

## Number Representation:

If  $N$  is  $(n+1)$  bit unsigned no.

std. form

$$N = \sum_{i=0}^{n-1} x_i 2^i$$

$x_i \Rightarrow$   $i^{\text{th}}$  binary bit of  $N$ .

$x_0 \Rightarrow$  LSB

$x_n \Rightarrow$  MSB

1. Decimal

2. Octal

3. Hexadecimal

4. Binary

5. BCD

6. 12 Decimal

7. Logarithmic

8. Residual Number System

9. CORDIC (Coordinate Rotation, Digital Computer)

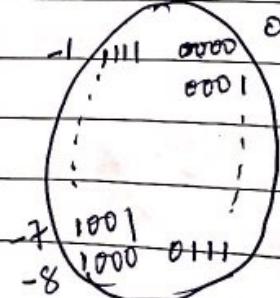
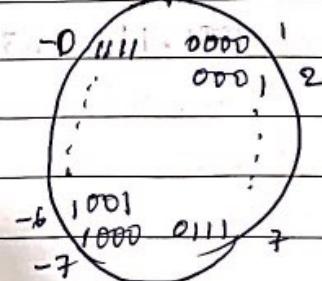
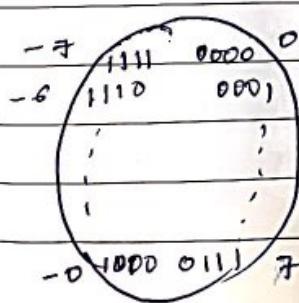
Binary

Signed, 1's complement

2's complement (adopted in hardware)



- 0 +ve  
- 1 -ve



Signed and 1st complement have 2 zeroes.  
Thus, 2's complement is used for hardware.

eg 4 bits no.	eg 8 bit
Signed $-7 \text{ to } +7$	Signed $-(2^7 - 1) \text{ to } 2^7 - 1$
Unsigned $0 \text{ to } 15$	Unsigned $0 \text{ to } 2^8 - 1$

Binary system:

Shifting is the key.

$2^2 | 2^1 | 2^0 | . | 2^{-1} | 2^{-2} | 2^{-3}$

Binary point.

To define fixed point we need 2 parameters:

- width of the number representation
- Binary point position within the number

$$00010.110 = 2.75$$

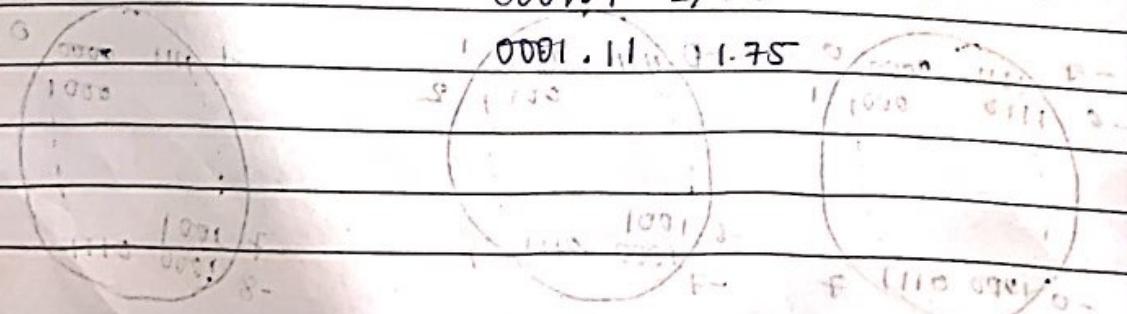
Negative no : 1's complement

Benefit - straight forward and efficient

Disadvantage - loss of range  $[-8, 7]$  instead  $[0, 15]$

Loss of precision, if we fix the  
bits after the binary point

$$0001.1 \Rightarrow 1.5$$



To avoid the compromise of range and precision with the fixed point (for scientific calculation purpose) floating point number system is preferred.

$$S \pi = 3.14159265$$

$$e = 2.711828$$

$$\pi = 31.4159265 \times 10^{-1}$$

$$e = 271.11828 \times 10^{-2}$$

normalized numbers.

standard numbers

32 bit fixed point number, what is the full range?

$$(2^{32}-1) \sim 10^9 \times 4$$

Precision?

$$2^{(32-1)} = 2^{(1-1)} = 2^1 = 2 \text{ bits}$$

16 bits

$$\text{Range: } 2^{16}-1$$

$$\text{Precision: } 2^{-16}$$

### Floating Point Number

32 bit representation of a float

Bit No.	Size	Field Name
31	1	Sign(s)
23-30	8	Exponent(E)
0-22	23	Mantissa(M).

$$f. No = -1^S \times (1 + \text{mantiss}) \times 2^E$$

Binary point

Overflow: Exponent is too large to be represented in Exp. field.

