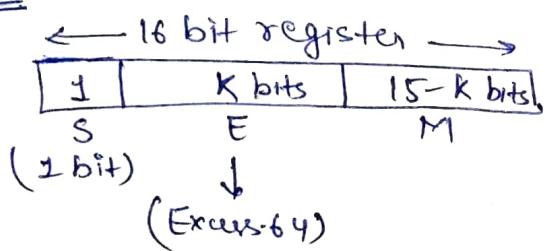
$$(-1)^S (1.M) \times 2^{E-\text{Bias}}$$

$$(-1)^0 (1.00111) \times 2^2$$

Example: Consider a 16-bit register of the following format is used to store a floating point number. Mantissa (M) is denoted as normalised signed magnitude fraction, Exponent(E) is expressed as excess-64 form. Base of system is 2.

- ① How many bits are allocated for fractional mantissa?

Sol<sup>n</sup>



With K bit max. negative no. possible =  $-2^{K-1}$

$$\text{Bias } 2^{K-1} = 64$$

$$\Rightarrow K = 7$$

for fractional mantissa = 8 bits

- ② What is the expression for decimal values

$$[(-1)^S \times (1.M) \times 2^{E-64}]$$

→ In general we use implicit norma.

- ③ What is the largest no. represented in base 10.

| S | E       | M        |
|---|---------|----------|
| 0 | 1111111 | 11111111 |

$$(-1)^0 \times 1.1111111 \times 2^{127-64}$$

$$= 1 \cdot \left( \frac{1}{2} + \frac{1}{4} + \dots + \frac{1}{2^8} \right)$$

$$= (1 + (0.1111111)) 2^{63}$$

$$= (1 + 1 - 2^{-8}) \times 2^{63}$$

$$= \underline{\underline{2^{64}-255}} \text{ Ans!}$$

- ④ What is the 16 bit pattern for  $(-7.5)_{10}$ .

$$= -(111 \cdot 1)_2 = \underline{\underline{-(1.111 \times 2^2)}}$$

|   |          |          |
|---|----------|----------|
| 1 | 1000010  | 11100000 |
| S | E<br>(7) | M<br>(8) |

$$E = 64 + 2 = 66$$

$$\underline{\underline{(-1)^1 \times 0.111 \times 2^{66-2}}}.$$

Example: Question is same as previous examples

- ⑤ What is the difference between 1<sup>st</sup> smallest (+ve) number and 2<sup>nd</sup> smallest (+ve) number

|       |       |       |
|-------|-------|-------|
| (-ve) | →     | (+ve) |
| ————— | ————— | 0     |

| S(1) | E(7)    | M(8)     |
|------|---------|----------|
| 0    | 0000000 | 00000000 |

← 16 →  
2<sup>16</sup> possible comb<sup>n</sup>

→ Numbers are densely distributed near '0' and sparsely distributed far from zero

1<sup>st</sup> smallest no. is,

| S | E       | M        |
|---|---------|----------|
| 0 | 0000000 | 00000000 |

$$(1 \cdot 0) \times (2)^{0-64}$$

$$= \underline{\underline{1 \cdot 0 \times 2^{-64}}}$$

2nd smallest no.

| S | E       | M          |
|---|---------|------------|
| 0 | 0000000 | 1000000001 |

$$(1.00000001) \times 2^{-64}$$

Difference

$$= (1.00000001 - 1.0) \times 2^{-64}$$

$$= \underline{2^{-72}}$$

Extension

1st highest positive no.:

$$(1.11111111) \times 2^{63}$$

2nd highest number

$$(1.11111110) \times 2^{63}$$

Difference = 255

So, difference is very small near 0 means densely packed and very far away from zero sparsely packed between two consecutive numbers.

⑥ What is the pattern denotes the value  $(-10.125)_b$

$$-(10000.001)_2$$

|   |         |         |           |
|---|---------|---------|-----------|
| 1 | 1000100 | 0000001 | → Pattern |
|---|---------|---------|-----------|

Implicit normalization

$$-\left(\frac{1.0000001}{M} \times \frac{2^4}{E}\right)$$

$$\underline{4+64=68} \rightarrow \text{to store in exponent}$$

Theory floating point representation

Floating point numbers stored in exponential and mantissa form.

→ Generally we represent mantissa as normalized signed magnitude fractions.

→ Normalization can be explicit or implicit.

→ Exponent denoted in biased form.

