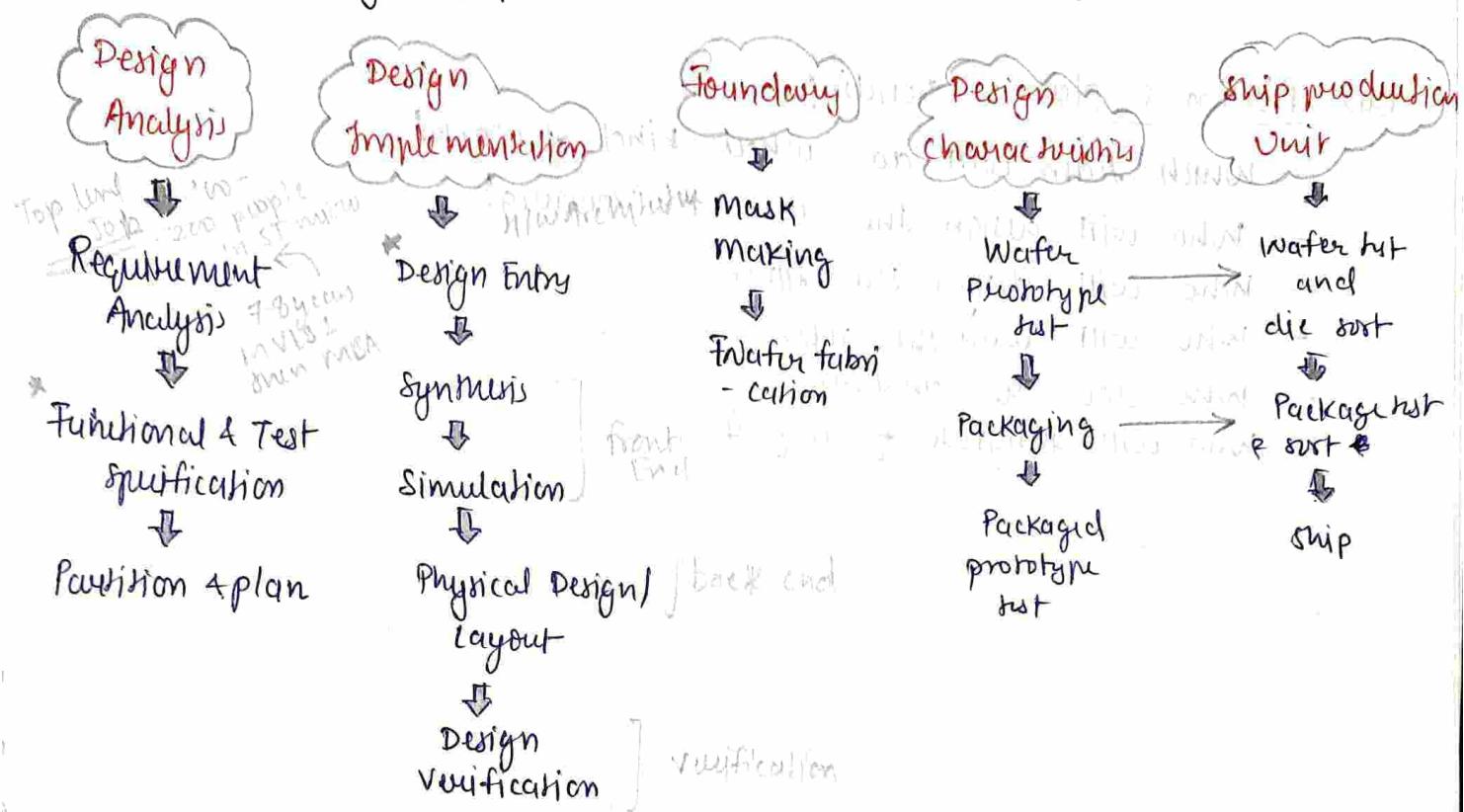


1. Introduction to VLSI Design

Video-1
Lecture 2

SOC Design steps



① Design Analysis:

- (i) Requirement Analysis :- Here the market Analyst will look for the consumer's interest at which price point they are ready to pay for the product.
 - (ii) Functional & Test Specs :- While making the face recognition of the iPhone we have to decide its functionality, it should be inferred, it should recognize the face, it should unlock the screen, it should unlock the screen if 80% or above face matches.
 - Along with the functionality, you have to decide how we are going to test the functionality.
 - We will be taking 10,000 samples from the camera & if any sample matches with the input face, it will unlock the screen.
- "functional specification" (and) "Test specification"**

** If the functional specification changes at any time of the design flow, then we have go back to the functional & start over again. and start over again.

(iii) Partition & plan : Deciding

- which team will do what kind of work?
- Who will design the layout for camera?
- Who will design the H/w.
- Who will write the software
- Who does the analysis
- Who will integrate all of it.

⇒ All these process can be done by 10 people.

- These kinds of job are provided by headquarter by:
① very close to customer
 - TI - US
 - micron - US
 - Samsung - Korea
 - ST microelectronics - Italy & France

② Design implementation

(i) Design Entry :- here we have to make a design of H/w which will do the calculation to ensure that face recognition happens accurately.

(ii) Synthesis :- here we decide like what components are required in here like adder, counter etc. multiplier, comparator etc. And I synthesize it & make the connections.

(iii) Simulation :- here we simulate our ckt & tell whether the design work is correct or not.
if yes we can emulate our ckt using FPGAs etc.

* Aging simulation model.

simulation → Emulation

(CAD) (FPGA)

* Monte carlo

(iv) Physical design & layout:

Tape out

Model core implementation.

(v) Design verification: Verifying the completely designed ckt. then we tape out the design & send it to foundry.

DRC and all.

③ Foundry:

the design is send to mask shop & then

- (i) mask making the masks are send to fabur worker
- (ii) Wafer fabrication

- Depending on the technology, there can be 39, 28, 47, 60 masks that the design would go through.

④ Design characterization:

- (i) Wafer prototype test

- (ii) Packaging

- (iii) Packaged prototype test

- We have 12 inch diameter wafer with ~~our~~, which ~~has~~ billions of dies (ie, IC used in our phone). We test if all the dies are working fine or not.

~~the design means~~

- While packaging, we are connecting wires from the pins of the package to the pads on your die.
- There could be some error in the connection above, therefore we do Packaged prototype test.

Q. What is meant by design characterization?

- If I have designed a ckt at 1.2 GHz then ~~I have~~ I have to check after its fabrication that for silicon at what clock speed is it working 1.1 GHz or 1.25 GHz.

⑤ Ship Modulation Unit

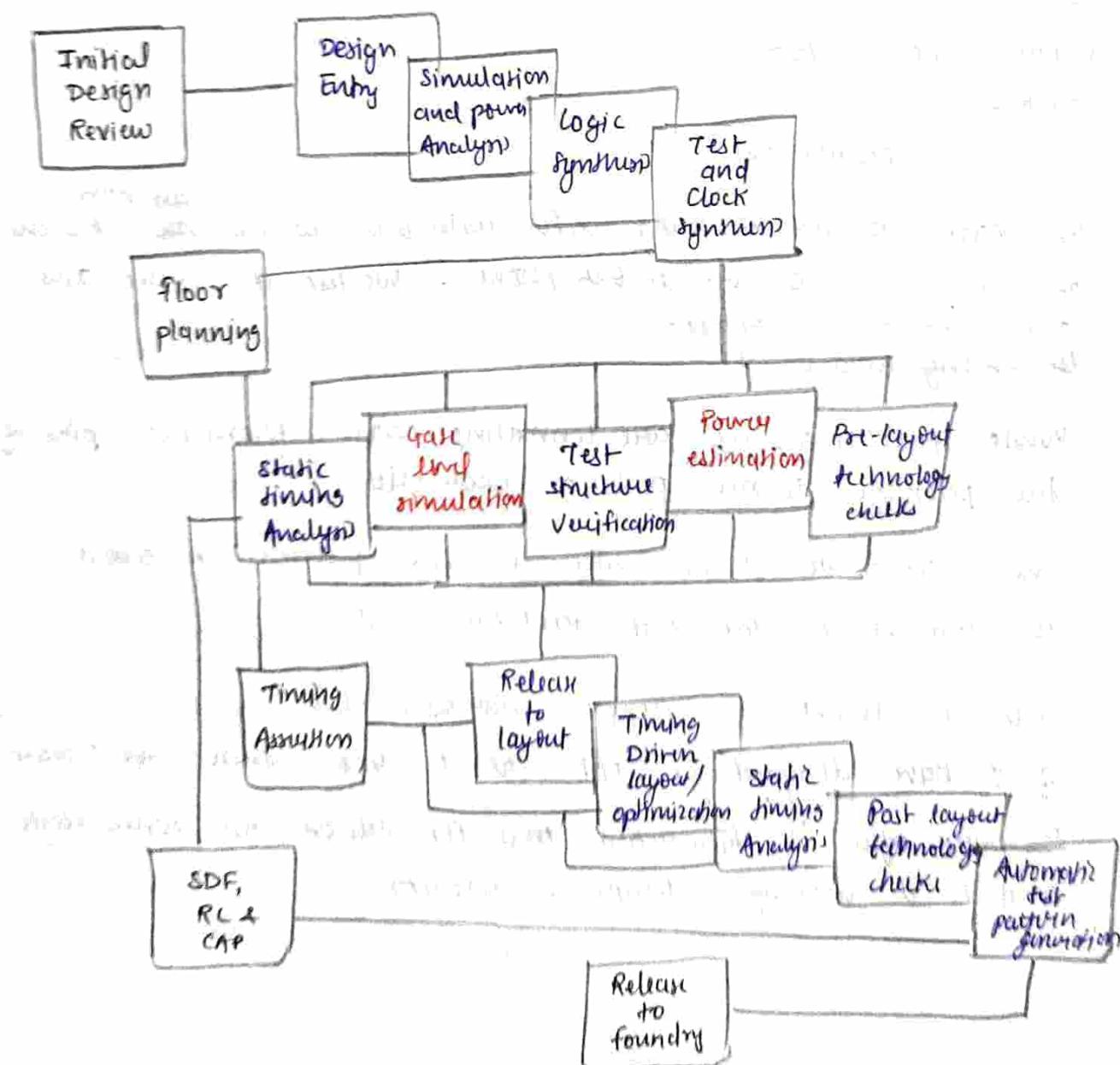
PPA \Rightarrow power - performance - Area.

- (i) Wafer test & die sort
- (ii) Package ^{test} & sort
- (iii) ship

ship is just an
 ↗
 IC.

- Testing the wafer & calculating its frequency & then working the die on the band of its freq & performances
- Die sort means when the dies are on the wafer.
- Package sort means when the dies are packaged into their ICs.

Design Implementation :-



- Initial Design ~~Review~~ :- Where you say I need this kind of a processor, I need this kind of memory, interface, DDR and so on.
- Design Entry :- RTL design, where you write a verilog code.
- Simulation and power analysis :- Now you check whether after simulation that power calculation is meeting your needs.
- Logic synthesis :- Converting an RTL verilog code to netlist to a gate level netlist. If we have design a adder logic in our design entry (RTL) now we make the (design time) logic behind that adder.
- Test & clock synthesis :- Suppose you have 1 million FF and every FF needs a clock signal to run. So can we use a single clock to run all the FF? No.
We will use buffers b/w the clocks & FFs, such that 10 buffers will drive the 10 FFs. So this part of driving clock buffers, such that none of the clock pins nodes has a load more than 100 pF.
If you have 2 FF & both are supposed to be run sequentially. ∵ equal no. of buffers must be put in path of both the FF. This is called as clock synthesis.
- Floor planning :- Here we decide where we are going to put memory, where the processor is going to be placed, so on.

~~skip module for now~~

- (i) wafer test & die sort
- (ii) Package test & sort
- (iii) ship

*⑦ static timing Analysis (STA)

- Here we are going to analyze the timing behaviour of the circuit.
- checking whether the circuit meeting the setup & hold time conditions.

⑧ Gate level simulation

- When we have converted the $A+B$ logic to ~~an~~ Adder circuit then simulating the ~~an~~ gate level netlist is done at this level.

⑨ Test structure verification

- Here we test the components individually.

⑩ ? Power consumption

- From the usual gates available how we estimate the power consumption.

⑪ Pre-layout checks

- We can check whether ~~or~~ more any single node that is driving ~~two~~ from two inputs.

— or is there any node which is floating.

— or is there any load at particular node gain which it can't bear to hold.

After ~~int~~ checks you can generate your netlist.

⑫ Timing constraint

⑬ Relax to layout

- After deciding where we are going to place all of our components we can now start placing it in the chip.

⑭ Timing driven layout optimization:

- If you know that output of one AND gate is going into this OR gate. Therefore we can put this AND & OR gates together. This is layout optimization.

⑮ STA in the layout:-

- Because ⁱⁿ layout we ~~can~~ ~~not~~ know how long the wire is & the ~~actual~~ delay of the clk must be known.

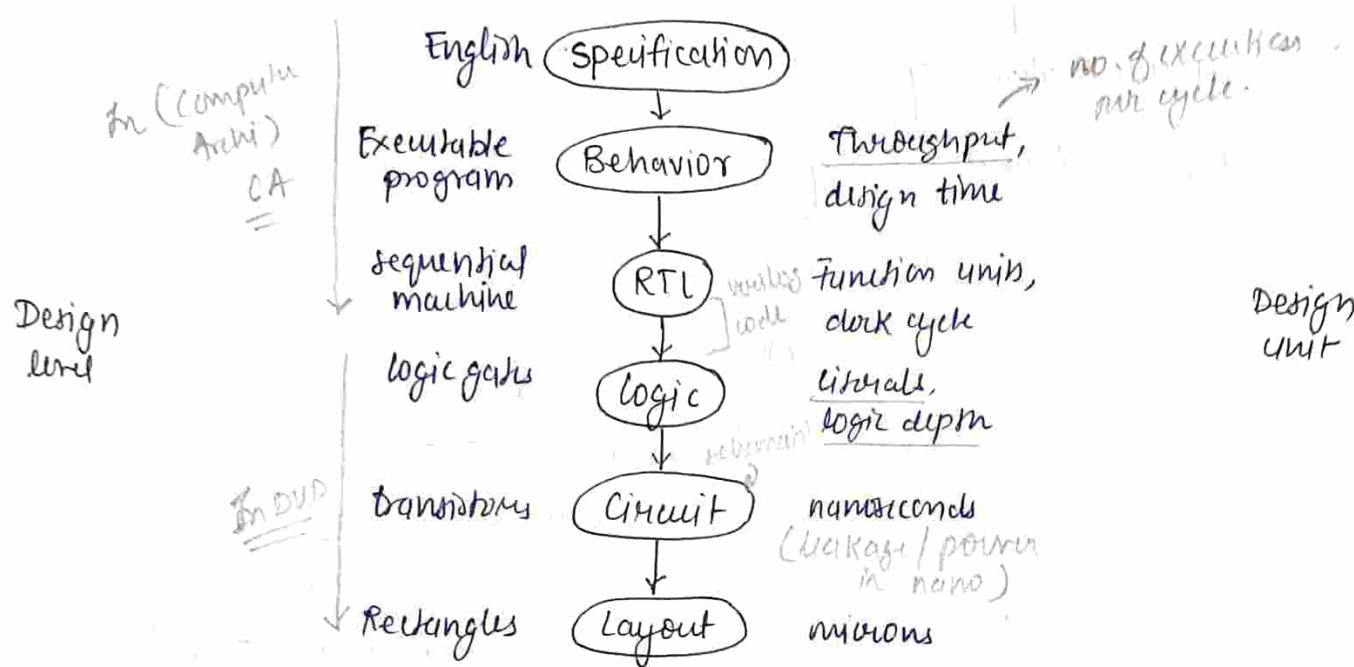
- Every polysilicon wire needs to be connected to source or drain because the charge due to plasma on the polysilicon must be sunk into the substrate.

Date
01-09-2022

class notes - 01/09/2022

• finer granularity

- STA → worst case analysis of clk wrt time
- Vector based analysis
- SDF → standard delayed format

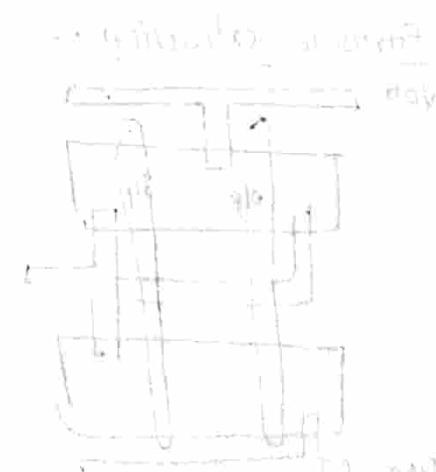
1. Design Abstraction

- .lib file
- zodiac
- High level synthesis (Behavioral) in C/C++ or verilog
then go to ~~RTL~~ RTL phase by doing high level synthesis tool.
- Layout ~~can~~ can be automatic but it may take more area
as compared to manual layout.
- Usually we do verification after every step of design abstraction.

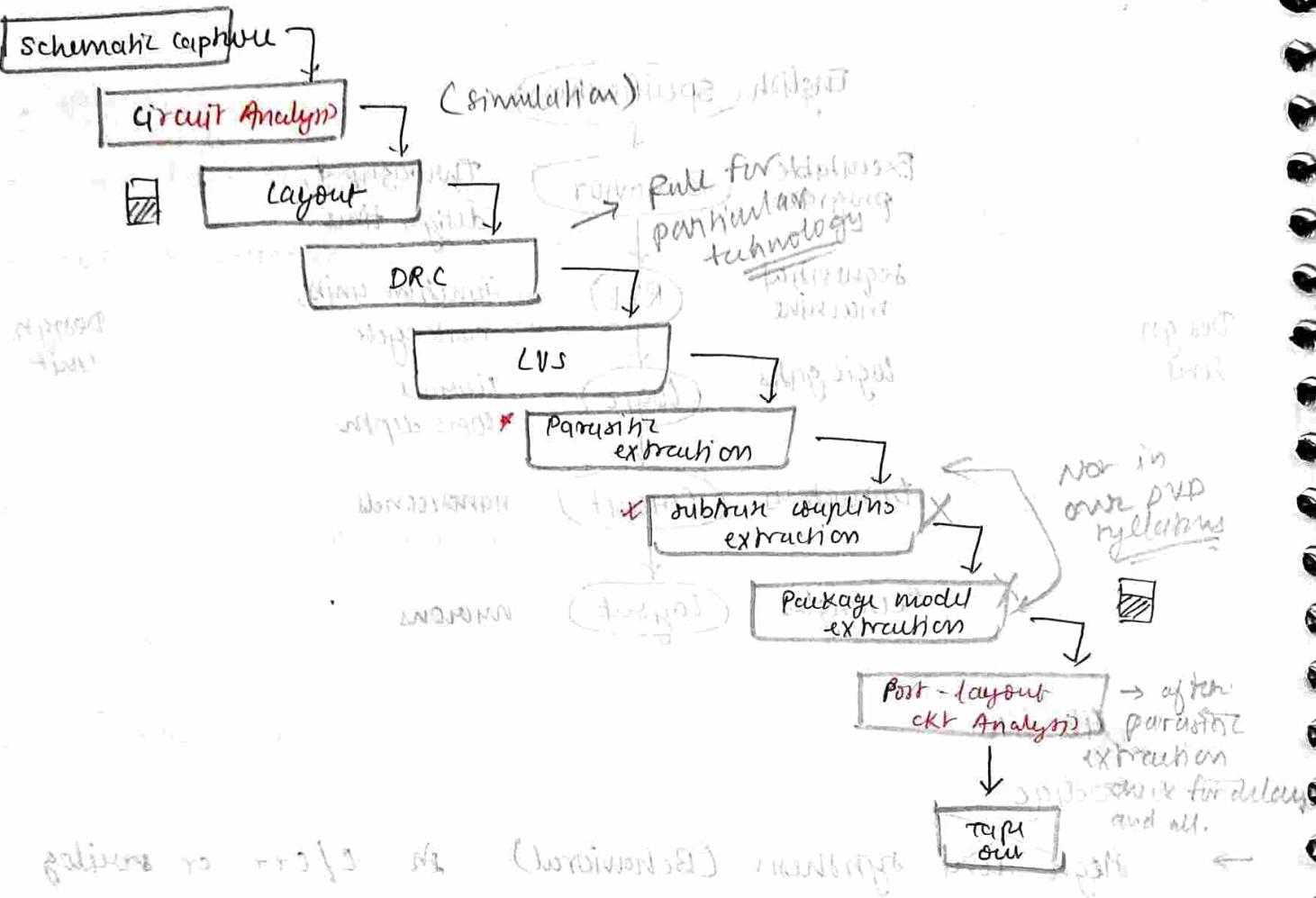
2. UVM Design Flow

UVM is designed to work top down.
The verification is done from bottom up.
Top down means you first define the
interface and then implement it.
Bottom up means you first implement
the interface and then verify it.

The interface should be defined in
such a way that it is testable, reusable
(compliant).

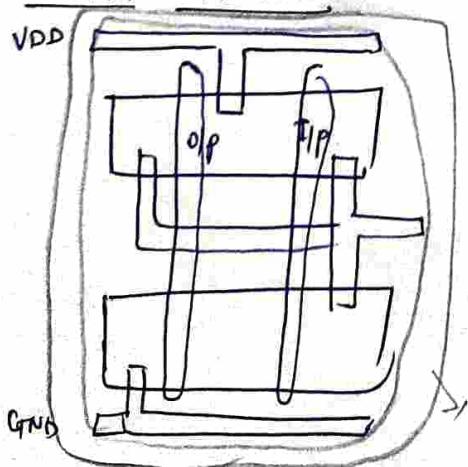


2. Library Design Flow :-



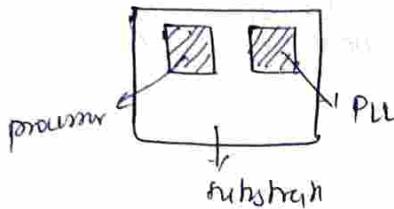
- Suppose we make a NAND gate layout then we have to maintain some minimum spacing b/w the wires, min area must be full filled & this conditions are DRC.
- LVS is used to check whether the ~~the~~ obtained layout matches with the designed schematic or not.

Parasitic extraction :-



- Metal will have some initial R & C value
- As we can see that output is connected to the metal hence some capacitive effect will be more (due to metal width) & I want to extract that (find the value of capacitance) capacitance & see its effect on my values
- We will do all the capacitance & resistance extraction. (means finding more guarded values)

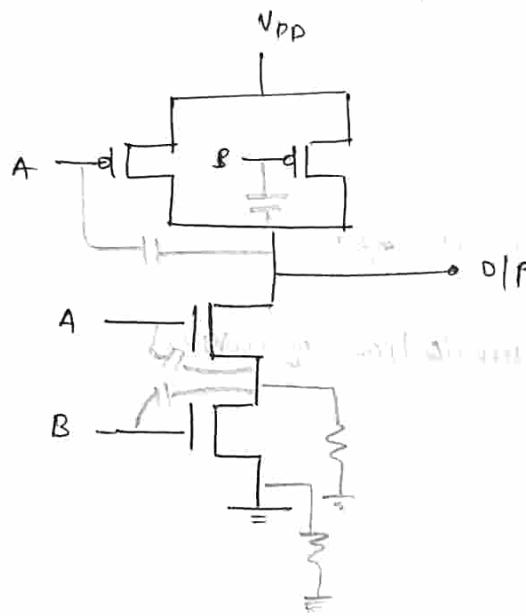
substrate coupling \rightarrow



Since processor & PLL are on same substrate, then the noise due to processor will effect the PLL. So there should not be any substrate noise.

- It will be like source to body capacitance.
- In DVD we will not observe this
- It can be removed by the ground ring.

If there is only RC in the ckt, it adds delay to the ckt.



If Post-layout simulation will have all the capacitance & then we simulate.

If we have to see whether the ~~R & C~~ present due to metal connection affecting our ckt or not.

- When you increase the size of your PMOS & NMOS then it can flow more current. Then this extra current can discharge the charges in RC in less time hence the delay will be reduced.

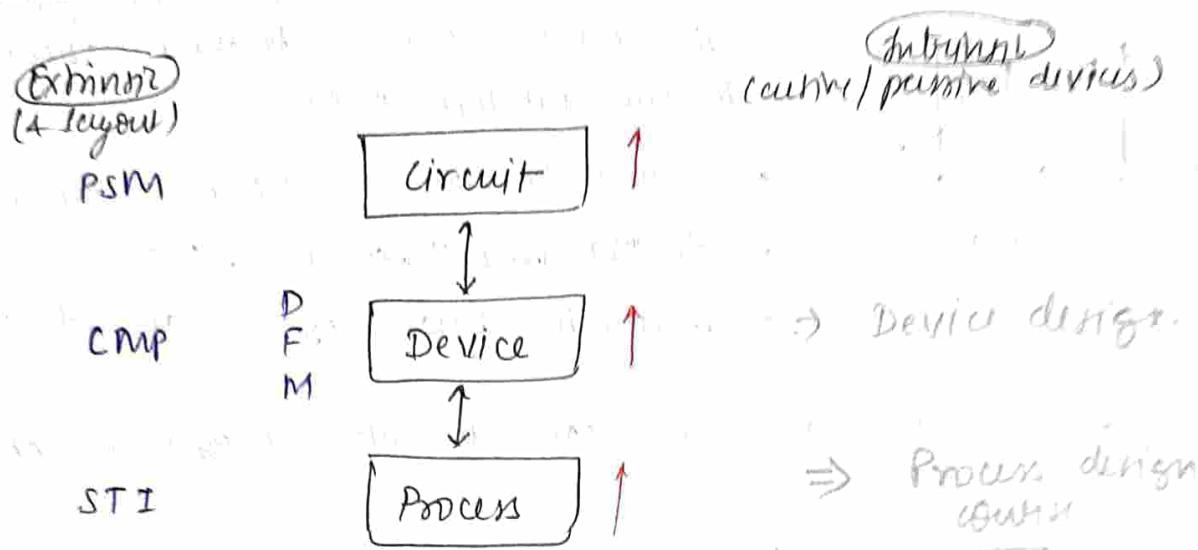
What is
the significance of doing
parasitic extraction?

\rightarrow By doing parasitic extraction in the metal wire we get to know about the $R \& C$ value & finally get the delay in the ckt now we can reduce the delay by scaling down $R \& C$ or we can decrease the delay by increasing the size of PMOS & NMOS.

The value of $R \& C$ depends on the length of the wire.

Inductance

3. Before ckt Analysis - Model Development



TCAD tools → Silvaco

Before ckt design - Device design

- silvaco → physics based simulation of devices
- library design
- SOC design
- IP design

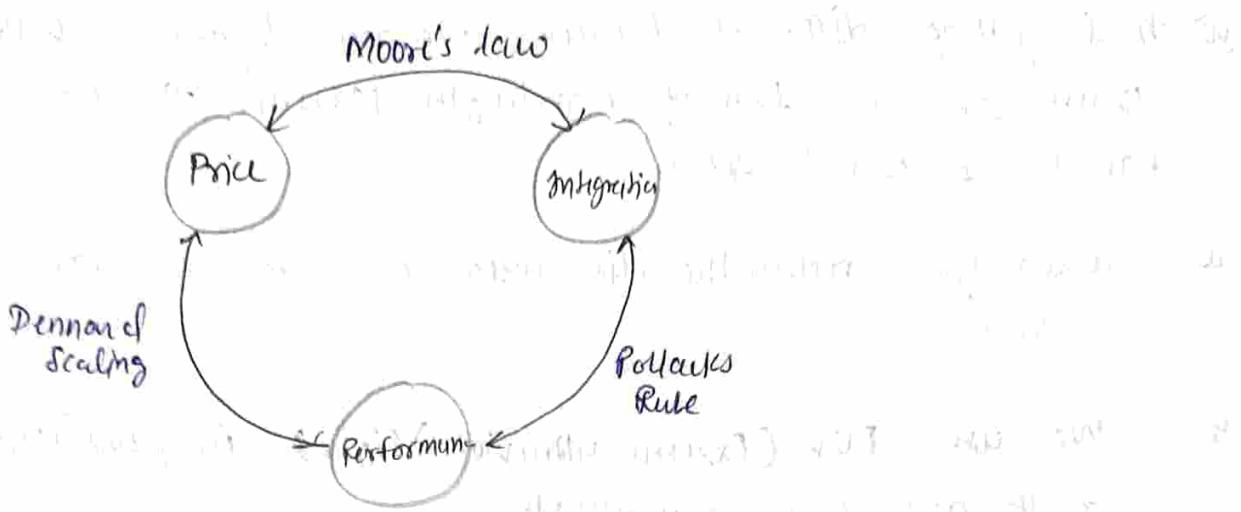
Key factors for Tremendous Growth of microelectronics

Moores' law : Demand scaling : Pollacks Rule

- Every 18 months the no. of ~~sister~~ units you can pack gets double in per cm^2 area.
- Every 18 months the area will be reduced to design the same ckt.
- So we can add more ckt on the same die. i.e. we can add more functionality.

Demand scaling : Area $\downarrow \Rightarrow$ price \uparrow

Pollacks Rule : Integration \Rightarrow performance \uparrow



{ Moore : time \uparrow , integration improves, area \downarrow

Dennard : area \downarrow , price \downarrow , performance \uparrow

Pollak's : ~~price \downarrow~~
integration improves, area \downarrow , performance \uparrow

- * In 1954 125 Rs for 1 call for 1 min (Eg)
- * Now the value of money is degraded in VLSI industry.
- * Moore law is just an observation.
- * Dennard gives the scaling rule for reducing the λ_{min} size.

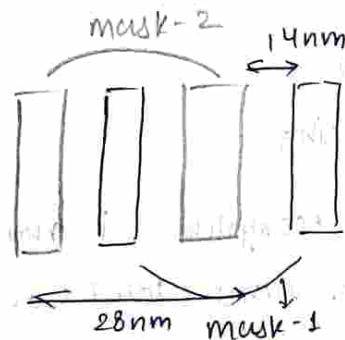
(*) till 20nm, price \rightarrow reducing
after that, price stops reducing.

- * In actual, the fabrication ecosystem is not able to handle finer geometries than 28 nm. Hence total no. of mask increases.
- * physical ^{upto} 45nm are designed only & after that like 10nm
it means not physical 10nm dimension but it some what 20-20% better than its previous version like 32nm.
 \Rightarrow physical λ_{min} is not changing, but somehow the performance of the CTR is improving by other things.
- * for finer geometries, the cost of making mask increased exponentially is from 180nm to 28nm

- It is getting difficult because we are trying to design 10nm using the light of wavelength 193nm. Therefore we have to do lots of optics.
- When you reduce the die area the yield of correct die increases.
- We use EUV (Extreme ultraviolet) light for photolithography as it has 13nm wavelength.
→ It will reduce the no. of mask you require.
- EUV is made by ASML company (only one)
wavelength (stock of ASML is getting higher + higher).
EUV > 100 W/hr
- ⇒ We are also adding some additional content in our chips which we can't scale. Like Analog CKT, Sensors in Automobile.

**

RET → Resolution enhancement techniques



⇒ no. of mask increases
price for chip ↑

↳ Dual patterning

2.5D design → image sensor + processor together

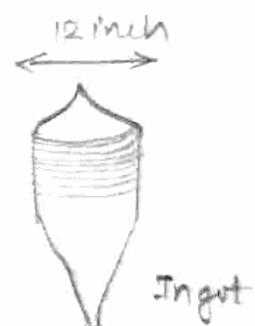
Crystal structure:

- Diff. planes of same material, say diagonal plane are dissimilar
 - diff. density of atoms
 - diff. topology of atoms
- Different EF on both planes
 - diff. electrical properties
- * Plane of silicon wafer decides its behaviour
 - has to be accounted for when wafers are being made.

Now such required plane can be fabricated?

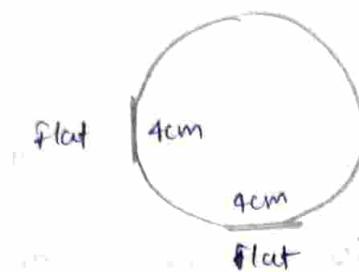
(i) Crystal from sand

- silicon crystal (seed) dipped into molten silicon.
- Gradually pulled out.
- molten silicon bonds to crystal
- solidifies according to crystal's orientation
- speed determines diameter of crystal - called Ingot.
- size of diameter can be: 6", 12", 8"

(ii) Ingot to Wafer

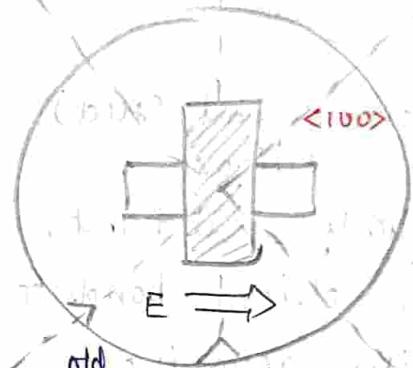
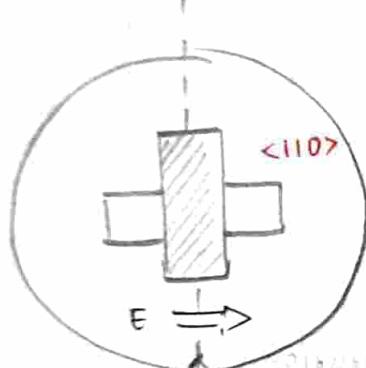
- Slicing - To get thin wafers for fabrication.
- Grinding
 - To make circular & of desired size.
- Lapping
 - To make of desired thickness.
- Polishing
 - To remove top layer that got dislocated.

- substrate doping
- Primary & secondary flat
 - P-type $<111>$ substrate is used - so no secondary flat.



n-type wafer - 2 flat regions
p-type wafer - 1 flat

(iii) Rotated substrate (wafer)



"standard substrate"

* Current is flowing in
 $<110>$ direction

"Rotated substrate"

* Current is flowing in
 $<100>$ direction

* After rotation, mobility of holes are improved.
i.e. to get the same current, we can have smaller device.

* hole mobility \rightarrow thousand significantly after $65nm$ tech.
 $\therefore p:n = 1:2$.

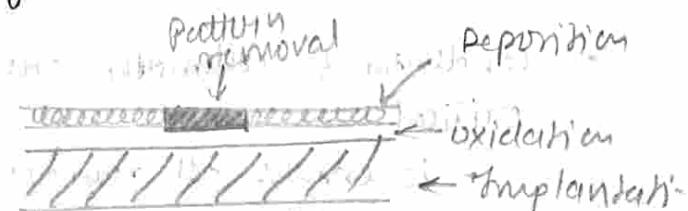
* e⁻ mobility remains same & you could save more area.

it Because of the change in effective mass of carrier, mobility is changing.

