

Module 2

MOS and BiCMOS Circuit Design Processes

- Methods of realizing circuit design in silicon
- The design process can be understood by means of stick diagrams and symbolic diagrams along with set of design rules.
- Design rules: is a communication link between designers specifying the requirements and the fabricator.

MOS Layers:

- MOS circuits are basically formed by 4 layers
 - Metal
 - Polysilicon
 - N diffusion
 - P diffusion
- Here all the 4 layers are isolated from each other through thick or thin oxide layer (i.e., silicon dioxide layer)
- The thin oxide (thinox) layer includes n-diffusion, p-diffusion and transistor channel.

Stick diagram:

- Stick diagrams are a means of capturing topography and layer information using simple diagrams.
- They convey layer information through color codes (or monochrome encoding).
- Acts as an interface between symbolic circuit and the actual layout.
- Stick diagrams do show all components/vias(contacts), relative placement of components and helps in planning and routing. It goes one step closer to layout.
- However they do not show exact placement of components, transistor sizes, length and width of wires also the boundaries. Thus we can say that it does not give any low level details.
- The color encodings chosen for different technologies is shown below.
- **Encodings for NMOS process:**

COLOR	STICK ENCODING	LAYERS	MASK LAYOUT ENCODING	CIF LAYER
GREEN		n-diffusion (n+ active) Thinox*		ND
RED		Polysilicon		NP
BLUE		Metal 1		NM
BLACK		Contact cut		NC
GRAY	NOT APPLICABLE	Overglass		NG
nMOS ONLY YELLOW		Implant		NI
nMOS ONLY BROWN		Buried contact		NB
FEATURE	FEATURE (STICK)	FEATURE (SYMBOL)	FEATURE (MASK)	
n-type enhancement mode transistor				
Transistor length to width ratio L:W should be shown.				
n-type depletion mode transistor nMOS only				
Source, drain and gate labelling will not normally be shown.				

Procedure to draw Stick Diagram:

Nmos Design Process.

- 1) Draw two metal lines/ power rails providing sufficient space to accommodate all transistors. i.e: Vdd & Vss.
- 2) Draw common n+ diffusion layer for all the transistors.
- 3) Provide Vdd and Vss contacts.
- 4) Draw polysilicon to cross n+ diffusion layer to form transistors.
- 5) Create buried contact for depletion transistor.
- 6) Provide input and output connection.

CMOS Design Process:

- Two type of transistors are used i.e: Nmos and Pmos, thus in stick diagram demarcation line is used to separate them.
- All Pmos transistors are placed above Demarcation line and Nmos transistors below demarcation line.
- While drawing stick Diagram
 1. Diffusion paths must not cross the demarcation line
 2. N-diffusion and P-diffusion wires must not join.
 3. Nmos and Pmos transistors are joined by Metal layer when it is required.
 4. Cross must be placed on Vdd and Vss which represent substrate and P-well connection respectively.

Encodings for CMOS process:

COLOR	STICK ENCODING	LAYERS	MASK LAYOUT ENCODING	GIF LAYER
GREEN	Encoding as in Color plate 1 (a)	n-diffusion (n ⁺ active) Thin _{ox}	Encoding as in Color plate 1 (a)	CAA or GNA
RED		Polysilicon		CPF
BLUE		Metal 1		CMF
BLACK		Contact cut		CC
GRAY		Overglass		COG
YELLOW (STICK)	green outline here for clarity	p-diffusion (p ⁺ active)	p ⁺ mask	CAA or CPA
YELLOW	Not shown on diagram	p ⁺ mask	either or	CPP
DARK BLUE OR PURPLE		Metal 2		CMS
BLACK		VIA		CVA
BROWN	Demarcation line p-well edge is shown as a demarcation line in stick diagrams	p-well		CPW
BLACK	X	V _{DD} or V _{SS} contact		CC
FEATURE	FEATURE (STICK)	FEATURE (SYMBOL)	FEATURE (MASK)	
n-type enhancement mode transistor (as in Color plate 1(s)) Transistor length to width ratio L:W may be shown.	Demarcation line L:W			
p-type enhancement mode transistor	L:W S G D Demarcation line	S G D	S G D G p ⁺ mask	
Note: p-type transistors are placed above and n-type below the demarcation line.				

Procedure to draw Stick Diagram:

- 1) Draw two metal lines/ power rails providing sufficient space to accommodate all transistors. i.e: Vdd & Vss.
- 2) Draw demarcation line in the middle of the two power lines.
- 3) Draw P+ diffusion above demarcation and N+ diffusion below demarcation
- 4) Draw polysilicon to represent Pmos and Nmos which represents gates of the transistor.
- 5) Connect source terminal of transistors to supply.
- 6) Drain terminals of transistor are connected using metal 1.
- 7) Place contact cuts wherever necessary.
- 8) Draw X which represents substrate and P-well contact on power lines.

Layout: describes actual layers and geometry on silicon substrate to implement a function(Expressions).

[Diffusion region where transistor can be formed is called active region, polysilicon serves as the gate of MOS transistor. L defines channel length and W represents width of channel/active region]

Design rules: are set of guidelines which specify minimum dimension and spacing allowed in layout drawing. Design rules also acts a communication link between circuit designers and process engineers during manufacturing phase.

Goal of design rule: is to achieve optimum yield. Yield = (No. of good chips on wafer)/(Total no. of chips on wafer).

Design rules are also called layout rules. If the circuit performance has to be increased then rules must be more aggressive. Else this leads to non-function of the circuit or yield reduction. There are two rules.

1. Micron Rule - Absolute Dimension rule, here all sizes and spacing are specified in micron. Here the circuit density is the important goal.
2. Lambda (λ) Based Rules - The Lambda based design rules are Proposed by Mead and Conway. Scalable design rules, here this design rule normalizes all geometric design rule by parameter lambda (λ) also called as scaling factor/feature size. In this all mask patterns are expressed as multiples of lambda.

Advantages of lambda based design rules:

1. The mask layout can be scaled down proportionally if the feature size of the fabrication process is reduced.
2. Design rules are conservative.
3. This rule enable technology changes and design reuse and reduced design cost.

Disadvantages:

1. Linear scaling cannot be extended and is limited over range of dimension (1-3 μm)
2. As rules are conservative, results in over dimension and density of design is less.

The Design rules can be conveniently set out in diagrammatic form as shown in fig. 1 for width and separation of conducting path. In fig. 2 shows the design rules associated with extensions and separations with transistor. Fig. 3 and 4 demonstrates the design rules for

contacts between layers. Table below also gives the layer and distance of separation dimensions.

Layer	Dimension
n-diffusion	2λ
p-diffusion	2λ
Polysilicon	2λ
Metal 1	3λ
Metal 2	4λ

Layer -Layer	Dimension
n-diffusion – n-diffusion	3λ
p-diffusion – p-diffusion	3λ
n/p diffusion - polysilicon	1λ
Poly-poly	2λ
Metal 1	3λ
Metal 2	4λ

Layer dimension

Distance of Separation

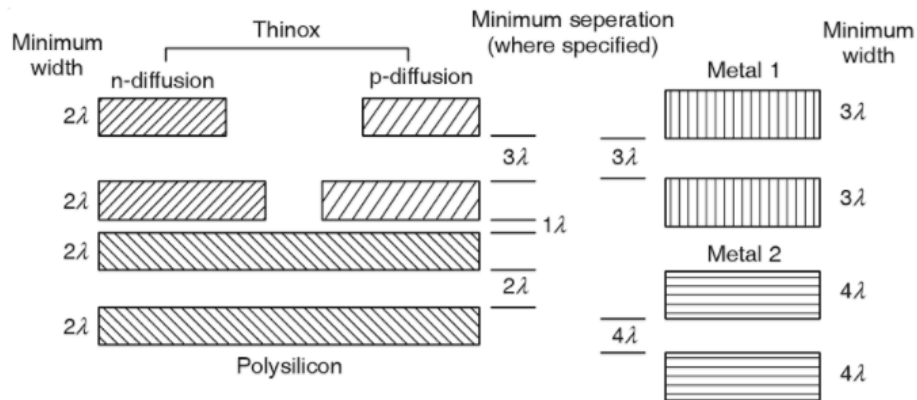


Fig. 1 Design rules for wires and separations (nMOS and CMOS)

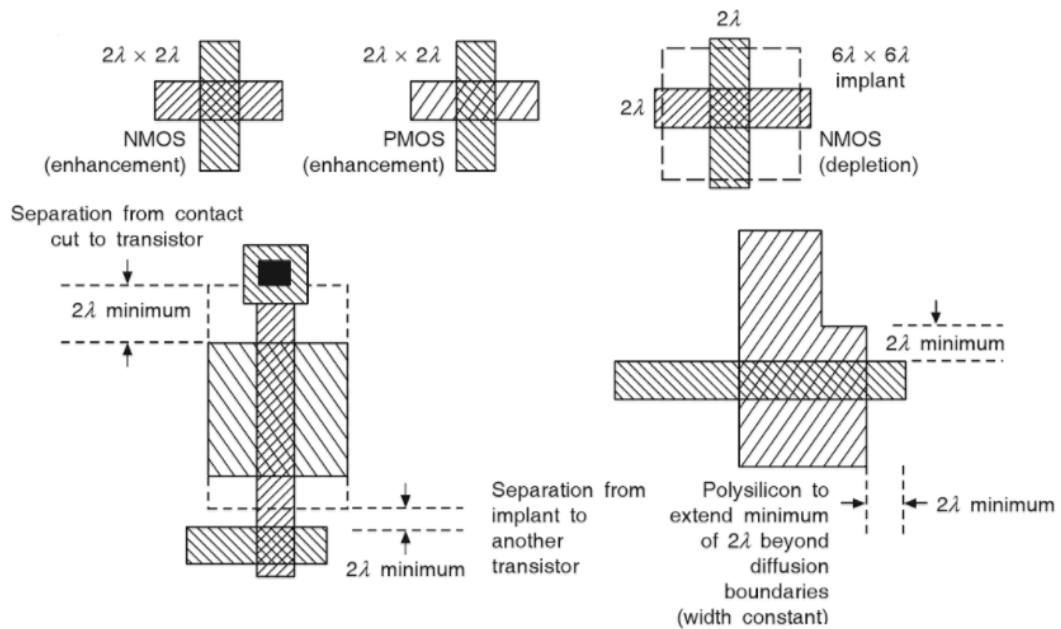


Fig. 2 Design rules for Transistors (nMOS, pMOS and CMOS)