→ Biased exponent is an unsigned number, which can represent in signed exponent of original numbers.

$$\text{True Exp.} = \text{Biased Exp.} - \text{Bias}$$

→ Representation

| S                          | E | M |
|----------------------------|---|---|
| $\leftarrow n \rightarrow$ |   |   |

$$S = 1 \text{ bit} \quad E = k \text{ bit}$$

$$M = (n-1-k)$$

$$\text{Biased Exp.} = 0 \leq E \leq 2^k - 1$$

$$\text{Bias} = 2^{k-1} \quad (\text{2's complement})$$

IEEE Standards

1985 → Single, double precision

2008 → half, quad single precision

1985 standards → More popular.

→ Base of system is 2.

→ Single Precision → 32 bits  
Double → 64 bits

→ Floating point no. can be represented as -

•) fractional form

•) 9mplt Normalized form

→ Certain mantissa and exponent combinations doesn't represent any number (Not a number - NAN)  
it's used to represent 0 and ∞.

### Single Precision

| S | E | M  |
|---|---|----|
| 1 | 8 | 23 |

← 32 bits →

bcz of 8 bits its (Excess-127)

| Sign. | E(8)              | M(23)                        | Value               |
|-------|-------------------|------------------------------|---------------------|
| 0/1   | 00000000<br>E=0   | 000---0<br>M=0               | ± 0                 |
| 0/1   | 11111111<br>E=255 | 000---0<br>M=0               | ± ∞                 |
| 0/1   | 1 ≤ E ≤ 255       | <u>M = xx--x</u><br>anything | Implicit Normalized |
| 0/1   | E=0               | M ≠ 0                        | Fractional form     |
|       | E=255             | M ≠ 0                        | NAN                 |

Because we can't represent 0 and ∞ in general so here we reserved some combinations for it, otherwise we will not be able to represent them.

→ In Implicit Normalization value is  $(-1)^s \times (1.M) \times 2^{E-127}$ .

→ Fractional form also known as denormalized form and value is

~~group~~  $(-1)^s \times (0.M) \times 2^{-126}$

Instead of  $2^{-127}$  here it is  $2^{-126}$  and it is exception no.

→ we can't used E=255, M≠0 as a combination.

### Double precision

| S | E  | M  |
|---|----|----|
| 1 | 11 | 52 |

Excess → 1023

| Sign. | E(11)         | M(52)   | Value               |
|-------|---------------|---------|---------------------|
| 0/1   | E=0           | M=0     | ± 0                 |
| 0/1   | E=2047        | M=0     | ± ∞                 |
| 0/1   | E ∈ [1, 2046] | M=xx--x | Implicit Normalized |
| 0/1   | E=0           | M ≠ 0   | Fractional form     |
|       | E=2047        | M ≠ 0   | NAN                 |

M=xx--x means, it can hold any combination.

→ In Implicit Normalization values are

$$(-1)^s \times (1.M) \times 2^{E-1023}$$

→ Fractional form

$$(-1)^s \times (0.M) \times 2^{-1022}$$

Exception.

Ex → Consider 32 bit register which stores floating point numbers in IEEE single precision format.

1) What is the value of number if 32 bits are given below.

|   |          |          |
|---|----------|----------|
| 0 | 10000011 | 1100---0 |
|   | 8        | 23       |

Sol<sup>n</sup>

$$(-1)^S (1.M) 2^{E-127}$$

$$(-1)^0 (1.11) 2^{131-127}$$

$$= (1.11) \times 2^4$$

$$= \underline{11100}_{10} \cdot (28)_{10}$$

2) What is the value if 32 bits are given as.

$$\begin{array}{c} 0 \quad 11111111 \quad 11000 \dots 01 \\ \text{S} \qquad \text{E} \qquad \text{M} \end{array}$$

$\rightarrow E = 255$  M ≠ 0 means  
it is NAN.

Q If we have 64 bit register to stores floating point no. in IEEE double precision format. What is the value of number given in 64 bits below

$$\begin{array}{c} \pm \quad \pm 1 \quad 52 \\ 1 \quad 10000000011 \quad 1101100 \dots 0 \end{array}$$

Implicit normalized form

$$(-1)^1 (1.11011) \times 2^{1027-1023}$$

$$-1 \times 1.11011 \times 2^4$$

$$-(11101.1)_2 = -\underline{(29.5)}_{10}$$

Ex: Represent  $(117)_{10}$  in IEEE 754 single and double precision formats

Sol<sup>n</sup>

$$(1110101)_2 = (117)_{10}$$

$$(1.110101 \times 2^6)$$

Mantissa going to be same in both formats

for single precision Bias = 127

$$E = 6 + 127 = 133$$

for double precision Bias = 1023

$$E = 1023 + 6 = 1029$$

Single precision,

|     |          |              |
|-----|----------|--------------|
| 0   | 10000101 | 11010100...0 |
| (2) | (8)      | (23)         |

double precision

|     |             |              |
|-----|-------------|--------------|
| 1   | 10000000101 | 11010100...0 |
| (1) | (32)        | (52)         |