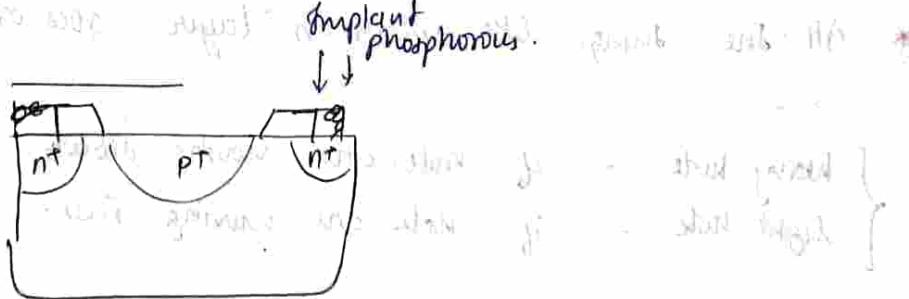
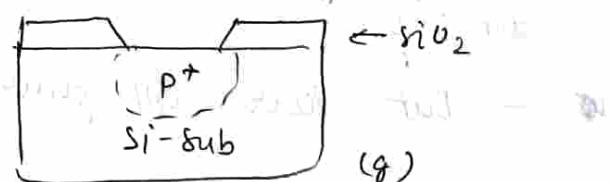
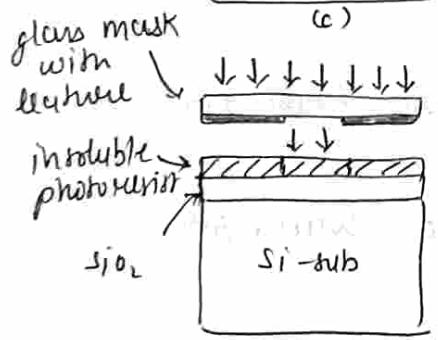
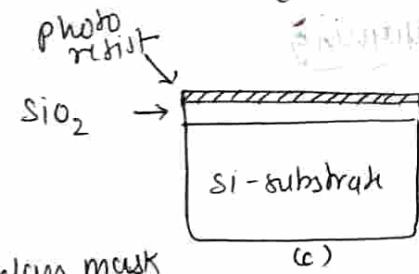
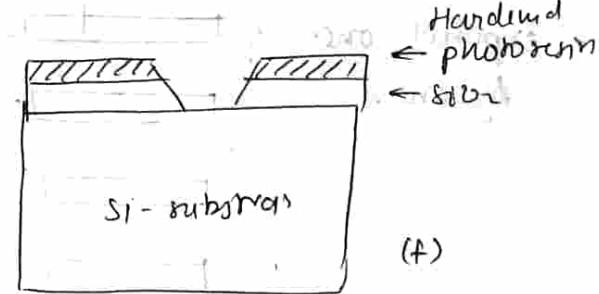
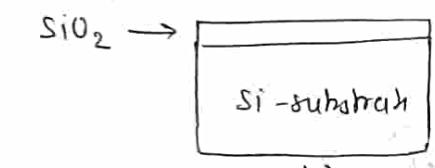
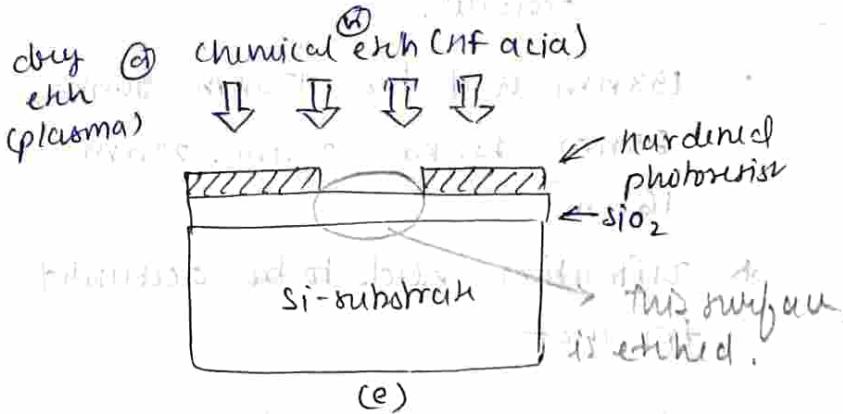
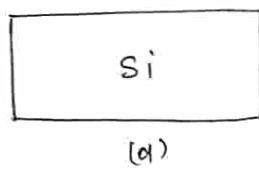
(iv) Process in Si Manufacturing

- Implantation
- Oxidation
- Deposition photo resist
- Patterned removal

Repeated operation of such steps
for diff. time durations results in
a functional chip:

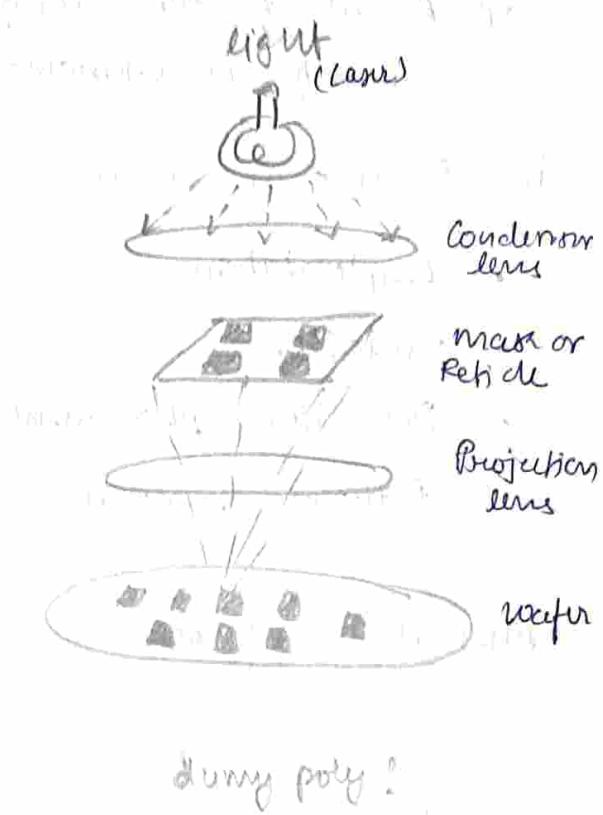


Photolithography ➡



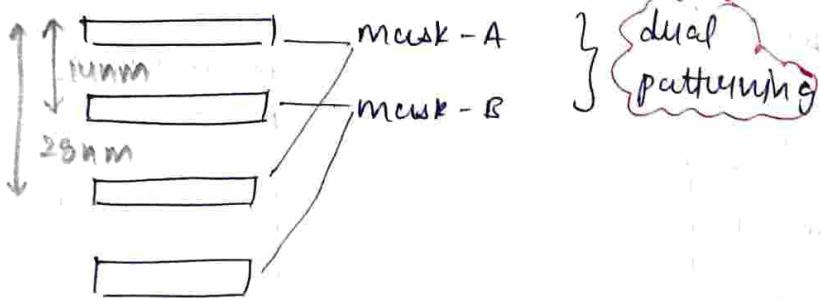
Photolithography Technique

- Pattern made on a mask
 - Positive mask ✓
 - Negative mask ✓
- Photomask exposed to light through the mask
- Condenser & collimator are critical
 - to ensure all rays are parallel.
- 193nm and for 180nm, 90nm, 65nm, 45nm, 32nm, 22nm, 16nm.
- * Diffraction needs to be accounted for that



dumy poly?

- expose one by one.



- * Tubelight generates a spectrum of light with large range of wavelength.
- * - But laser will generate light of single wavelength.

- * All the things like insulation layer goes on the ~~top~~ surface of the Si.

{ heavy hole - if holes are moving slowly
light hole - if holes are moving fast.

193 nm - the feasibility

- Diffraction for present range of feature sizes.
- Loss of accuracy & precision
- $\pm 1-5 \text{ nm}$ is easily in scope of error
- Optical proximity correction
 - to reduce ill-effects of diffraction
 - to diffraction to enhance resolution
- Change in medium of light to change effective wavelength
 - So it is 193 nm days
- Increasing lithography cost in every technology
 - increase focus on "first time right"
- EUV is necessary for price scaling at lower than 20 nm.
- + During the fabrication of a transistor, the actual metal contact is less ^{than} as it was supposed to be. Less area means same current will flow from the less area hence resistance will increase.

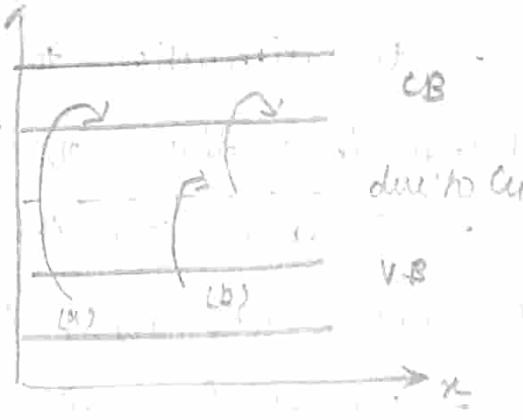
Process in Silicon Manufacturing

- Oxidation
 - Gate dielectric
 - Inter metal dielectric
 - Masking
- Deposition
 - Gate
 - Interconnects (metals)
- Patterned removal
 - Etching
- Implantation
 - Dopant to modify electrical properties of device.

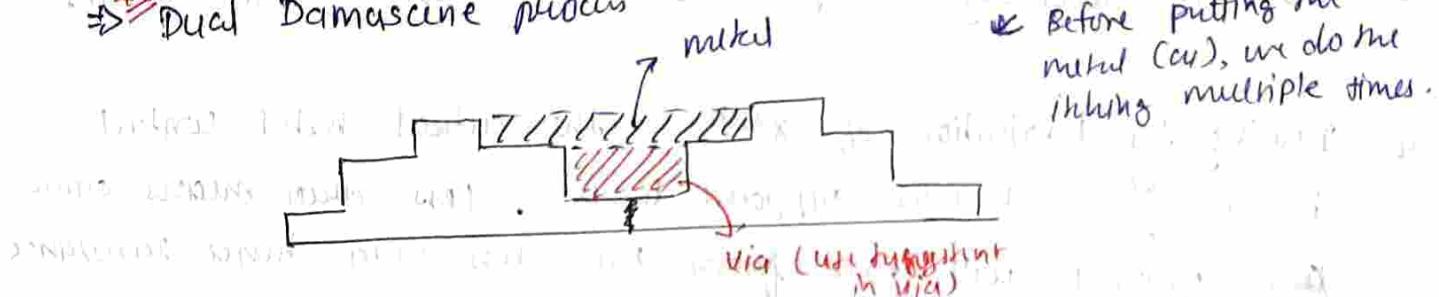
CMP \rightarrow chemico-mechanical polishing.

other - dry etching \rightarrow OIA
(conformal method)

- Copper is more ductile than Aluminium
- Copper has ~~less~~^{less} resistance
- Nowadays Cu is used in IC fabrication.
- earlier it was Al.
- When copper, gold & silver when they interact with Si, they make energy states deep inside the band gap
- But the constant ~~is~~ from silicon is ~~by~~ by tungsten because Cu can't come in contact with Si.

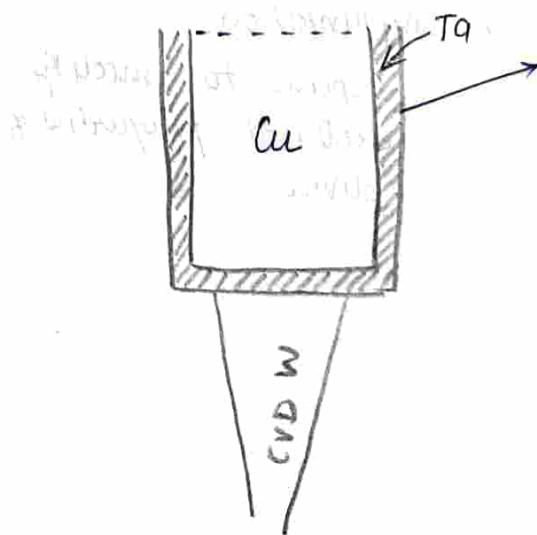


⇒ Dual Damascene process



Because of Cu we can now go to 10-12 layers of metal.

Copper wires



- . For 65 nm technology, the barrier width is say 5nm.
- . But for finer geometries ~~is~~ if we keep the barrier width same then the effective resistance ↑ significantly.
- . so for advanced technologies we are continuously looking for advanced barrier materials.

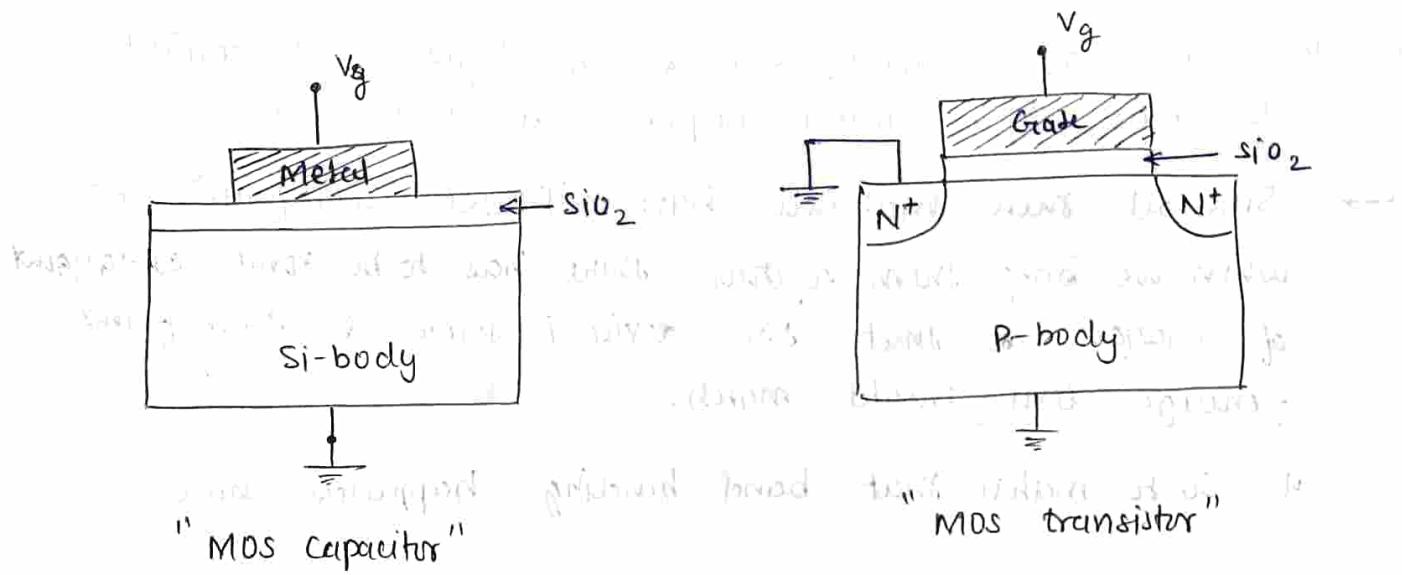
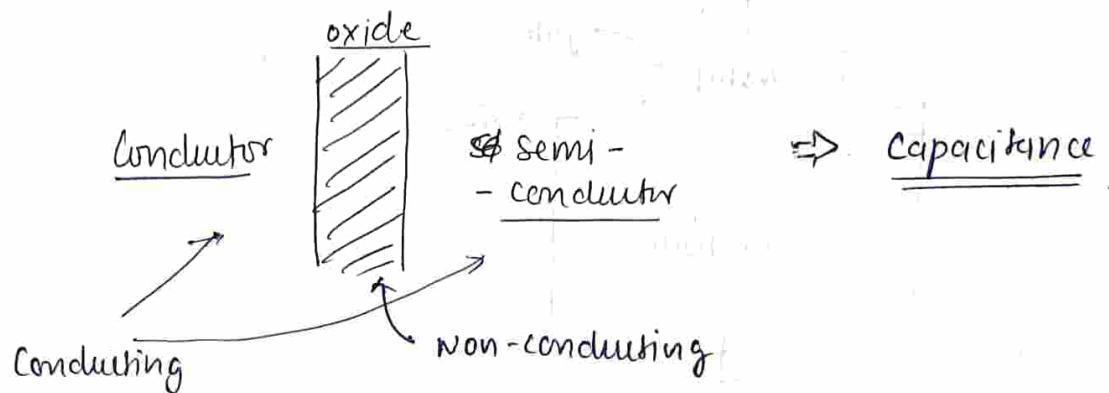
ALD → atomic layer deposition
(atom-by-atom)

Variability

- Device operating in range of nanometer.
- 65nm technology node
 - Gate oxide only 1.2nm (12\AA) thick
 - * ~~is~~ only 2-3 molecule of SiO_2
 - * Additional molecule leads to 25% change in performance.
 - 60nm gate length (CLP process)
 - * etching for few additional milliseconds may result in reduction in length by 5nm.
 - * change in performance (speed) by 8%
 - * change in leakage by more than 30%.
- 22 nm technology node?
 - Gate oxide replaced by high K-dielectric & metal gate (work-function management)
 - 18nm gate length (CLP process)
 - * sensitivity to variation increases further.
- clean Room
 - liquid cooling
 - self Aligned polysilicon.

2. VLSI Devices

Metal Oxide S/c (MOS)



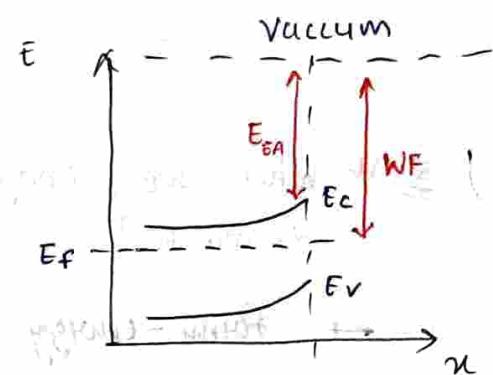
Electron affinity: It is the energy required by the e^- to escape from the material.

Electron affinity: Energy required by e^- to move from CB to vacuum.

- it doesn't change with doping
- it is material constant (it don't change)

Work funcⁿ: Energy required by e^- to move from from Si energy level to vacuum

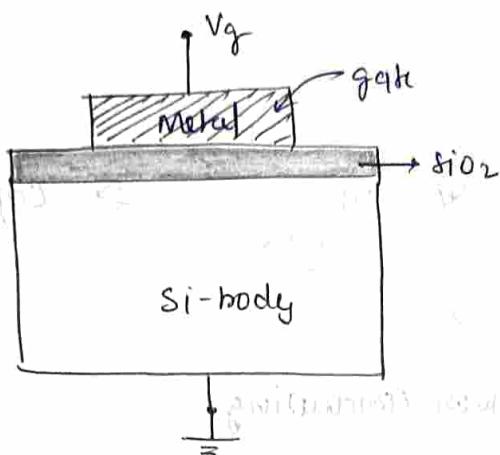
- it changes wrt doping
- it depends on surface termination (crystal face, surface chemistry).



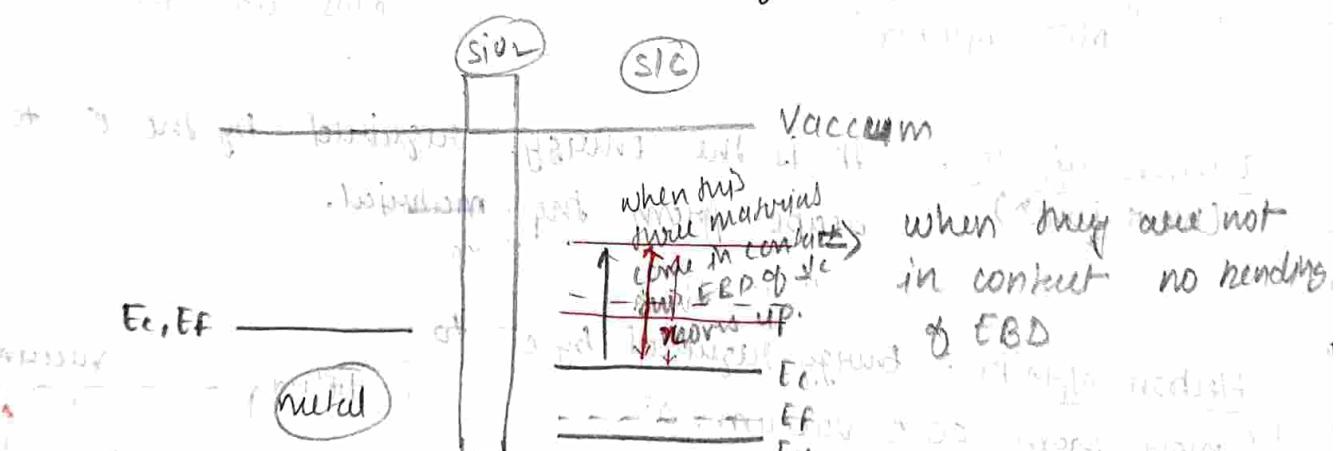
"energy band diagram of MOSCAP"

QUESTION 4

MOS Capacitors



- Q When we bring metal, oxide & SiC together in contact with each other, what happens at their interface?
- Since all these materials have different workfuncⁿ, so when we bring them together there has to be some arrangement of charges so that one device is stable & their fermi-energy level should match. (Why?)
- + So to match that band bending happens here.

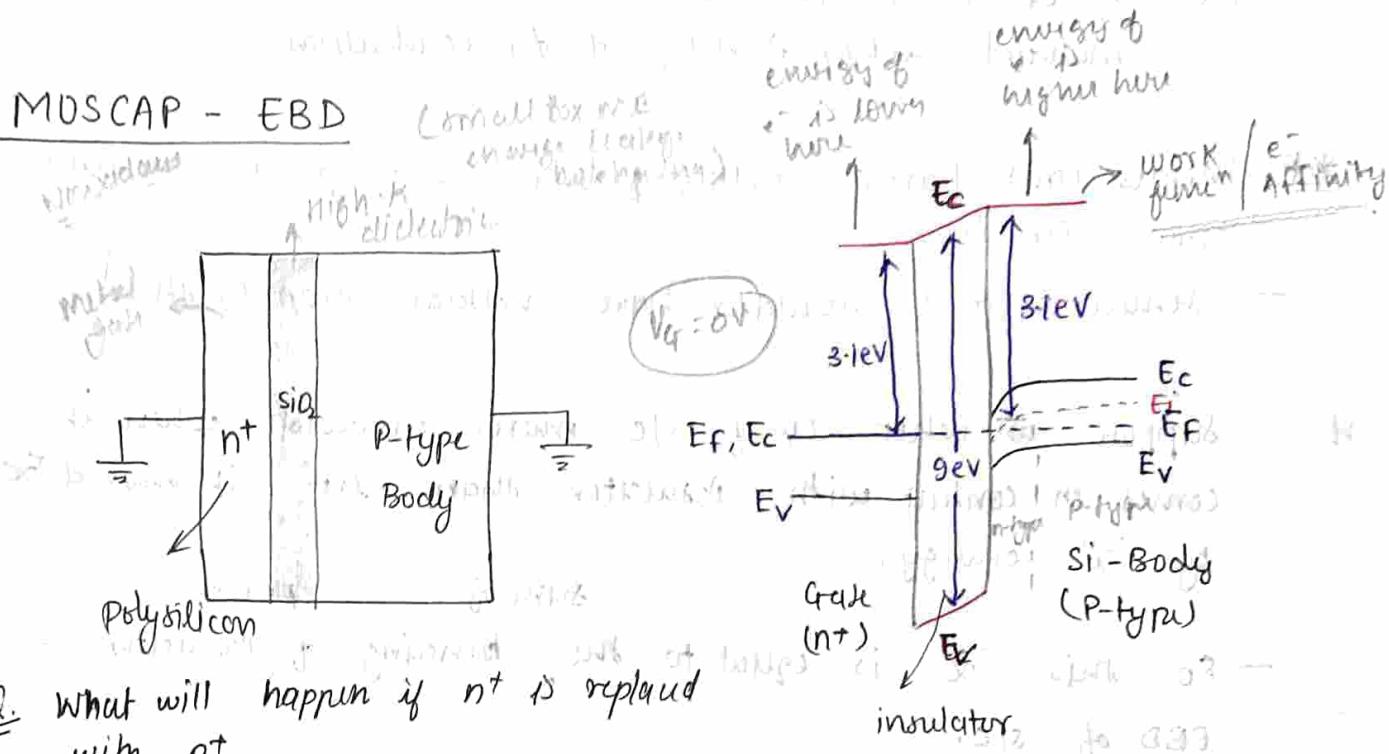


- Q When we bring this materials together then what needs to match?
- Fermi-energy level

When we talk about capacitance b/w two metals & capacitance b/w metal - oxide - s/c, then both are very different.

- because in this case, when diff. materials come in contact then fermi-level of materials needs to match & for that band-binding occurs.

MOSCAP - EBD



Q. What will happen if n⁺ is replaced with p⁺.

* n⁺ - poly-silicon : In this region, the amount of doping is very high so that in poly-region E_F & E_C both are at same level.

* on increasing the doping in n-type material the fermi level moves toward CB.

$$\uparrow E_{Fn} = E_C - \sqrt{KT} \ln\left(\frac{N_c}{N_D}\right) \rightarrow \text{position of } E_F \text{ for } n\text{-type s/c.}$$

no. of dopants in
n-type material.

⇒ Fermi-Energy level :- (E_F)

- It tell us the prob. of finding an e⁻ at a particular level.
- It tells that 50% of e⁻s are above E_F & 50% of e⁻s are below E_F.

at T > 0K

- At T=0K, e⁻s don't lie above E_F.

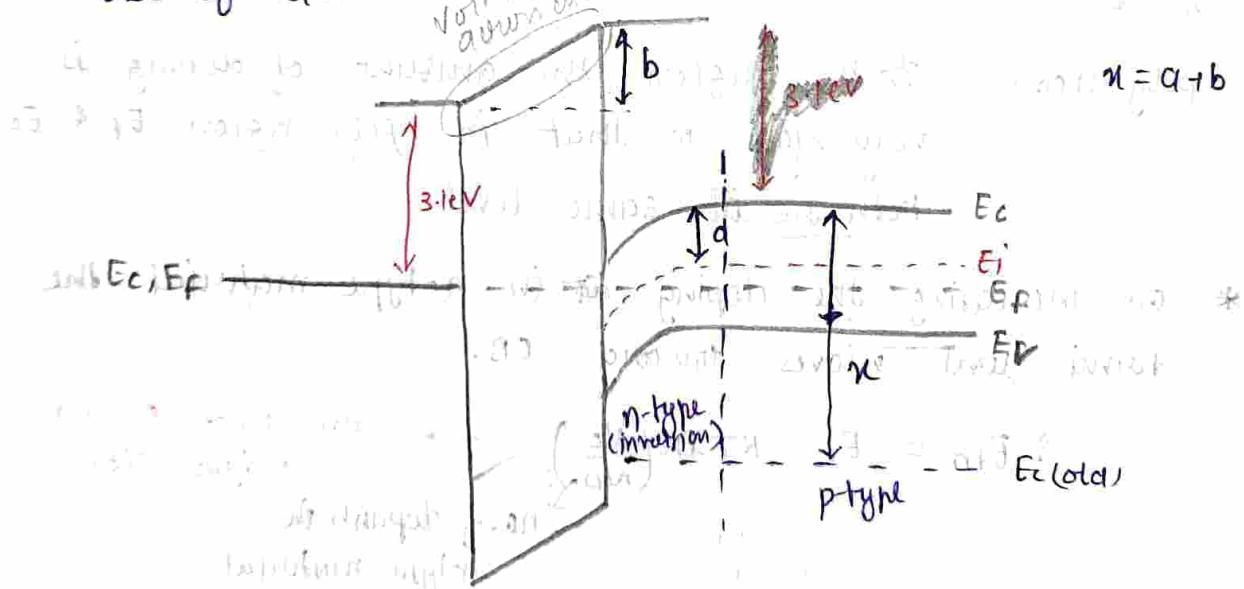
- Q: What does the fermi level in b/w VB & CB represents?
- (i) If EF near VB represents, it is a p-type material & prob. of finding an e^- in VB is higher.
- (ii) If EF is near CB, it means a n-type material & prob. of finding an e^- in CB is higher.
- (iii) If EF is in b/w VB & CB, then it means a intrinsic material, which is not good for conduction.

* Metals can't have a voltage drop across it as the current starts flowing.

- Semiconductor & insulator have voltage drop in them.

W: Suppose when CB of s/c moves upward when it comes in contact with insulator then let it move by χ energy.

- So this χ is equal to the binding of insulator + EBD of s/c.



W: Both the bending is occurred due to potential difference b/w the materials.

* It moved upward in such a way that $\text{workfunc} = 3.1 \text{ eV}$ (remains const) A EF matches with metal.

* Why both polyniticon + Body have 3-level of work func?

→ Because poly is also a kind of sic hence both have same work func & hence one can easily change the work func of the material.

→ ~~disadv~~ These days we use high-k dielectric material, whose workfunc can be easily changed by V_T of PMOS & NMOS can be varied independently.

* The slope of bending EBD in oxide region, deviates with the increase in thickness of oxide-layer.

[for thinner oxide-layer it will lead to tunneling.]

* If we replace dielectric material with high-k dielectric material then Band Bending will also deviate EBDing will change].

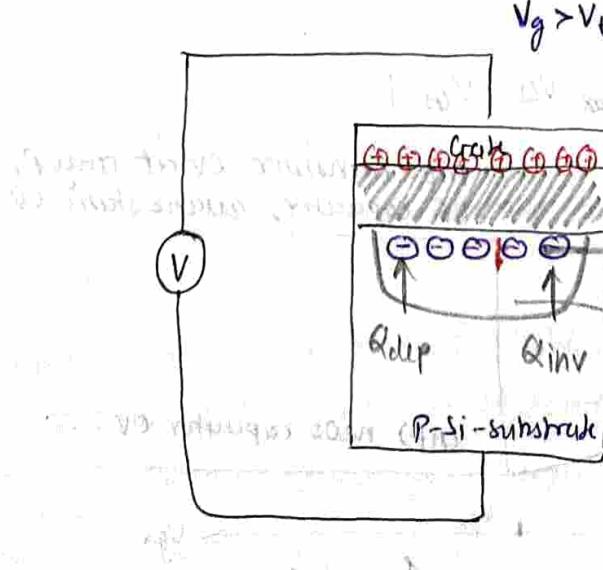
Bending is func of dielectric const.

Q. Why we use high-k dielectric? (21:30 - 22:30)

Q. What is the significance of work func if e⁻ can conduct from CB?

→ Band Bending determines the threshold voltage.

MOSCAP - Characterization :

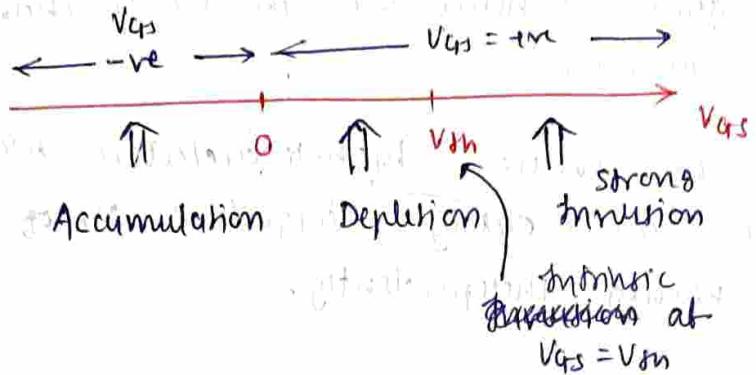


at low freq change
in charge happen
over here

- When we apply & vary voltage V, we observe the change in charge, capacitance & flow of current. (i.e. curr Vg)

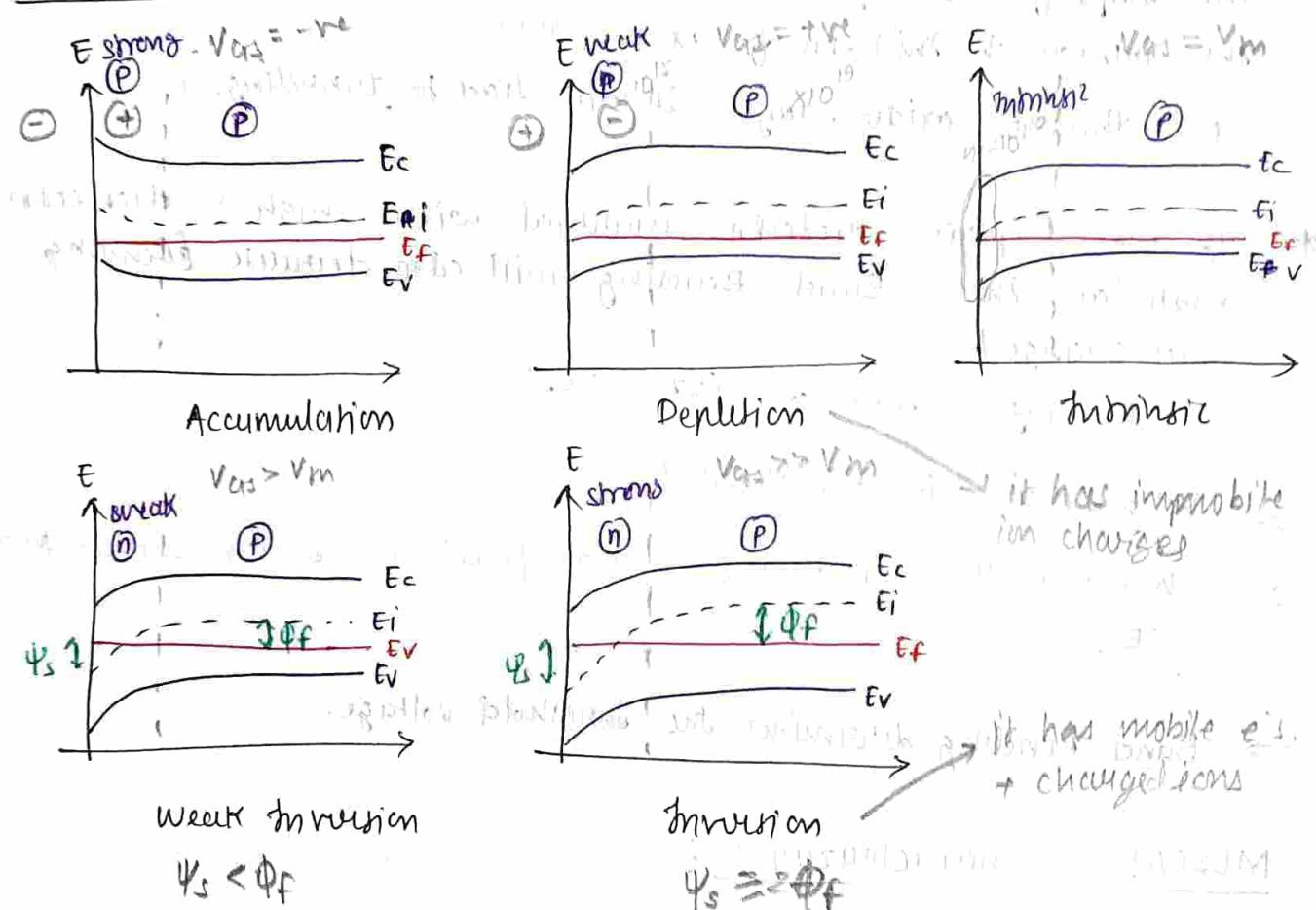
at high freq.
change in charge
happens over here.