Transistor design rules

- Minimum dimension of transistor is $2\lambda \times 2\lambda$ – overlapping of diffusion and poly
- Poly and diffusion both must extend beyond the boundary of transistor at least by 2λ

- Implant for depletion mode transistor is $6\lambda \times 6\lambda$ i.e., implant must extend boundary of transistor by at least 2λ in all direction.
- From the boundary/ implant of one transistor, the next transistor should maintain min distance of 2λ
- The distance from contact cut to transistor should be at least 2λ

Metal contact – contact between metal 1 to polysilicon OR metal 1 to diffusion (active region) is called metal contact. This is shown in fig. 3

- A $2\lambda \times 2\lambda$ cut centered on $4\lambda \times 4\lambda$ superimposed area is used to connect layers
- In case of multiple contacts the distance between adjacent contacts should be 2λ

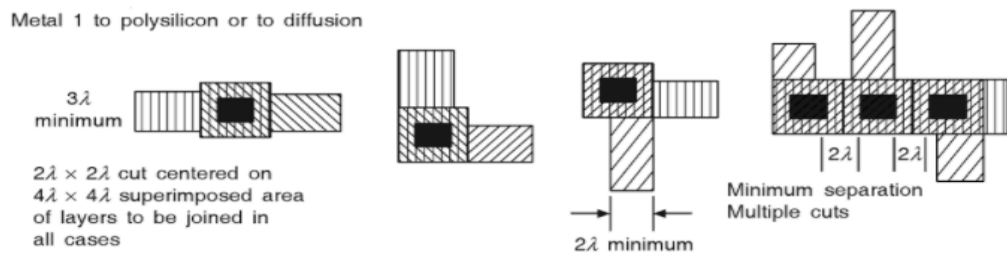


Fig. 3. Contacts (nMOS and CMOS)

Via contact – the contact between metal 1 and metal 2 is called via contact as shown in fig. 4.

- A $2\lambda \times 2\lambda$ cut centered on $4\lambda \times 4\lambda$ superimposed area is used to connect layers
- To connect metal 2 with diffusion via and cut both are used

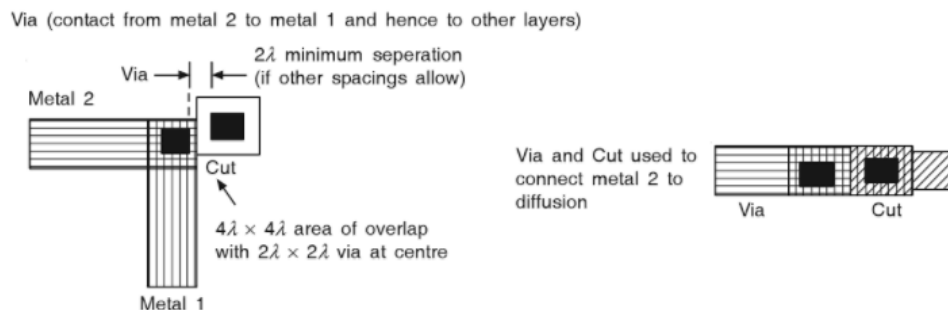


Fig. 4 contacts

Contact Cuts:

- Electrical connection between layers can be done using special structures 'contact cuts'.
- There are 3 approaches for contacts between polysilicon and diffusion in nMOS circuits. They are
 1. Polysilicon to metal and then to diffusion
 2. Buried contact - polysilicon to diffusion
 3. Butting contact - polysilicon to diffusion using metal
- ✓ Among the three buried contact is most used as it gives economy in space and reliable contact.
- Buried contact is distinguished feature in nMOS for connection between poly and diffusion and this is most widely used than butting contact.

Buried Contact (nMOS):

- Layer is joined over the area of $2\lambda \times 2\lambda$ with buried contact cut extending by 1λ in all directions except in the diffusion path. It extends by 2λ in order to avoid formation of unwanted transistors.
 - The contact cut shown in broken line indicates the region where thinox is removed on the silicon wafer and polysilicon gets deposited on wafer.
 - When impurities are added, it diffuses into poly and also to diffusion region within the contact area. This provides satisfactory contact between poly and diffusion as shown in fig 5.
- In CMOS poly to diffusion connections are made through metal. The process of making connection between metal and either of 2 layers (poly or diffusion) is by buried contact.

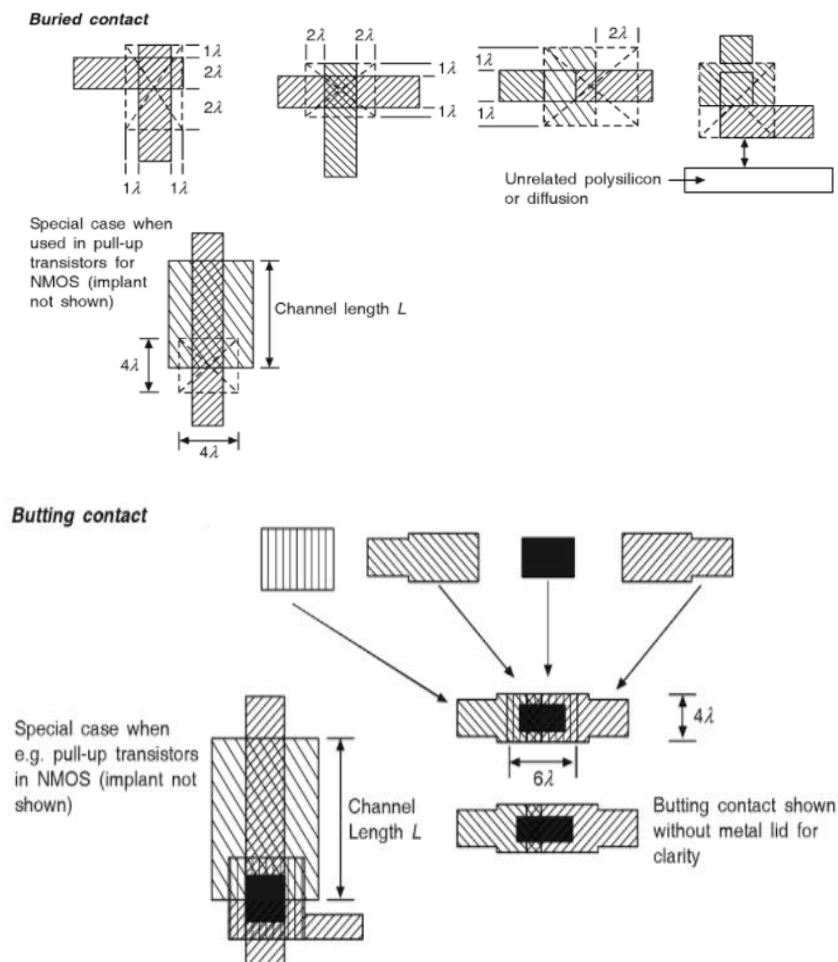


Fig. 5 Buried and butting contacts only for nMOS

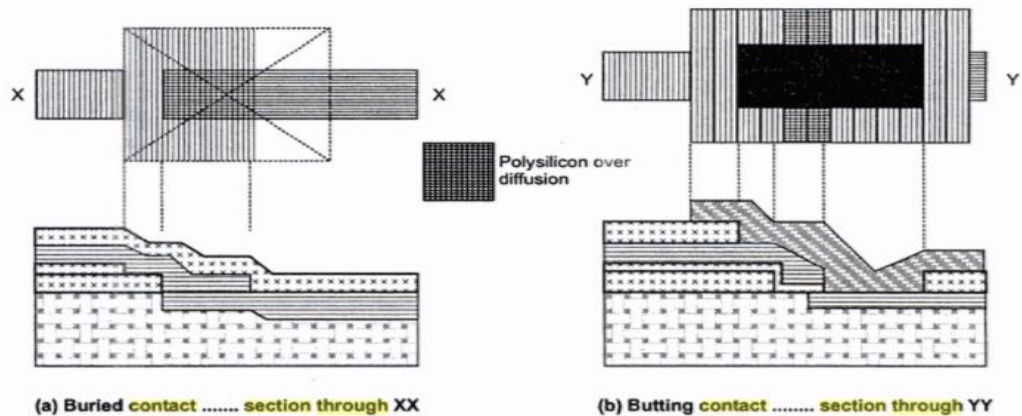


Fig 5. Cross section through contact structures

Butting contact

- Butting contact process is complicated and done when two layers do not overlap. Contact cut of $2\lambda \times 2\lambda$ is made until each of layers is joined. The layers are held in such a way that these two contacts become continuous.
- The poly and diffusion outlines overlap and thinox under poly acts as mask during diffusion process. Finally contact between two butting layers is done by a metal. This can be seen in fig. 5 cross-sectional view.

Double Metal MOS Process Rules:

- If to process considered till now introduction of second metal layer will boosts the design capabilities. It gives more freedom. Ex. this will be helpful for power rail (Vdd and Vss/Gnd) distribution and also for clock.
- This process is called Double Metal MOS Process
- This technique involves connecting metal 1 and metal 2 contacts called 'via'. This is shown in fig. 4 and fig. 6

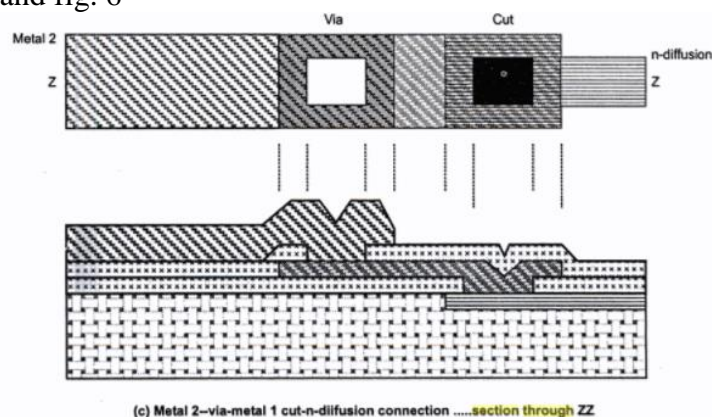


Fig. 6 cross section of via contact structure

- The 2nd metal layer is coarser than 1st metal layer (conventional) and the isolation layer between the 2 is usually thicker than normal.
- To distinguish contacts between 1st and 2nd metal layer they are called as 'vias' rather than contact cut
- In stick diagram representation its color code is dark blue or purple.