Underflow: no. is too small to be represented in Exp. field.

The standard expresses the exponent as a biased notation. (This is to accommodate larger binary nos. and -ve no. representable)

$$1 \times 2^{-1} = [0] \boxed{11111111} \boxed{000\underset{8}{\dots}000} = 1.1111111 \times 2^{-1}$$

$$1 \times 2^{+1} = [0] \boxed{00000001} \boxed{000\underset{8}{\dots}000} = 1.0000001 \times 2^1$$

2's  
comple  
ment

with the biased notation the FB number is equal to expressed as

Effectively  $\rightarrow$  FB FP =  $(-1)^S \times (1+M) \times 2^{E-\text{bias}}$

it is 33 bit no. where bias = 127.

-1 exponent will always be represented as

$$-1 + 127 = 126$$

$$+1 + 127 = 128$$

$$0 + 127 = 127$$

This is a normalized number with 1 intrinsic to the number

(2.3) example

1

13

(3.3) floating point

8

00-00

(4.1) decimal

8.2

48-0

$$3.2 \times 10^{13} \times 10^{(13-13)} = 3.2 \times 10^{13}$$

representing it

in memory and in memory is represented in binary

using left int

in memory and in memory is represented in binary

using left int

# IEEE 754 - 1985 standard

a representation

## Single Precision (32 bit)

Double Precision (64 bit).

Extended bit (80 bit)

S | Exponent | Fraction or Mantissa

Single Precision : 1 8 23

Double Precision: | 1 | .1 | 52 |

$$x = (-1)^s \times (1 + \text{fraction}) \times 2^{\frac{e}{\text{Exponent-bias}}}$$

## Floating point number

$$S = 0 \quad +ve$$

$$S = 1 - ve^{-2\pi x^2} \text{Im}(\zeta A)$$

Normalize significant : eq  $1.23 \times 10^{-21}$

$$1.0 \times 10^0 \rightarrow 5 \times 10^0$$

$1.0 \leq |z_{\text{significant}}| \leq 2.0$

~~0.111... = 1/9~~

$$\frac{1}{2} \cdot \frac{1}{2} = 0.25$$

$$0.5 + 0.25 + 0.125 + \dots = 1.875$$

0.999... =

⇒  $\text{X} \otimes \text{S} \oplus \text{S} \otimes \text{X}$  为  $\text{B}(H)$  的子代数

Significand is fraction with "1." restored.  
always has a leading 1 bit. So, no need  
to represent it explicitly.

Exponent : excess representation

actual exponent + bias

This ensures that the exponent  
is unsigned.

single precision = 127

double precision = +1023 1023

Single Precision range:

Exponents: 00000000 11111111

$\downarrow$   $\leftarrow$   $x \times (\text{sign}) \times 2^{\text{exp}} = x$

smallest value : 00000001

of exponent

Actual expo =  $127 - 127$   
 $= 1 - \text{bias}$

$= 1 - 127$

$= -126$

Fraction =  $0.1 \times 0.0000 \dots \times 2^{-126}$

Significand = 1.0

smallest value :  $\pm 1.0 \times 2^{-126} \approx 1.2 \times 10^{-38}$

Largest value : expo = 11111110 = 254

actual expo =  $254 - 127$   
 $= 127$

Fraction = 11...1111

significand  $\approx 2.0 \pm 2.0 \times 2^{127} \approx 3.4 \times 10^{38}$

Double Precision: 64 bits | 128 bits | 112 bits

smallest value: actual expo = 1022

$$\text{Fraction} = 0. \overline{0001001}$$

Significand = 1.0

$$\pm 1.0 \times 10^{-1022} \approx \pm 2.2 \times 10^{-808}$$

Largest value: actual expo =  $10^{23}$  111...110

$$\text{fraction} = \frac{111}{111+222+333}$$

significand  $\cong 2.0$

$$\pm 2.0 \times 10^{23} \simeq \pm 1.8 \times 10^{308}$$

Precision is all fractional bits are significant.

Single Precision approx:  $2^{-23}$

$$\text{Equivalent : } 23 \log_{10} 2 = 23 \times 0.3$$

$\approx$  6 decimal digits of precision

~~1.03~~ Double Precision approx.  $10^{-52}$

$$\text{Equivalent} \therefore 52 \log_{\frac{1}{2}} 2 = 52 \times 0.3$$

$\approx$  16 decimal digits of precision

eq - 0.75

8

$$S = 1.$$

$$0.77 = 0.75 \quad (\cos 3 + i \sin 3) = \text{Fraction} = 000\ldots0001$$

$$0.11 \rightarrow 1.1 \times 2^{-1}$$

$$\text{exponent} = -1 + 127$$

$$\frac{1}{2} \sin(\alpha + \beta) = 0.76856 \approx 76.856^\circ$$

$$-0.75 = (-1)^1 \times (1.1) \times 2^{-1}$$

$$= 0.75 = (-1)^1 \times (1:1) \times 2^{-1} \quad 0111$$

Cost = 3.00

$$-0.75 = (-1)^1 \times (1.1) \times 2^{-1}$$

وَهُنَّ مُنْذَرٌ

అంబుల్ కు

eg

$$1 \quad 1000\ 0001 \quad | \quad 01000\ 0000$$

↓129

 $s = 1$  $\text{fraction} = 0.1000\ldots$  $\text{significand} = 1.0100\ldots$  $\text{exponent} = 129 - 127 = 2 \pm$ 