MOSCAP - Region of Operation

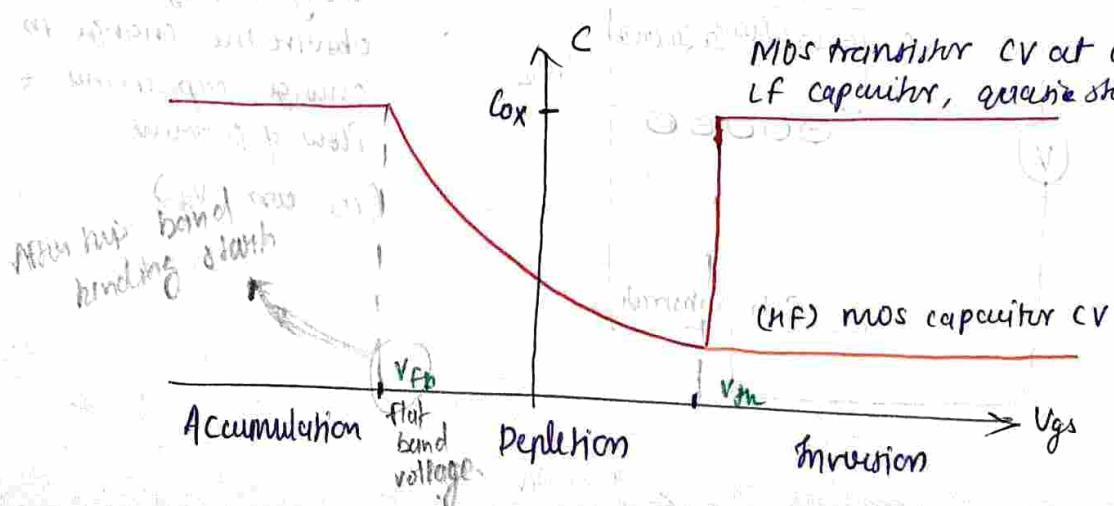


for NMOS

For NMOS



MOSCAP - Transistor for c/s [C_{ox} vs V_{GS}]



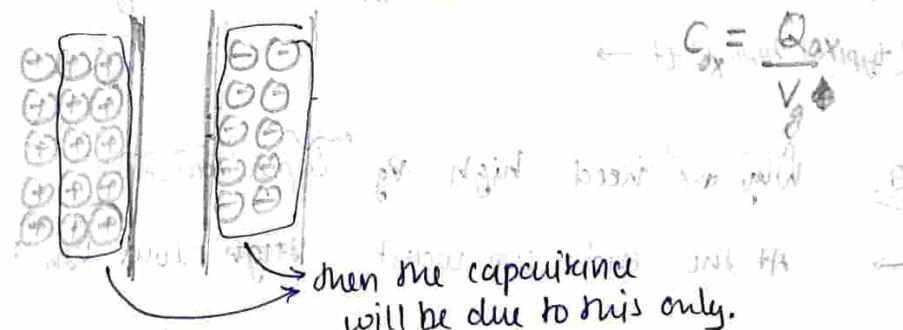
At low freq.: In accumulation region, capacitance is high because the change in charge is happening across the oxide only. (i.e. -ve charge) $V \uparrow$ (in-re), $Q \uparrow$ (in-re), $\Rightarrow C \rightarrow \text{constant}$

- As $V_{GS} \uparrow$, in depletion region the depletion layer starts to form hence net capacitance starts decreasing. (due to depletion region effective length of insulator is increased hence $C \downarrow$)
 - Depletion layer means, there is no available charges. Therefore any change in the charge would happen at the boundary of depletion layer. $V \uparrow$, $Q \uparrow$ (very small), $\Rightarrow C \uparrow$

$$Q = CV$$

$$C = \frac{Q}{V}$$

- When inversion layer forms, the capacitance \uparrow which reduces due to the addition of additional charges at the interface of silicon & SiO_2 . $V \uparrow$, $Q \uparrow$, $C \uparrow$ $\Rightarrow C \uparrow = \text{const}$ (because V is low voltage)
 - \Rightarrow Above condition (rectangular) is for low frequency.
- After inversion layer is formed, when we increase the V_{GS} further then value of capacitance remains same because on increasing V_{GS} the charge near gate is only increasing & now is no change in the charge amount near Si-SiO_2 interface.



- At high freq: At very high freq, inversion layer does not have any influence on capacitance because it can't change any charges in the inversion layer wrt to change in gate voltage. $V \uparrow$

- Means inversion is formed here, only depletion is here & depletion width can't go beyond a limit hence for high freq. also, capacitance value gets saturated.

$$\text{At high freq. } V = \text{const}$$

$$(capacitance) \propto \frac{1}{V} \rightarrow \propto$$

$$\text{constant}$$

* for MOSFET (CMOS Transistor), the value of capacitance will remains ~~same~~ except as that for the LF capacitance of MOSCAP.

* the band bending starts after applying applied flat Band voltage (V_{fb}). ~~so when the are applied~~

Q: Why do we use high ϵ_r dielectrics?

$$C = \frac{\epsilon_0 A}{d} \quad \text{with } C \propto \epsilon_r \quad (\text{C} \uparrow)$$

as $d = t_{ox} \downarrow$, "tunnelling occurs" \rightarrow problem.

\therefore we keep $t_{ox} \rightarrow \text{high}$ & $\epsilon_0 \rightarrow \text{high}$ ϵ_r value.

$C \uparrow$, coupling $\equiv Q (Ch)$ \uparrow

* if cap. \uparrow , then coupling of charges \uparrow , hence more charge on inversion will form for low value of gate voltage.

(slow-slow) SS \rightarrow $V_m \uparrow$ \rightarrow $Q = C V_m$

(fast-fast) FF \rightarrow low $V_m \rightarrow$ inv. layer.

(typical-typical) tt \rightarrow

Q: Why we need high ϵ_r capacitance?

$$\downarrow V_g = \frac{Q}{\epsilon_r} \uparrow$$

\rightarrow At the end we want ~~high~~ low V_m :

$$\text{and, } V_m = \phi_{ms} + 2\phi_f + \frac{Q_{dep}}{C_{ox}}$$

$$\phi_{ms} \rightarrow \phi_m - \phi_s$$

$$\phi_f = \frac{kT}{2} \ln \left(\frac{N_{sub}}{n_i} \right)$$

$$Q_{dep} = \sqrt{4\pi \epsilon_s (\phi_f) N_{sub}}$$

$$t_{ox} \rightarrow \text{SiO}_2 \text{ (oxide layer)} \Rightarrow \text{for } t_{ox} = 20\text{ nm} \quad C_{ox} = 17.25 \text{ pF}/\mu\text{m}^2$$

MOSFET

linear \Rightarrow
Region

$$V_{DS} \leq V_{GS} - V_m$$

$$I_D = K_n \frac{W}{L} \left[(V_{GS} - V_m) V_{DS} - \frac{V_{DS}^2}{2} \right]$$

where, $K_n = 4nC_{ox} = \frac{4n\epsilon_0}{t_{ox}}$ \rightarrow process transconductive parameter

Saturation \Rightarrow
Mode

$$I_D = \frac{K_n}{2} \frac{W}{L} (V_{GS} - V_m)^2 (1 + \lambda V_{DS})$$

\uparrow channel length modulation.

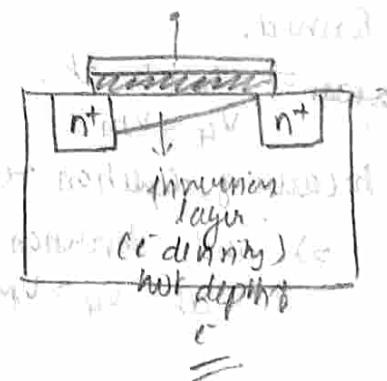
Q. For $V_{DS} = V_{GS} - V_m$, what will be the region of operation?
 → linear

Q. In saturation region, ~~what~~ will the charges in inversion layer also stops inverting as we increase V_{DS} at const. V_{GS} .

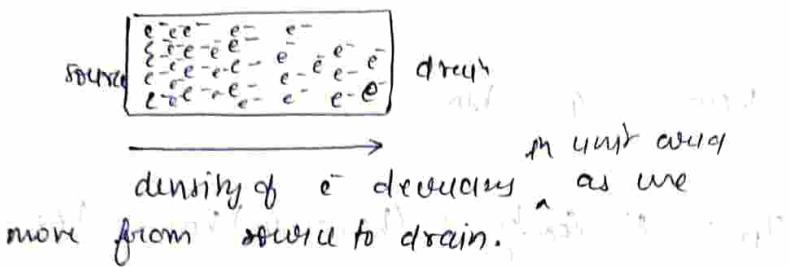
Q. For $V_{DS} > V_m$, will the charge in inversion layer change or remains same.

Q. What is the ~~depth~~ depth of inversion layer?
 → inversion is always at the surface, there is no physical depth in inversion layer

Q. What is the prob. of finding an e^- in inversion layer?



Q. How to visualize the pinch-off in the sheet of paper (contd.)?



Q. Why won't I_D become zero after pinch off as there is no e^- near the drain region?

* When you increase V_{DS} , the depletion region b/w n^+ -source/drain and substrate will increase as it forms a diode b/w them called a parasitic.

- But the $W_{d,max}$ for depletion region at the ~~top~~ surface of SiO_2 -substrate will not increase.

Inversion layer
Depletion region b/w SiO_2 -substrate, max. \Rightarrow no width/depth
Depletion region b/w n^+ -substrate, its width ↑
am I right?

"Depletion Region"
"Inversion Layer"

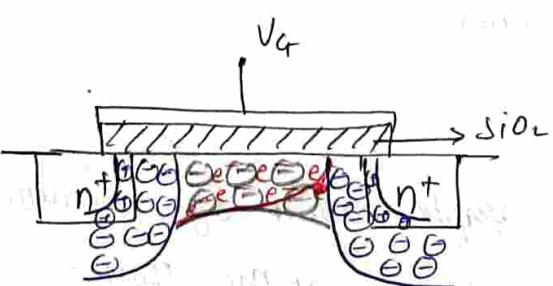
Velocity saturation

As $V_G \uparrow$, depletion region b/w SiO_2 -substrate is also formed.

~~∴~~ $I_D \neq 0$ at $V_G < V_m$ (threshold conduction)

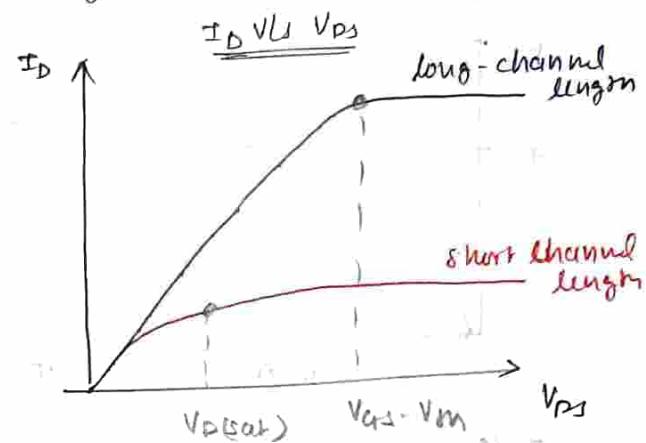
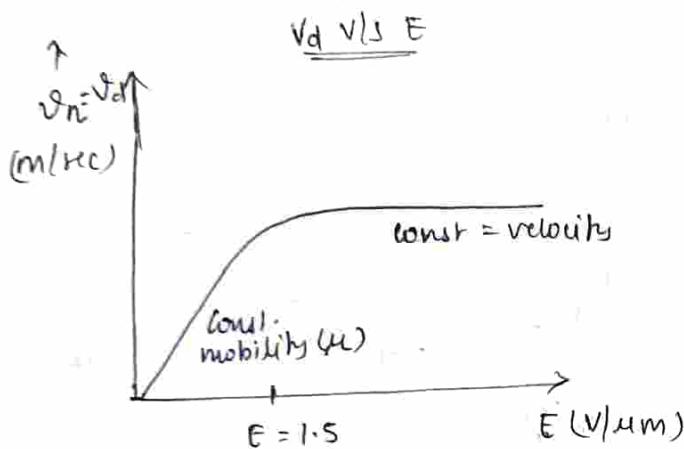
because of depletion region b/w n^+ -p-type

\Rightarrow and inversion layer is formed at $V_G > V_m$.



Velocity Saturation

- As $V_{DS} \uparrow$, velocity of $e^- \uparrow$.
★ and pinch-off occurs as we can't observe the e^- (as they are moving very fast at high V_{DS}) near the drain region.
- But after some V_{DS} , the v_n starts saturating because of the collision of e^- as they are moving at high speed.



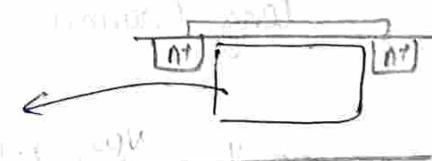
it for long-channel \rightarrow current saturates at $V_{DS} = V_{DS(sat)}$
or for short-channel \rightarrow current saturates at $V_{DS} = V_{DS(sat)}$

$$V_d = -\frac{e}{m^*} E \quad , \quad \mu = \frac{-e}{m^*}$$

$$\mu = \frac{V_d}{E}$$

Q: What is the short-channel exactly?

\rightarrow when n^{+} region no longer remains as rectangular as it has to be then we start to observe short channel effect.



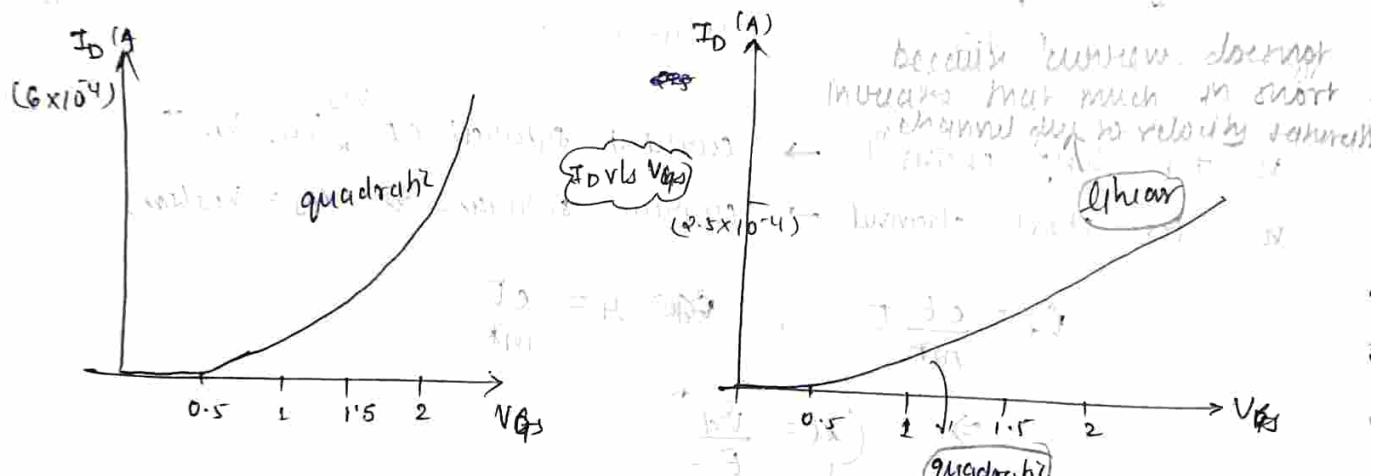
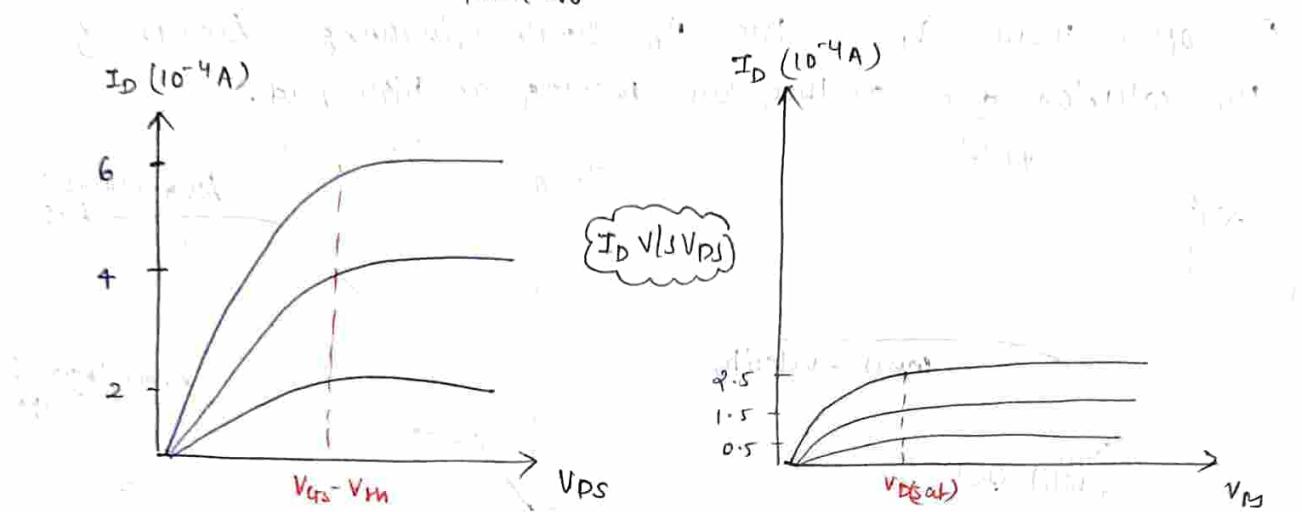
when depth of depletion region is much compared to the length of channel then we start observing short-channel effect.

* for long-channel you can dope the substrate more & the depletion width will remain the same (How?) it will not induce halo implant
- But for short-channel we can't dope the substrate.

halo implant
super halo implant

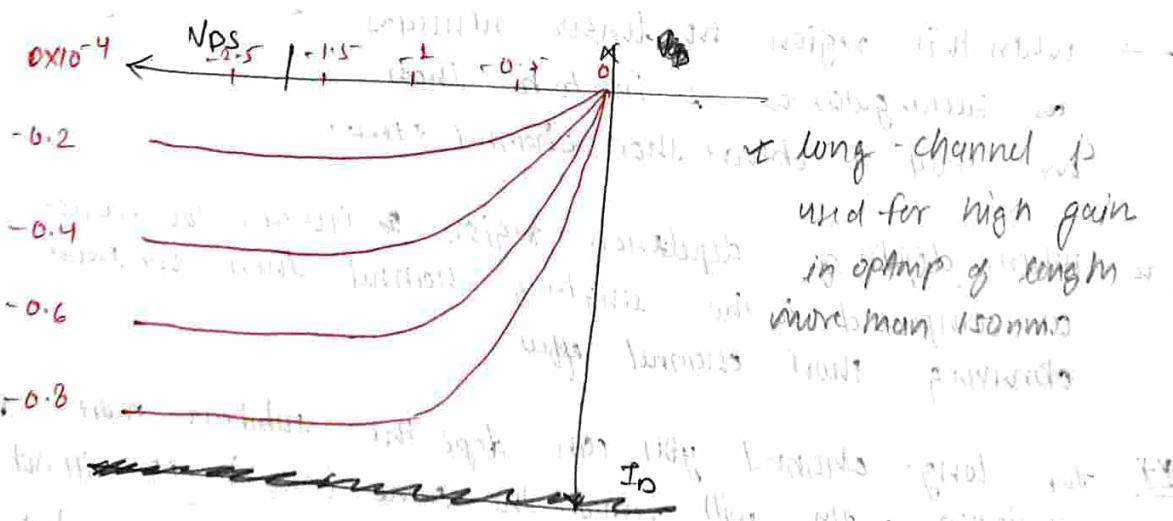
for short-channel you will reach the saturation at lower value of ext EF.

while talking about short-channel we only discuss $V_D(sat)$ & not V_m phn off



"long Channel"

"Short channel"

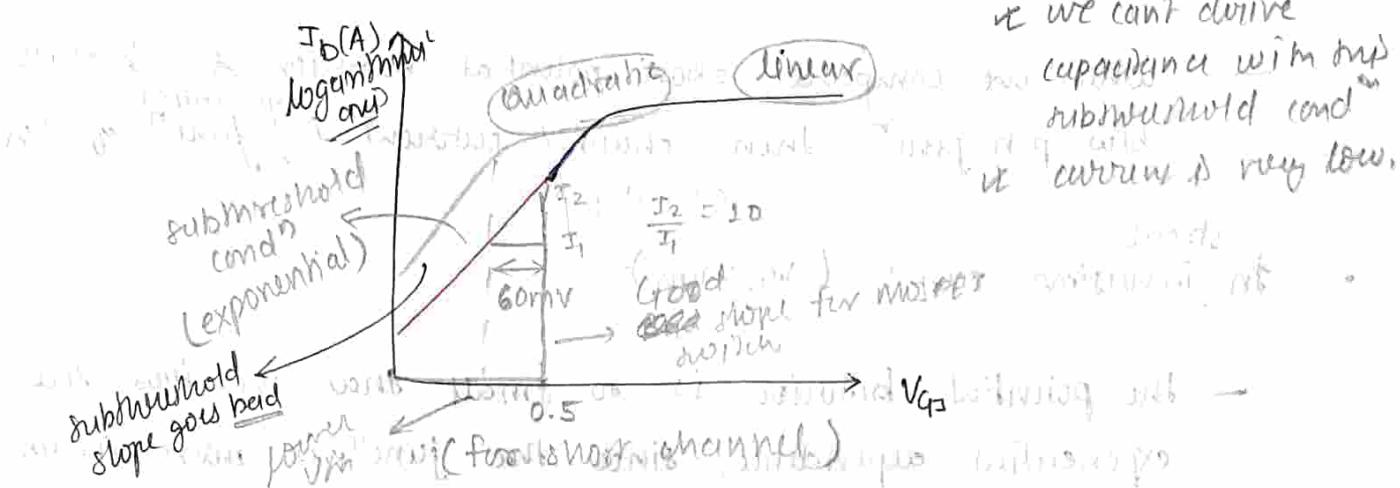


Charging time constant $\tau = C_{GS} \cdot R_{DS(on)}$

for JNP for JNP Subthreshold Conduction (doubt)

Neamen Pg - 444

- Current from source to drain when $V_{GS} < V_m$.



- In reality, $V_{GS} \leq V_m$, $I_D \neq 0$, but I_D is finite & I_D depends on V_{GS} exponentially as $I_D \propto e^{\frac{V_{GS}}{nKT}}$ due to non-ideality factor.

- V_{GS} must decrease by approx 80mV for one decade change in I_D .

$$\text{Eg. } I_{D1} = 1 \mu\text{A} \text{ at } V_{GS} = V_m$$

$$\therefore I_{D2} = 18 \text{ mA} \text{ at } V_{GS} = 0 \text{ V}$$

- due to subthreshold condⁿ significant amount of power dissipation.

Subthreshold slope \rightarrow relationship between V_{GS} and I_D in subthreshold region.

s is ΔV_{GS} for $I_{D2}/I_{D1} = 10$ i.e. increase in I_D by 10 times.

$$s = n \left(\frac{kT}{q} \right) \ln(10)$$

} Typical value for s :
 $60 \dots 100 \text{ mV/decade}$
 $V_0 = 0 \text{ V}$

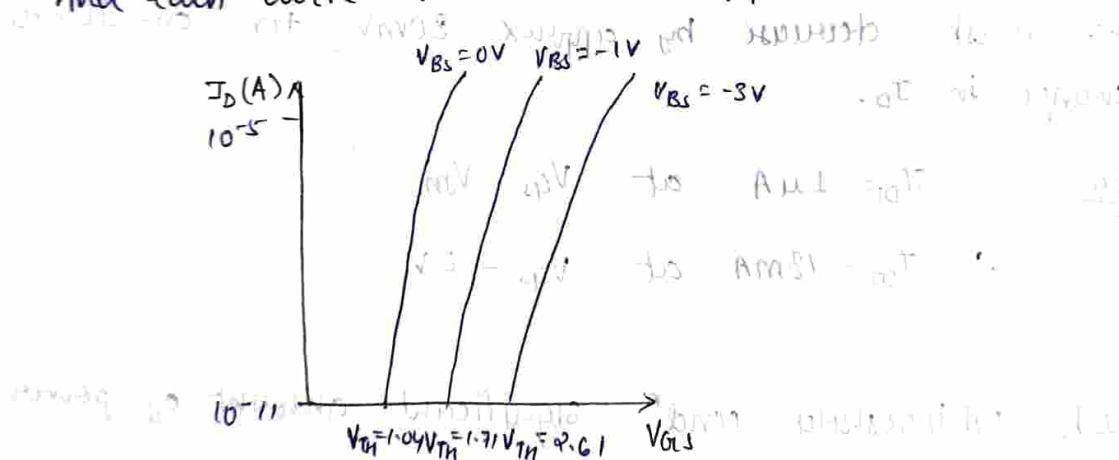
$V_{GS} = V_0 + s \ln(I/I_0)$ for subthreshold condition.

$$V_0 = \frac{kT}{q} \ln(10)$$

$$= \frac{1.38 \times 10^{-23} \times 290}{1.6 \times 10^{-19}} \text{ V} = 25 \text{ mV}$$

- Accumulation & weak inversion case: ($V_{GS} < V_m$)
 - To generate current, e^- must overcome the potential barrier b/w n+ source & channel.
 - when we compare above potential barrier & barrier b/w p-n junction, then channel current \propto ^{exponential} "junc" of V_{GS} .
- ^{strong} In inversion case: ($V_{GS} > V_m$)
 - the potential barrier is so small that we loose the exponential dependence, since the junction is more like an ohmic contact.
- Below graph shows the exponential behavior of the subthreshold current of serial's body-to-source voltages.

- And each curve has different V_m values.



* Ideally, a change in gate voltage (V_g) of 60mV produces a ~~and~~ change in drain current (I_D) as $I_{D2} = 10^{0.3} I_D$

E.g. $V_m = 0.3$ (for low voltage operation)
 $V_{GS} = 0V$

Drain current decreases by a factor of $10^{0.3/60\text{mV}}$
 $\approx 10^{3.75}$
 $= 5.62 \times 10^3$.

- When we bring "m-o-si" together, to make E_F equal, a potential drop is created due to which bending of wave occurs.

* All the behavior is at the surface of s/c.

- If you deviate the device size, then applied voltage must be divided. (e.g. n -poly silicon)
- Hence, type of resistance pattern is anti-symmetrical

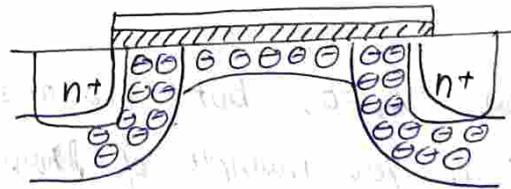
Eg. 2 for $I_{D1} = 1\text{mA}$, $V_{G1} = V_m$

and for $I_{D2} = 18\text{nA}$ when $V_{G2} = 0\text{V}$, (off)

Note: If mos device conducts at $V_{G1} < V_m$, then how we define threshold voltage (V_m)?

- On logarithmic scale, the weak & strong inversion c/s and consider their inversion voltages as the threshold volt.

for $V_{G1} < V_m$



Q. ** What happens when subthreshold slope goes bad (minimum have lower subthreshold current). What should be done to avoid this?

→ We want subthreshold slope such that within 60mV you should be able to deviate current by one decade (@ $70\text{...}100\text{mV}$)

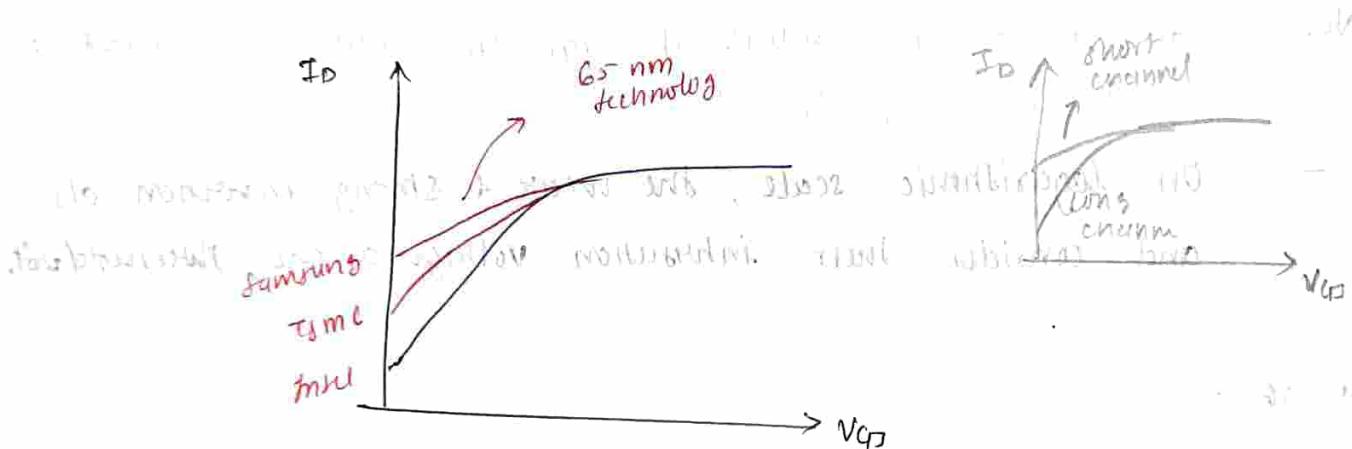
$$\text{i.e., } \Delta V_{GS} = 60\text{mV} \Rightarrow \frac{I_2}{I_1} = 10$$

- Q. MOSFET is off at $V_{GS} < V_m$. But because of subthreshold cond'n how will it be in off state.
 ✓ sub-threshold \rightarrow used for high gain in Analog CMOS
 ✓ Normal cond'n \rightarrow for switching.

Note: \rightarrow leakage increases as you decrease the V_T
 \rightarrow leakage also increases when technology is bad & $V_m = \text{same}$.
 (short channel effect, punch through)

* A technology now also called subthreshold slope, for diff. technology V_m will change.

* When you look at device, you don't only look at the V_{TH} but also at the subthreshold leakage.



- leakage is '0' when $V_{GS} = 0$, but increases to its full value when V_{GS} is a few multiple of thermal voltage V_T (e.g. $V_{GS} > 50\text{mV}$)

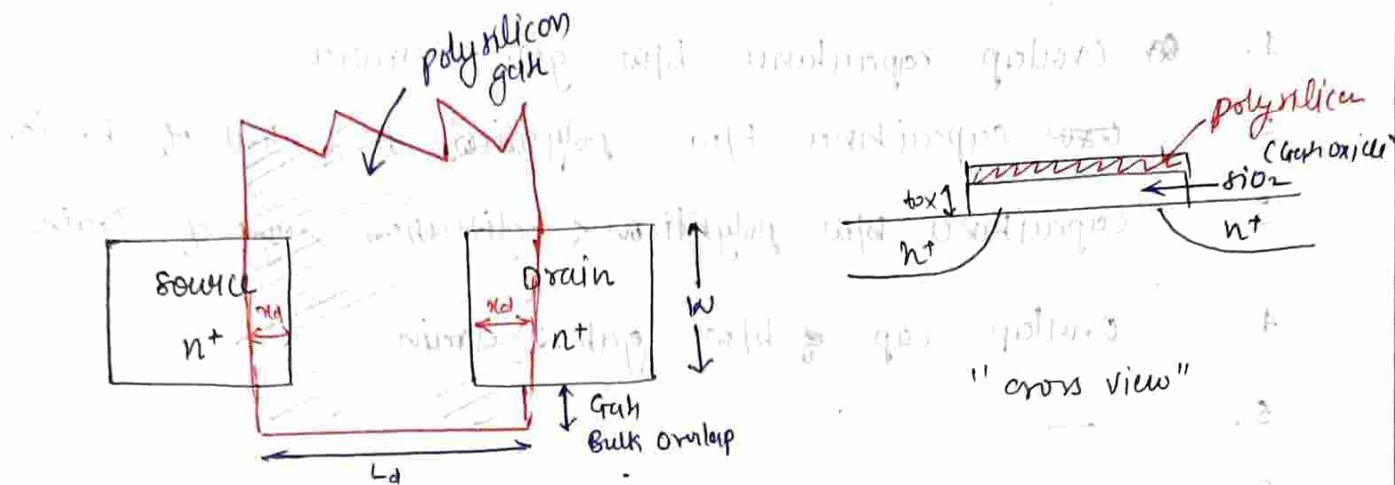
Q. what is subthreshold cond'n current dependence on V_{GS} ?