The steps of fabrication process is as follows:

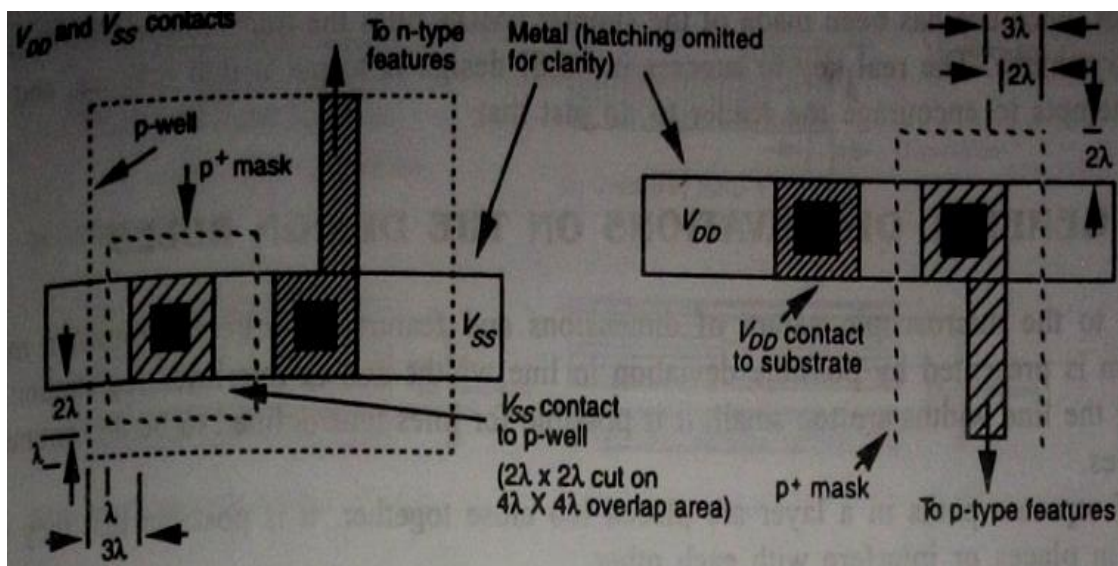
1. Using chemical vapor deposition oxide layer under 1st metal layer is deposited.
2. using same method oxide between 2 metal layers are formed.
3. Selected areas of oxide are removed by using plasma etching. The etching process is done under high vertical ion bombardment to get high and uniform etching.

The layout strategy used with double metal process is summarized as below

1. Second metal layer is usually used for global power railings and clock lines
 2. First metal layer is used for local power distribution and signals.
 3. The layout of the two metals are such that are mutually orthogonal wherever possible
- Similar to double metal process, other process allows second poly layer. The process steps are similar to previously described process.
 - The first polysilicon (poly 1) layer is deposited and patterned on this a second thinox (thin oxide) layer is grown. On this the second polysilicon (poly 2) layer is deposited and patterned. Thus 2nd thinox isolated the poly layers.
 - Presence of poly 2 provides greater flexibility in interconnections and allows transistors to be formed by intersection of poly 2 and diffusion.

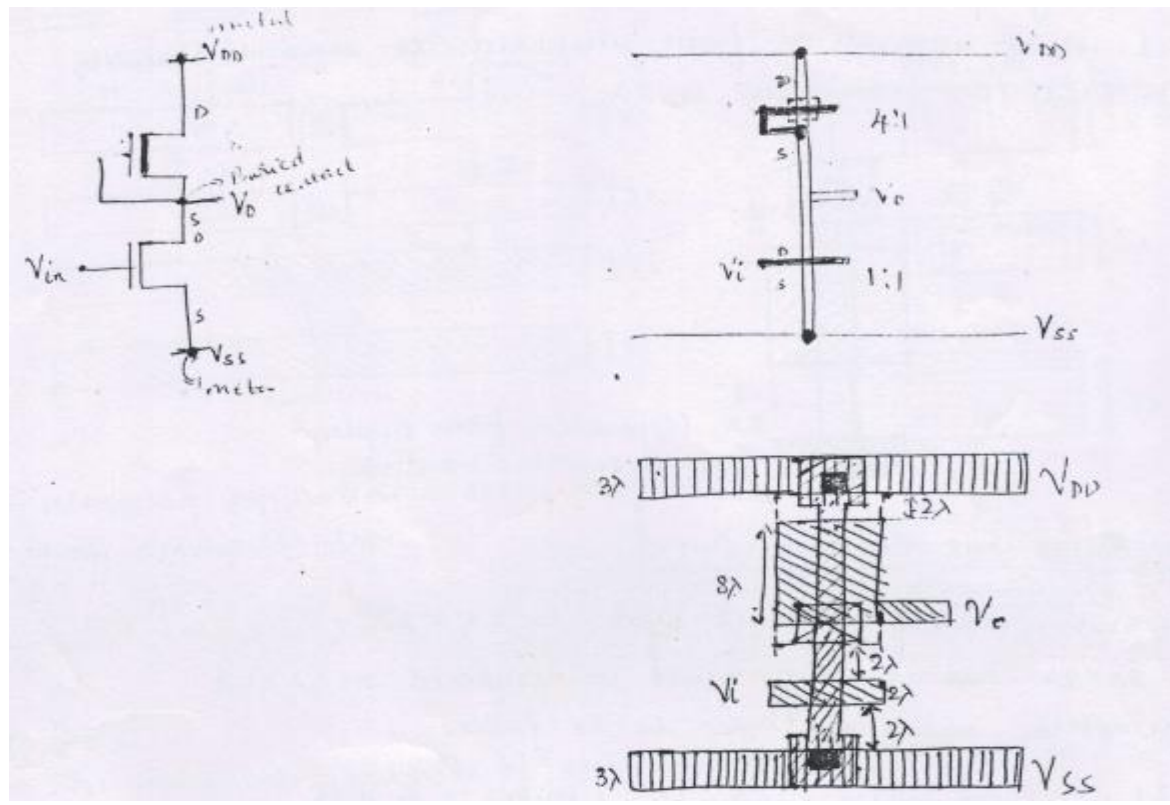
CMOS Lambda-based Design Rules:

- Comparing to Nmos fabrication process, CMOS fabrication is more complex.
- Extending the Nmos design rules, Noting exclusion of butting contact and buried contact rules.
- Additional rules associated with CMOS process concerned with unique feature p-well CMOS, i.e: p-well and P+ Mask and Substrate contact.

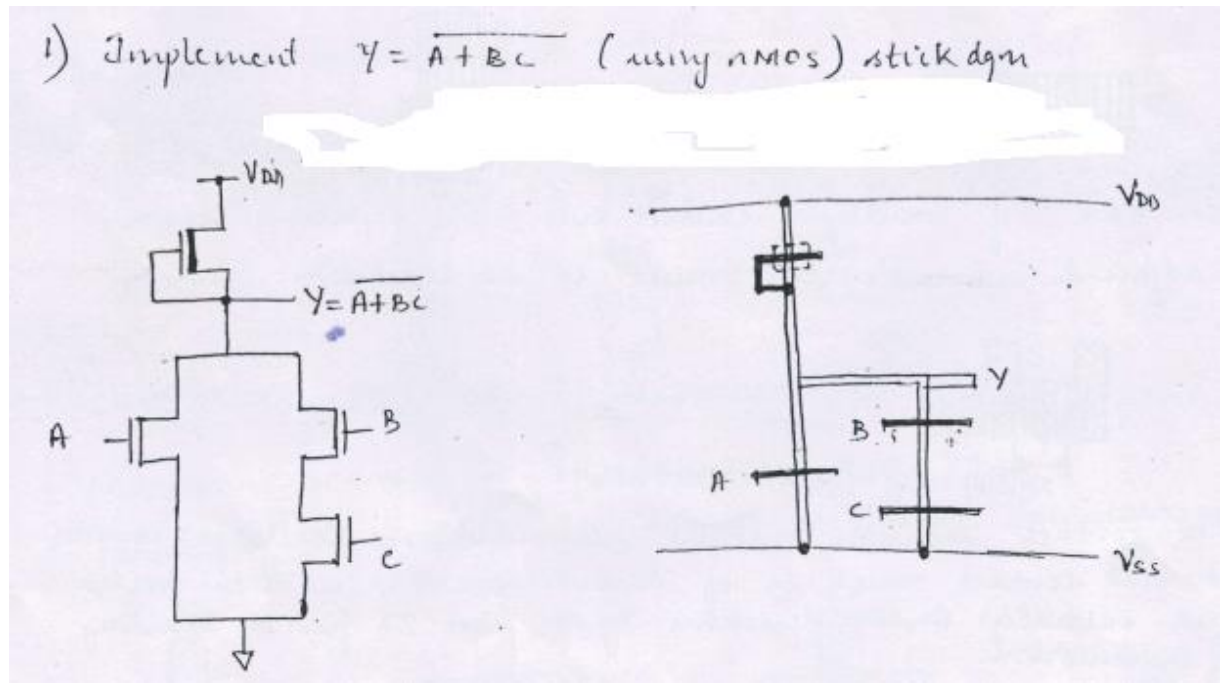


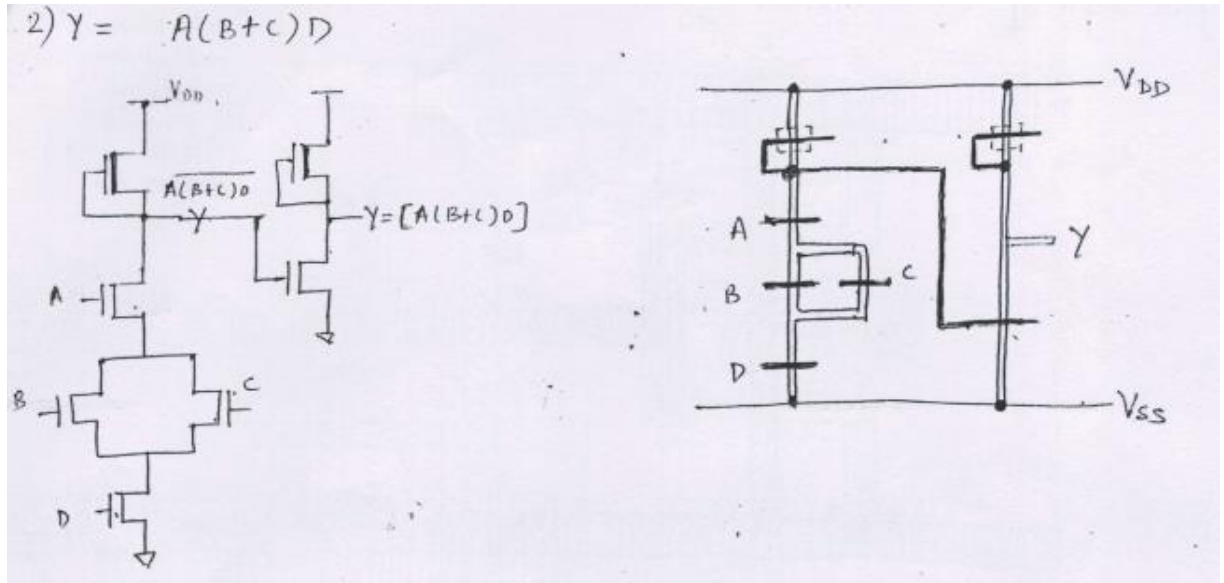
Problems on stick diagram and layouts.