$$\begin{aligned} & (-1)^1 \times 1.25 \times 2^2 \\ & -1.25 \times 4 \\ & -5.0 \end{aligned}$$

↓  
Excess 7 bias

Adder:

Arithmetic operations relating to input  
and output is

$$x + y + \text{cin} = 2^n \text{cout} + z$$

For n bit adder:

Inputs:

$$x = (x_{n-1}, x_{n-2}, \dots, x_0) \quad x_i \in \{0, 1\}$$

$$y = (y_{n-1}, y_{n-2}, \dots, y_0) \quad y_i \in \{0, 1\}$$

$$\text{cin} = \{0, 1\}$$

Outputs:

$$z = (z_{n-1}, \dots, z_0) \quad z_j \in \{0, 1\}$$

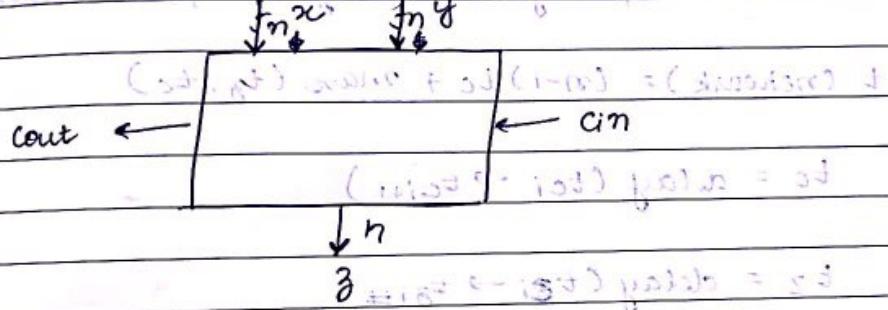
Functions:  $z = (x+y+\text{cin}) \bmod 2^n$ 

$$\text{cout} = \begin{cases} 1 & \text{if } x+y+\text{cin} > 2^n \\ 0 & \text{otherwise} \end{cases}$$

if  $n=5$ 

Iteration	$x_i$	$y_i$	$c_{in}$	$z_i$ ( $(x_i + y_i + c_{in}) \text{ mod } 2$ )	$c_{out}$
1	11	12	1	$(1+12+1) \text{ mod } 2 = 24$	0
2	15				

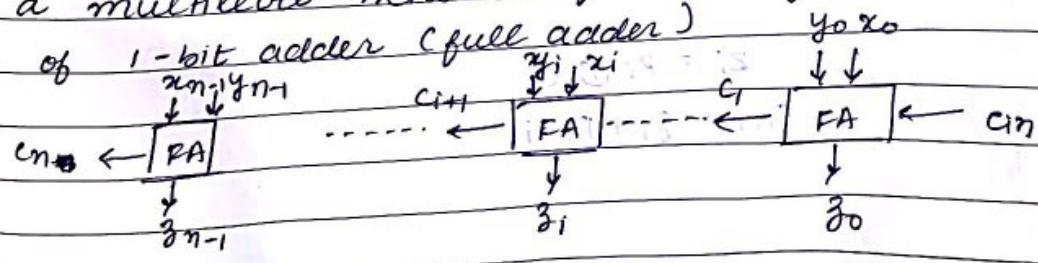
Implementation of ripple carry adder:



Ripple carry Adder:

i/p :  $x_i, y_i, c_i \in \{0, 1\}_2$ .o/p :  $z_i, c_{out} \in \{0, 1\}_2$ Function :  $z_i = (x_i + y_i + c_i) \text{ mod } 2$ 

$$c_{i+1} = \begin{cases} 1 & \text{if } x_i + y_i + c_i \geq 2 \\ 0 & \text{otherwise} \end{cases}$$

For  $n$  bit ripple carry adder module, if it is a multilevel network formed by correct connection of 1-bit adder (full adder) $n$  bit adder (RCA) $x + y$

Critical Path :

Begins in  $c_m$ , traverses all the module, terminates in  $c_{out}$  or  $z_{n-1}$ .