\rightarrow Good slope \rightarrow sharp a/s (sharp slope of I_D vs V_{GS}) before V_m .
 Bad slope \rightarrow exponential & graph of I_D vs V_{GS}

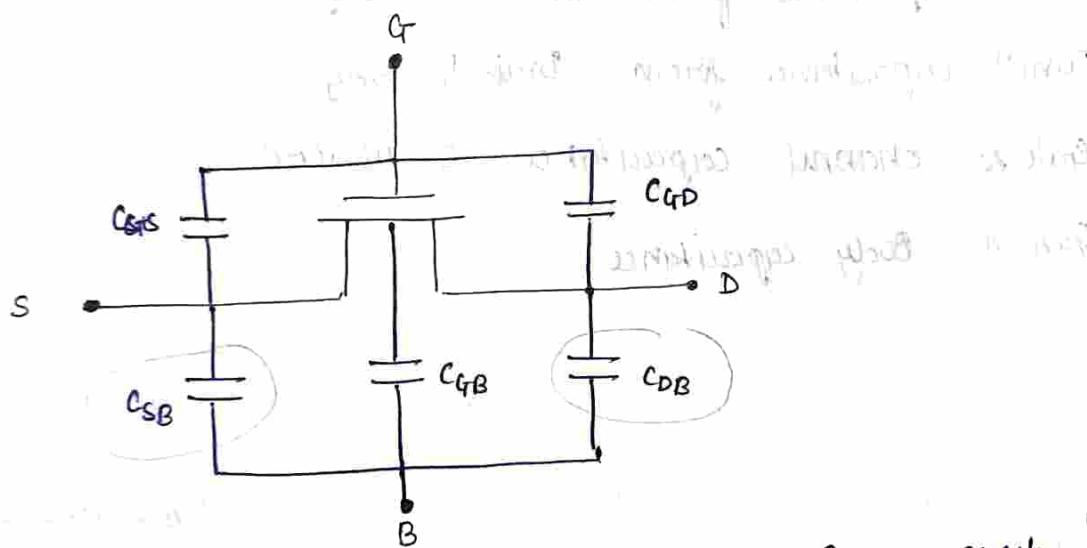
lect-7

MOSFET Capacitance

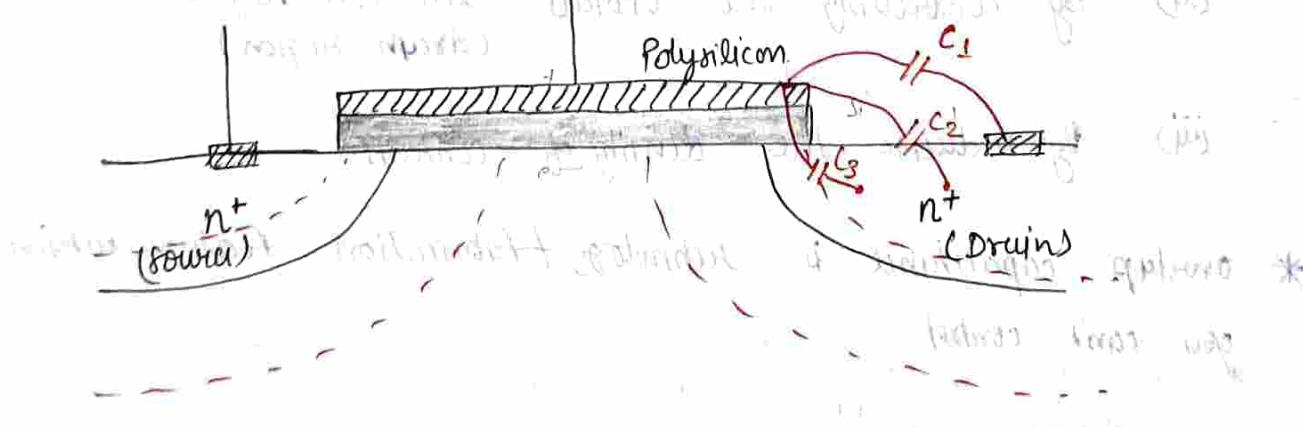
Rubacy Pg. 109



"Top view"



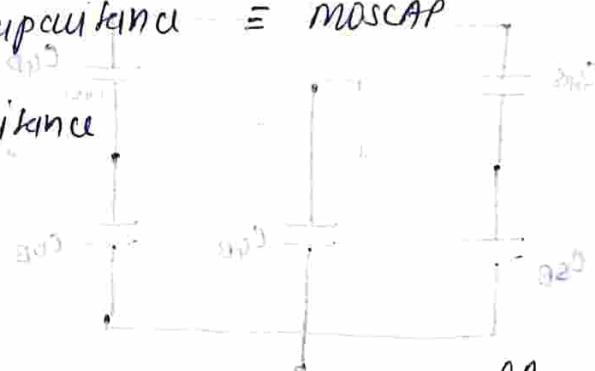
(3) $C_{GD} \rightarrow$ overlap capacitance of $G\&D$,
 (1) + poly to contact coupling
 (2) + poly to diffusion layer coupling.



Body.

Types of capacitance available are:-

1. ~~overlap~~ capacitance b/w gate + source
2. ~~gate~~ capacitance b/w poly/silicon + contact of drain.
3. capacitance b/w polysilicon + N⁺ diffusion layer of Drain
4. overlap cap. of b/w gate & drain
5. —
6. —
7. Tunⁿ capacitance from source to Body C_{SB}
8. Tunⁿ capacitance from Drain to Body C_{DB}
9. Gate to channel capacitance = MOSCAP
10. Gate to Body capacitance



(As a designer)

Q. How do we reduce gate to drain cap. C_{GD} ?

-
- (i) By controlling the spacing b/w contact + poly. (They are placed in P&R)
 - (ii) By controlling the overall diffusion region (drain region)
 - (iii) By reducing the density of contacts.

* overlap capacitance is technology / fabrication feature, which you can't control

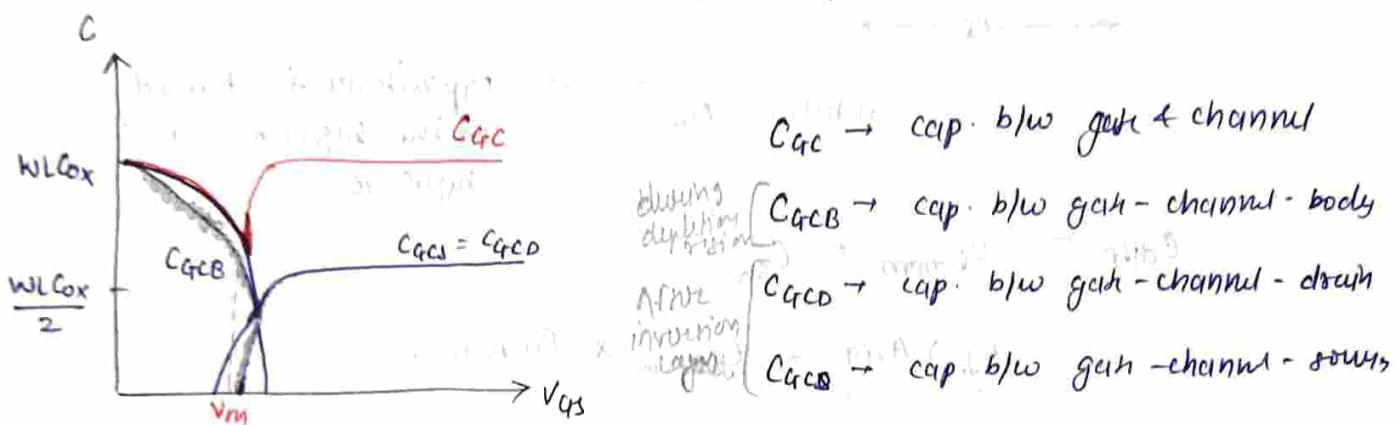
Q. Why do we have multiple contacts over drain region?

→ To reduce resistance & increase current of drain junction which is high voltage diode.

MOSFET - Gate capacitance

$$C_{gate} = \frac{\epsilon_{ox}}{t_{ox}} WL$$

$$\epsilon_{ox} = \frac{\epsilon_0}{t_{ox}}$$



- During depletion region, width of depletion region ↑ as we ↑ V_{GS} hence C_{GCB} ↓.
- During inversion region, a capacitance is formed b/w G-C-S & G-C-D because of inversion layer formation.

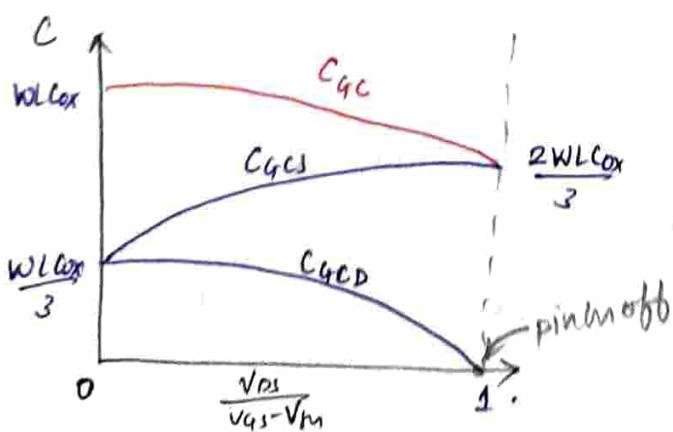
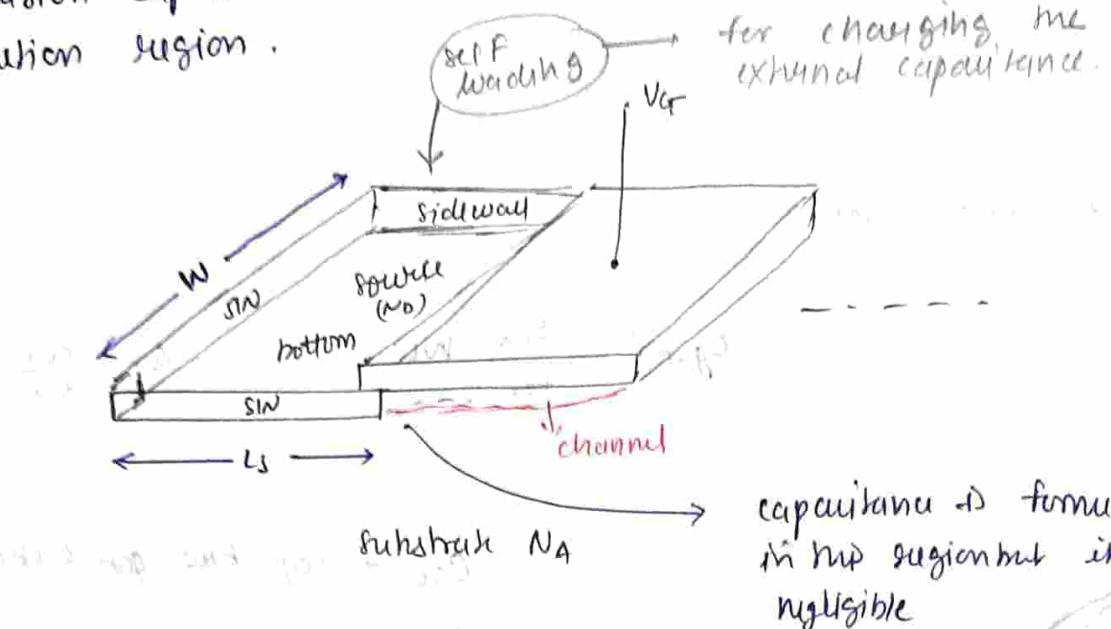


fig.2 Capacitance as a function of degree of saturation

- ~~mean way less~~
- C_{GCS} ↑ & $\cancel{C_{GCD}}$ ↑ because pinchoff has occurred & drain is not connected to the channel on ↑ V_{DS} , it is going towards pinchoff & after some time pinchoff occurs & there will be no C_{GCD} capacitance.

MOSFET - Diffusion capacitance

- Diffusion capacitance is formed whenever there is diffusion depletion region.



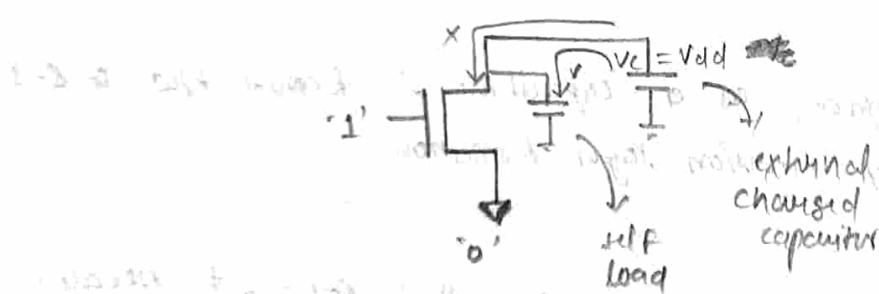
$$C_{\text{diff}} = C_{\text{bottom}} + C_{\text{sw}}$$

$$= C_j \times \text{Area} + C_{jsw} \times \text{Perimeter}$$

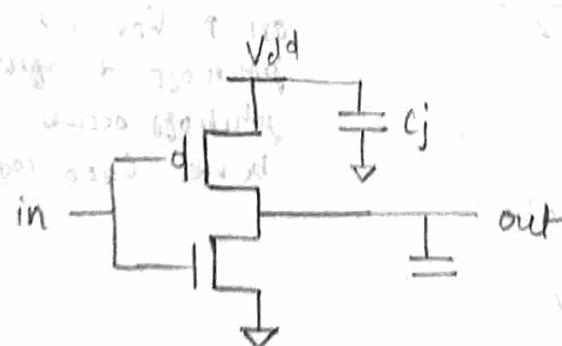
$$= C_j L_S W + C_{jsw} (2L_S + W)$$

⇒ Diff capacitance acts as self load when you want

into discharge external capacitance.



[cap. due to
 C_{GS}] \Rightarrow diffusion capacitance



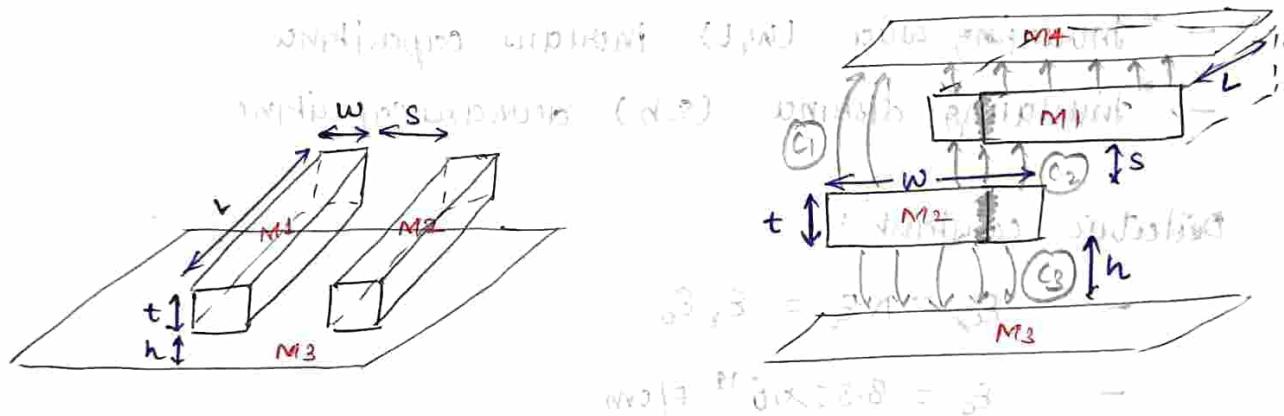
VLSI Devices - Wires

- For making interconnects in integrated circuits we use wires
+ wires are generally of Cu or Al.
— Nowadays Cu is mostly used.

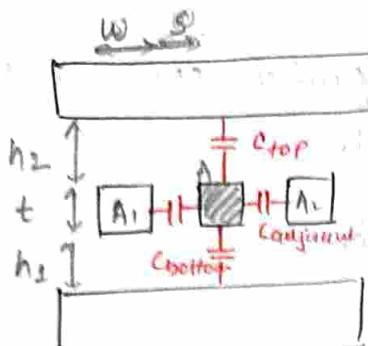
(VDF)

1. Wire Geometry and Capacitance

- Pitch = $w + s$ Old process had $AR \ll 1$
- Aspect ratio $\tau : AR = \frac{t}{w}$ Modern process have $AR \approx 2$
(Ratio of thickness & width)
Aspect ratios pitch & width determine low dimensional —



- Between diff. metal layers there is ~~oxide~~ = oxide layer of SiO_2 . because of which diff. metal layers shows capacitance to each other.
- Different metal layers are connected with via when we want to transfer data b/w them hence no capacitance is formed.
- But if M₁ → carries CLK pulse
M₂ → carries data signal
so we don't connect them with vias hence a capacitance is observed b/w them because of SiO_2 b/w them.



$$C_{\text{total}} = C_{\text{top}} + C_{\text{mid}} + C_{\text{bot}}$$

\Rightarrow How to deviate net cap. of shaded box in middle.

\rightarrow (i) Invert distance of A1 & A2 and A1+A2.

(ii) By not connecting the wire "way" there will be no effect of A1 capacitance.

Capacitance Trend:

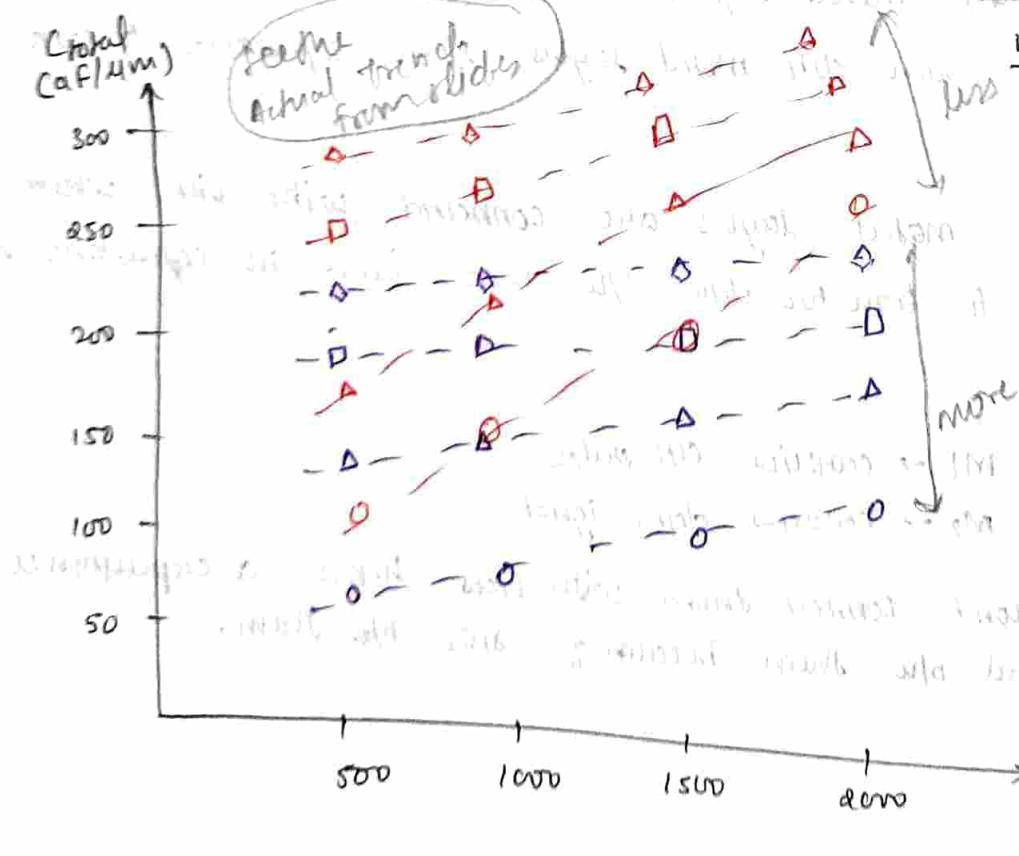
- Parallel plate equation $C = \frac{\epsilon_0 A}{d}$
 - Wires are not parallel but may obey trend.
 - Increasing area (w, t) increases capacitance $C \propto w, t$
 - Increasing distance (s, h) decreases capacitance $C \propto \frac{1}{s, h}$

Dielectric constant:

$$\epsilon_{ox} = K \epsilon_0 = \epsilon_r \epsilon_0$$

$$\epsilon_0 = 8.85 \times 10^{-12} \text{ F/cm}$$

$$K = \epsilon_{si} = 3.9 \quad (\text{for } SiO_2)$$



M1, M3 plane

- $\diamond \rightarrow s = 320$
- $\square \rightarrow s = 480$
- $\triangle \rightarrow s = 640$
- $\times \rightarrow s = \infty$

Trench

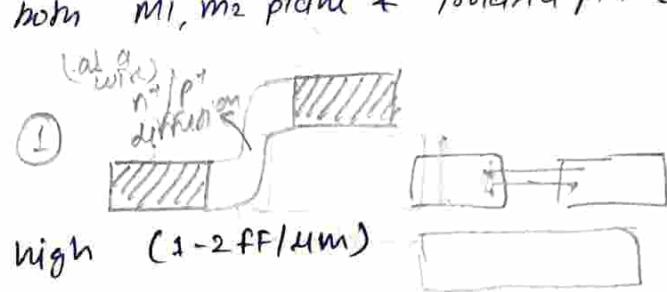
- $\diamond \rightarrow s = 320$
- $\square \rightarrow s = 480$
- $\triangle \rightarrow s = 640$
- $\times \rightarrow s = \infty$

"Width of metal"

- (a) for less 'w' width of wire, * typical dnm wire ham
≈ 0.2 fF/ μ m
- (b) for isolated metal:
- (i) For isolated metal, for constant 'w' → width of wire, when we increase 's' → distance b/w adjacent metal, then change in capacitance is very high.
 - But for M₁, M₂ plane the change is low because of adjacent metal capacitance being high.
 - (ii) for isolated metal, w → increasing, s → same then increase in capacitance is more in M₁, M₂ plane as compared to isolated plane.
 - (iii) for small value of w, the capacitance don't have ~~much~~^{met} much both M₁, M₂ plane & isolated plane.

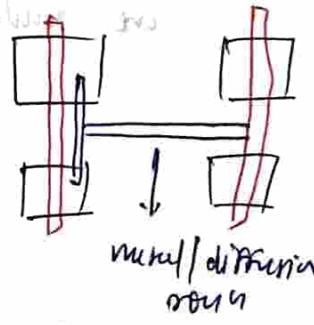
2. Diffusion & Polysilicon

- Diffusion capacitance is very high (1-2 fF/ μ m)
 - similar to gate capacitance
 - diffusion also have high resistance
 - Avoid using diffusion runners for wires
[We don't use diffusion to make own capacitance]
 - (?) But if we use diffusion then we put lots of contacts over it, put lots of wire over it, metal over it so that resistance reduces



Eg. Suppose we want to connect two inverters & we have two option first through metal route & another through diffusion route.

- we prefer metal route because it provides less capacitance but sometimes metal is not available for data transmission ∴ we can diffusion runners. (Another reason which can transmit data)



Type of wires:

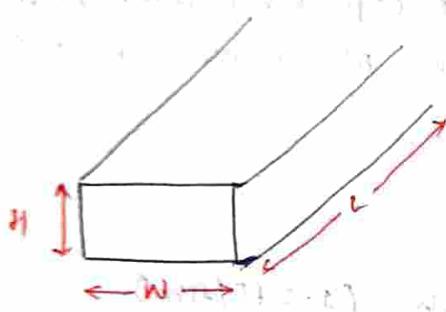
- metal wires
- Poly silicon wire (for very short distance)
- diffusion routes

Q: What is active region, for transmitting signals?

→ Active region \equiv OD \equiv diffusion region

\equiv Region of source + drain \equiv drain region

3. Wire : Resistance



$$R = \frac{\rho L}{W}$$

↳ sheet resistance (R_s) [n/square]

* As a designer we don't have control over the length, height & width of the wire, it's a technology feature.

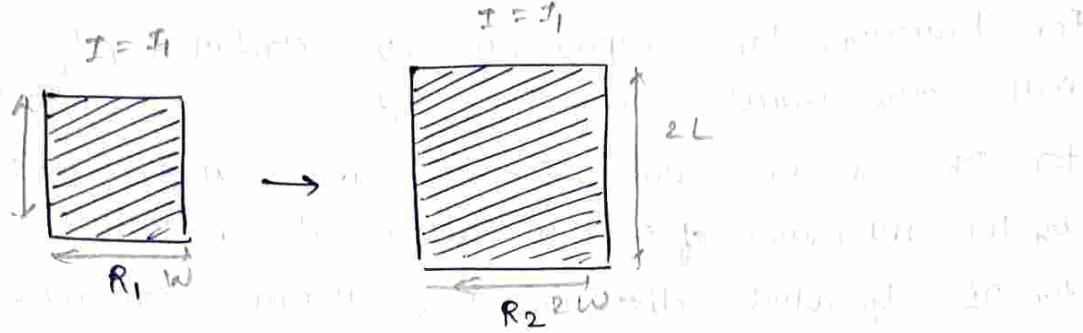
$\frac{\rho}{H} \rightarrow$ constant, \therefore total resistance depends on $\frac{L}{W}$.

and $\frac{L}{W}$ says that the no. of squares that we can make in direction of flow of current.

* Unit of sheet resistance $R_s = \frac{\rho}{H} \rightarrow$ n/square

\therefore For calculating the no. of resistance of material we multiple sheet resistance with no. of squares.

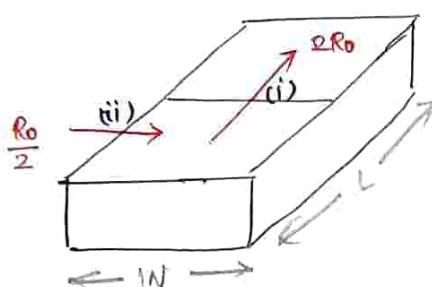




Unit of sheet resistance (R_0) \rightarrow ohm/square

$$\textcircled{a} \quad R_0$$

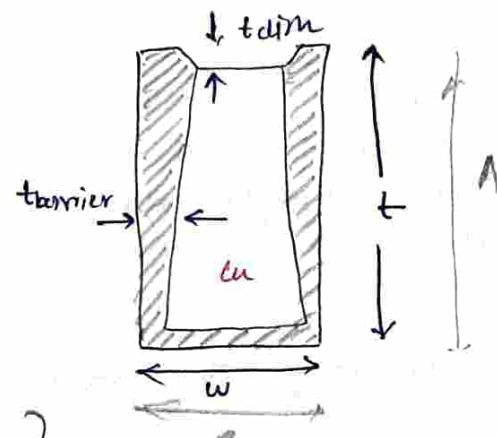
- $I = I_0$ will simply multiply the no. of sheets/squares in the direction of flow of current, with no sheet resistance to get total resistance in the direction of current.



$$\begin{aligned}
 \text{i)} \quad R &= R_0 \frac{L}{w} = R_0 \cdot (2) \\
 &\stackrel{\text{fixed}}{\downarrow} \\
 &= 2R_0 \\
 \text{ii)} \quad R &= R_0 \left(\frac{w}{L} \right) = R_0 \left(\frac{1}{2} \right) \\
 &= \frac{R_0}{2}
 \end{aligned}$$

Cu-wire me barrier

- Cu-wire diffusion barrier has high resistance
- Cu is also prone to dishing during polishing
- Effective resistance is higher



$$\left\{
 \begin{aligned}
 R &= \frac{\rho}{(t - t_{\text{dish}} - t_{\text{barrier}})} \cdot \frac{l}{(w - 2t_{\text{barrier}})} \\
 A & \qquad \qquad \qquad B
 \end{aligned}
 \right\}$$

Q: For decreasing the capacitance b/w contact & poly, we can make our contact away from poly (as $d \uparrow, C \downarrow$). But for DRC it says that for OD also - it has to be covered by the minimum of outside the contact, therefore we shift the OD by which diffusion & contact capacitance comes into the picture. So is there a trade off we can't do this?

- By decreasing the diffusion region area overall capacitance decreases because the cap. due to poly-contact does not get affected as it is always placed at the min-DRC distance hence not affecting the cap. due to poly-contact.

$$(\frac{1}{2}) \pi r^2 = \frac{\pi r^2}{2} = R \quad (i)$$

$$(\frac{1}{2}) \pi r^2 = (\frac{\pi r^2}{2}) = R \quad (ii)$$



* Tab poly ke upper se metal jata hai to capacitance add hote hai but unke Karan phleki hi bahut se capacitance hua Rukta tha.

Q: Elmore delay numericals on Mid-term & End-term.
→ solve from Weste & Quinton Banks.

Q: For deviating the capacitance b/w contact & poly, we can make our contact away from poly (as $d \uparrow, C \downarrow$). But for DRC it says that for OD also - it has to be covered by the minimum of outside the contact, therefore we shift the OD by which diffusion \uparrow & diffusion capacitance comes into the picture. So is there a trade off we can't do this?

- By deviating the diffusion region over overall capacitance decreases because the cap. due to poly-contact does not get affected as it is always placed at the min. DRC distance hence not affecting the cap. due to poly-contact.

$$C_{diff} = \frac{\epsilon_0 A}{d} = 9 \quad (i)$$

$$\epsilon_0 = 2$$

$$\left(\frac{1}{2}\right) \times 9 = \left(\frac{W}{2} \times 9\right) = 9 \quad (ii)$$

$$\frac{W}{2} =$$



With different materials (e.g., SiO₂, Al, TiN, etc.)

different effect

depends on many factors

primarily film

height is important factor

** Jab poly ke upper se metal jata hai to capacitance add ho ja hei but use karen phleki hi bahut se capacitance hua karta tha.

Q: Elmore delay numerical on Mid-sem & End-sem
→ solve from Nestle & Question Bank.

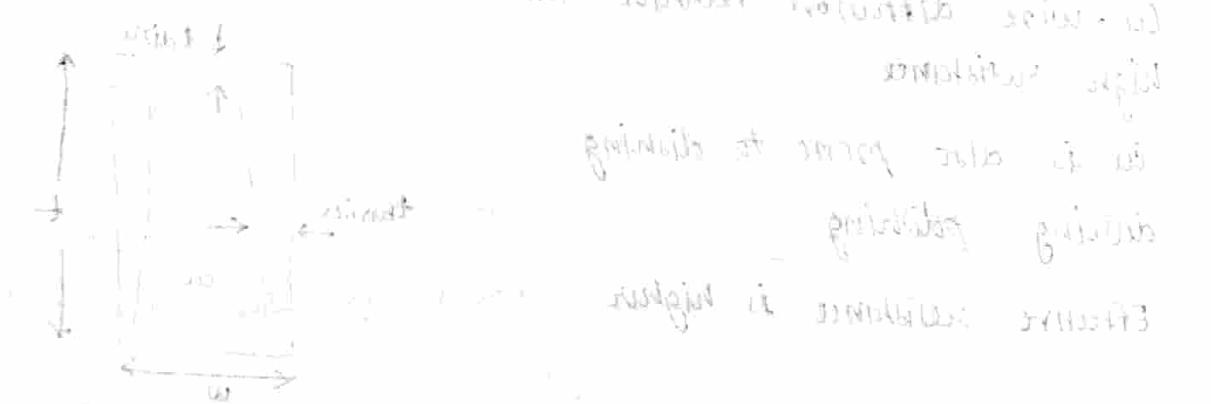
Q. For decreasing the capacitance b/w contact & poly, we can make our contact away from poly (as $d \uparrow, c \downarrow$). But for DRC it says that for OD also - it has to be covered by the minimum of outside the contact, therefore we shift the OD by which diffusion & diffusion capacitance comes into the picture. So is there a trade off we can't do this?

- By decreasing the diffusion region area overall capacitance decreases because the cap. due to poly-contact does not get affected as it is always butted at the min-DRC distance hence not affecting the cap. due to poly-contact.

$$C_{\text{diff}} = \frac{\epsilon_0 \epsilon_r}{d} = k \quad (\text{eq})$$

$$(1) C_{\text{diff}} = \frac{W}{2} \cdot k = k \quad (\text{eq})$$

$$\frac{k}{2} =$$



* Jab poly ke upper se matal jala how to capacitance add hote hui but use Karun phelwadi hi bahut lesser capacitance hua karta tha.

Q. Elmore delay numericals on Mid-jim & End-jim.
→ solve from Weste & Runtian Banerji.

Standard cell layout Design

Layout of NAND gate

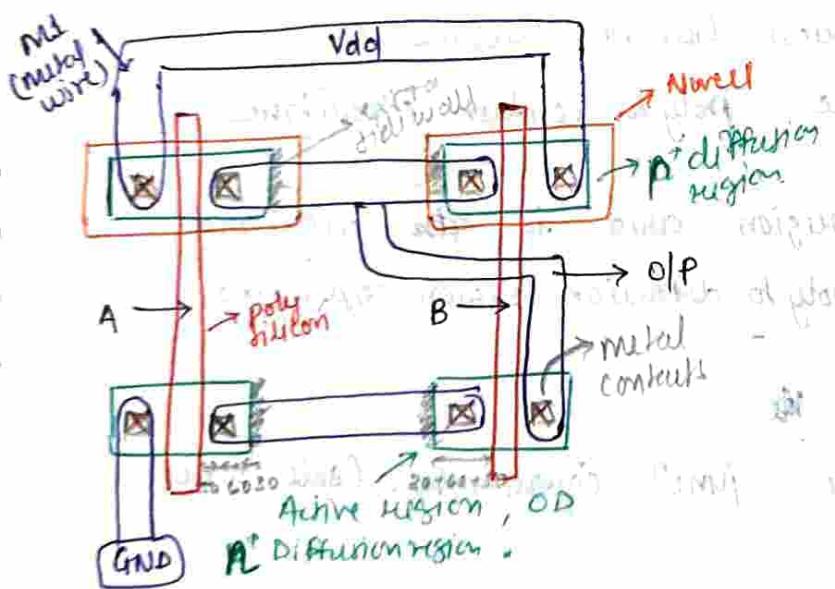
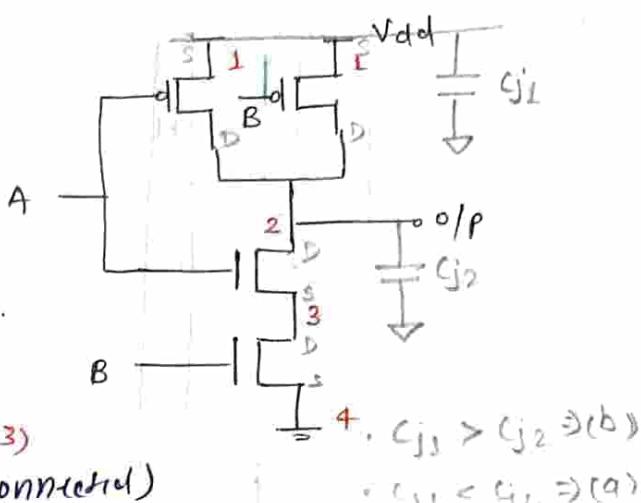
- For PMOS, one of the diffusion region from both PMOS are connected to Vdd. (1)

- And other diffusion region of both PMOS are connected to o/p. (2)

- For NMOS, one of the diffusion region is shorted (or directly connected) (3)

- and one of NMOS is connected to o/p. (2)

- and another of NMOS is connected to GND. (4)



* If we make NAND gate using this method then overall capacitance will increase in IFA (ii)

(i) No. of metal contacts = 8

- Poly to contact capacitance
is due to having 8 metal contacts.

(iii) Now it has extra side walls because of which junction capacitance b/w diffusion region & substrate increases.