Nmos Inverter

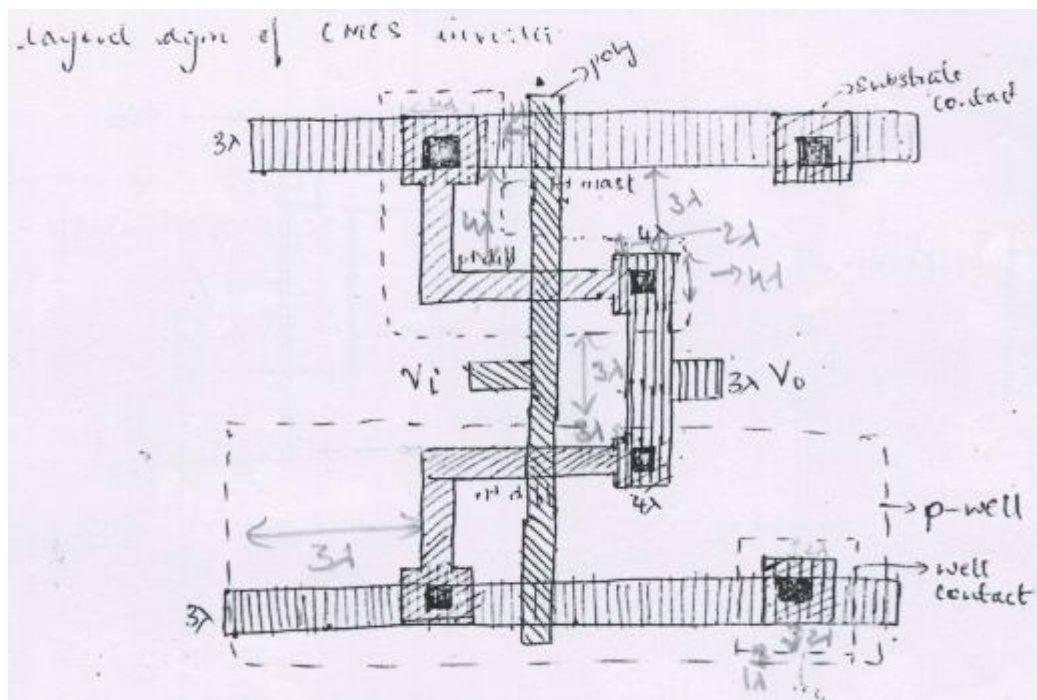
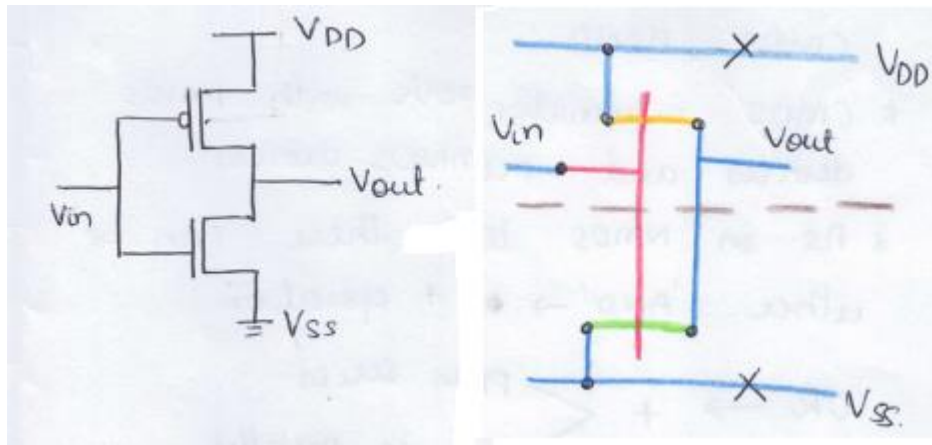


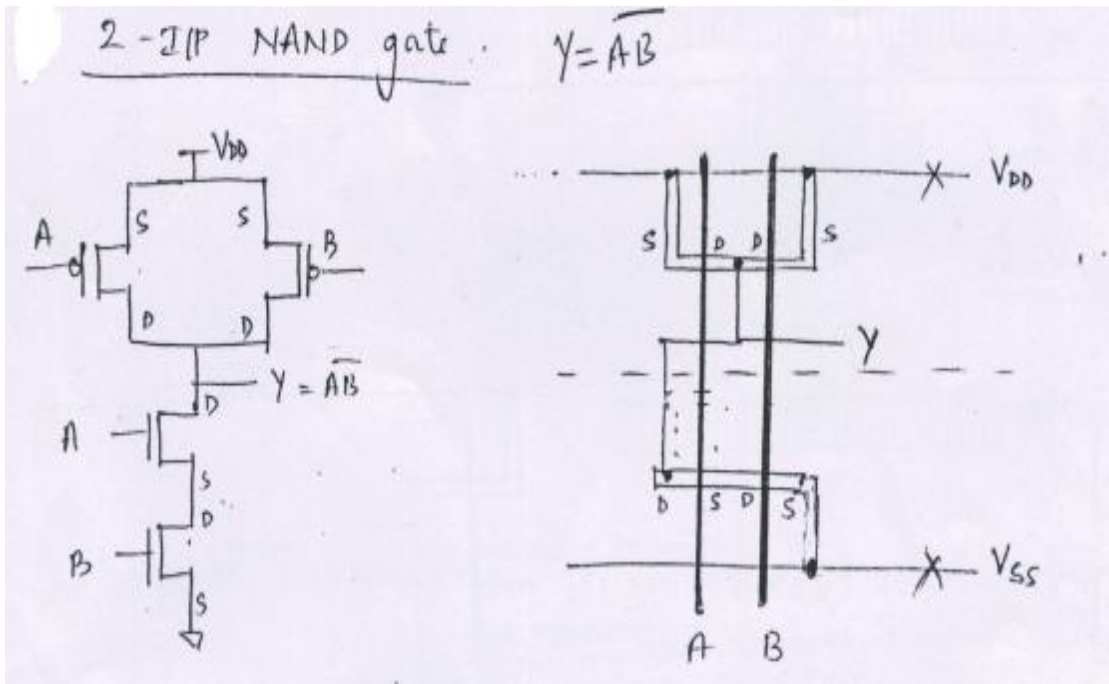
1) Implement $Y = \overline{A+BC}$ (using nmos) stick dgm