Delay :

Worst case propagation delay is approximately:

$$t_{\text{network}} = (n-1)t_c + \max(t_g, t_c)$$

$$t_c = \text{delay } (t_{ci} \rightarrow t_{ci+1})$$

$$t_g = \text{delay } (t_{gi} \rightarrow t_{gi+1})$$

$$z_i = x_i \oplus y_i \oplus c_i$$

$$x_i y_i' c_i' + x_i' y_i c_i + x_i' y_i' c_i +$$

is 1 when  $x_i + y_i = 2$ , i.e.  $(x_i y_i = 1)$ ,

or when  $x_i + y_i = 1$  and  $c_i = 1$ .

$$(x_i \oplus y_i) \cdot c_i = 1$$

$$c_{i+1} = x_i y_i + \bar{x}_i (x_i \oplus y_i), c_i = 0$$

Variable

Propagate:  $p_i = x_i y_i$  (and we can not)

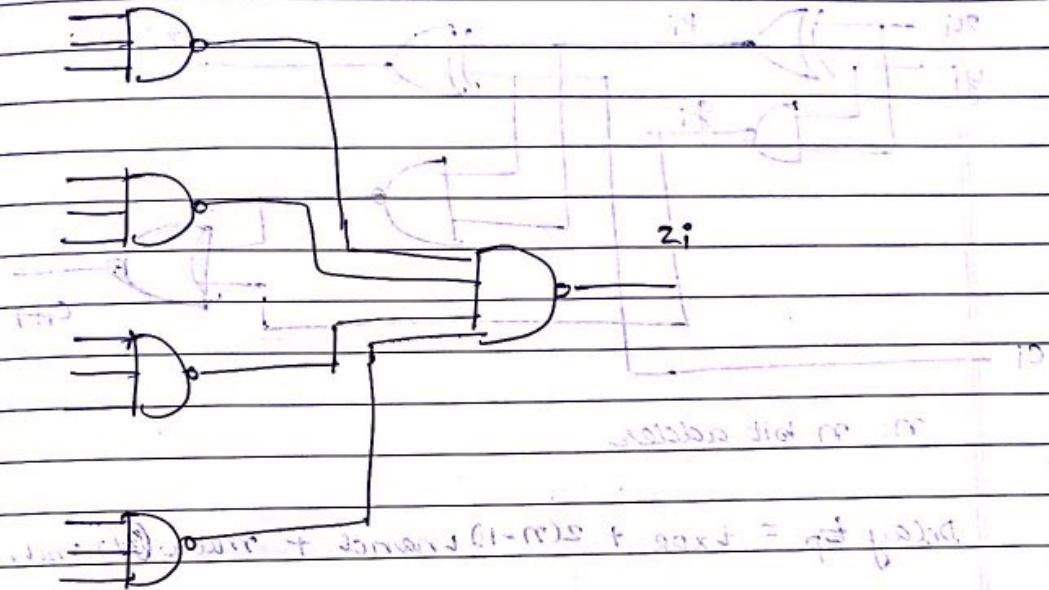
Generate:  $g_i = x_i y_i$  (and we can)