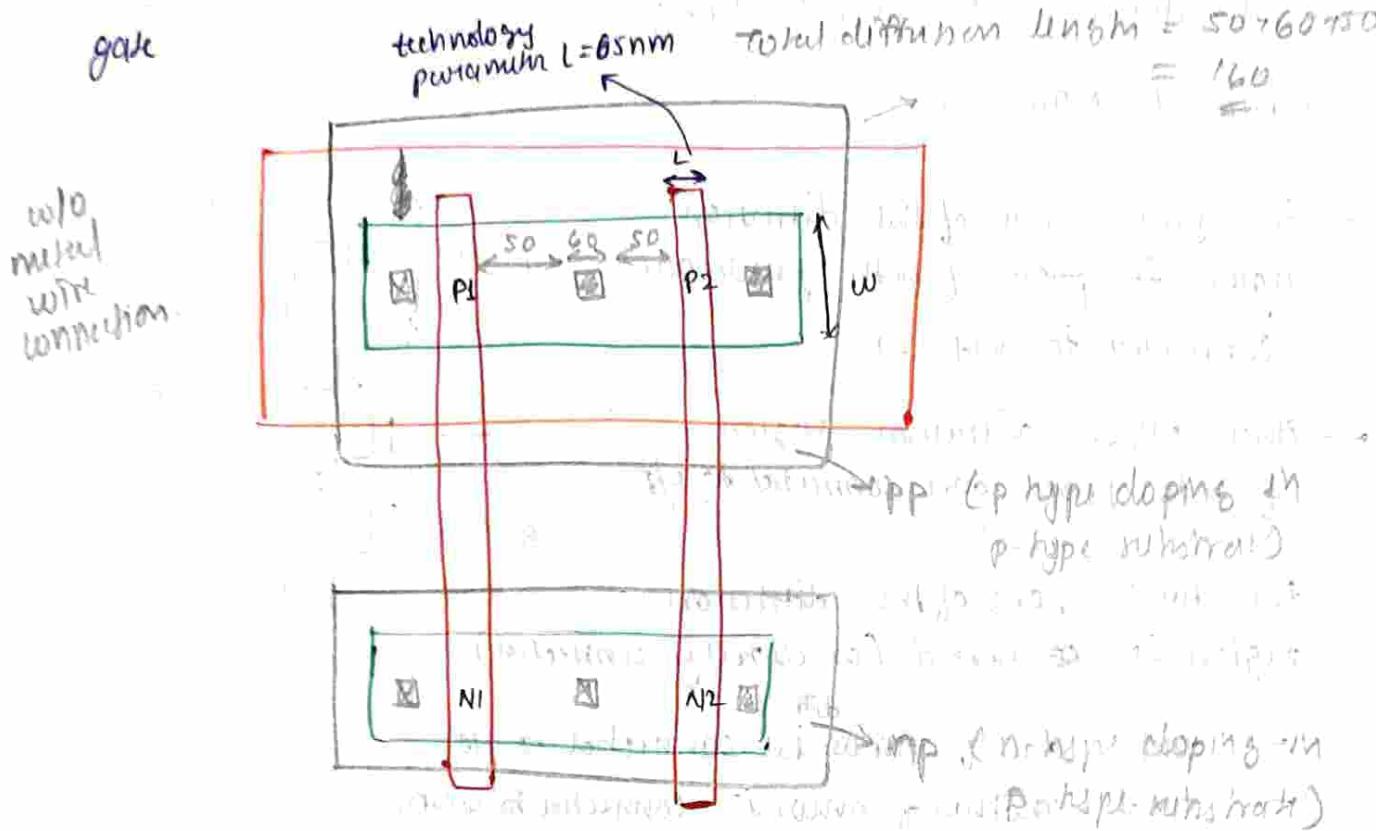
(ii) Total 8 source & drain capacitance

- Overall area for the active diffusion region is increased

- Capacitance due to metal contact & diffusion region increases.

* Any of the poly region can be circuit any can be source.

- Therefore we use the below standard layout for NMOS gate



(i) Here no. of metal contacts has been reduced to 6.

- Hence decreasing the poly to contact capacitance.

(ii) Also the diffusion region area is decreased

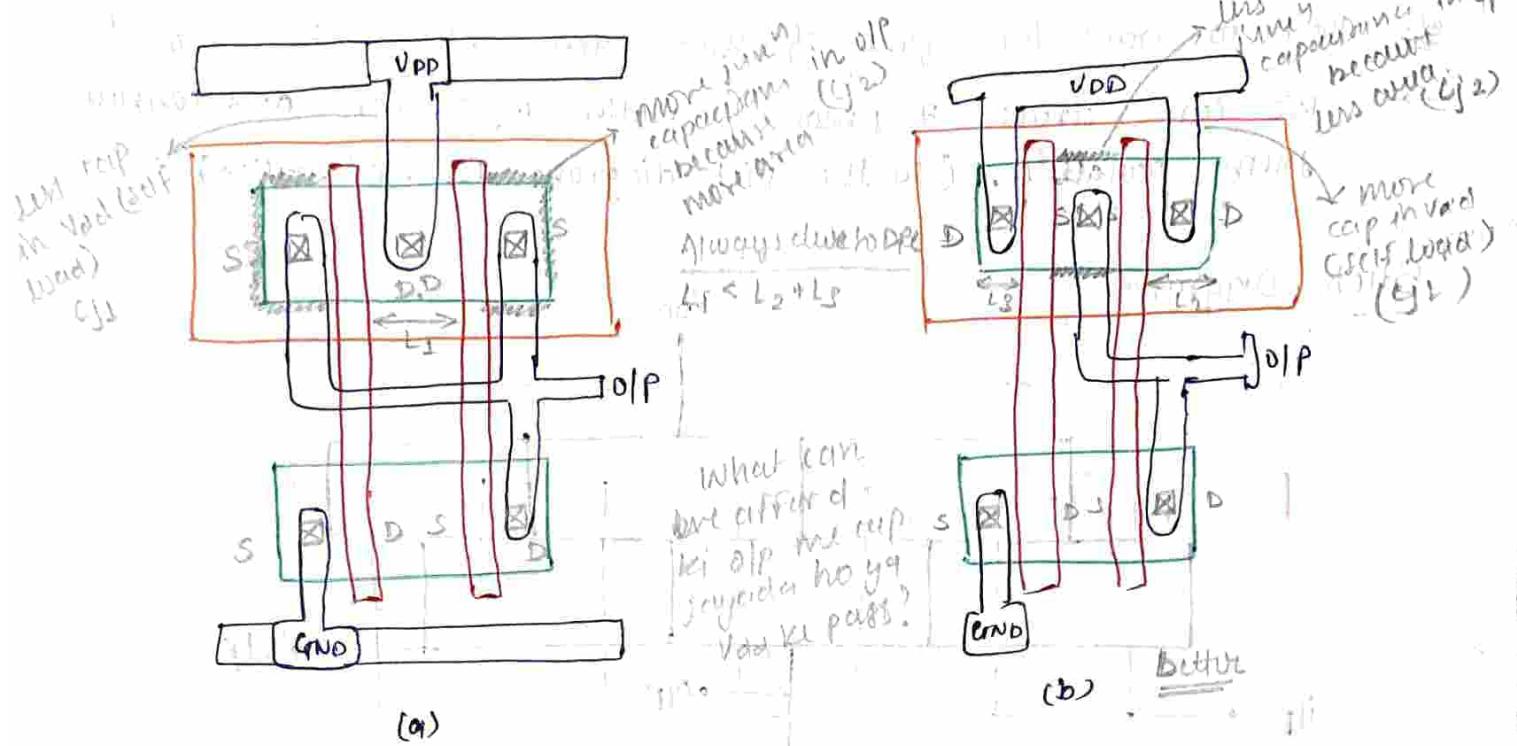
- Decreasing the poly to diffusion region capacitance

(iii) No extra side walls

- Hence decreases the "junc" capacitance. (Side wall capacitance)

Q. 9. In the above layout which one should be op & which one should be Vdd for pmos. (Tell all the metal connections).

→ Electrically there is no difference but we have to see if we

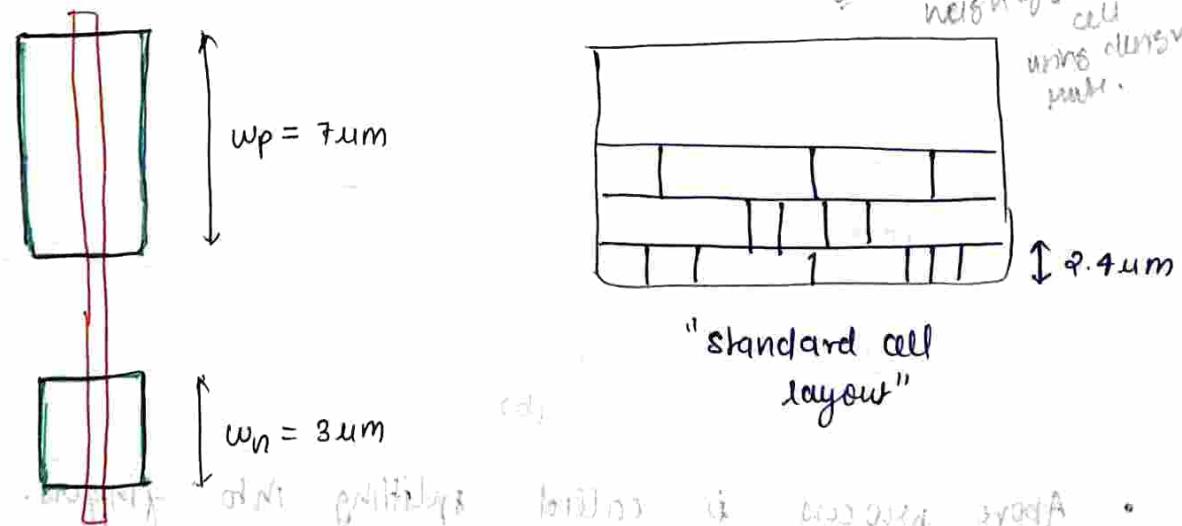


- To me go by fig. (a) then more metals will be used in the o/p resulting in increase in capacitance at the o/p.

Point 2

While drawing the each layer in our layout we have to follow design rules that (i) minimum width of contact metal & poly (ii) min area of ND, NDLL, PP, NP must be satisfied (iii) spacing b/w the contact & poly should be more.

- Why we use NTAP & PTAP?
- How to calculate height of standard cell using design rule.

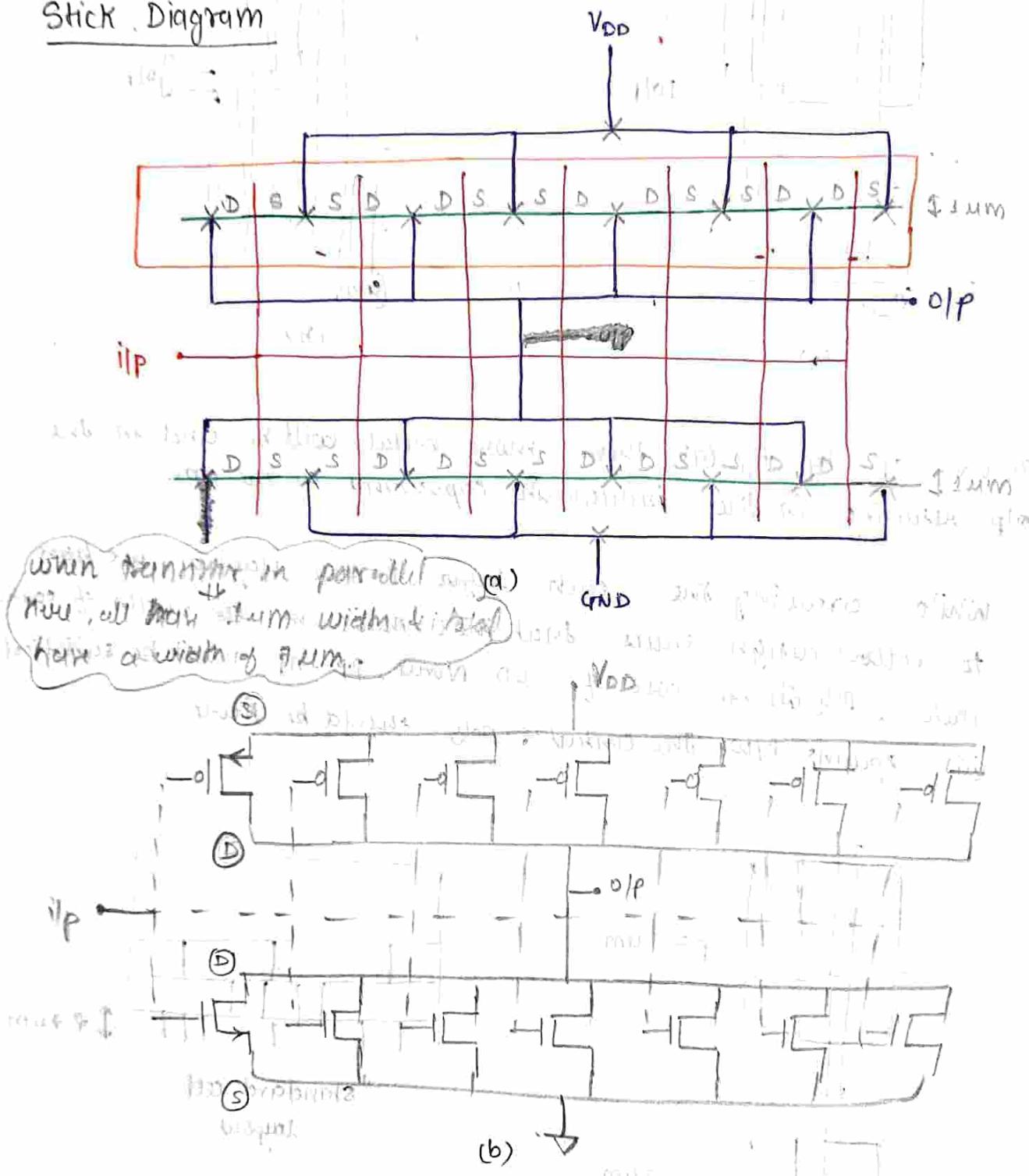


a. $w_p = 7 \mu m$ & $w_n = 3 \mu m$ in our CMOS if but the height constraint per cell is $2.4 \mu m$ then how do we draw this cell under $2.4 \mu m$?

→ explanation is shown in stick diagram.

- Q. How to draw the pmos of $w_p = 7\mu m$ into a height of $1\mu m$.
 → We can draw 7 pmos of width $w_p = 1\mu m$ and connect them parallelly. (see the stick diagram for explanation)

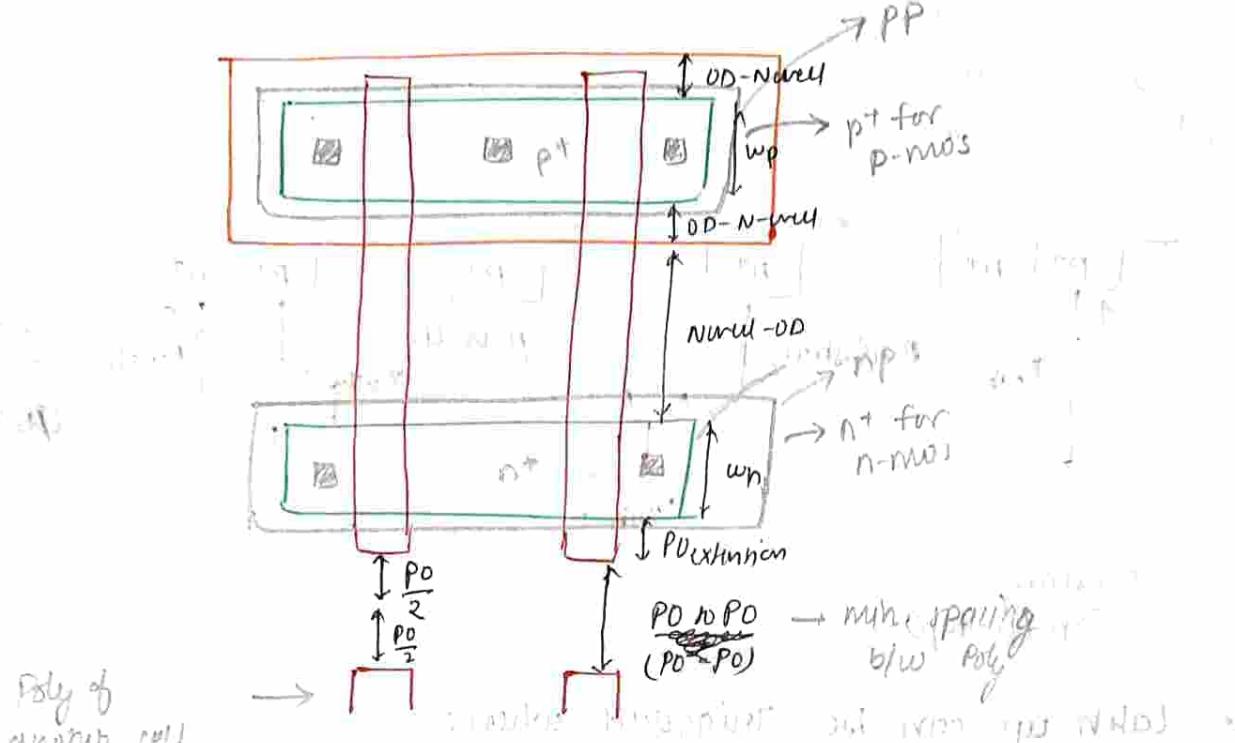
Stick Diagram



Above process is called splitting into fingers.

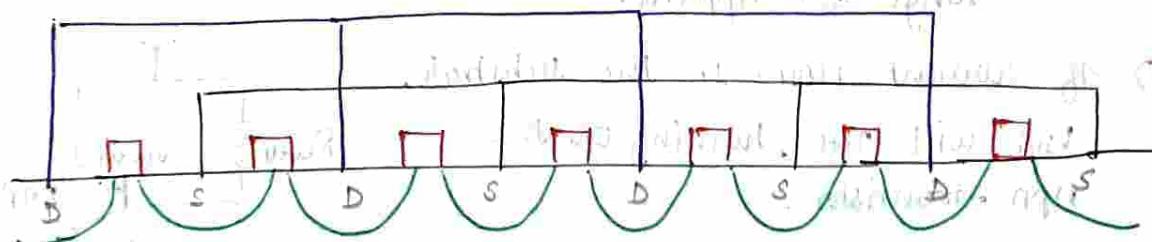
When we do short fingers then area increases because we have to constraint the height.

Q. How calculate the height of given CMOS?

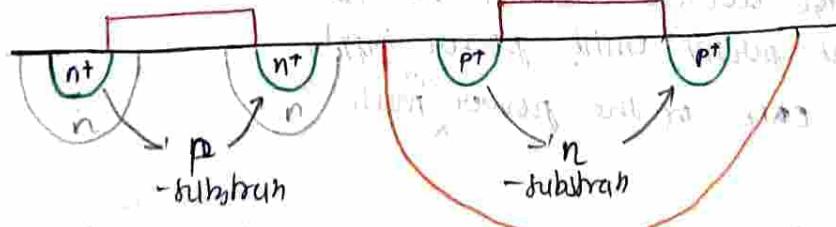


$$\text{Height} = (\text{OD-Nwell}) + (\text{wp}) + (\text{OD-NW}) + (\text{NW-OD})$$

$$+ (\text{wn}) + (\text{POext}) + (\text{PO-PO})$$

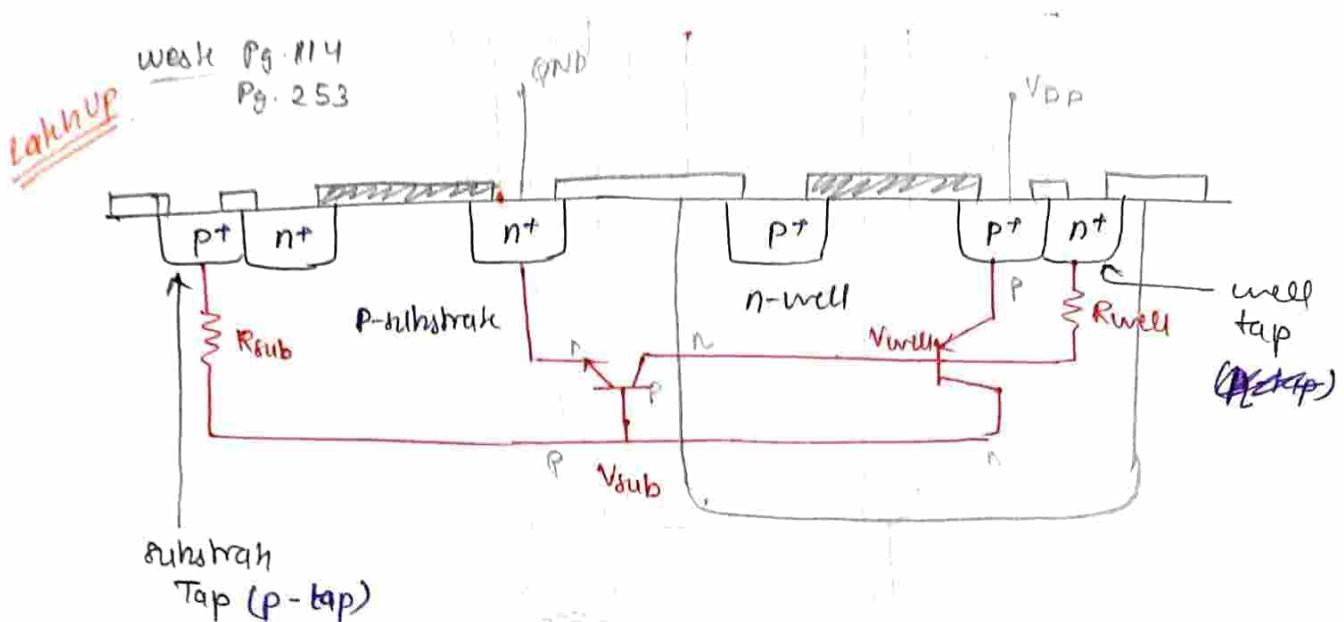


Part 3 Q. What is leech up? And how do we avoid it?



- n^+p^+ & p^+n^+ kind of BST & CL is observed here when resistance of substrate is high. In designs with narrow A and low k, leeching occurs.

- So a CKT is made with a FB of npn + pnp transistors. [Latch-Up]

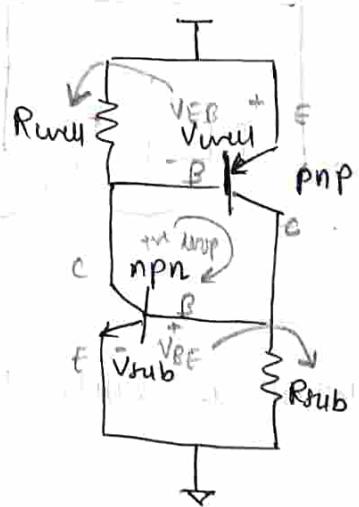


- Latch up can be triggered when:
 - when transient current flows through the substrate during normal chip power up.
 - when external voltages outside the normal operating range are applied.

- ① If current flows in the substrate,

V_{sub} will rise, turning ON the npn-transistor.

- This pulls down the V_{well} & turning ON the pnp-transistor.
- The pnp-transistor in turn raises V_{sub} , initializing the FB loop with large current flowing b/w V_{DD} & GND that persist until power supply is turned OFF or the power wire melts.



- * Latch up can be prevented by reducing R_{sub} & R_{well} .

- * A common approach is to place 'tap' adjacent to every source connected to V_{DD} & GND.

The following guidelines:-

- (i) Every well should have atleast one tap.
 - because the n-well sheet resistance can be several k Ω per square, it is necessary to ground the well thoroughly by providing a sufficient no. of well taps.
 - This will prevent excessive voltage drops due to well currents.
- (ii) All substrate + well taps should connect directly to the appropriate supply in metal.
- (iii) A tap should be placed for every 5-10 transistors, or more often in 'spare' areas.
- (iv) nmos xistor should be clustered together near GND & pmos xistor should be clustered together near V_{DD}.

Q: What is latch up & how do we avoid them? (Conclusion)

PTAP & NTAP are used for?

- the feedback loop due to npn & pnp xistor.
- It can be avoided by reducing the resistances
- ^{multiple} PTAP & NTAPS (substrate + well taps) are used to reduce R_{sub} & R_{well}.

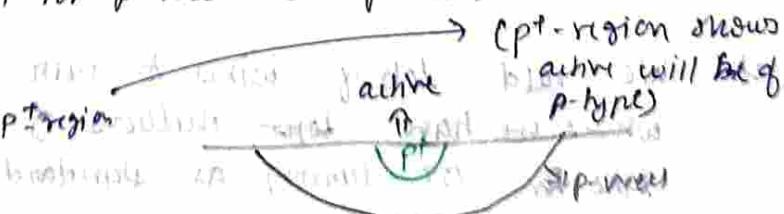
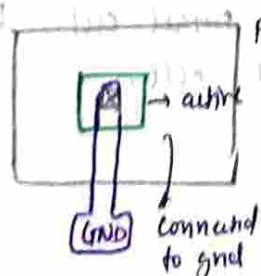
* After LVS we do substrate coupling extraction & cadence.

- LVS will show error if you don't put PTAP & n-well in layout

tap → substrate connection → contacts are sometimes else

→ 'vias' in layout are different from 'taps'

Q: how do you ~~connect~~ make a tap, if you want to make the substrate connection for p-well. (for p-MOS).



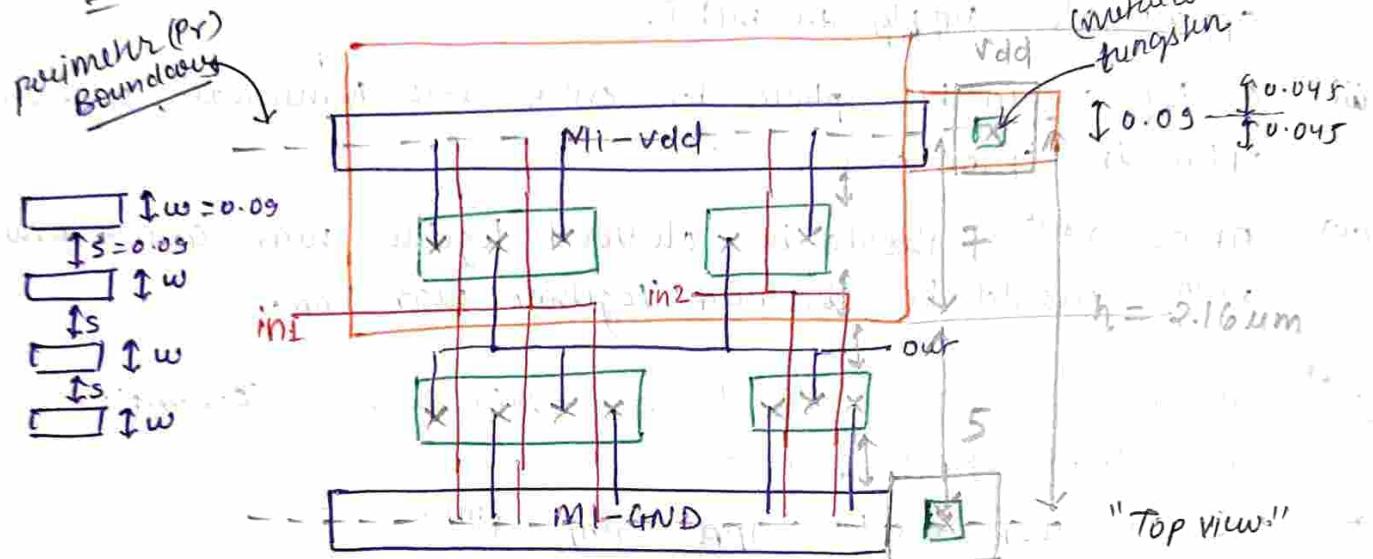
∴ There is no depletion region b/w them

* nt + pt regions shows us what kind of doping goes where they are not actually n+ layers, they are n+ implants. (why?)

Q. How do we make interconnects?

Standard Cell - A basic unit of cell library consisting of:

Part 1 Q. how the implants of n+ standard cells are made?



Q. How do we decide the height of 'h'.

→ the minimum width of M1-metal = 0.09 um] DR rule
the min. spacing b/w M1 & M1 = 0.09 m] DR rule

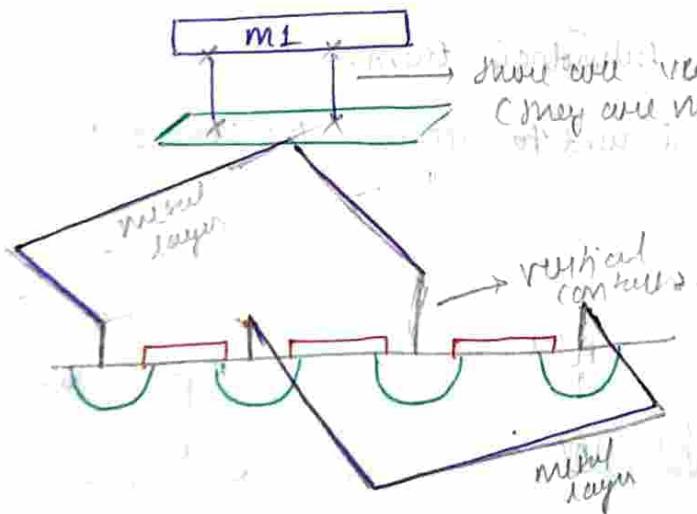
Hence we can use using 12T (12 track) library
where $T = 0.18 \mu\text{m}$, then we can place only
12 metal layer (M₂, M₃, ..., & others)

above our cell with the min. spacing b/w them

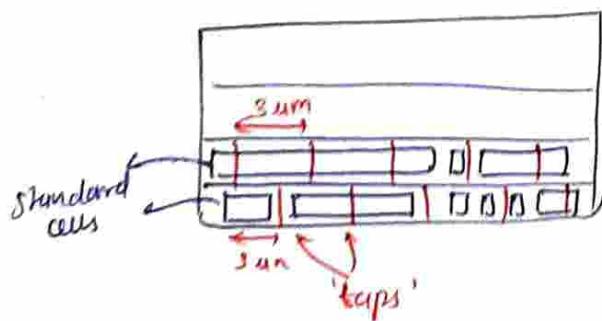
$$h \Rightarrow 12T = 12 \times 0.18 = 2.16 \mu\text{m}$$

We need lots of wires to run across the whole chip
where we have tens millions of small cell hence we
take 12T library as standard cell.

- We only take half DRC of the metal layers above & below the cell.
 - As metals layers (VDD, GND) are shared with other cell above it & below it.
 - all above & below will be rotated 180° for connection of VDD & GND.
- DRC is only b/w M₁-M₁, P_O-P_O, above & below.
 - we can pass M₁ above the P_O or vice-versa.
- Metal contacts are tungstens.
 - via is copper from dual damascene process
- When we make any vertical connection b/w metal M₁ & active region, as M₁ is running above active region. Is that by via?
 - No it is not a via.



- * Minimum spacing b/w the taps should be ~~34~~ 34 μm.
- How do we place the 'taps' in our chip?



- During placement of M₁'s standard cell in our chip we also put the 'taps' with minimum spacing.
 - So no need to make tap for every cell.

* For no-cutting region one tap should be far away from 3μm (min. DRC) distance away.

Q: If we have many no. of cells b/w, so do we have to share the n-well for all the stand. cells, as 'taps' are placed around at 3μm distance only.

→ Yes, the n-well will also be from right to left.

- therefore it is advised to make np & pp vertically parallel to each other.

* 'taps' are also called as 'straps'.

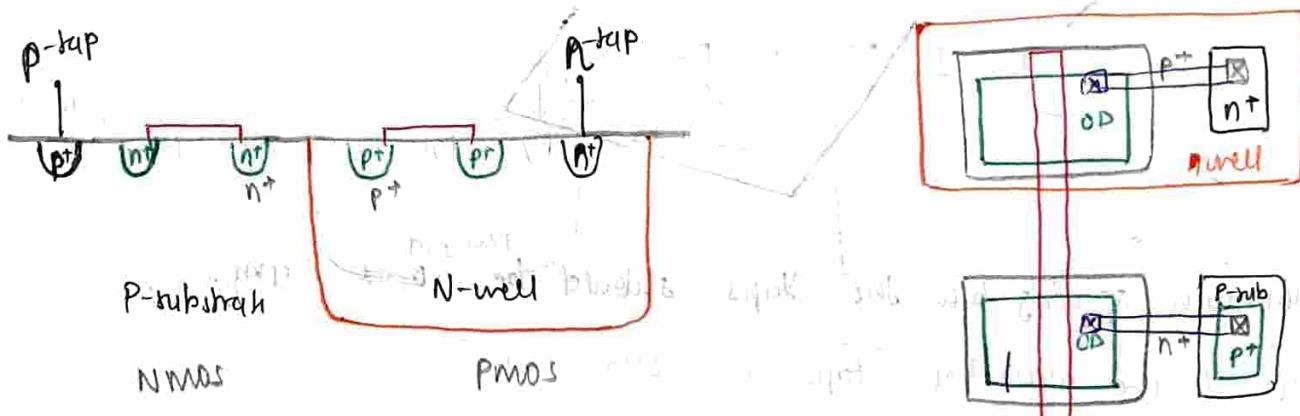
Q Why we have DRCs?

→ Agar hum min DRC de rakhna width lete hoi to technology wale logo ke koi dikkaat nahi hoi, dikkaat goliya hogi kuki functionality of the device may change.

- Agar min ~~width~~ spacing or Kam spacing b/w two rings may get short circuited.

• DRC is defined by technology team.

* V_T implant layer is used to deviate the value of V_T.



4. Challenges in Cu-wire:

- ~~Copper~~ Diffusion barrier in Cu-wire has high resistance.
- Cu is also prone to ~~afford~~ by dishing during polishing.

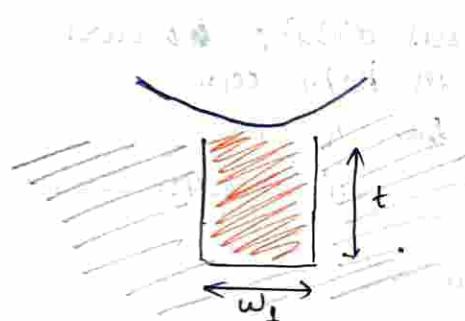
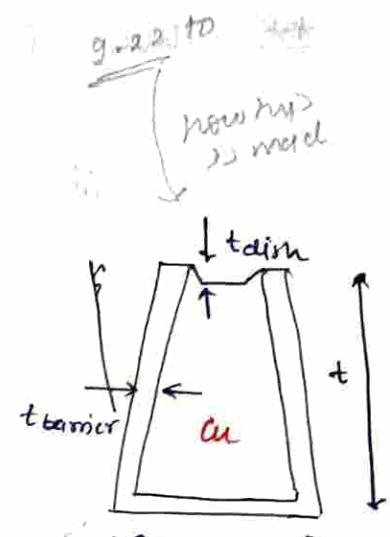
VDF

Q: What is dishing in Cu-wire?

→ polishing → flattening out the surface

- Due to this polishing a mechanical force is also applied which press the Cu-wire & don't affect the Barrier.

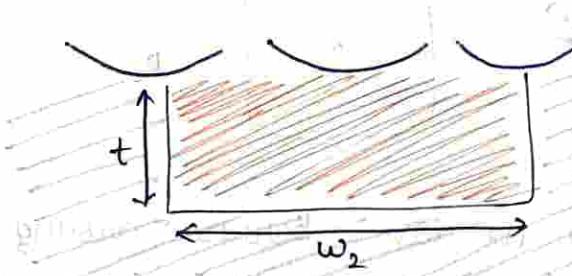
→ thus changing the thickness of Cu is called dishing.



for less width w_1 ,

less dishing \Rightarrow observed

"Mores
dishing
view"



for high ~~less~~ width w_2

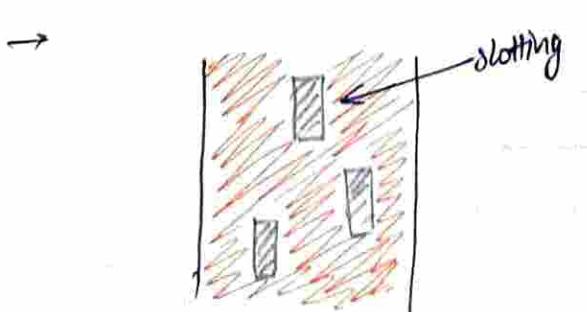
more dishing \Rightarrow observed.

but becomes
sharp.
due to dishing.