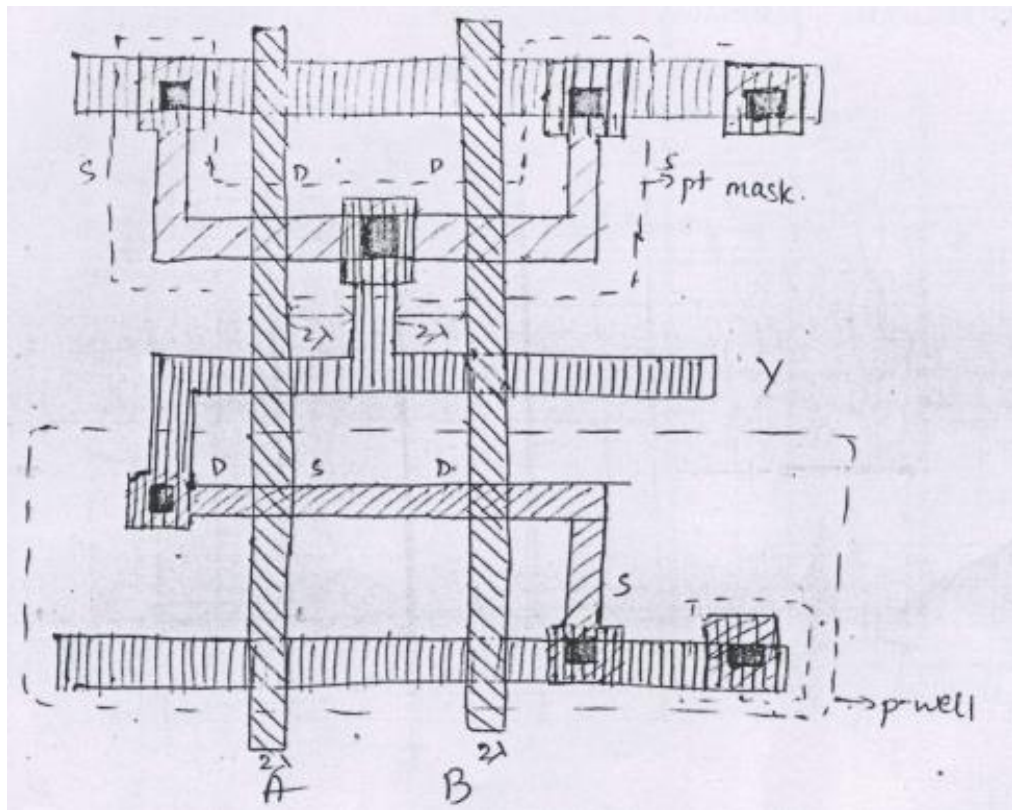


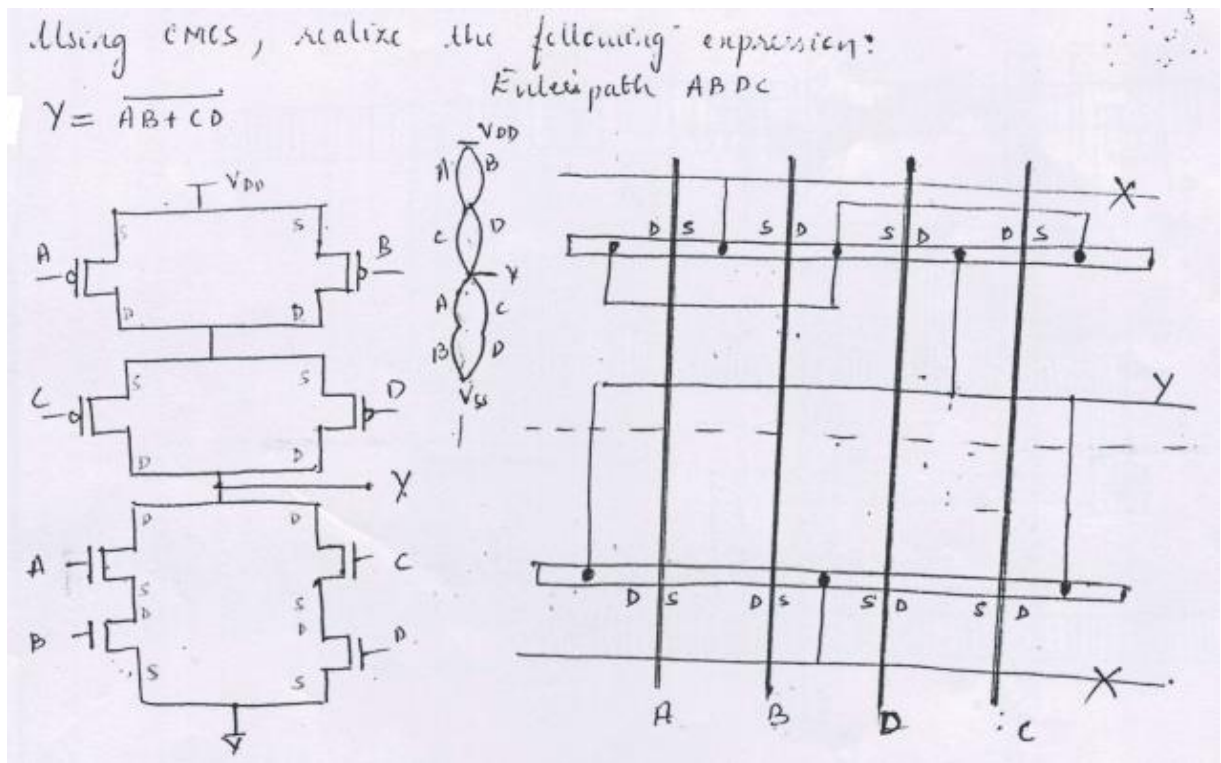
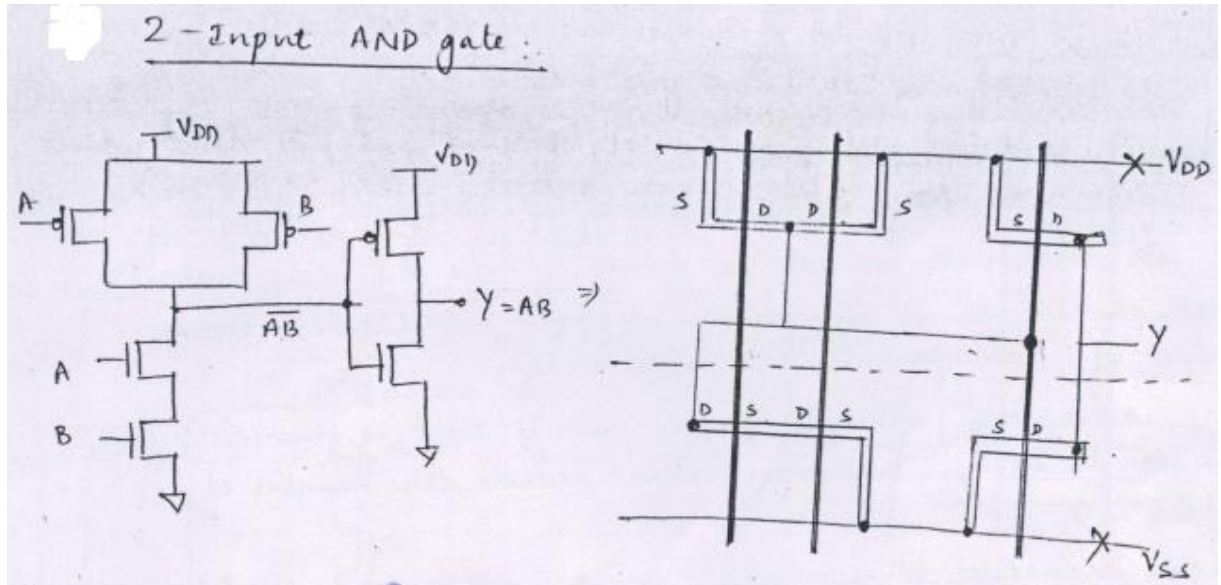
CMOS Inverter



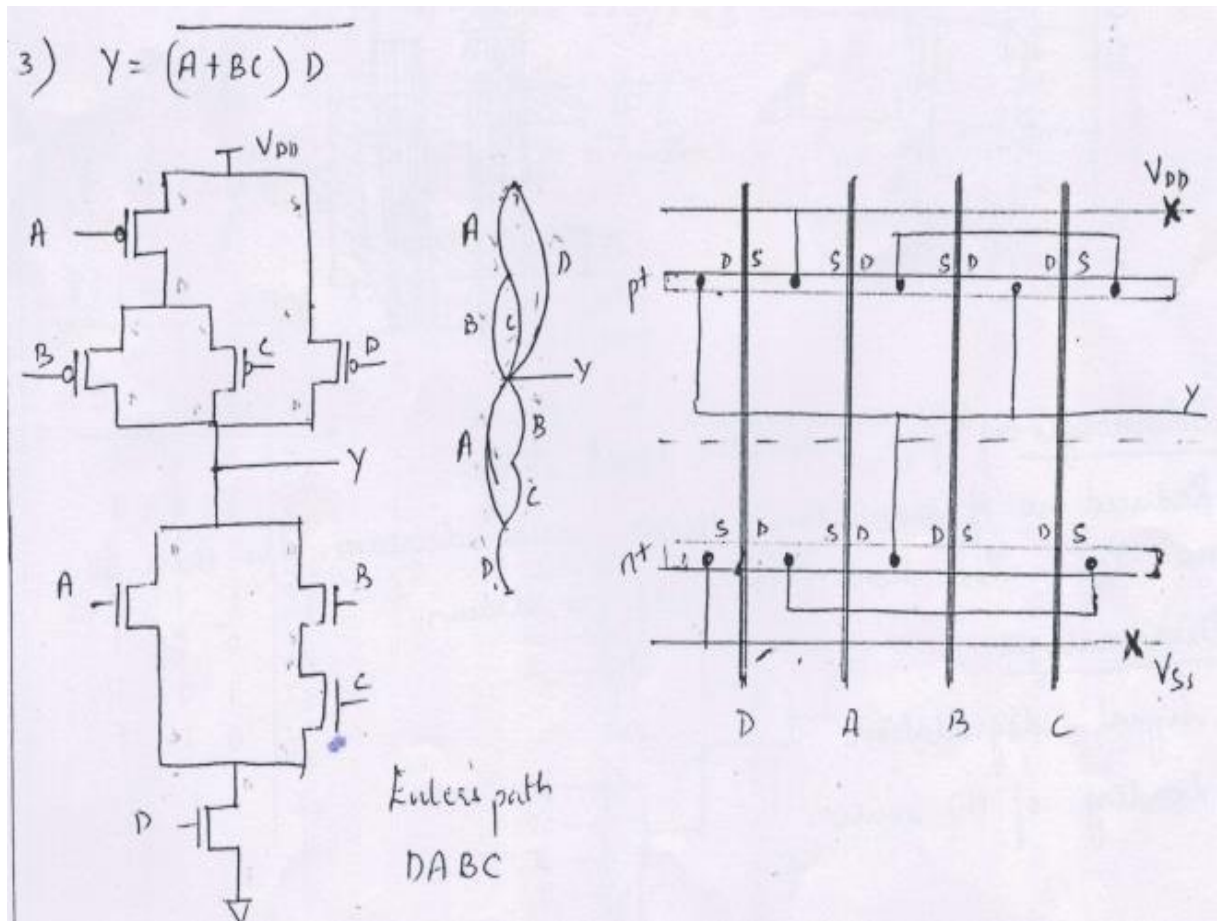
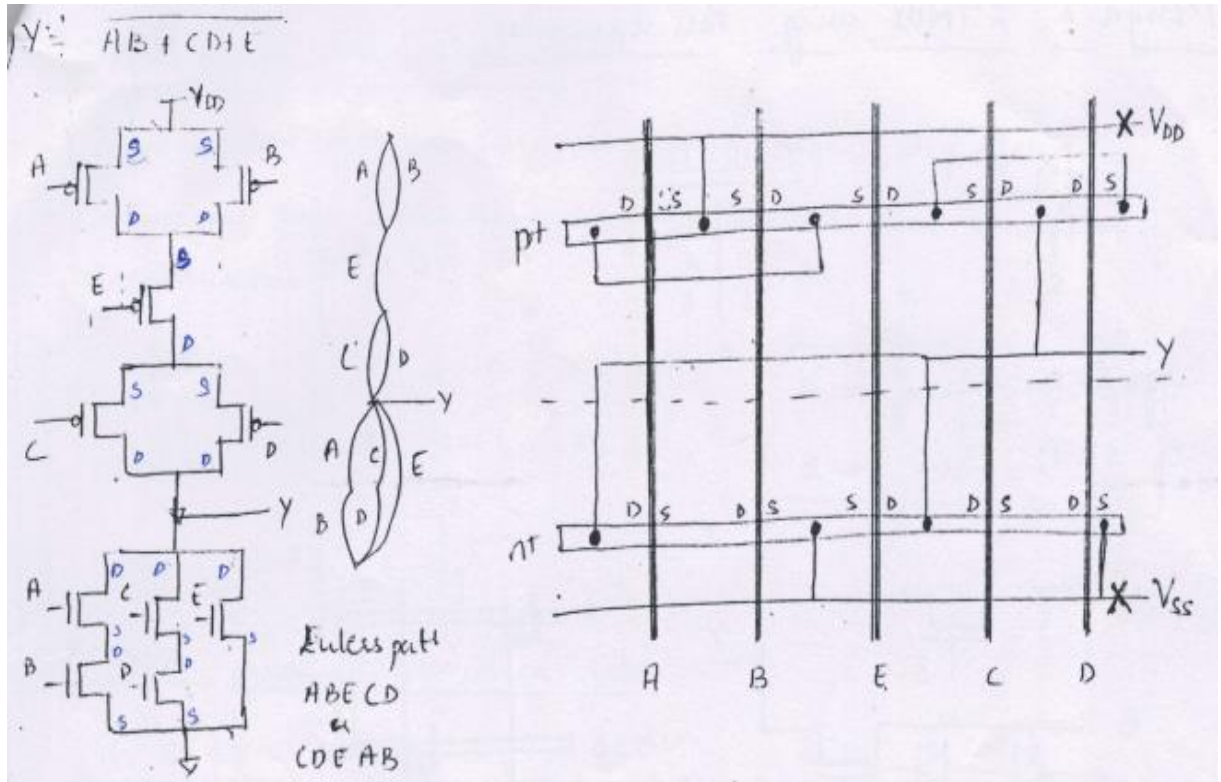


Diffusion lines running horizontal and polysilicon lines in vertical direction.