$$z_i = p_i \oplus c_i$$

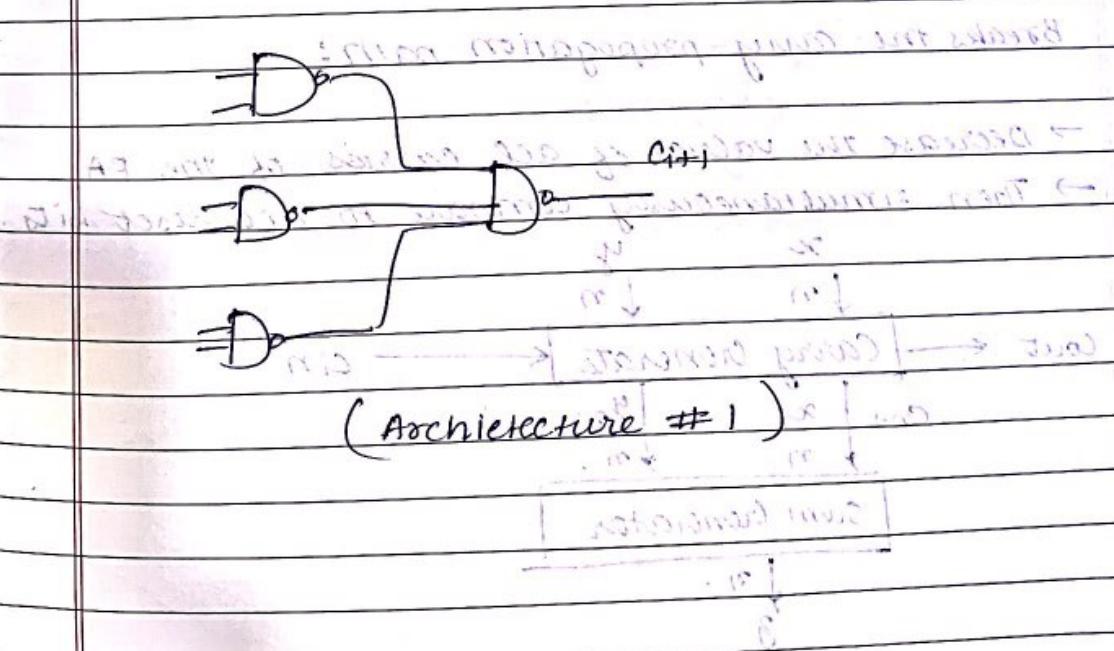
$$c_{i+1} = g_i + p_i c_i$$

$$g_i = x_i y_i$$

$$p_i = x_i y_i$$



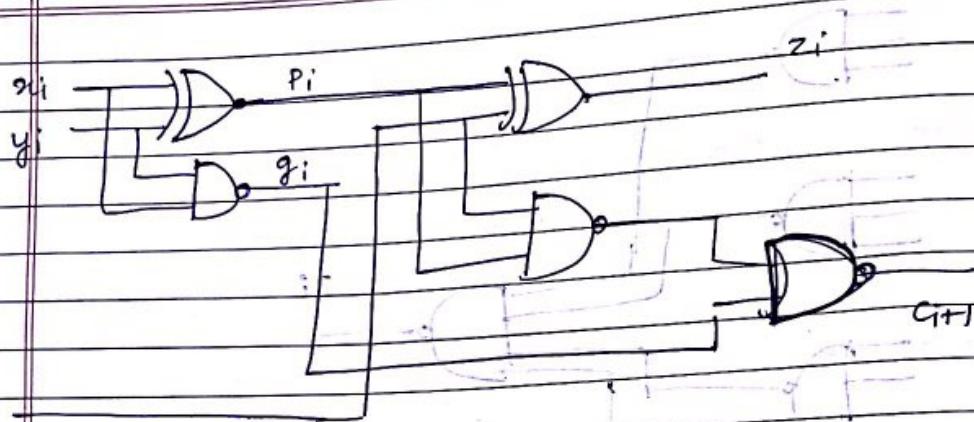
$$C_{i+1} = x_i y_i + x_i C_i + \text{previous state result}$$



1011 + 10 = 1101

1101 - 10 = 1011

1011 + 10 = 1101



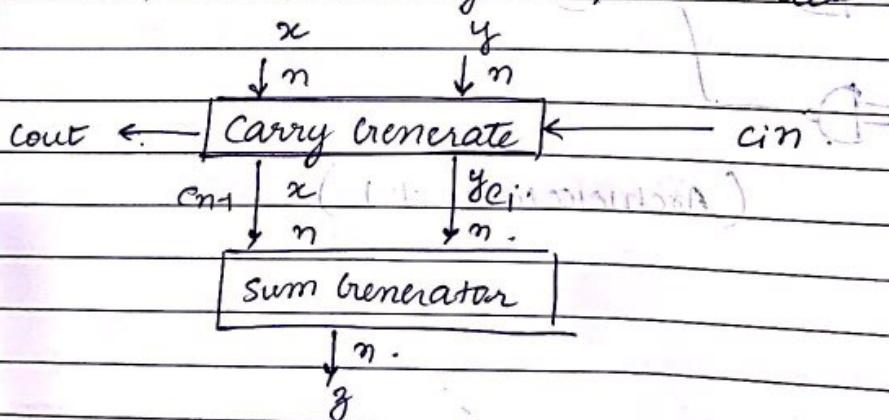
n: n bit adder

$$\text{Delay } t_p = t_{\text{XOR}} + 2(n-1)t_{\text{nand}} + \max(t_{\text{nand}}, t_{\text{XOR}})$$

Carry Look Ahead Adder:  $t_{\text{PLA}} = t_{\text{XOR}}$

Breaks the carry propagation path:

- Decrease the values of all carries of the PA.
- Then simultaneously compare to all reset bits.