Q: What is the problem due to dishing?

→ We wanted less resistance hence we increase the width 'w' but due to dishing the effective area is reduced hence increasing the Resistance

Q: How to decrease this dishing?



"Top View"

• if we can somehow increase the mechanical strength of Cu, then dishing can be reduced.

— therefore we can slot in b/w the Cu-wire, this slots are of SiO_2 .

— This also reduces the Cu-area, $\therefore R \uparrow$ but it is better than reducing the width of Cu.

(i) Reduce dishing by decreasing width but it will increase resistance more.

(ii) Reduce dishing by slotting. It will also increase resistance but not that much.

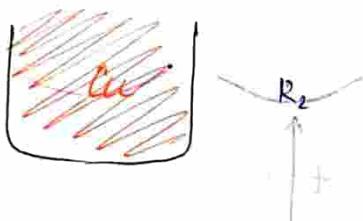
* There is DRC to all Cu-wire, slot, barrier etc.

(iii)



slotting

(ii)



less dishing occurs in both cases but $R_1 < R_2$
∴ (iii) is better than (ii)

Q: Why do we have barrier around Cu-wire?

protective and insulating effect

if thickness more

4.85

Wid Pg 220 (6.3.1)

Rubber Pg. 249 (6.3.5)

(6.4.2)

Thickness not mentioned (6.4.3)

(6.4.4)

RC Delay

protective and insulating effect

(Wid) paper back

cross talk (6.3.3)

(6.4.3)

Thickness not mentioned (6.4.3)

Thickness not mentioned (6.4.4)

polishing material

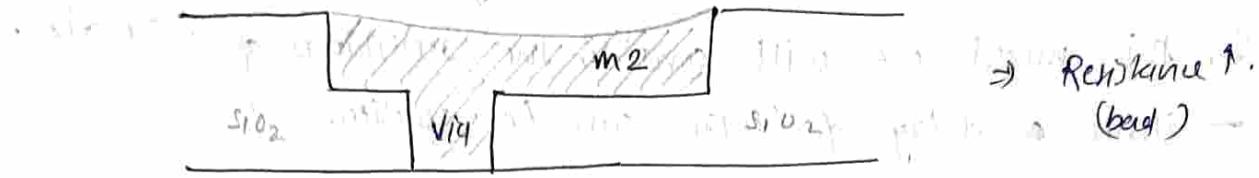


diffused in metal resistivity
- 0.14 for the initial diameter

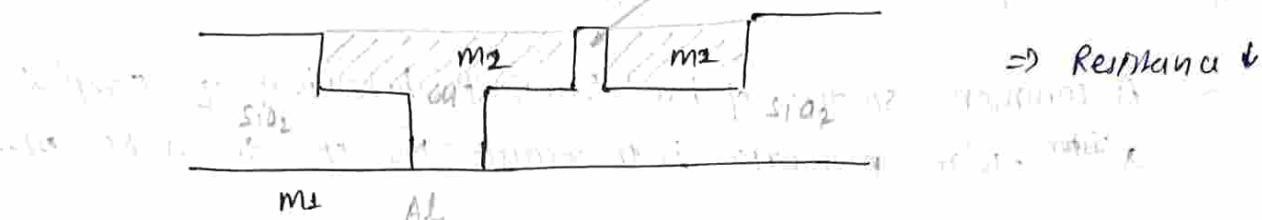
TOP, diffusion and diffusion rate

diffusion length of the metal

no problem



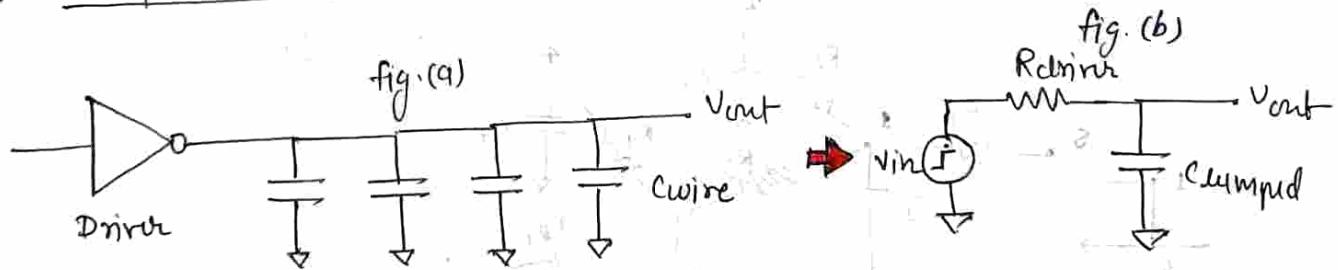
• ~~slotting~~ → slotting will result in a much higher self-inductance due to air gap. \rightarrow Aperture slotting is a better way to increase coupling.



- We use low-K dielectric polymer material for low capacitance.
- We use high-K dielectric MgO as gate insulator material.
- * Miller coupling.

5.4 RC Delay

• Robotic 4.4.2
i) Lumped Model "Lumped capacitor Model"



- In this model we only consider the capacitive component of wire, and to lump the distributed capacitance into a single capacitor.
 - In this model the wire still represents an equipotential region, and the wire itself does not introduce any delay.
 - Distributed v/s lumped cap. model of wire [fig.(a)-(b)]
- C_{lumped} = L × C_{wire} if the driver is modelled as a voltage source + a source resistance = R_{driver}.
- (→ length of wire)
- C_{wire} → cap. per unit length

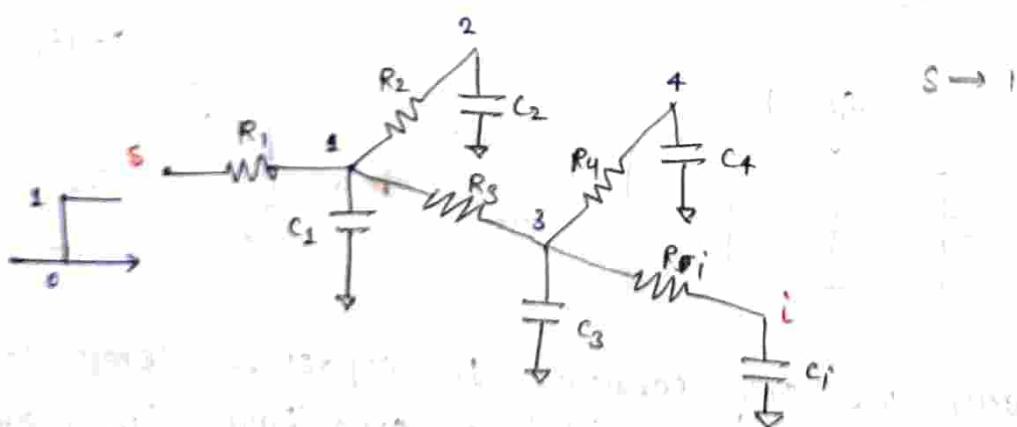
(ii) Lumped RC model [Elmore Delay]

Paraboy
4.4.2

- In this model, we will consider the resistance of wire also.
- Hence a delay of wire can be studied.
- In this model, we lump the total wire resistance of each wire segment into a single R & similarly combines the global capacitance in a single capacitor C .
- A common study of the transient behaviour of complex X_{bus} -wire networks is to reduce the ckt to an RC n/w.
- A large complex RC n/w can be solved by two methods:
 - i) Elmore Delay method
 - ii) Distributed RC delay method

Elmore Delay : [only good for short length n/w]

Consider a less complex RC n/w for understanding Elmore delay.



Properties of above RC-tree:

- i) The n/w has single i/p node.
- ii) All the capacitors b/w a node & the ground.
- iii) the n/w doesn't contain any sensitive loop (which makes it a tree).

- The total Resistance along path 's' to 'i' is, path Resistance (R_{ji})
 - The path resistance b/w one source node 's' & node 'i' is,
- The resistance shared among the path from the root node 's' to node 'k' & 'i', is shared Path Resistance (R_{ik})

* Assume 'N' nodes of the N/w is initially discharged to GND, and that a step ilp is applied at node 's' at $t=0$.

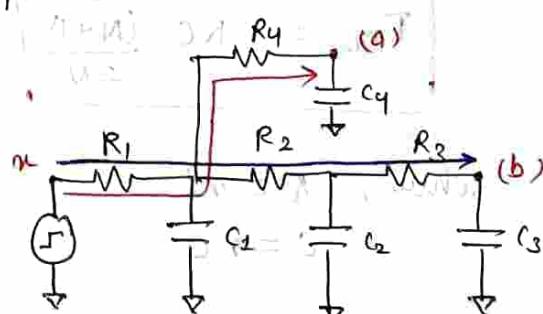
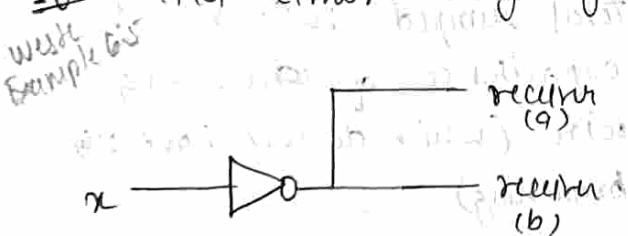
- the Elmore delay at node 'i' is,

$$T_{Di} = \sum_{K=1, K \neq i}^N C_K R_{iK}$$

shared path R_{ik}

$$\therefore T_{Di} = R_1 C_1 + R_1 C_2 + (R_1 + R_3) C_3 + (R_1 + R_3) C_4 + (R_1 + R_3 + R_4) C_5$$

Eg. Find Elmore Delay from ilp 'x' to each receiver (a) & (b)



$$\therefore T_{D(a)} = R_1 C_1 + R_1 C_2 + R_1 C_3 + (R_1 + R_4) C_4$$

$$T_{D(b)} = R_1 C_1 + (R_2 + R_3) C_2 + (R_1 + R_2 + R_3) C_3 + (R_1 + R_2 + R_3 + R_4) C_4$$

initial voltage $V_0 = 2V$ is required.

Example: Time constant of Resistive-capacitive wire

Time constant \rightarrow RC delay.

A wire of total length L is partitioned into N -identical segments, each with a length of L/N . The resistance & capacitance of each segment are hence given by $\frac{rL}{N}$ & $\frac{cL}{N}$ resp.

Compute the dominant time constant of the wire.

$$\begin{aligned} \rightarrow T_{DN} &= \frac{rL}{N} \cdot \frac{cL}{N} + \frac{2rL}{N} \cdot \frac{cL}{N} + \frac{3rL}{N} \cdot \frac{cL}{N} \\ &\quad + \dots + \frac{(N-1)rL}{N} \cdot \frac{cL}{N} \\ &= \left(\frac{L}{N}\right)^2 [rc + 2rc + 3rc + \dots + Nrc] \\ &= (rc) \left(\frac{L}{N}\right)^2 [1 + 2 + 3 + \dots + N] \\ &= rc \cdot \frac{L^2}{N^2} \cdot \frac{N(N+1)}{2} \\ &= rc L^2 \frac{(N+1)}{2N} = (rc)(cL) \frac{(N+1)}{2N} \end{aligned}$$

$$T_{DN} = RC \frac{(N+1)}{2N}$$

\rightarrow Total lumped resistance & capacitance of wire or long wire (which does not have any branching)

where, $R = rL$

$C = rc$

for very large value of N , this model asymptotically approaches the distributed 'rc' line, which gets simplified to

$$T_{DN} = \frac{RC}{2} = \frac{rcL^2}{2} \rightarrow \text{"delay for distributed rc"}$$

$$T_{DN} = RC \rightarrow \text{"delay for distributed lumped rc"}$$

Conclusion: Time delay with wire is always less than that of step response.

1. The delay of the wire is a quadratic funcⁿ of its length.

$$T_{\text{on}} \propto L^2$$

2. The delay of the distributed rc-line is half of the delay that would have been predicted by RC model.

$$T_{\text{distribution}} = \frac{1}{2} T_{\text{lumped}}$$

- * In a switch model, the resistors are replaced by their equivalent, linearized on-resistance. Then delay is calculated by RC model.

(iii) Distributed rc-line

- . We will not use this method for our delay calculation.
- + The step-response of lumped & distributed RC N/W

<u>Voltage range</u>	<u>Lumped RC N/W</u>	<u>Distributed RC N/W</u>
0 - 50% (t_p)	$0.69 RC$	$0.38 RC$
0 - 63% (T)	RC	$0.5 RC$
10% - 90% (t_{tr})	$2.2 RC$	$0.9 RC$
0 - 90%	$2.8 RC$	$1.0 RC$

(Note that 1) Amplitude doesn't change in response of transmission line. It is same.

2) Response of distributed rc-line is faster than that of lumped.

3) Step response of distributed rc-line is faster than that of lumped.

Ques: When should be use lumped ^{RC} model & when to use lumped capacitive model.

→ (i) RC delay should be considered when

$$t_{PRC} \gg t_{propagation}$$

driving gate
delay

delay for RC model

$$t_{PRC} \approx L_{int} \Rightarrow \sqrt{\frac{t_{propagation}}{0.38 R C}}$$

depends on

- sizing of driving gate
- choosing interconnect material.

(ii) RC delay should be considered when min (full) time

at the flip flop is smaller than $R C$, the min (full) time of the line.

$$t_{min} < R C$$

- with $R + C$ the total resistance & cap. of wire

* When this cond'n is not met, the changes in signal is slower than the propagation ~~to~~ delay of the wire, and a lumped capacitive model suffices.

Note: (i) 3-segment π -model is more accurate (is upto 30%)

(ii) But we commonly use single-segment π -model

* This π -model is nothing but the lumped RC model.

A lumped RC model is intermediate b/w lumped cap. & distributed RC.

Q. Now does simulator works?

→ Simulator works on matrix.

Matrix size increases with more no. of elements or nodes.

It

- Simulation time is less for π -model than distributed RC.
• ~~discrete~~ [lumped RC]

- Simulation time is ~~less~~ more for π -3 model than π -1 model.

Q.* "Contacts" & "vias" have completely different resistance profile than the ~~vias~~ "wires" (why?)

- If you want to decrease the resistance of wire you should its ~~area~~.
- If you want to decrease the resistance of via you don't have any control. Therefore you make the multiple resistances in parallel.

* Cost of fabricating via is more than metals. (Why?)

DFM → Design for manufacturability

⇒ which checks the no. of contacts in our design.

- if increasing the no. of vias not affecting the area of the chip then we can do that.

Q.* How many contacts we can put in our design or mould?

→ If we instantiate a pmos (pull) from the cadence of width $w_p = 10\mu m$ then it would have a number of contacts then it is good or bad?

- Increasing the no. of contact \rightarrow resistance.

- Increasing the no. of contact \rightarrow gate-contact capacitance

⇒ So we keep the no. of contacts such that our requirement is filled & both g-cap. & resistance is at good value.

• $\Delta f = \frac{1}{2\pi} \cdot \frac{1}{d} \cdot \frac{1}{2\pi} \cdot \frac{1}{d} = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{2\pi} \cdot \frac{1}{d} \cdot \frac{1}{2\pi} \cdot \frac{1}{d} = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{2\pi} \cdot \frac{1}{d} \cdot \frac{1}{2\pi} \cdot \frac{1}{d} = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{2\pi} \cdot \frac{1}{d} \cdot \frac{1}{2\pi} \cdot \frac{1}{d} = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{2\pi} \cdot \frac{1}{d} \cdot \frac{1}{2\pi} \cdot \frac{1}{d} = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• **6. Crosstalk effect:** $\text{werte Pg. 222 (6.3.3)}$

• $\text{werte Pg. 222 (6.3.3)}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

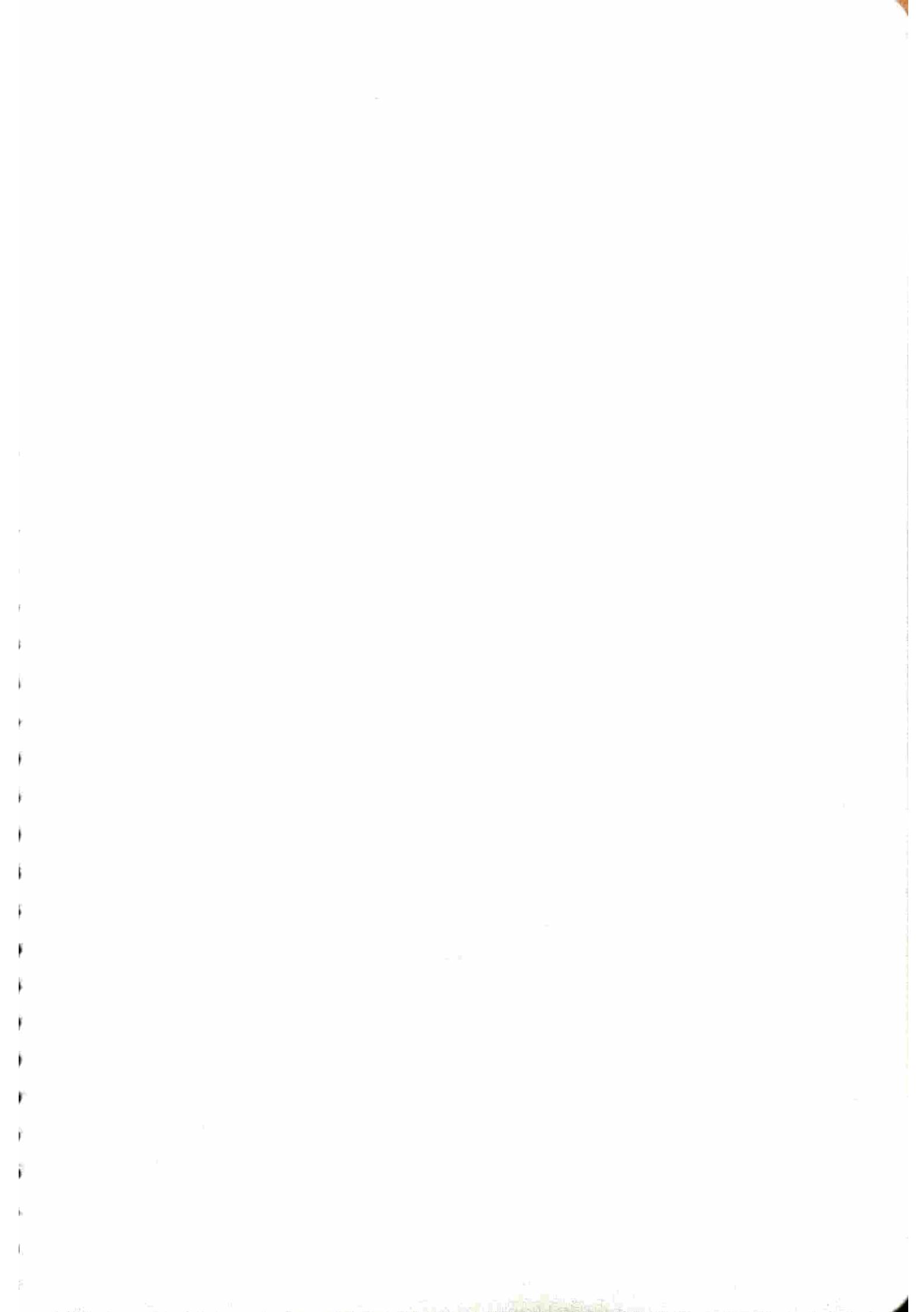
• $\Delta f = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

• $\Delta f = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$

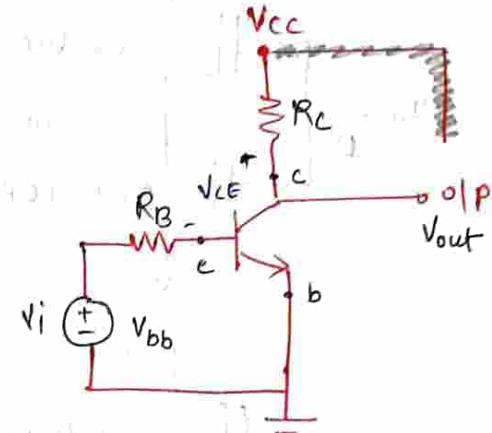
• $\Delta f = \frac{1}{4\pi^2 d^2}$ \rightarrow $\Delta f = \frac{1}{4\pi^2 d^2}$



3. Design of an Inverter

BJT Inverter

- When the input voltage is low, the transistor does not conduct, there is no current through R_C , and the collector voltage ($V_{O/P}$) is pulled up to V_{CC} .

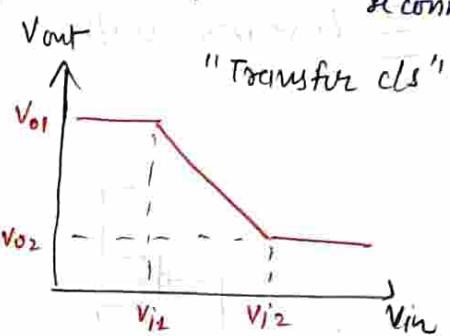


- When the input voltage is increased, the transistor begins to conduct, the voltage drop across R_C starts increasing, and the output voltage falls.

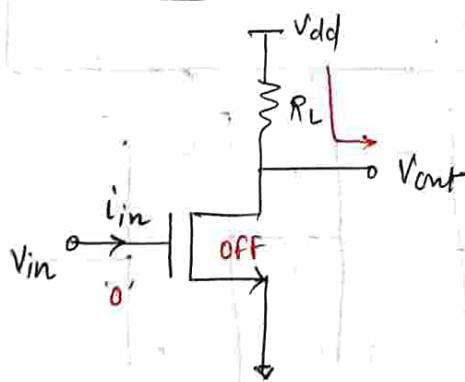
* when V_{bb} is less
 $\therefore V_{eb} < 0.7V$
 $\Rightarrow O/P = \text{high}$

- When the input voltage is high enough to drive the transistor into saturation, there is a small fixed drop ($V_{CE} = 0.1$ to $0.2V$) across the transistor, & output voltage saturates to its low value.

* when $V_{bb} \rightarrow \text{high}$
 $\therefore I_c$ is high &
 $O/P = \text{low(GND)}$

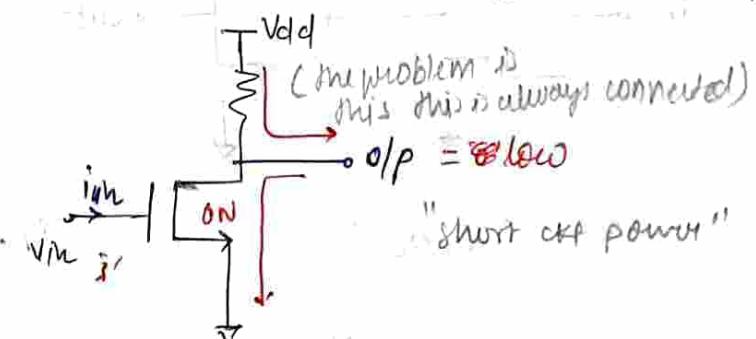


MOS Inverter



- if $V_{in} = 0$, MOSFET is off

$$\therefore V_{out} = V_{dd} \quad '1'$$

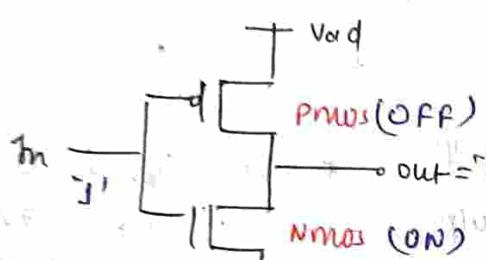
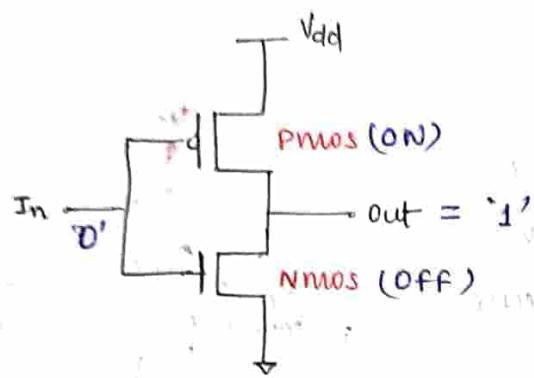


- if $V_{in} = 1$,

MOSFET is ON

$$\therefore V_{out} = \text{GND} \quad '0'$$

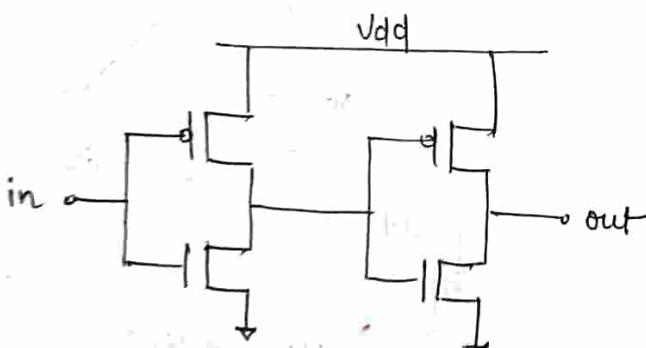
CMOS Inverter



why CMOS?

- At steady state, no power consumption.
- Better noise margin

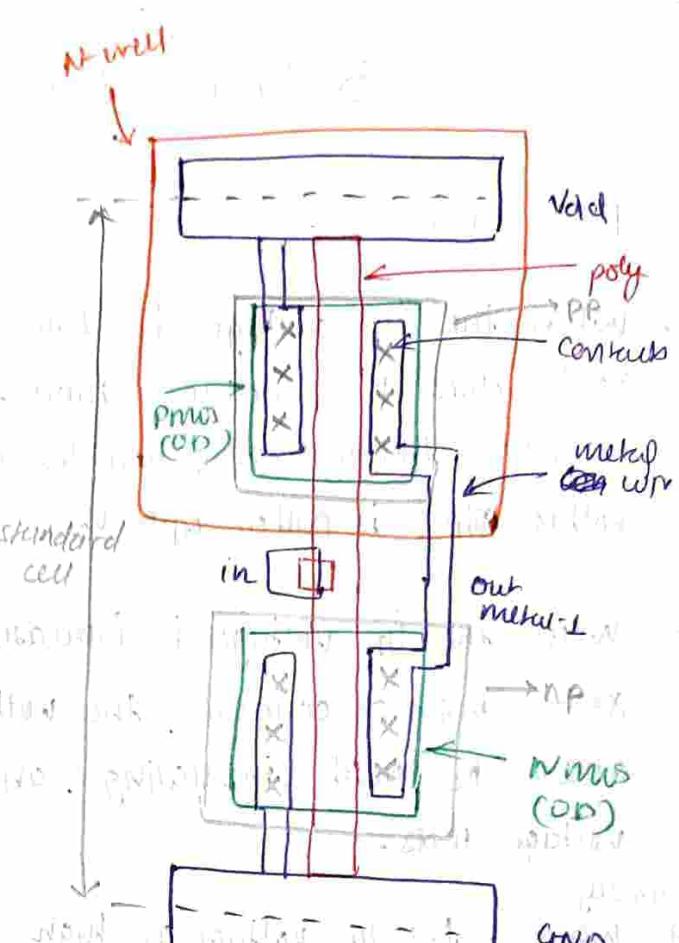
2nd inverter with 'straps' (taps)



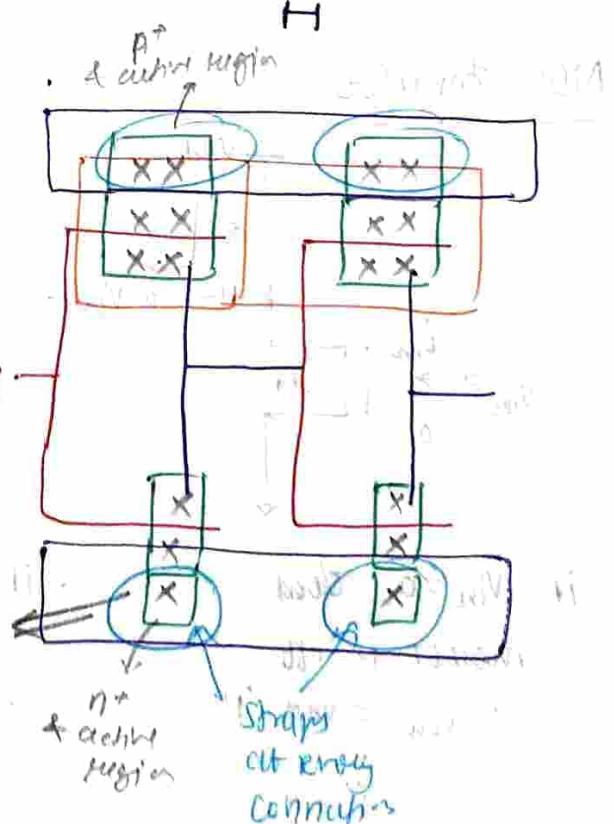
"Buffer"

- We can also use this straps later at the required place also.

- * Straps are required only at 32um distance.

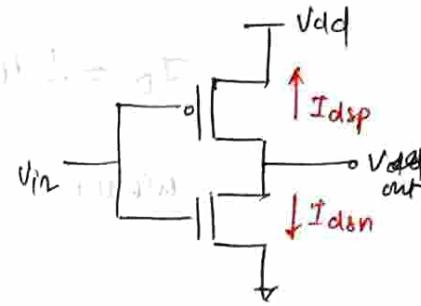


$1.2 \mu\text{m} = 2A$



DC Response

- DC Response → If my input is held constant at a particular voltage level at what voltage my output be held at.
 - ↳ V_{out} v/s V_{in}. (V_{th} for gate)
 - for CMOS: when $V_{in} = 0$, $V_{out} = V_{dd}$
when $V_{in} = V_{dd}$, $V_{out} = 0$
 - In b/w, the V_{out} depends on transistor size & current.
 - By KCL, $|I_{dsn}| = |I_{dsp}|$
 - We could solve equations but graphical solution give more insight.
- Note → DC response leads to static delay → transient response



MOSFET: Region of operation

NMOS

Cutoff

$$V_{gsn} < V_{tn}$$

$$V_{in} < V_{tn}$$

Linear

$$V_{gsn} > V_{tn}$$

$$V_{in} > V_{tn}$$

$$V_{dsn} < V_{gsn} - V_{tn}$$

$$V_{out} < V_{in} - V_{tn}$$

$$\begin{cases} V_{gsn} = V_{in} \\ V_{dsn} = V_{out} \end{cases}$$

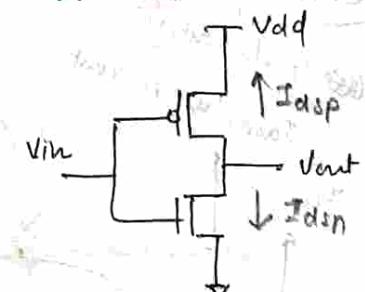
Saturation

$$V_{gsn} > V_m$$

$$V_{in} > V_m$$

$$V_{dsn} > V_{gsn} - V_{tn}$$

$$V_{out} > V_{in} - V_{tn}$$



PMOS

Cutoff

$$V_{gsp} > V_{tp}$$

$$V_{in} > V_{DD} + V_{tp}$$

Linear

$$V_{gsp} < V_{tp}$$

$$V_{in} < V_{DD} + V_{tp}$$

$$V_{dsp} > V_{gsp} - V_{tp}$$

$$V_{out} > V_{in} - V_{tp}$$

$$\begin{cases} V_{gsp} = V_{in} - V_{DD} \\ V_{dsp} = V_{out} - V_{DD} \end{cases}$$

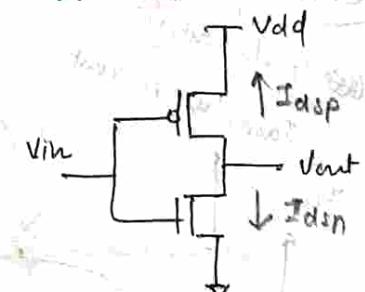
Saturation

$$V_{gsp} < V_{tp}$$

$$V_{in} < V_{DD} + V_{tp}$$

$$V_{dsp} < V_{gsp} - V_{tp}$$

$$V_{out} < V_{in} - V_{tp}$$



NMOS - PMOS : I-V c/s

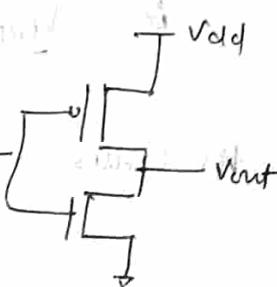
We are assuming that $B_n = B_p$ by making the PMOS wider than the NMOS.

$$I_D = \frac{1}{2} \mu_n C_{ox} \frac{W}{L} (V_{GS} - V_m)^2$$

$$\text{where, } B_n = \mu_n C_{ox} \frac{W_n}{L_n}$$

$$B_p = \mu_p C_{ox} \frac{W_p}{L_p}$$

in saturation

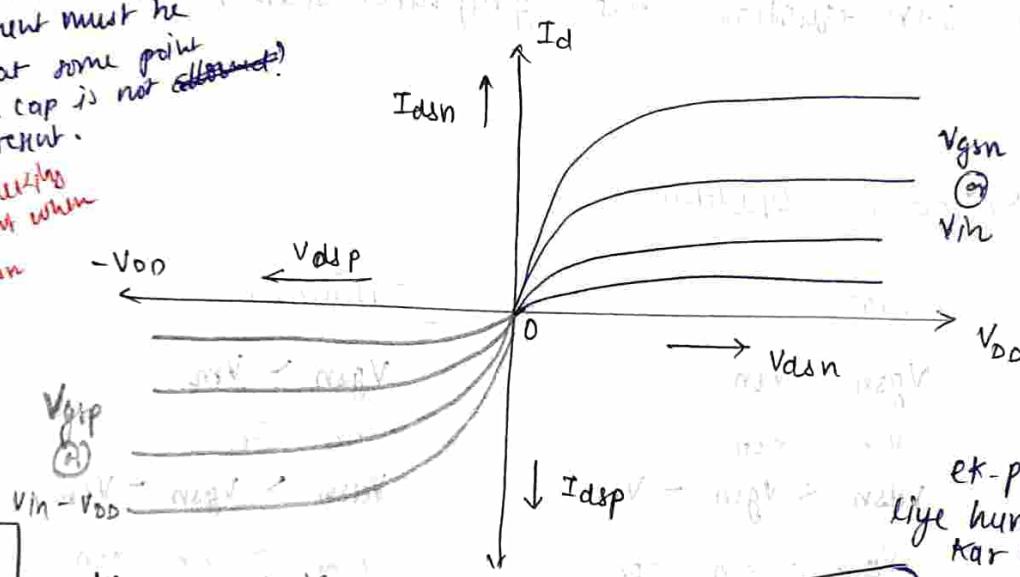


Since, $\mu_n \neq \mu_p$ or $\mu_n > \mu_p$

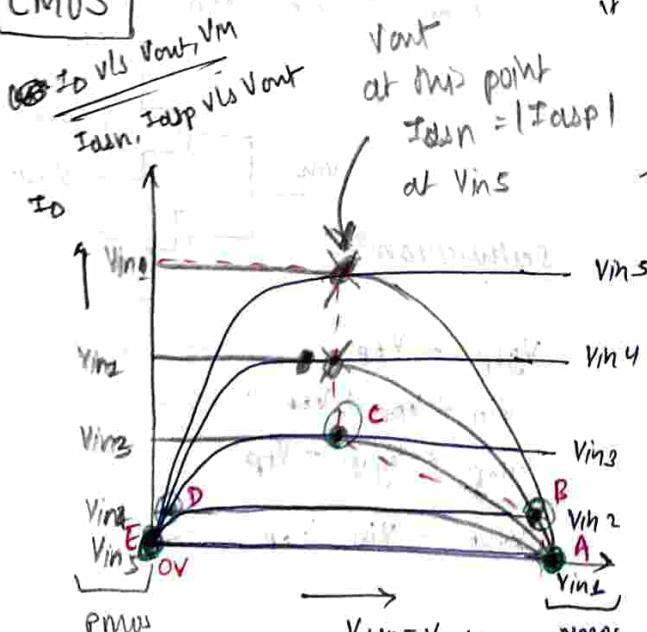
\therefore To make $B_n = B_p$ we make wider PMOS $\therefore W_p > W_n$.