Gate 2016: Let x be the no. of distinct. — —, then  $x-y$  is —

Sol<sup>n</sup>

Using 16 bits we get  $2^{16}$  comb<sup>n</sup>.  
In 2's comp. we denoted each comb<sup>n</sup> uniquely so,  $2^{16}$  combinations are possible.  
In signed mag. representation we waste two comb<sup>n</sup> for 0. So we consider  $2^{16}-1$  numbers distinctly.

$$\text{So, } x-y = 2^{16} - (2^{16}-1) = \underline{\underline{1}}$$

Binary Multiplication (partial sum method)

$$\begin{array}{r}
 & 0101 \rightarrow 5 \\
 \times & 0110 \rightarrow 6 \\
 \hline
 & 0000 \\
 & 0101 \times \\
 & 0101 \times \times \\
 0000 & 0 \times \times \times \\
 \hline
 & 0011110
 \end{array}
 \quad 5 \times 6 = 30$$

$$= \underline{\underline{(11110)}}_2 = (30)_1$$

To multiply 2, n bits numbers, in worst case it will take  $2^n$  bits to store answer.

Ex → To multiply 2, 4 bits no. we need register of size 8 bits

$$\begin{array}{r}
 0101 \\
 \times 0110 \\
 \hline
 0000 \\
 0101 \\
 0100 \\
 0000 \\
 \hline
 00011110
 \end{array}$$

$$\begin{array}{r}
 01011010101010 \\
 \text{shift} \gg 4 \\
 010110110101010 \\
 0101 \\
 01011111101010 \\
 \text{shift} \gg 1 \\
 010111111101010 \\
 + 0000 \\
 0101111111101010 \\
 \text{shift} \gg 1 \\
 0101011111101010
 \end{array}$$

→ After adding into Most significant bits shift the value by 1 digit.

No of Adds require = No. of 1's in multipliers

No of shift require = No. of bits in multiplier.

### Idea Behind Booth's Algo.

$$\begin{array}{r}
 011110 = (2^{4+1} - 2^1) \\
 2^5 \underline{2^4 2^3 2^2 2^1 2^0} \\
 \hline
 \end{array}$$

$$0110 = (2^{2+1} - 2^1) \\
 \underline{2^3 2^2 2^1 2^0}$$

General

$$0111 \dots 1 \dots 0 = (2^{j+1} - 2^i)$$

$\overset{j}{\underset{i}{\dots}}$

$2^j$  represent MSB and  $2^i$  represent LSB in sequence of 1's in binary no.

$$\begin{array}{r}
 M \times Q \\
 \downarrow \quad + \\
 \text{Multiplicand} \quad \text{Multiplier}
 \end{array}$$

if  $Q = 011110$  means it will take 6 shifts and 4 addition will be performed.

$$\begin{aligned}
 & \Rightarrow M \times (2^5 - 2^1) \\
 & = (2^5)M + 2^1(-M)
 \end{aligned}$$

Multiply any no with  $2^n$  means  $n$  times left shifts.

so in above computation total shift =  $2^5 + 2^1$   
 $= 6$  shift

But instead of 4 addition we will perform only 1 additions.

### Example

$$\begin{array}{r}
 01011 \rightarrow M \\
 \times 01110 : (2^4 - 2^1) \rightarrow Q
 \end{array}$$

$$(2^4 - 2^1)M$$

$$= 2^4 M + 2^1(-M) \quad \text{--- (1)}$$

$$\text{Here } f(M) = 10101$$

$$(2^1)xni = \underline{101010} \quad (\text{Left shift by } \frac{1}{2})$$

$$(2^4)xM = \underline{010110000}$$

$$\begin{array}{r}
 \boxed{111}101010 \\
 \times 010110000 \\
 \hline
 \textcircled{1}010011010
 \end{array}$$

we fill it with 1's to maintain the significance of (-ve) sign no.  
 $\hookrightarrow (154)_{10}$

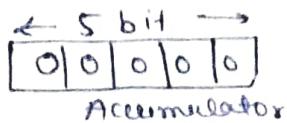
discard it bcoz we are using 2's comp. arithmetic

If don't understand here watch other [YT video](#)

we are explaining the concept by same example discussed above.

$$M = 01011$$

$$-M = 10101$$



| A     | Q     | $Q_{-1}$ |             |
|-------|-------|----------|-------------|
| 00000 | 01110 | 0        |             |
| 00000 | 00111 | 0        | Right shift |

|       |       |   |            |
|-------|-------|---|------------|
| 10101 | 00111 | 0 | A - M      |
| 11010 | 10011 | 1 | Rightshift |
| 11101 | 01001 | 1 | >> shift   |
| 11110 | 10100 | 1 | >> shift   |

↳ obtained pattern

Now add M to to obtained pattern

$$\begin{array}{r}
 1111 \ 01010 \\
 0101 \ 10000 \\
 \hline
 010011010
 \end{array}
 \xrightarrow{\text{A} + \text{M}}$$

discard it       $\hookrightarrow (154)_{10}$

~~Note:~~ Do atleast 5 iteration on Booth's Algorithm.