$$c_{i+1} = g_i + P_i c_i$$

$$g_i = x_i y_i$$

$$P_i = x_i \oplus y_i$$

A 10 A

$$e_1 = g_0 + p_0 c_0$$

$$c_2 = g_1 + p_1 c_1$$

$$= g_1 + g_0 p_1 + p_1 p_0 c_0.$$

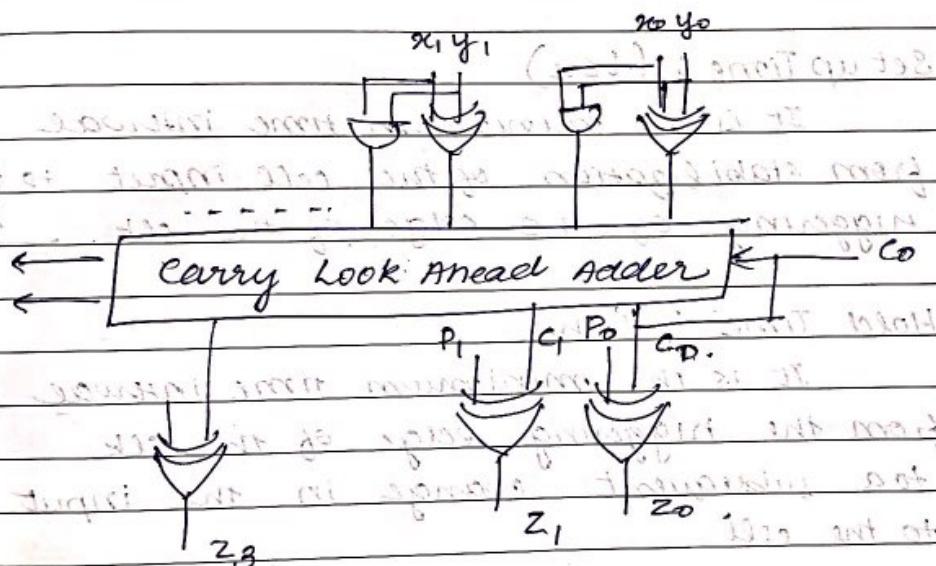
$$c_3 = g_2 + p_2 c_2$$

$$= g_2 + g_1 p_2 + g_0 p_1 p_2 + p_0 p_1 p_2 c_0$$

$$(g_0 p_1 p_2 + g_1 p_2 + g_0 p_1) + g_2 + p_0 c_0 = z_0$$

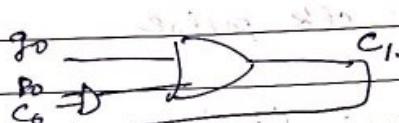
$$+ \text{and } z_i = x_i + y_i + c_i \text{ for } i = 0, 1, 2$$

$$\text{so } z_i = p_i + c_i$$

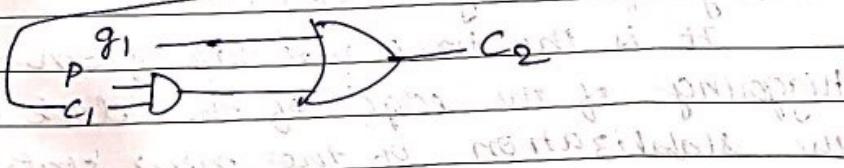


(cont.) : previous slides

$$c_1 = g_0 + p_0 c_0$$



$$c_2 = g_1 + p_1 c_1$$



4 CLA

CLA-4

$$C_4 = G_0 + P_0 C_0$$

$$C_8 = G_1 + P_1 G_0 + P_1 P_0 C_0$$

$$C_{12} = G_2 + G_1 P_2 + G_0 P_1 P_2 + P_0 P_1 P_2 C_0$$

$$C_{16} = G_3 + P_3 G_2 + P_3 P_2 G_1 + P_3 P_2 P_1 G_0 + P_3 P_2 P_1 P_0 C_0$$

Set up Time : (t<sub>su</sub>)

It is the minimum time interval from stabilization of the cell input to the triggering of the edge of the clk.

Hold Time : (t<sub>h</sub>)

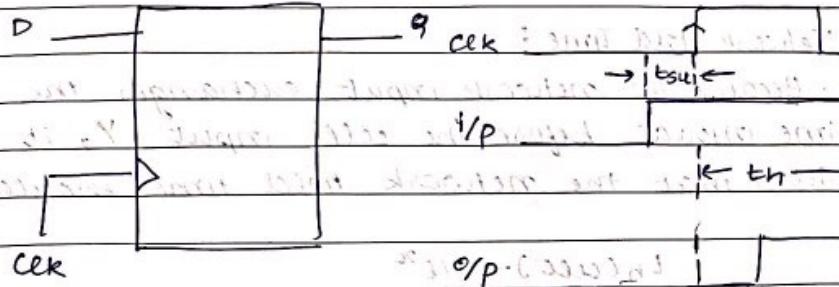
It is the minimum time interval from the triggering edge of the clk to a subsequent change in the input to the cell.