* Current must be same at some point when cap is not present.

We are discussing at the point when $I_{dsn} = I_{dsp}$

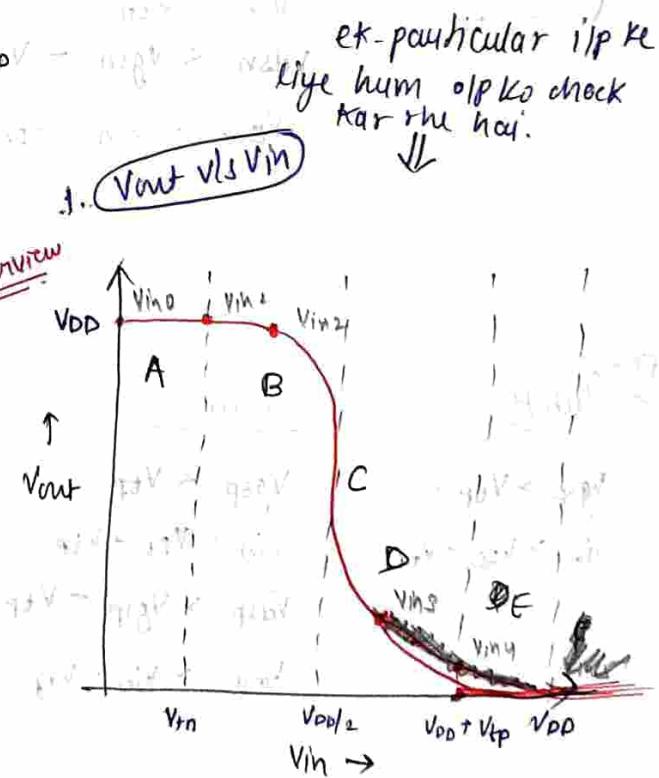


CMOS

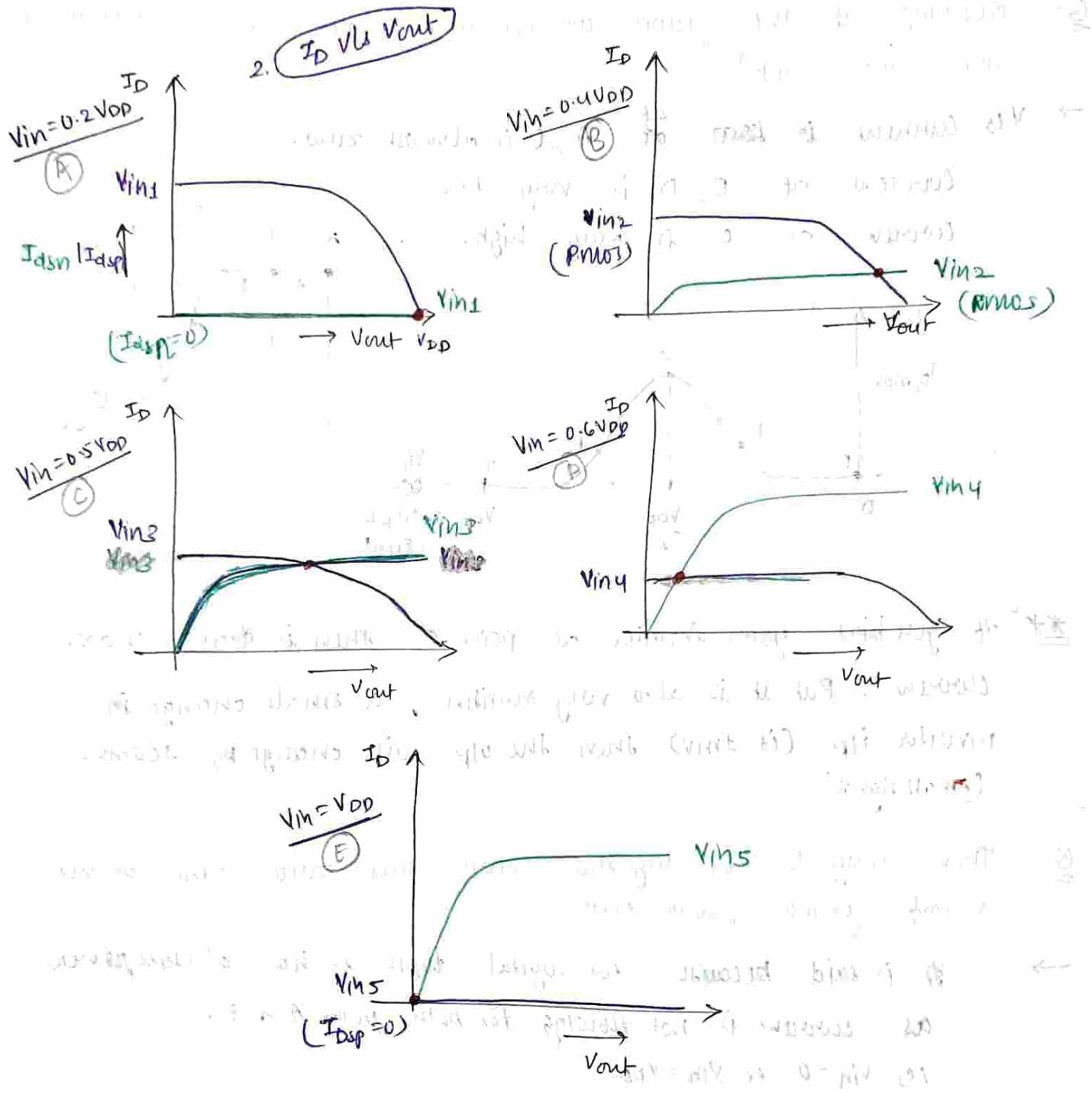


$$V_{ds} = V_{out} - V_{dd}$$

"load line Analysis"



"DC trans for curv of CMOS"



Region

NMOS $V_{in} = 0.2V_{DD}$ A

cutoff

linear

p MOS $V_{in} = 0.4V_{DD}$ B

saturation

linear

NMOS $V_{in} = 0.5V_{DD}$ C

saturation

saturation

p MOS $V_{in} = 0.6V_{DD}$ D

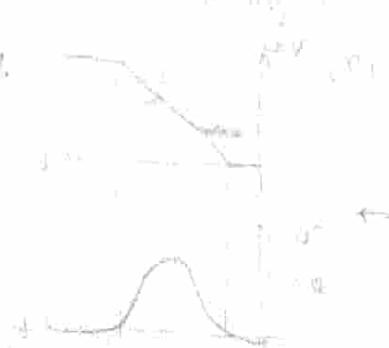
linear

saturation

NMOS $V_{in} = V_{DD}$ E

linear

cutoff



Q: Looking at the graph can we also predict the current (I_D) from the graph?

→ Yes, current at A, E is almost zero.

Current at B, D is very low.

Current at C is very high.

N-linear P-saturation

N-saturation N-saturation

N-saturation P-linear

V_{in}

V_{DD} Current Count

N-saturation P-linear

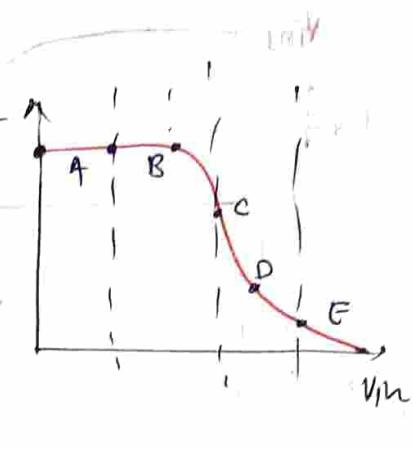
V_{DD}/2

-10mV

N-linear P-linear

E

(small signal)



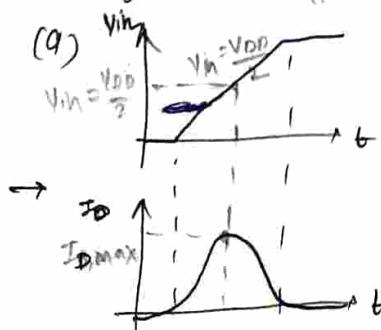
** If you bias your inverter at point C, then it draws max. current. But it is also very sensitive, a small change in inverter input ($\pm 1\text{mV}$) then the output will change by 100mV .

Q: Then why do we say that CMOS has zero static power & only dynamic power exists.

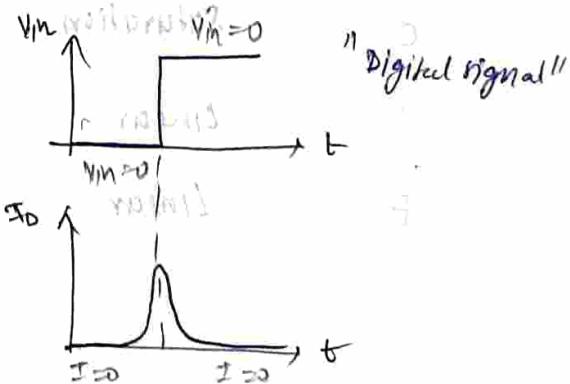
→ It is said because for digital logic it has '0' static power as current is not flowing for both points A + E.
i.e., $V_{in} = 0$ or $V_{in} = V_{DD}$.

Also for analog it has some static power when small signal voltage ΔV_{in} (small change in input voltage) is applied when operating in region C ($V_{in} > V_{DD}/2$). (soft-nonsaturation region)

Q: What would be the value of I_D over time if V_{in} vs time is given.



Small signal
Analog



- * So, no current is obtained for digital signal hence, static power = 0
- Any other signal with $V_{in} \neq 0$, it has some static power.
- And, for small signal (V_{in}) $I_D \neq 0$.

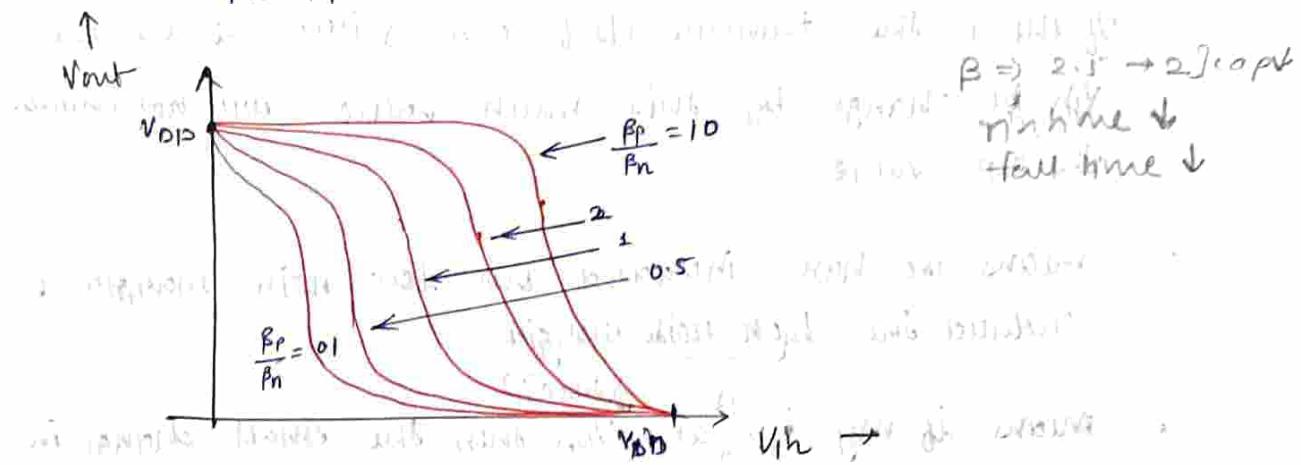
m(b)

- * For very less sum and full time of V_{in} , we don't pay so much of attention in region 'c', we directly move from E to A.

Part 2

B-Ratio [With pg 90 (2.5.2)]

- What would happen if $\beta_n > \beta_p$?
- The transfer cts of V_{out} vs V_{in} (i.e., CMOS) will move towards left.
- for $\beta_n > \beta_p \rightarrow$ curve shift left
- $\beta_n < \beta_p \rightarrow$ curve shift right



- If $B_p/B_n \neq 1$, switching point (i.e., region 'c') will move from $V_{DD}/2 \rightarrow$ called skewed gate

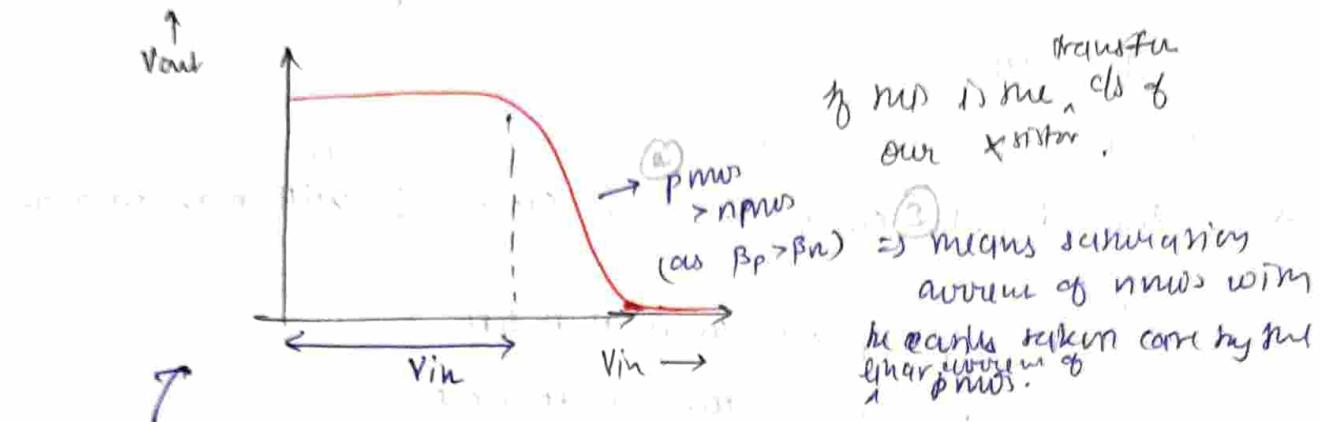
- Other gates will also have similar kind of V_{out} vs V_{in} cts
- not same because their B-ratios can also be different.

- The method that we have used to draw the transfer cts of CMOS inverter, we can also use this method to draw cts of other gate also. [for inverting gates only]
- for non-inverting gates the transfer cts may change.

Q: As a designer, what is the relevance of the change in the β -value?

→ By changing the β -value, you change the noise margin of your inverter.

• At times I don't want my inverter to switch.



① If nmos is the transfer func of our inverter & let me V_{in} be changed by this much value will not change the o/p value

② means we have increased our low-noise margin & reduced the high noise margin.

• means if my o/p is at V_{DD} then the small change in V_{in} in my o/p will go to 0.

E: 8:02 now $\beta_p > \beta_n$ but why is it shifting towards a right?

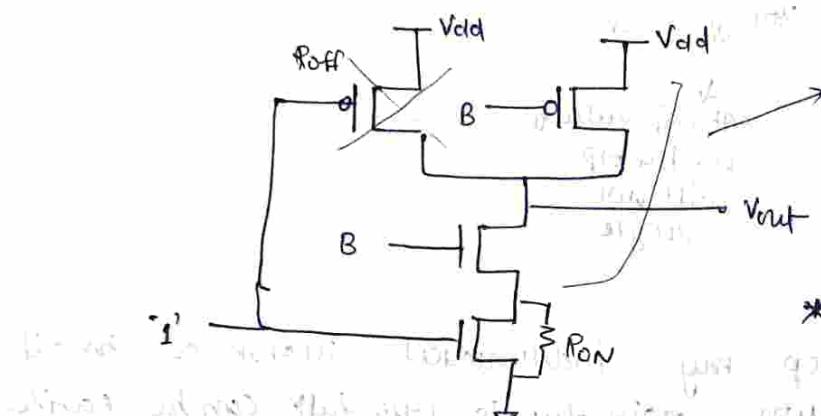
→ if $\beta_p > \beta_n$ then the saturation current of nmos is much less than the linear current of pmos.

$I_{dn(\text{sat})} < I_{dp(\text{linear})}$
So even the nmos is trying to sink as much current as it can, the linear current of pmos is so high that it kind of doesn't allow the o/p cap. to discharge.

Q. What is pull up & pull down now?

Q. Now the other gates will collapse into the inverter c/s?

→ Let us assume a NAND gate



Therefore only this c/s remains and it is now act as like NAND gate & collapsed into an inverter.

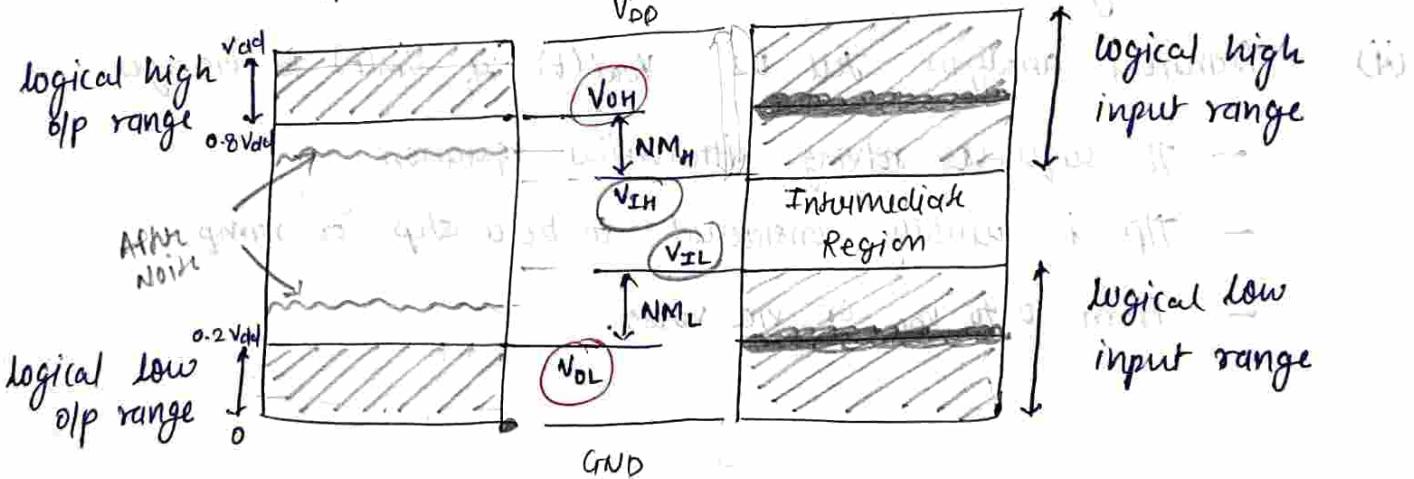
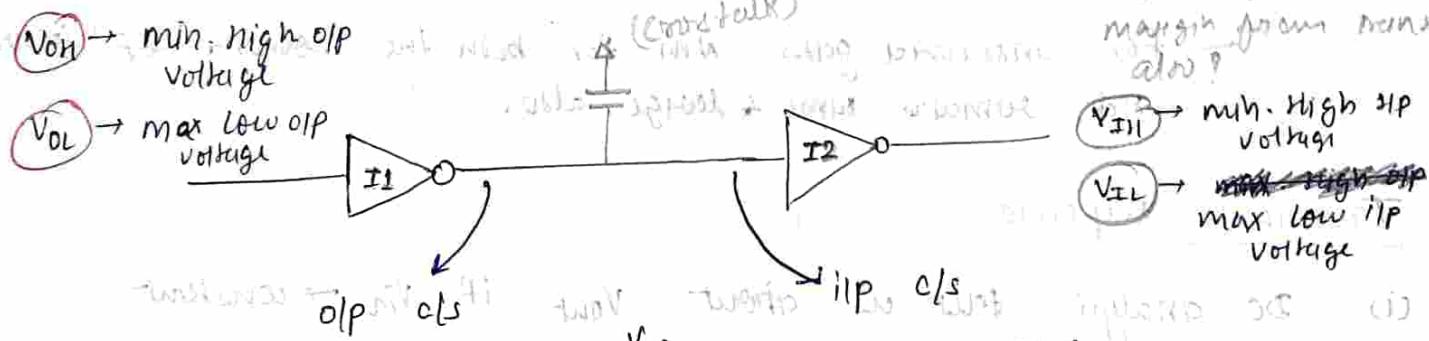
* So our main job is to select some i/p in other inverter gates to get the c/s of inverter.

Part - 2

Noise Margin

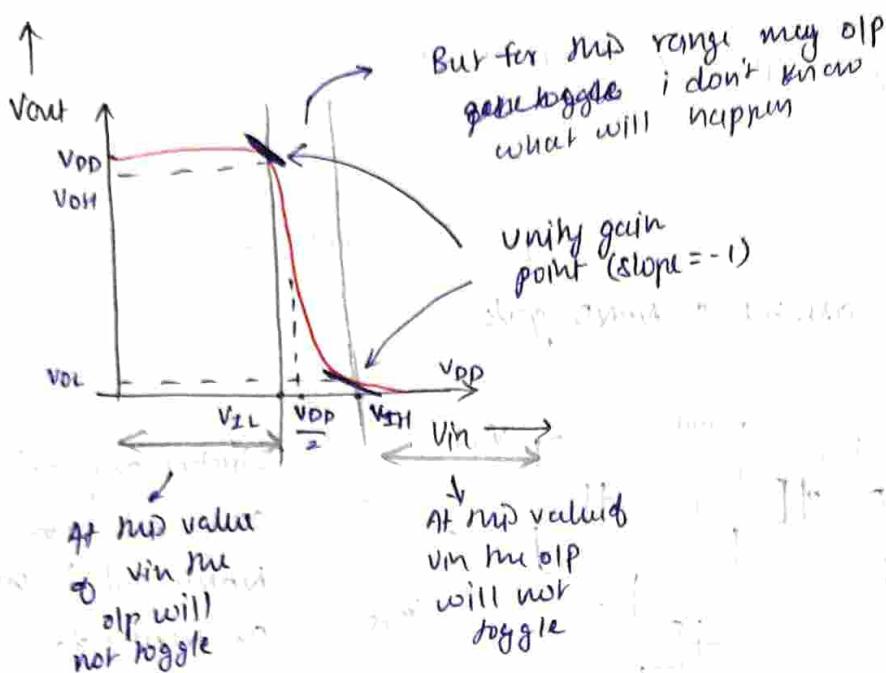
(Weste - Pg. 94) After reading this, it uses noise margin to do DC analysis.

If we get noise margin from inverter alone?



- As the o/p from I1 goes out due to some crosstalk noise added to the o/p voltage.
- so the noise margin ($NM_H + NM_L$) should be such that the voltage coming from I1 should be easily read by the i/p of the I2 inverter.

4C



- so we will try to keep my intermediate region as small as possible so maximum noise due to overshoot can be easily handled.

- for ~~skewed~~ skewed gate NM of one side ~~more~~ increases than NM of other side decreases
- for unskewed gates NM for both the inverters ~~is~~ higher side remains same & large also.

Transient Response

Wade Pg. 443 (4.2)

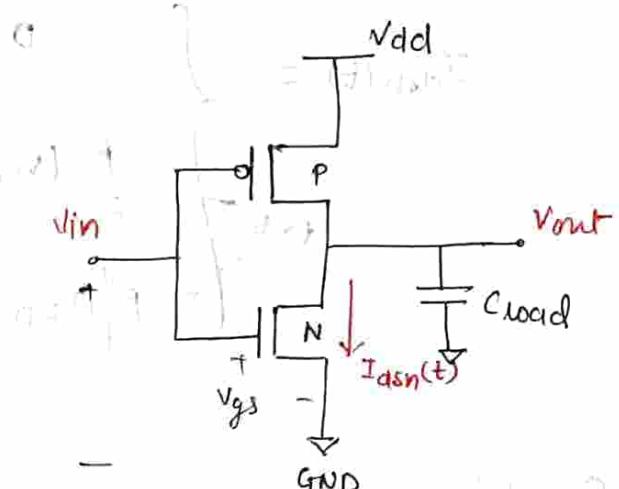
- DC analysis tells us about V_{out} if $V_{in} \rightarrow \text{constant}$
- Transient analysis tell us $V_{out}(t)$ if $V_{in}(t) \rightarrow \text{changes}$
 - It requires solving differential equation
 - t_{lp} is usually considered to be a step or ramp
 - From 0 to V_{DD} or vice-versa.

1. Delay Models :

inverter step response:

$$\cdot V_{in(t)} = \begin{cases} 0 & t < t_0 \\ V_{DD} & t \geq t_0 \end{cases}$$

$$= u(t - t_0) V_{DD}$$



Vout(t)

$$\cdot V_{out}(t) = \begin{cases} V_{DD} & t < t_0 \\ ? & t \geq t_0 \end{cases}$$

$$\underline{V_{gs} = V_{in}}$$

$$\therefore \frac{dV_{out}(t)}{dt} = -\frac{I_{dsn}(t)}{C_{load}}$$

t < t₀:
charging
of cap.

[NMOS - OFF, PMOS - ON]

$$I_{dsn} = \begin{cases} 0 & V_{ds} < V_{th} \\ \frac{\beta}{2} (V_{gs} - V_{th})^2 & V_{ds} > V_{th} \end{cases}$$

$$\beta [(V_{gs} - V_{th}) V_{ds} - \frac{V_{ds}^2}{2}]$$

$$V_{gs} = 0V \checkmark$$

$$V_{out} = V_{DD}$$

$$V_{ds} = V_{DD} \checkmark$$

$$V_{ds} < V_{th}$$

$$V_{DS} > V_{GS} - V_{th}$$

$$V_{DS} < V_{GS} - V_{th}$$

— ① —

initial
state

t > t₀:
discharging
of cap.

[NMOS - ON, PMOS - OFF] \rightarrow $V_{gs} = V_{DD}$

$$+ V_{ds} = V_{out}$$

$$V_{ds} = V_{DD} > V_{GS} - V_{th}$$

∴ NMOS is in saturation region.

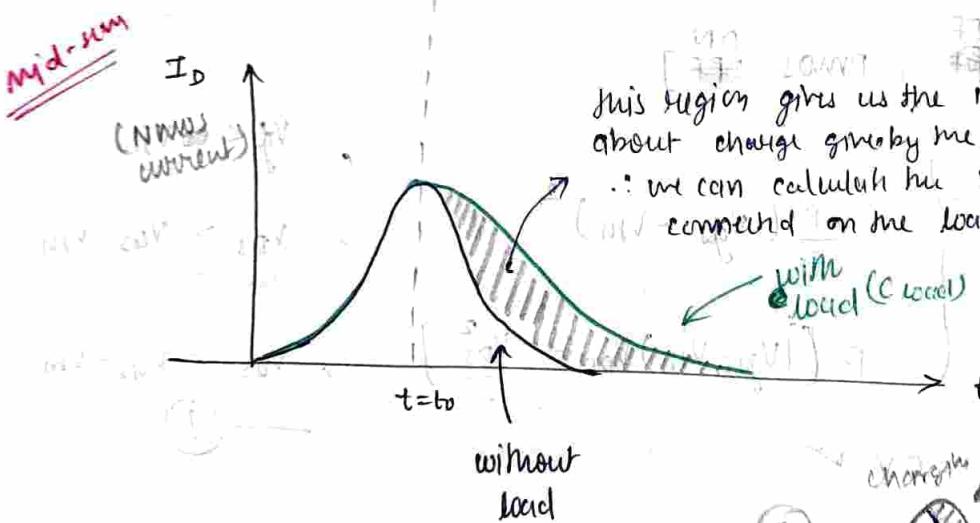
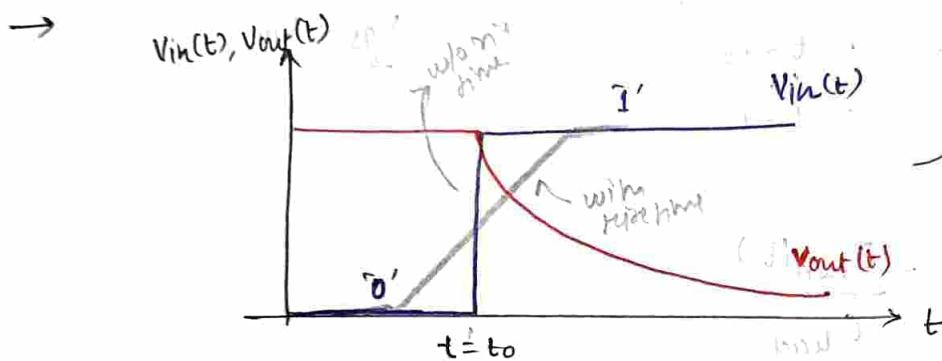
when $t \geq t_0$, $V_{ds} = V_{out}$ falls below $V_{DD} - V_{th} = 1.4V$

∴ NMOS comes in linear region.

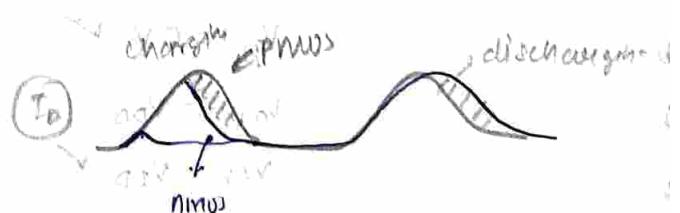
$t > t_0$ in eqn. ②.

$$I_{\text{out}}(t) = \begin{cases} 0 & t < t_0 \\ \frac{\beta}{2} \frac{(V_{DD} - V_m)^2}{V_{AS}} & V_{\text{out}} > V_{DD} - V_m \\ \frac{\beta}{2} \left[\frac{(V_{DD} - V_m)}{V_{AS}} V_{\text{out}} - \frac{V_{\text{out}}^2}{2} \right] & V_{\text{out}} < V_{DD} - V_m \end{cases}$$

Q: What will be the impact of C_{load} on current waveform?



This region gives us the information about charge given by the capacitor. ∴ we can calculate the value of cap. connected on the load.



without load : $I_{\text{out}} = I_{\text{up}}$

with load : $I_{\text{out}} \neq I_{\text{up}}$

no load

* when $V_{in} = 0V$ '0' Nmos - OFF PMOS - ON from V_{DD} ∴ I_{up} flows to V_{DD} (erwach)

$V_{in} = 1V$ '1' Nmos - ON PMOS - OFF

∴ I_{out} flows from V_{DD} to GND (erwach)

without load

* $V_{in} = 0 \Rightarrow NMOS - OFF, PMOS - ON$:

∴ I_{dsn} flows from V_{DD} to Cloud & charges the cap.

$V_{in} = 1, NMOS - ON, PMOS - OFF$:

∴ I_{dsn} flows from Cloud to GND

due to RC delay the current transition

takes more time discharging the cap.

Delay Definitions

(i) $t_{pdR} = t_{PLH} \rightarrow$ Rise propagation delay

- When o/p switches from low to high, the time difference b/w the instant when o/p is at 50% of its max value & when the o/p is at 150% of its max value.

(ii) $t_{pdf} = t_{PHL} \rightarrow$ fall propagation delay

- When o/p switches from high to low, the time difference b/w the instant when o/p is

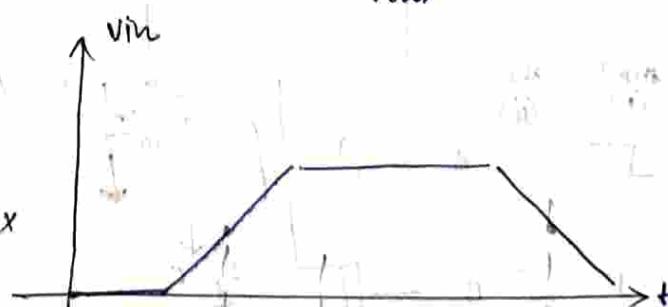
(iii) $t_{pd} \rightarrow$ average propagation delay

$$t_{pd} = (t_{pdR} + t_{pdf})/2$$

Buffer
 $V_{out} = V_{in}$

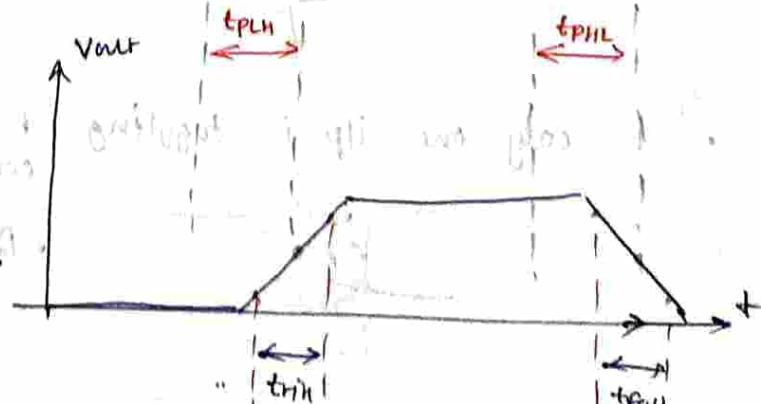
(iv) Rise time $\rightarrow (t_r)$

- The time taken by the signal to rise from 10% of the max value to 90% of max value.



(v) full time $\rightarrow (t_f)$

- The time taken by the signal to fall from 90% of max value to 10% of the max value.