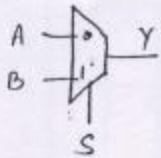




- When more number of inputs available, Euler's path is determined to know gate ordering.
- Advantage of using Euler's path is to that a common diffusion line can be used which reduces number of contact cuts.
- Uninterrupted path in both pull-up and pull down network represents optimized gate ordering which helps in drawing layout without breaking the diffusion layer.



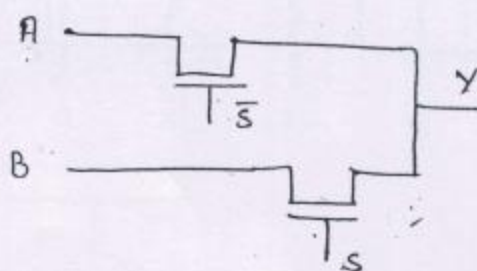
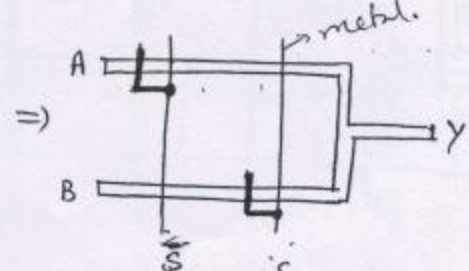
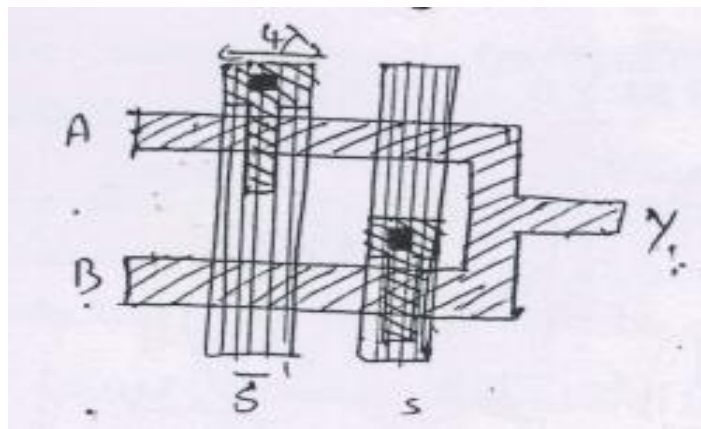
Design a 2:1 MUX using Pass transistors



S	Y
0	A
1	B

$Y = \bar{S}A + SB$

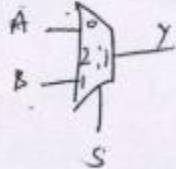
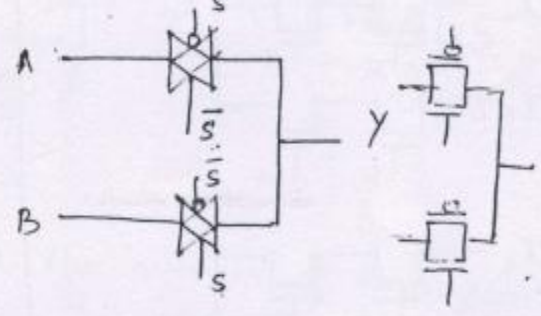
if $S=0$, then $Y=A$
 if $S=1$, then $Y=B$

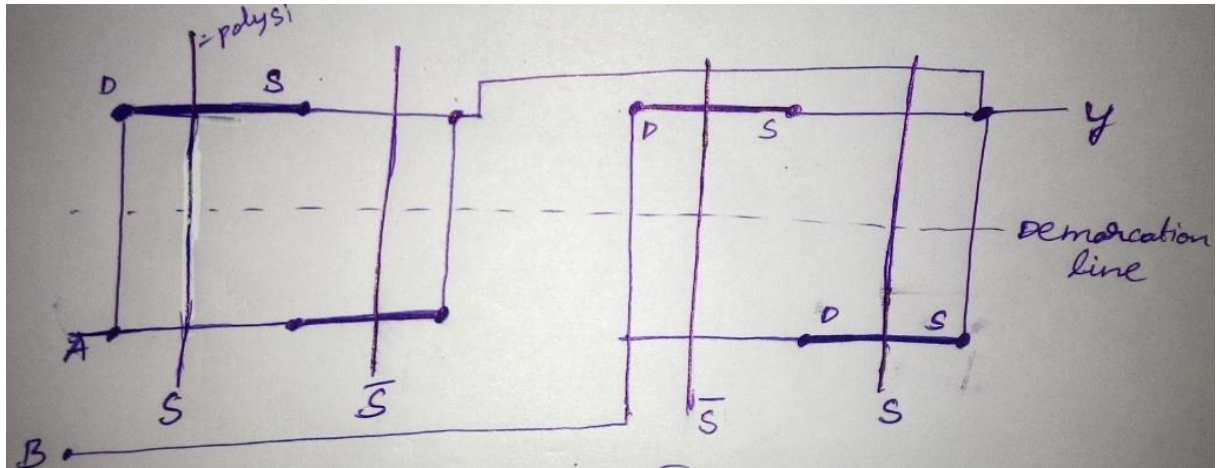




2:1 MUX using transmission gates

S	Y
0	A
1	B

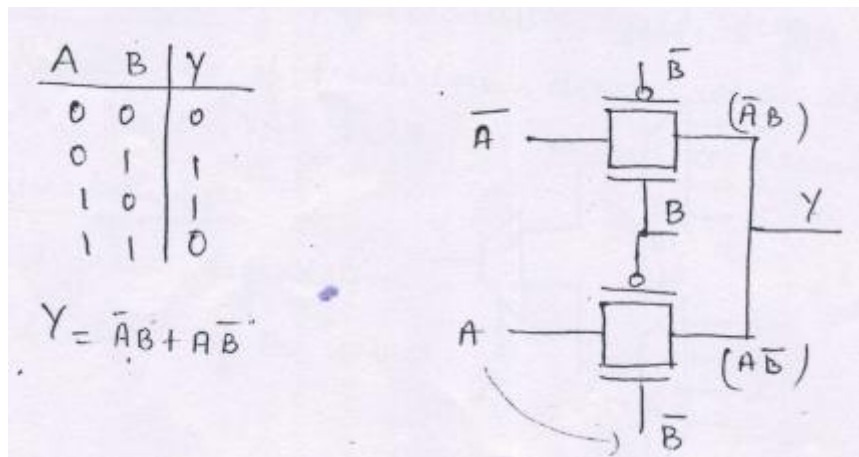
if $S=0$, $Y=A$
 if $S=1$, $Y=B$.

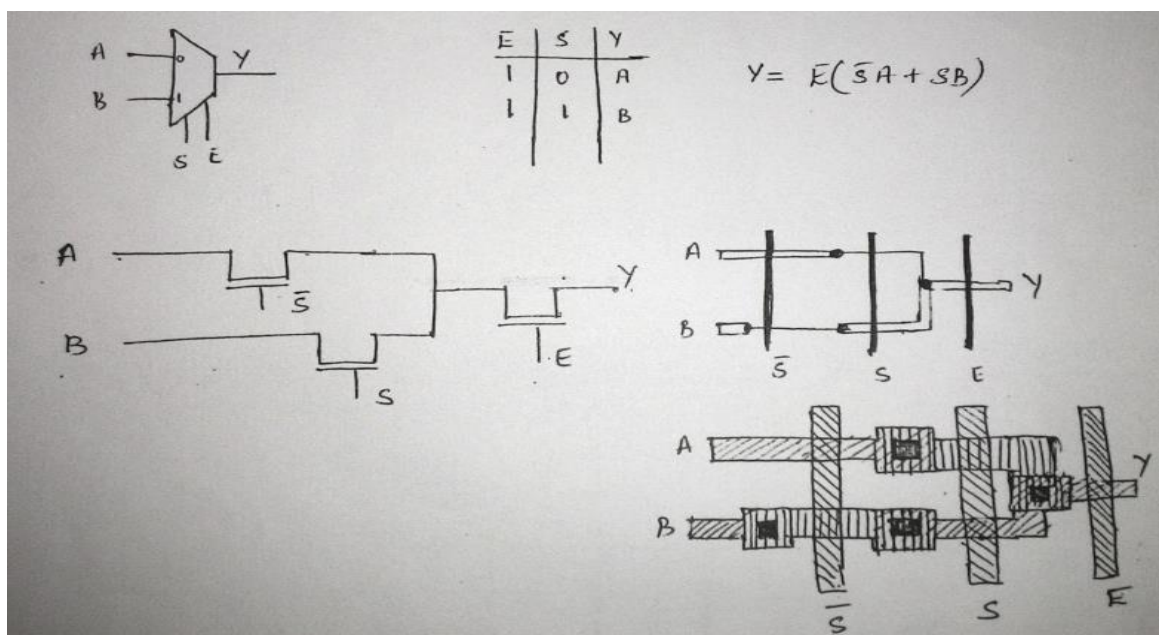


Stick diagram of 2:1 MUX using transmission gates.

Two input XOR gate realization using transmission gates.