Pulse width: (t<sub>w</sub>)

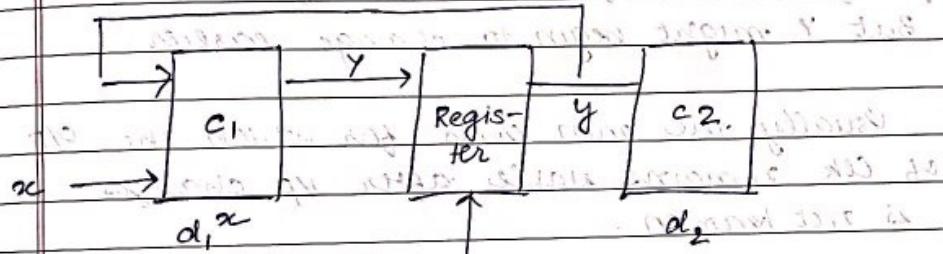
It is the minimum width of the synchronizing clk pulse.

Propagation Delay: (t<sub>p</sub>)

It is the time interval from triggering of the edge of the clock to the stabilization of the new state output.



It is used to store data until the clock edge  
as well as to receive data.



$d_1, x$  = delay of the combinational network C1.  
(Input + Output)

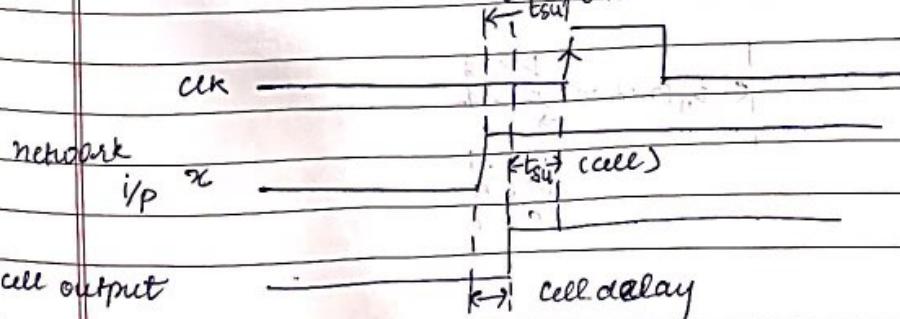
$d_1, y$  = worst input  $x$  and to the input from  
the state register  $y$

1. Network setup time : It is the time i/o to  
the cell to be stable at  $t^+$

$$t_{su(\text{net})} = d_1^x + t_{su(\text{cell})}$$

$t_{su}^*$  (network)

$t_{su}^*$  (cell)



Network Hold Time:

Because the network input exchanges the time interval before the cell input  $Y$ , it would seem that the network hold time should be

$$t_h(\text{cell}) - d_{\text{ex}}^x$$

However  $d_{\text{ex}}^x$  is the max delay between the stabilization of  $x$  and that of  $Y$ .

But  $Y$  might begin to change earlier.

Usually the min terms for which the O/P of CLK remains stable after IP changes is not known.

The network hold time is conservatively equal to the hold time of the cell.

$$t_h(\text{network}) = t_h(\text{cell})$$

With respect to  $X$  at  $t=0$  the output goes to  $Y_{10}$

At  $t=t_h(\text{cell})$  the output is

CLK

all inputs are at  $0$  & output goes to  $Y_{10}$

At  $t=t_h(\text{network})$  the output is

IP  $x$

$\leftarrow t_h(\text{cell})$

$$+ 30 = \text{final}$$

Consequently

we get  $Y_{10}$

O/P

$\leftarrow t_h(\text{network})$

$$+ 30 = \text{final}$$

Consequently

we get  $Y_1$

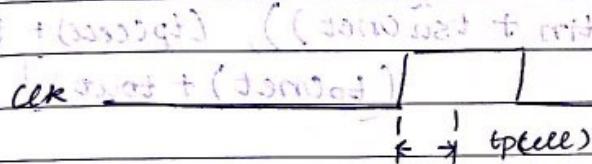
Final O/P

## Network Propagation Delay:

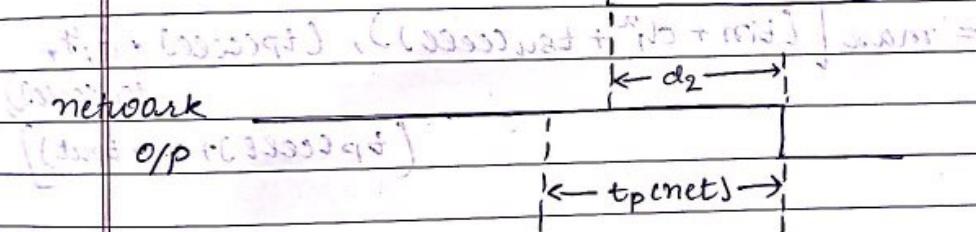
$C_2$  has a delay  $= d_2$

The netw. stabilizes  $d_2$  time units after a change in the o/p of the cell.