### Operations perform

| LSB(Q) | $Q_{-1}$ |                  |
|--------|----------|------------------|
| 0      | 0        | >> shift         |
| 1      | 0        | A - M / >> shift |
| 1      | 1        | >> shift         |
| 0      | 1        | A + M / >> shift |

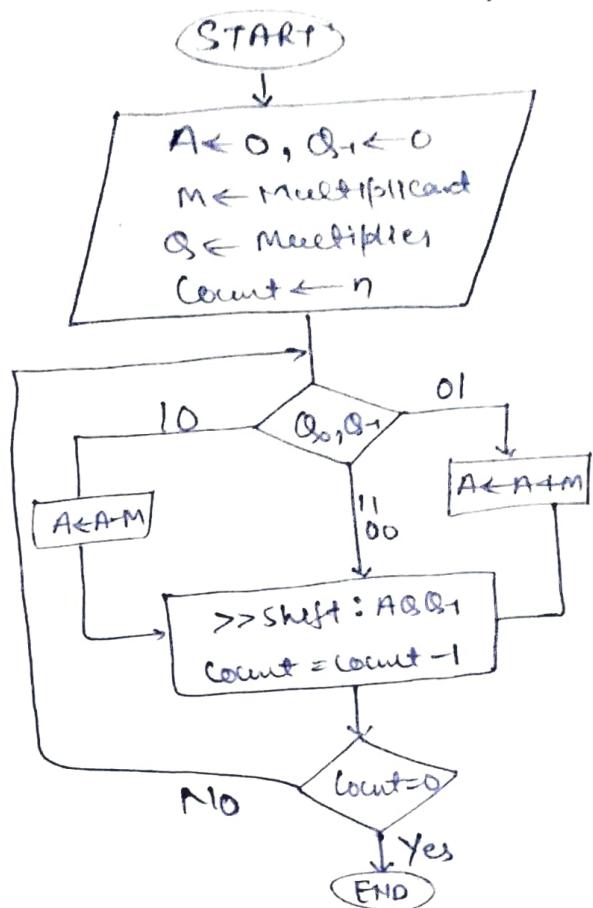
→ No. of steps = no. of bits in M

### Booth Algorithm flow

#### Diagram

See on Rus-top

Cont -



Advantages and disadvantages of Booth's

→ In general the number of operations required will be reduced.

→ Ex: 01110 → 3 ADD, 3 SOB and 6 shifts

for each sequence of consecutive 1's we perform 1 ADD, 1 SOB

Ex:

01110011110110 → 3 ADD, 3 SOB

→ It can do signed multiplications

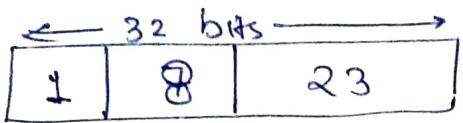
→ If negative no. are involved, it give better performance

Disadv.: If block contains only single 1.

Ex: 010101010101

## Example on fractional form

What is the range of floating point number represented using fractional form of IEEE 754 single precision?



for fractional form  $E=0, M \neq 0$ .

$$(-1)^S (0.M) \times 2^{-126}$$

for (+ve)  $S=0$ , (-ve)  $S=1$ , so other than M everything fix.

$$M = \underline{0 \ 0 \ 0 \ \dots \ 1} \quad (\text{To represent smallest no. } \rightarrow 0)$$

$$0.M = 0.000\dots 1$$

$$\begin{aligned} (-1)^0 (0.M) 2^{-126} &= 2^{-23} \times 2^{-126} \\ &= \underline{2^{-149}} \\ &\text{Smallest (+ve) fraction.} \end{aligned}$$

Largest possible fraction

$$\begin{aligned} (-1)^1 \times (0.111\dots 1) \times 2^{-126} &= (2^{-1} + 2^{-2} + 2^{-3} + \dots + 2^{-23}) 2^{-126} \\ &= \frac{1}{2} \left( 1 - \left(\frac{1}{2}\right)^{23} \right) \\ &= \underline{(1 - 2^{-23}) \times 2^{-126}} \end{aligned}$$

Done

## Qmb. Points regarding the Number System