Two way selector with enable

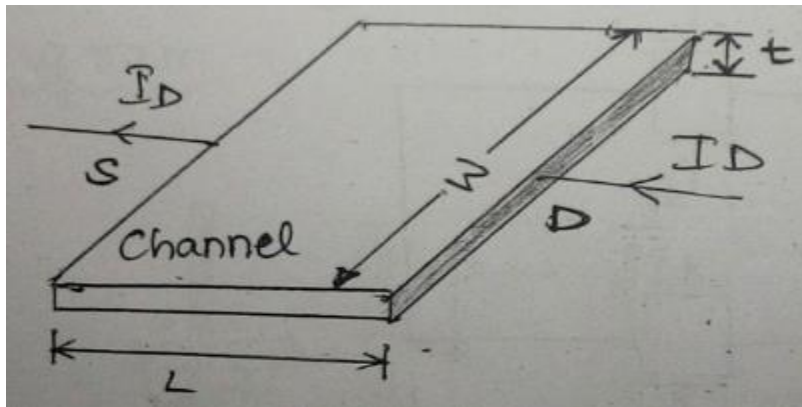


Basic circuit concepts

- In MOS technology, Active devices are dealt with some measurement.
- Wiring up of circuits is done through various conductive layers which is produced by MOS Processing.
- Therefore it is necessary to be aware of resistive and capacitive characteristics of each layers.
- For evaluating the effects of wiring, input and output capacitances, sheet resistance and standard unit capacitances are used.
- Further delay associated with wiring, inverters are evaluated by the term delay unit τ .

Sheet Resistance R_s

Consider a transistor with a channel having resistivity ρ , width W , thickness t and length between source and drain is L .



Resistance of the channel between drain and source is expressed as.

$$R_{DS} = \frac{\rho \cdot \text{Length}}{\text{Area of cross section}} = \frac{\rho \cdot L}{t \cdot W}$$

$$R_{DS} = R_s \frac{L}{W}$$

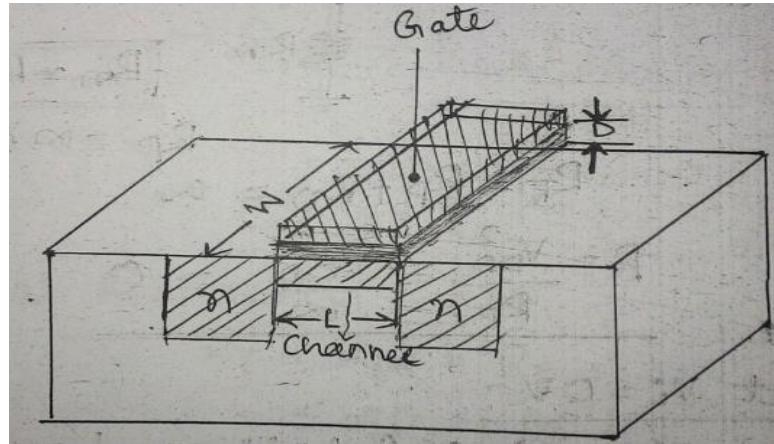
Where $R_s = \frac{\rho}{t}$ is a constant and it is called sheet resistance.

From the above equation, sheet resistance can be defined as resistance of the channel whose length and width are equal.

R_s is completely independent of area square. Ex: $1\mu\text{m}$ per side square slab of material has exactly same resistance as 1cm per side square slab of same material if thickness is same.

Area capacitance

In between gate and channel exists a capacitance and it is called gate capacitance and denoted by C_g .



From the above diagram.

$$C_g = \frac{\epsilon_0 \epsilon_r A}{D}$$

A is area of the channel or surface area of the gate

$$\frac{C_g}{A} = \frac{\epsilon_0 \epsilon_r}{D} \text{ pF}/(\mu\text{m})^2$$

$$C_A = \frac{\epsilon_0 \epsilon_r}{D}$$

ϵ_0 = permittivity of free space = 8.854×10^{-12} F/m.

ϵ_r = relative permittivity of a given material

D = thickness of SiO_2 constant for a given technology.

Area capacitance is defined as capacitance per unit area at the gate of transistor and denoted by C_A .

Standard unit of capacitance ($\square C_g$)

The standard unit of capacitance is defined as the capacitance at the gate of 1:1 transistor.

Ex: consider a 1:1 transistor where $L = 2\lambda$ and $W = 2\lambda$.

Gate area of transistor = $L \times W$

$$A = 2\lambda \times 2\lambda = (2\lambda)^2$$

Actual capacitance at the gate of transistor $C = C_A \cdot A = C_A \cdot (2\lambda)^2$

$$C = 4 \times 10^{-4} \text{ pF}/(\mu\text{m})^2 \cdot (2\lambda)^2$$

Consider $5 \mu\text{m}$ technology, i.e: $2\lambda = 5 \mu\text{m}$

$$C = 4 \times 10^{-4} \text{ pF}/(\mu\text{m})^2 \cdot (5 \mu\text{m})^2 = 4 \times 10^{-4} \times 25 = 0.01 \text{ pF}$$

$$C = 1 \square C_g$$

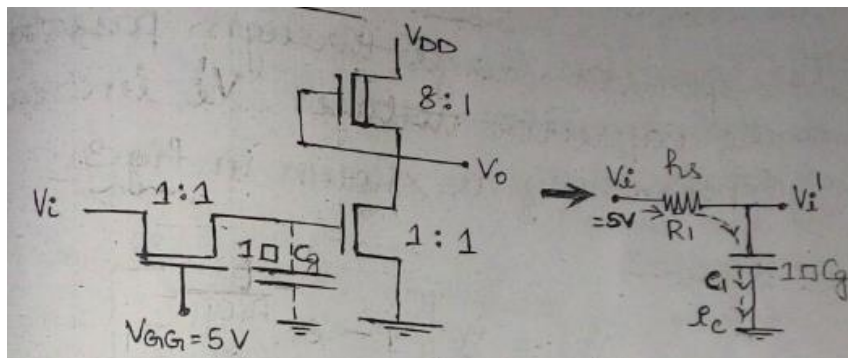
Standard Delay Unit (τ)

Time delay is measured in terms of standard unit τ .

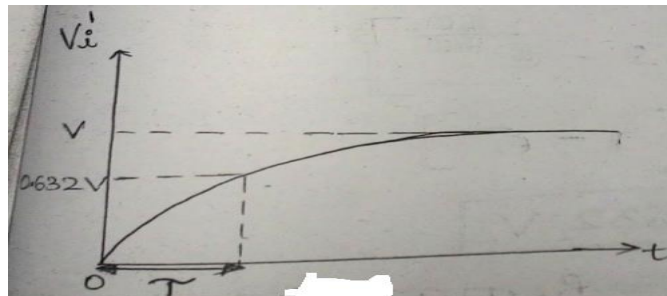
It is defined as product of R_s and C_g . i.e: $\tau = R_s \cdot \square C_g$

Measurement of τ .

Consider nmos driven by pass transistor shown in below figure and the dimensions are indicated. Pass transistor is ON for given gate voltage V_{GG} . Pass transistor is represented by R_s due to its equal length and width. Pull down transistor of inverter is represented by capacitance $\square C_g$. Since pull down has minimum dimensions.



τ is defined as time taken by capacitor to charge from 0 to 63.2% of maximum value as shown in below figure.

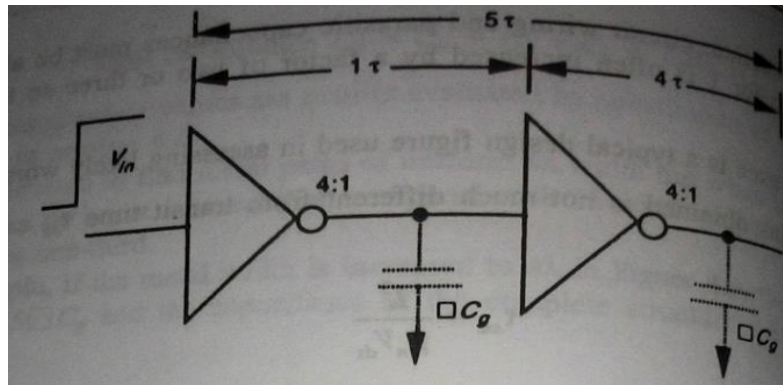


Inverter Delays

Consider basic 4:1 nmos inverter. To achieve 4:1 Z_{pu} to Z_{pd} ratio, R_{pu} will be $4R_{pd}$. Clearly resistance R_{pu} value is $R_{pu} = 4 R_S = 40K\Omega$. Meanwhile R_{pd} value is $10K\Omega$.

Delay associated with inverter depends on ON and OFF condition of transistors.

Consider a pair of cascaded inverter, delay in this pair will be constant irrespective of sense of logic level transition. The overall delay of nmos inverter is $\tau + 4\tau = 5\tau$. Shown in below figure.

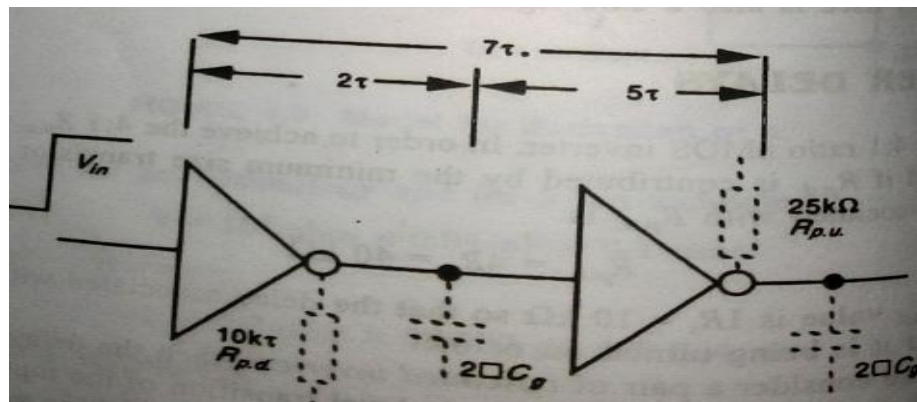


In general term delay through nmos inverter pair is given by $T_d = (1 + Z_{pu}/Z_{pd}) \tau$

So single 4:1 inverter exhibits asymmetric delays, delay in turning on τ (capacitor discharging condition) and delay in turning off is 4τ (capacitor charging condition). Asymmetry becomes worse for inverter with 8:1 ratio.

For CMOS inverter, nmos rules no longer applies, but we need to consider natural asymmetry of equal size pull up and pull down transistors.

Gate capacitance is double compare to nmos inverter since input is connected to both transistors and delay associated with pair of minimum size inverters is shown in below figure.