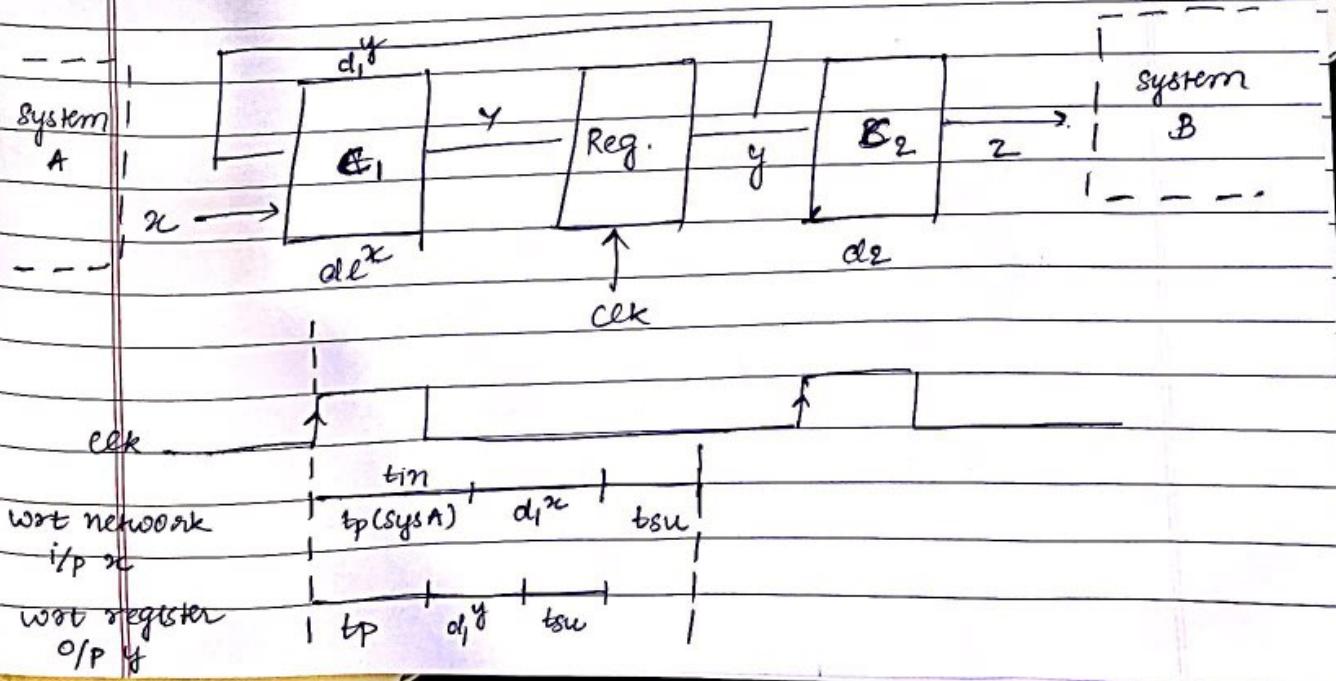
$$tp(\text{netw}) = tp(\text{cell}) + d_2$$

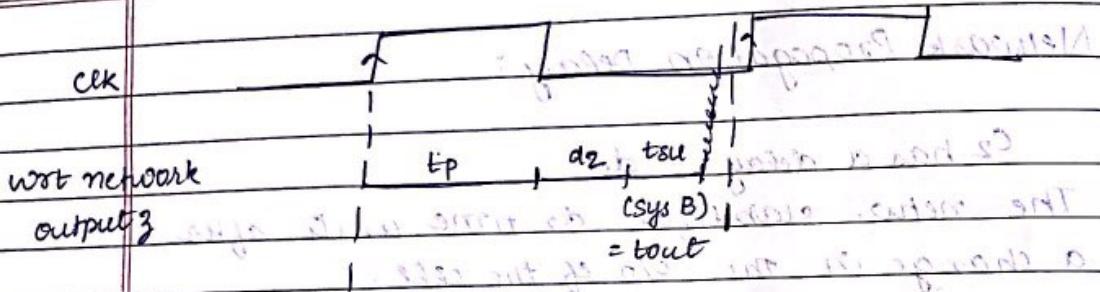


cell o/p :  $(\text{initial} \rightarrow \text{final})$  [same = direct]



## Maximum Clock Frequency:





Worst case path: Maximum clk freq (1<sup>st</sup>)

$$T_{\min} = \max \left[ (t_{\text{in}} + t_{\text{su}}^x), (t_{\text{pcell}} + t_{\text{su}}^y), (t_{\text{pcell}} + t_{\text{out}}) \right]$$

Assumption:  $t_{\text{hcell}} \leq t_{\text{pcell}}$

$$= \max \left[ (t_{\text{in}} + d_1^x + t_{\text{su}}^x), (t_{\text{pcell}} + d_1^y + t_{\text{su}}^y), (t_{\text{pcell}} + d_2 + t_{\text{out}}) \right]$$