Asymmetry of resistance is eliminated by increasing the width of p- device channel by factor of two or three, but gate capacitance increases by the same factor.

Driving large capacitive loads

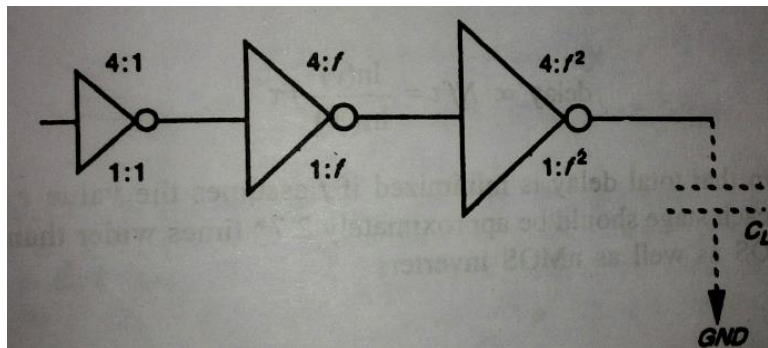
- A large capacitive loads problem arises when a signal to be transmitted from On chip to Off chip destinations.
- Off chip capacitance is generally higher than On chip $\square C_g$. And it is denoted by C_L .

$$C_L \geq 10^4 \square C_g$$

- A capacitance of this order to be driven through low resistance otherwise long delays will occur.

Cascaded Inverters as drivers

- Inverters to drive large capacitive loads resistance associated with pull up and pull down transistors to be low.
- Low resistance values of Z_{pu} and Z_{pd} implies low L:W ratio or channel width must be made wider to reduce channel resistance but consequently inverter occupies large area.
- Gate area $L \times W$ is more significant and large capacitance present at input which slows down rate of change of voltage at input.
- Remedy to use N cascade inverter is by maintaining L to a minimum feature size and width of each successive stage is increased by factor f as shown in below figure.



- With increase in width factor increases capacitive load at input side and area occupied by the inverter also increases.
- The rate of width increase influence on number of stages to be cascaded to drive particular C_L value.
- Total delay associated with nmos pair is 5τ and cmos pair is 7τ .

Let $y = \frac{C_L}{\square C_g} = f^N$, f and N are interdependent.

To determine value of f to minimize overall delay for given y

$$\ln(y) = N \ln(f)$$

$$N = \frac{\ln(y)}{\ln(f)}$$

$$\begin{aligned} \text{For } N \text{ even, total delay} &= \frac{N}{2} 5 f \tau = 2.5 f \tau \text{ (nmos) or} \\ &= \frac{N}{2} 7 f \tau = 3.5 f \tau \text{ (cmos)} \end{aligned}$$

$$\text{In all cases, delay} \propto N f \tau = \frac{\ln(y)}{\ln(f)} f \tau$$

- Total delay is minimized if f assumes the value e . i.e: each stage is approximately 2.7 times wider than its predecessor and it is applicable for both cmos and nmos inverters.

Thus assuming $f = e$, we have

$$\text{Number of stages } N = \ln(y)$$

And overall delay t_d

$$N \text{ even: } t_d = 2.5 N \tau \text{ (NMOS) or } t_d = 3.5 N \tau \text{ (CMOS)}$$

$$N \text{ odd: } t_d = [2.5 (N-1)+1]e \tau \text{ (NMOS) or } t_d = [3.5 (N-1)+1]e \tau \text{ (CMOS)}$$

For ΔV_{in} which indicates logic 0 to 1 transition of V_{in} .

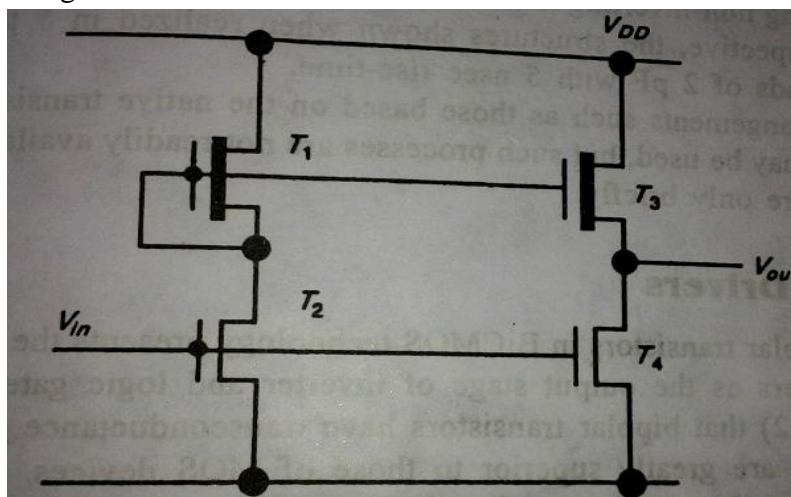
$$t_d = [2.5 (N-1)+4]e \tau \text{ (NMOS) or } t_d = [3.5 (N-1)+5]e \tau \text{ (CMOS)}$$

For ΔV_{in} which indicates logic 1 to 0 transition of V_{in} .

Super buffers

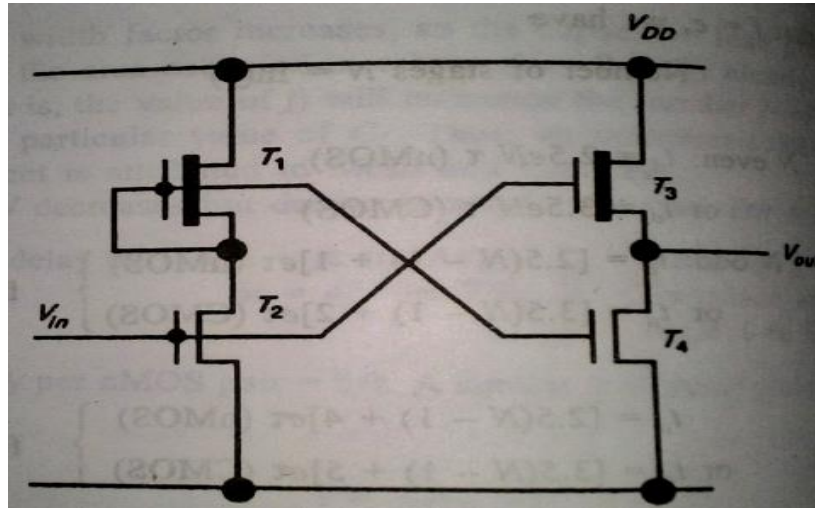
- Asymmetry of conventional inverter gives rise to significant delay problems when used to drive large capacitive loads.

Common approach used in nmos inverter is to use super buffers an inverting type nmos super buffer is shown in figure.



- Consider input $V_{in} = 1$, the inverter formed by T_1 and T_2 is turned On and thus gate of T_3 is pulled down to zero volts with small delay. So T_3 is in cut off and T_4 is turned On and output is pulled down.

- When $V_{in} = 0$, gate of T_3 is allowed to rise to V_{DD} . Thus T_4 turned Off, T_3 is made to conduct with V_{DD} on its gate. The voltage applied to gate is twice the average voltage of conventional nmos inverter.
- Doubling effective V_{gs} will increase current, thus reduces the delay in charging capacitor at output, so symmetry is achieved.
- The Non-inverting type nmos inverter is shown in below figure.



BICMOS Inverter

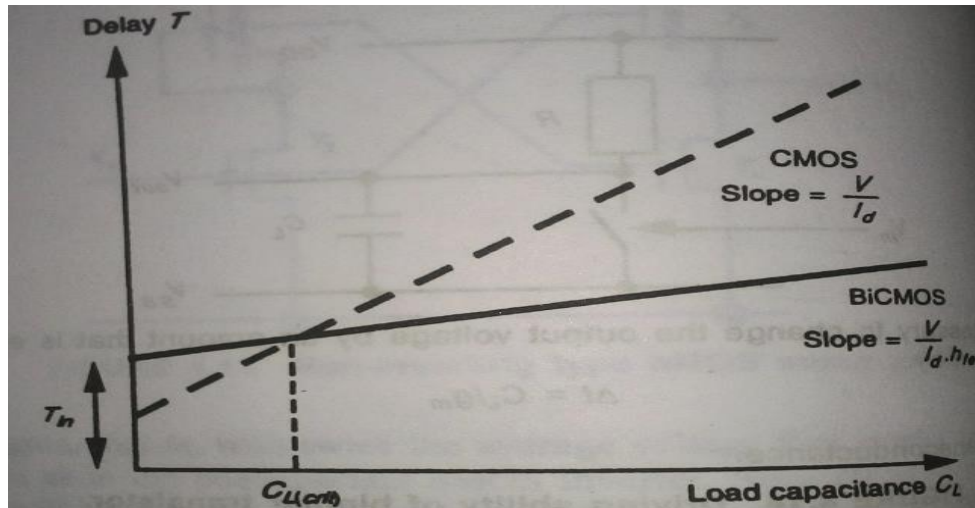
- Bipolar transistor availability in Bicmos technology presents possibility of using bipolar transistors as drivers at the output stage of the inverter.
- Transconductance and current/area characteristics are superior than MOS devices. so it has high current driving capability.
- Bipolar transistor has exponential dependence of output current on base emitter voltage which means transistor can operate with small input voltage swing compared to MOS transistors and switches large current.
- So the bipolar transistors have better switching performance results in small input voltage swing and switch large current.
- Switching performance of transistor driving capacitive load can be seen from simple model.
- The time required to change output voltage V_{out} by an amount equal to input voltage is given by

$$\Delta t = \frac{C_L}{g_m}$$

Where g_m is trans conductance of bipolar transistor. As g_m increases Δt decreases.

- Bipolar transistor delay has 2 main components T_{in} and T_L .
- T_{in} is the initial time required to charge the B-E junction of the transistor. It is time taken to charge the input gate capacitance.

- T_L is time taken to charge the output load capacitance C_L . This value is less for bipolar by factor of h_{fe} .
- As BJT has higher T_{in} , T_L is small and because of this faster charging takes place and helps in reducing the delay.
- Combined effect of T_{in} & T_L is in in graph. There is C_{Lcrit} critical load capacitance below which BICMOS driver is shown than CMOS driver.



- Delay of BICMOS is described by $T = T_{in} + (V/I_d)(1/h_{fe})C_L$.
- Delay for BICMOS inverter is reduced by a factor of h_{fe} when compared with CMOS inverter.