(vi) $t_{cdR} \rightarrow$ Rising contamination delay

- minimum time from input rising to o/p rising at 50% of max value of o/p.

(vii) $t_{cdF} \rightarrow$ falling contamination delay

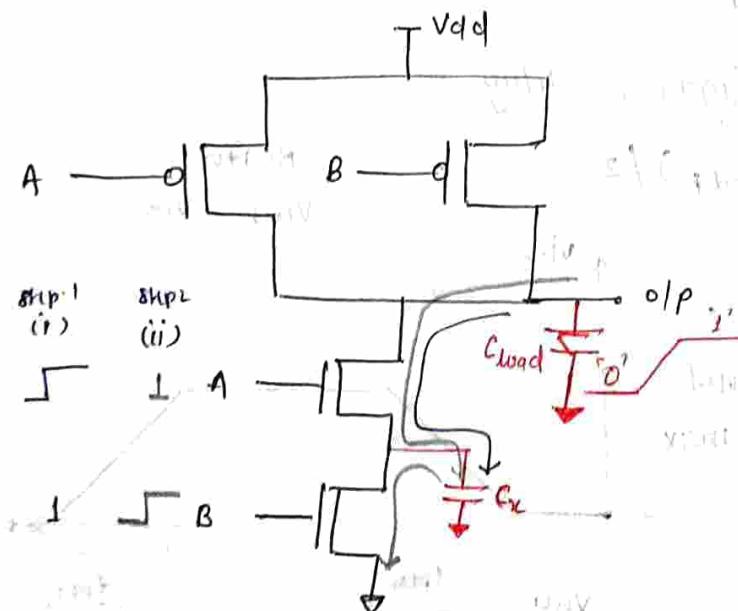
- minimum time from o/p falling/rising to o/p falling at 50% of max value of o/p.

(viii) $t_{cd} \rightarrow$ average contamination delay

$$t_{cd} = (t_{cdR} + t_{cdF}) / 2$$

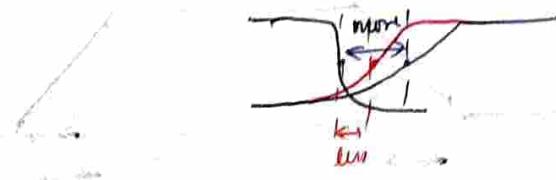
Note: When ilp changes, the o/p will change its old value for atleast the contamination delay & take on its new value in at most the propagation delay.

Eg: Understanding t_{cd} with example of 2NAND gate



A	B	\bar{AB}
0	0	1
0	1	1
1	0	1
1	1	0

- if only one ilp is toggling & other is constant then contamination delay is min



- But when both A & B are toggling together then contamination delay is less.

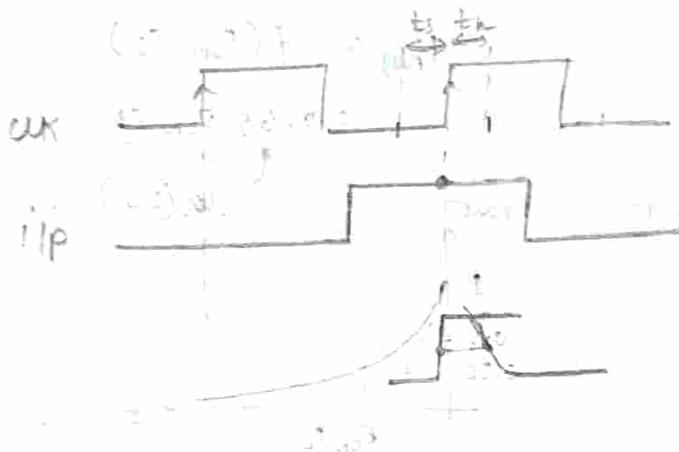
Q. Why contamination delay is important?

→ Contamination delay is important to validate hold time.
it Hold time error will appear if you don't consider contamination delay.

- Hold time is characterized by fastest delay. (t_{ch})
- Setup time is characterized by slowest delay. (t_{pd})

Setup time - the min. amount of time before the clock's active edge that the data must be stable for it to be latched correctly.

Hold time - the min. amount of time after the clock's active edge during which data must be stable.



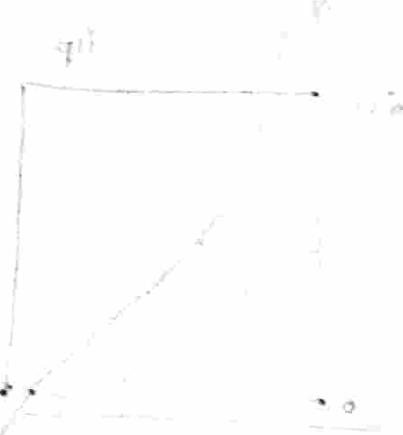
$$o/p = \overline{i/p}$$

pipelining.

throughput

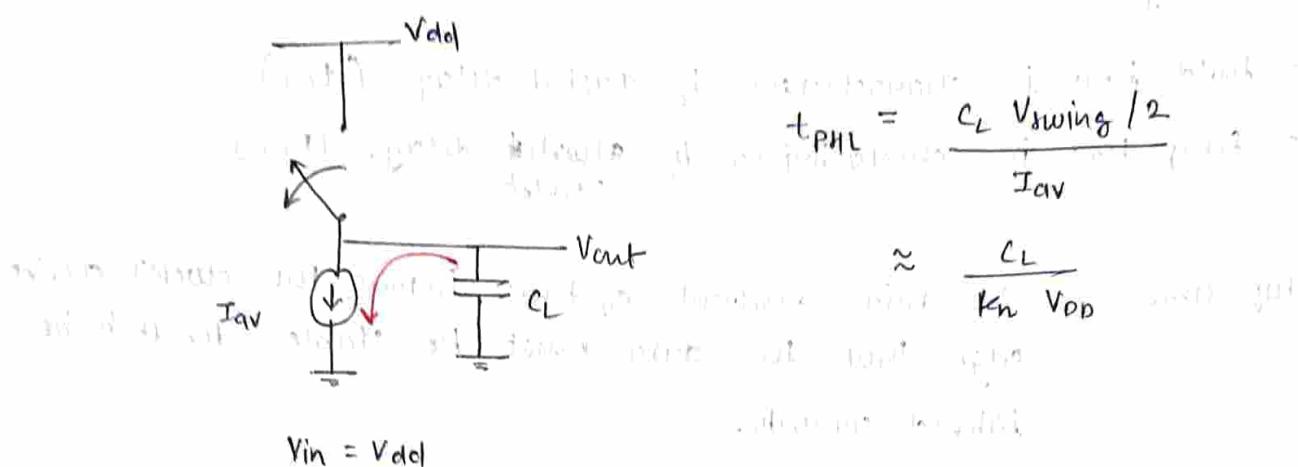
total = $t_p + t_{ch}$

A (current output) B (next output)

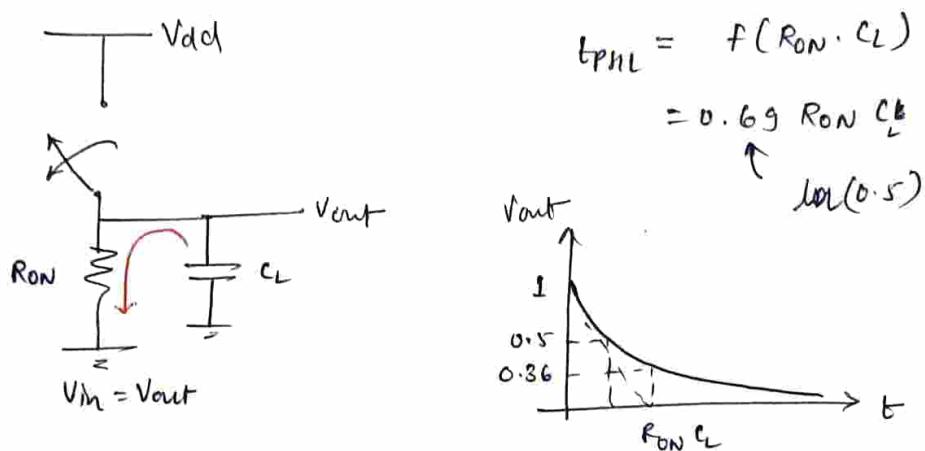


Delay models:

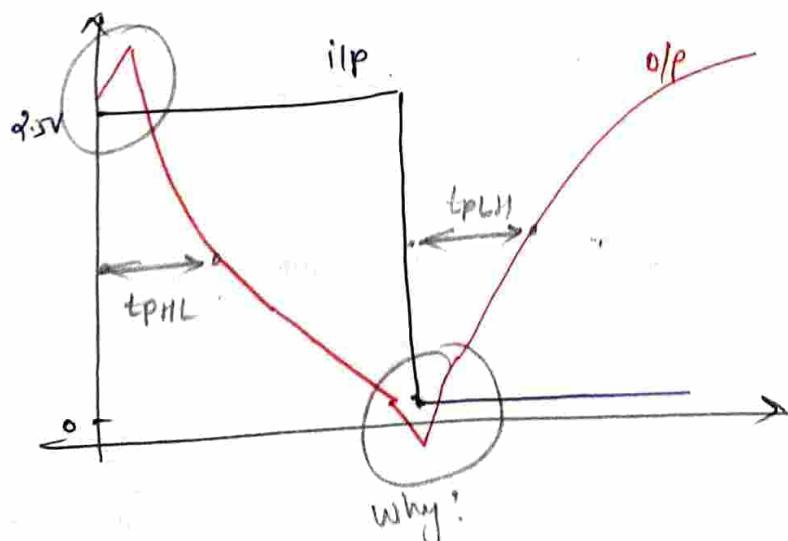
1. CMOS Inverter Propagation delay approach 1



2. CMOS Inverter Propagation delay approach 2 [Simpler]



Transient Response (of CMOS Inverters)

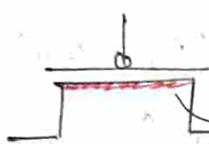


Q: Why the peak in V_{out} is observed?

→ Due to Miller capacitance

- As the gain i/p goes from 'two to one' there is coupling with source & drain
 - the source is gnd or vdd which doesn't change much
 - but drain is influenced by the change in i/p gain voltage
 - hence a little spike is observed.

Q: ⇒



more di diffusion layer here already but when you drive it ~~is~~ turned off then his charge has to go somewhere.

- then half of the charge goes to drain region & half goes to the ~~drain~~ source region.

Also,

$$t_p = 0.69 C_L \left[R_{eqn} + R_{eqp} \right]$$

Q: now to reduce the ~~switches~~ it

- Add extra cap. in parallel to the diffusion cap.

$$\Delta Q = C \Delta V$$

on $\uparrow C$, $\Delta V \downarrow$

$$\text{where, } R_{eq} = \frac{1}{V_{DD}/2} \int_{V_{DD}/2}^{V_{DD}} \frac{v}{I_{DSAT}(1+\lambda v)} dv$$

$$= \frac{3}{4} \frac{V_{DD}}{I_{DSAT}} \left(1 - \frac{7}{9} \lambda V_{DD} \right)$$

- Q: Why we don't want ~~gates~~?
 - ~~gates~~ may lead to breakdown.

with,

$$I_{DSAT} = k' \frac{W}{L} \left[(V_{DD} - V_T) V_{DSAT} - \frac{V_{DSAT}^2}{2} \right]$$

How to keep the delay less?

1. Keep load capacitance small.
2. Keep R_{eq} small
3. Increase V_{DD} , $R_{eq} \downarrow$, delay \downarrow
but at V_{DD} near subthreshold region, delay \uparrow exponential

1. By keep the capacitance small:-

- By $\downarrow C_L$.

- By reducing the self-loading capacitance.

- capacitance due to the biasing drain & source

2. By reducing the Row of the device (Nmos or pmos)

* On increasing the width of the device ($w_p + w_n$)

- (i) diff. cap. \uparrow , self-loading cap. \uparrow

- $\therefore C_{load}$ (overall) \uparrow

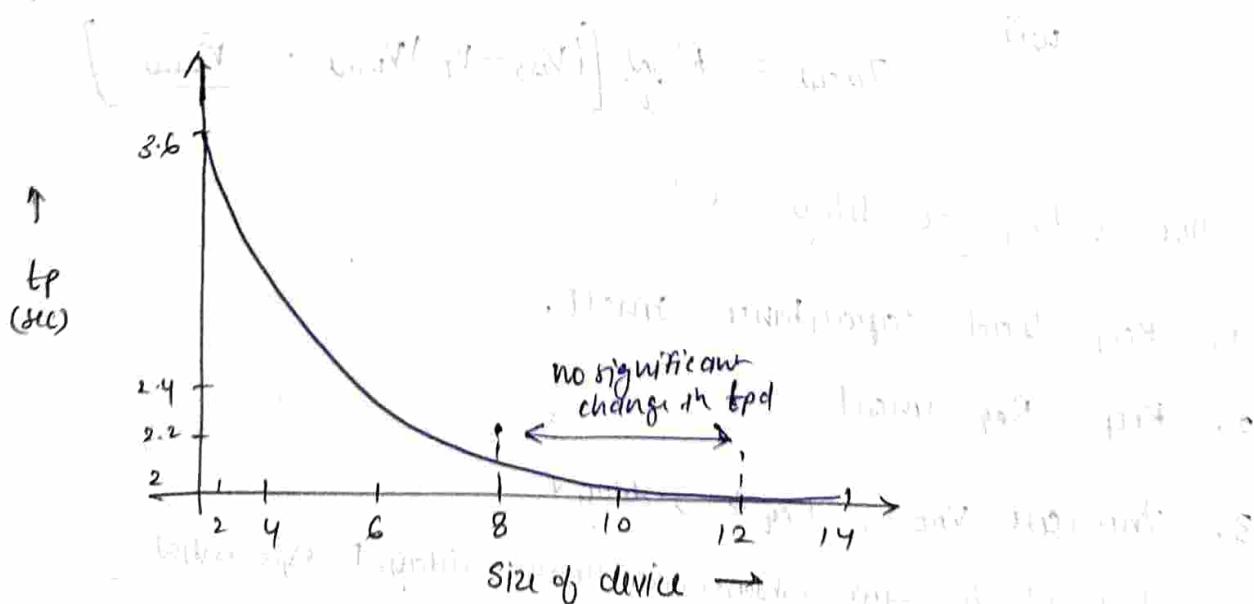
(ii) $R \propto \frac{1}{w} \therefore w \uparrow, R \downarrow$

\Rightarrow But effectively [on $\uparrow w$] $t_p \downarrow$

But on further increasing the value of w , the change in w is not significant.

Q. How large the logic gate should be to achieve least delay? (Ans?)

3. By increasing the VDD, $R_{eq} \downarrow$ } from formula
on $\uparrow V_{DD}$, $I_{D, sat} \uparrow \Rightarrow R_{eq} \downarrow$ } on previous page.

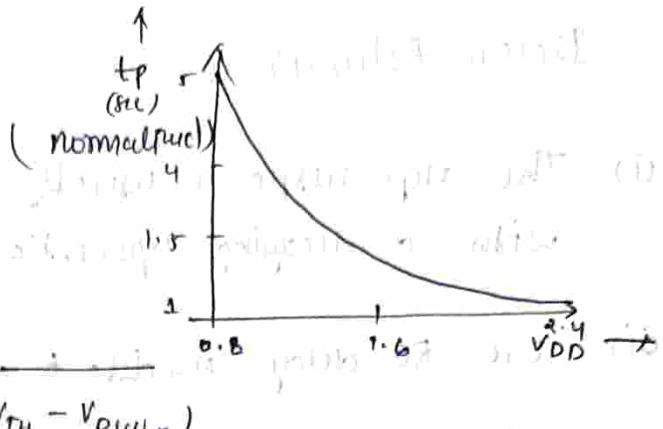


Delay as a funcⁿ of V_{DD} :

$$t_{PHL} = 0.693 \frac{3}{4} \frac{C_L V_{DD}}{I_{sat_n}}$$

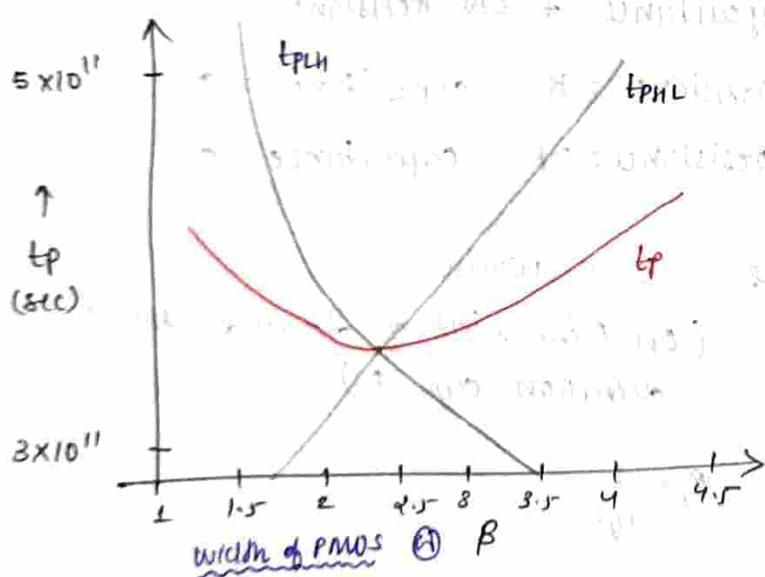
$$= 0.52$$

$$\frac{C_L V_{DD}}{(W/L)_n K_n (V_{DD} - V_{TH} - V_{DSAT_n})}$$



- Q: What are ways of reducing rise & fall time?
- min delay & fall delay (t_{PD}) increase for calculating delays.
 - rise & fall time are only calculating powers.
 - * the way we are reducing t_{PD} we can reduce t_r & t_f also.
 - load cap. $\uparrow \Rightarrow t_r$ & t_f \downarrow .

NMOS - PMOS : B ratios



∴ β -ratio, $(\frac{w_p}{w_n})$, should lie b/w 1.5 to 2.6.

(for 65nm technology)

- ③ Overall propagation delay first \downarrow as PMOS size \uparrow & time taken to charge the cap \downarrow & after further \uparrow in PMOS size due to significant effect of self loading t_p \uparrow as time taken to charge the cap load cap increases.

Q: On changing $\frac{w_p}{w_n}$

- how cap changes
- how delay changes
- for 1 cap, how $\frac{w_p}{w_n}$ will change
- $\frac{w_p}{w_n} \rightarrow \beta$ ratio of PMOS & NMOS

- ① as $\beta \uparrow$, size of PMOS \uparrow , more current will flow to the cap to charge the cap \downarrow . $\therefore t_{PHL} \downarrow$.
- ② as $\beta \uparrow$, size of PMOS \uparrow but NMOS size \downarrow less current is required to charge the cap. $\therefore t_{PHL} \uparrow$.

Delay Estimation

Pg. 146 (4.3)
Pg. 155 (4.4)

- (i) The step response usually looks like a 1st order RC response with a decaying exponential.
- (ii) Use RC delay models to estimate delay
 - $C = \text{total capacitance on o/p node}$
 - Use effective resistance R
 - So that $t_{pd} = RC$
- (iii) characterize x_{min} by finding their effective R .
 - Depends on avg current as gate switches
- (iv) Too inaccurate to predict current at any given time
 - But good enough to predict RC delay.

RC Delay model

- RC equivalent circuit of MOSFET is:
 - (i) Ideal switch + capacitance + ON Resistance
 - (ii) One nmos has, resistance = R , capacitance = C
 - (iii) One pmos has, resistance = $2R$, capacitance = C

$$* R \propto \frac{1}{W} \quad & C \propto \text{width}$$

(On nmos width of drain & source diffusion cap. \uparrow)

$$* R_n \propto \frac{1}{\mu_n} \quad & R_p \propto \frac{1}{\mu_p}$$

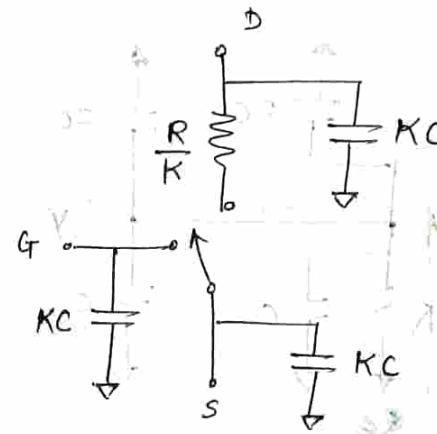
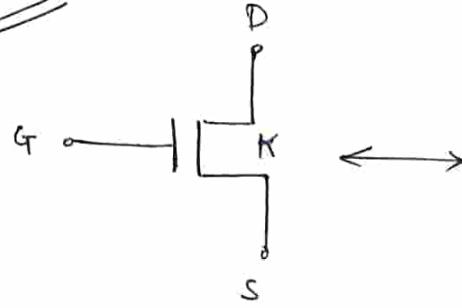
$$\therefore \text{let } \mu_n = 2\mu_p \quad (\mu_n > \mu_p) \\ (\because R_n < R_p)$$

$$\text{if } R_n = R$$

(parallel connection)

$$\therefore R_p = 2R$$

NMOS

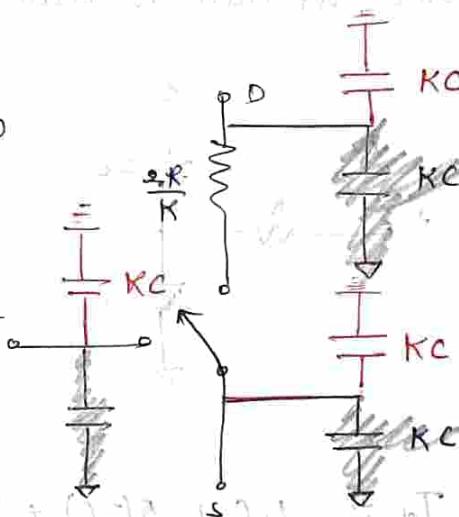
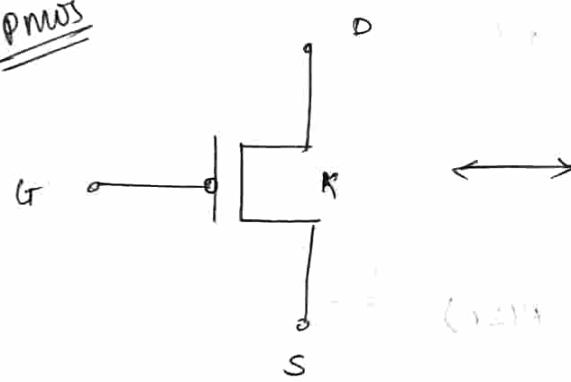


where K represents the k -times the width of the unit device size

$$\text{if } w_{n(\text{old})} = \alpha w$$

$$\therefore w_{n(\text{new})} = K \cdot w$$

PMOS



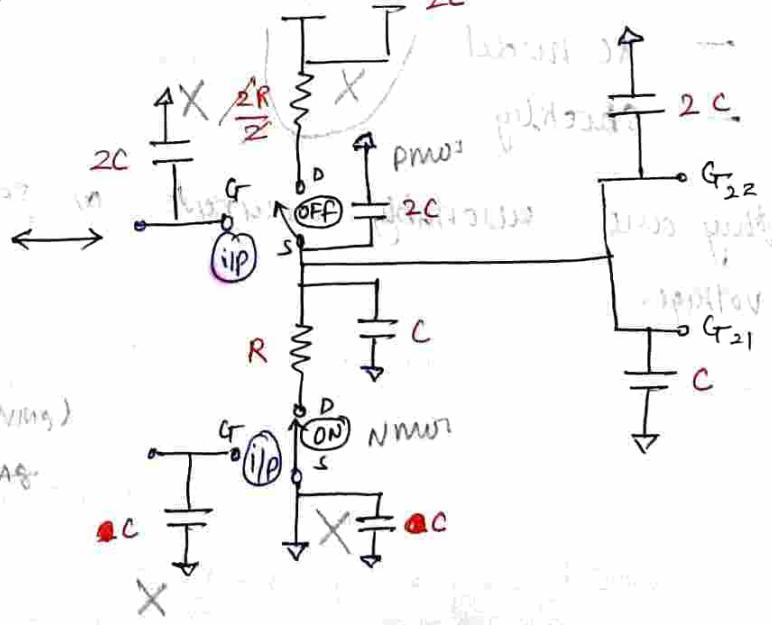
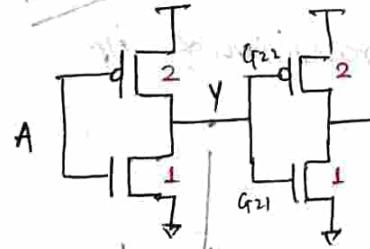
So make the inverter with same delay for both NMOS + PMOS we take $w_p = 2w_n$ (ie, $K=2$ for PMOS) (ie, $K=1$ for NMOS)

$$\frac{R_p}{K_n} = 2$$

Fanout Delay estimation

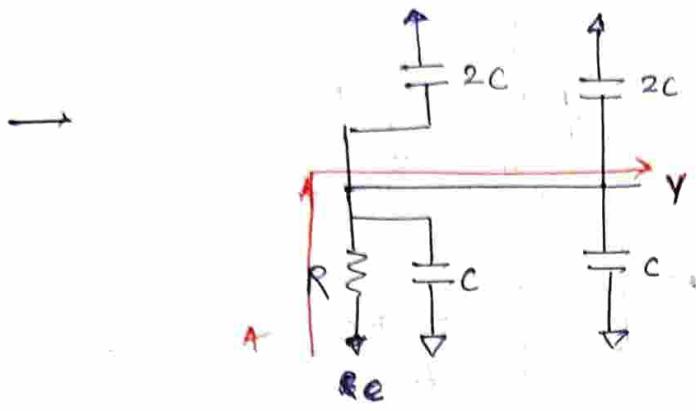
Estimate the delay of a fanout \rightarrow structure?
(when i/p goes from '0' to '1')

for fanout = 1

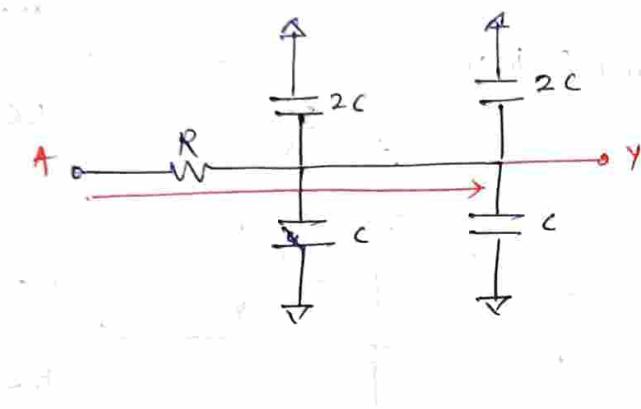


Fanout = 1
means No. of inverter (driving)
a minimum (drive) \rightarrow driving

driving.



Now we have to calculate delay for path AY



$$\begin{aligned}
 T_d &= RC + R(2C) + RC + R(2C) \\
 &= R(C + 2C + C + 2C) \\
 &= R(6C)
 \end{aligned}$$

** When try to calculate the delay in the path AY

we can use multiple Delay models.

- Spice formulation
- RC model
- Shockley

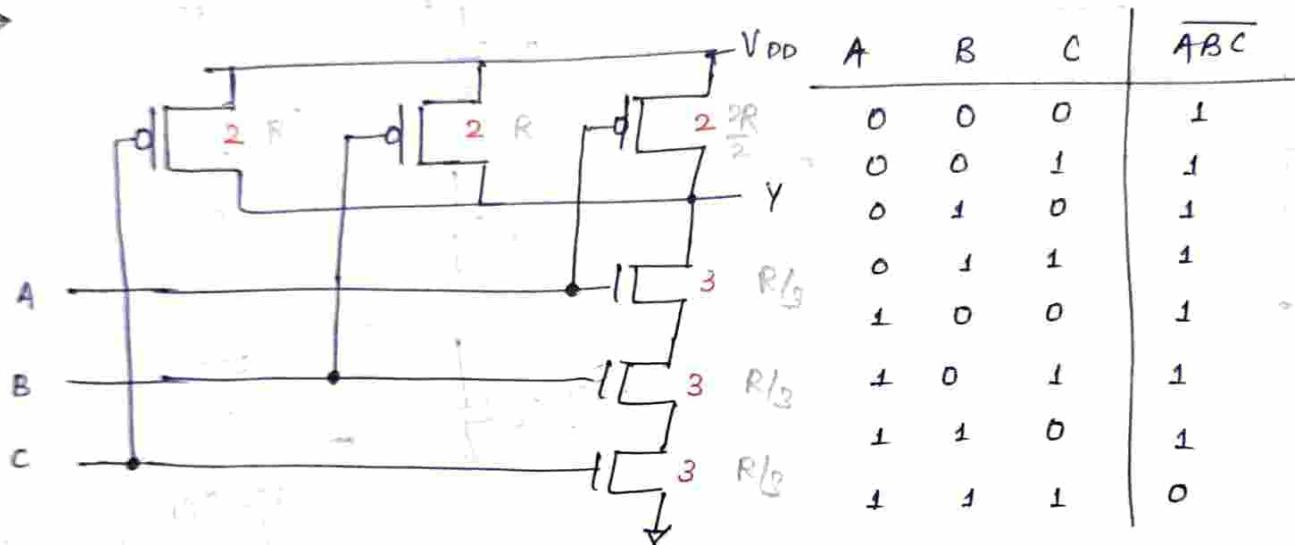
They are reasonably accurate at 50% transition of input voltage.

Example: "3-input NAND gate" work example \rightarrow 4.2
 Calculating worst case state transition & full delay (IP parameter) \rightarrow likely value.

- Q. Sketch a 3-input NAND with transistor width chosen to achieve effective rise & fall resistances equal to a unit inverter (R). Annotate the gate with its gain & diffusion capacitance. Assume all diffusion nodes are connected. Sketch equivalent circuit for: (Also calculate delay)

- (i) falling o/p transition (Best case + worst case) \rightarrow discharging
- (ii) rising o/p transition (Best case + worst case) \rightarrow charging

\rightarrow



Q. How to take the value of size of PMOS & NMOS (K_p & K_n)?

\rightarrow the size of PMOS & NMOS should be such that the (charging)
 (discharging)

$$R_{ON} \text{ for NMOS} = R_{ON} \text{ for PMOS} = R_{ON} \text{ (for the worst case)} \quad \text{(discharging)}$$

$$I_D \text{ for NMOS} = I_D \text{ for PMOS} \quad \text{(for the worst case)}$$

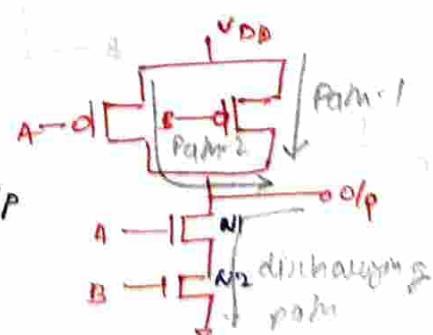
- (i) If we want to charge a load at o/p 'Y' then we can charge it through any path-1 or path-2 (for the worst case)
 i.e., not considering both path-1 & path-2 together

\therefore Both the paths should have resistance = R .

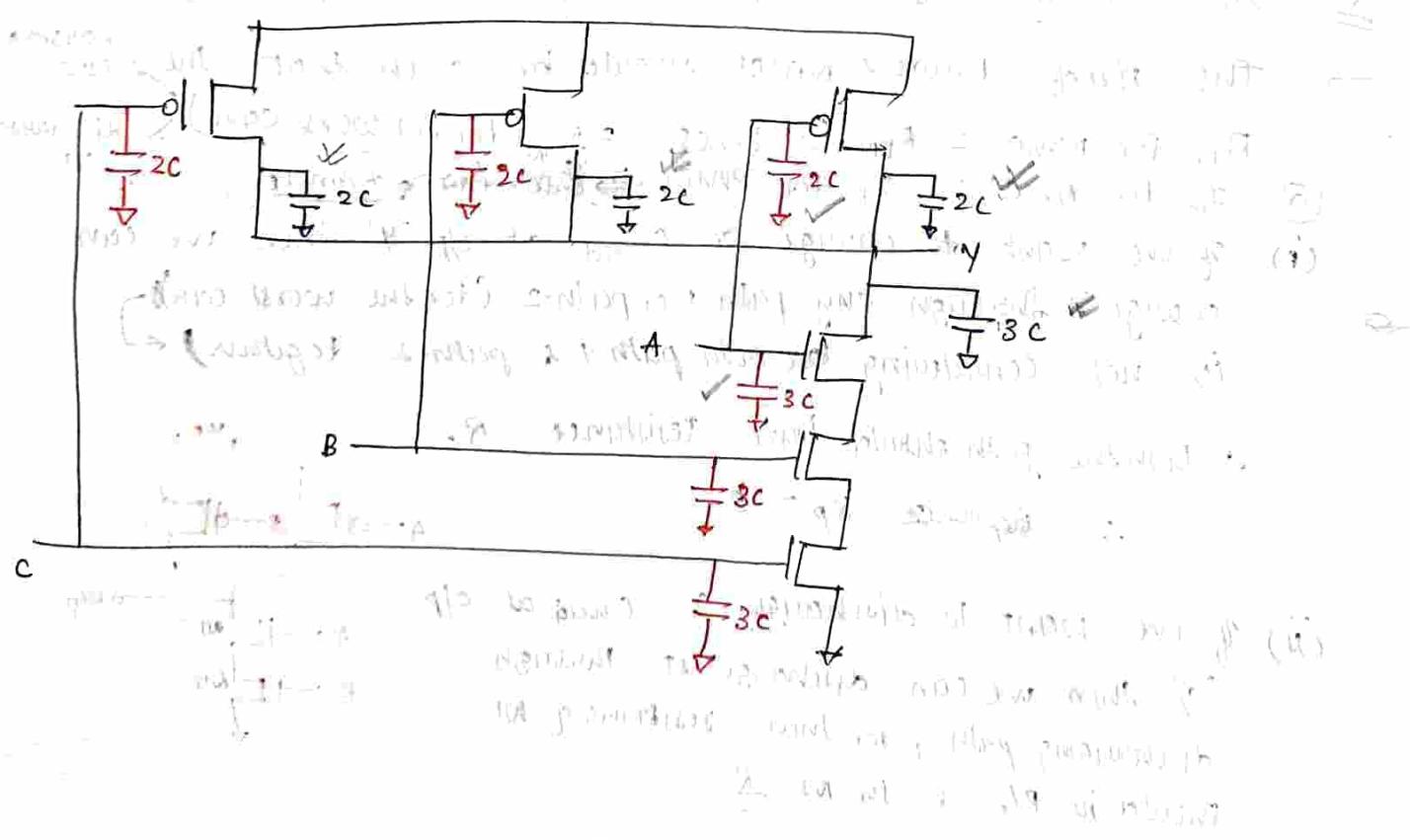
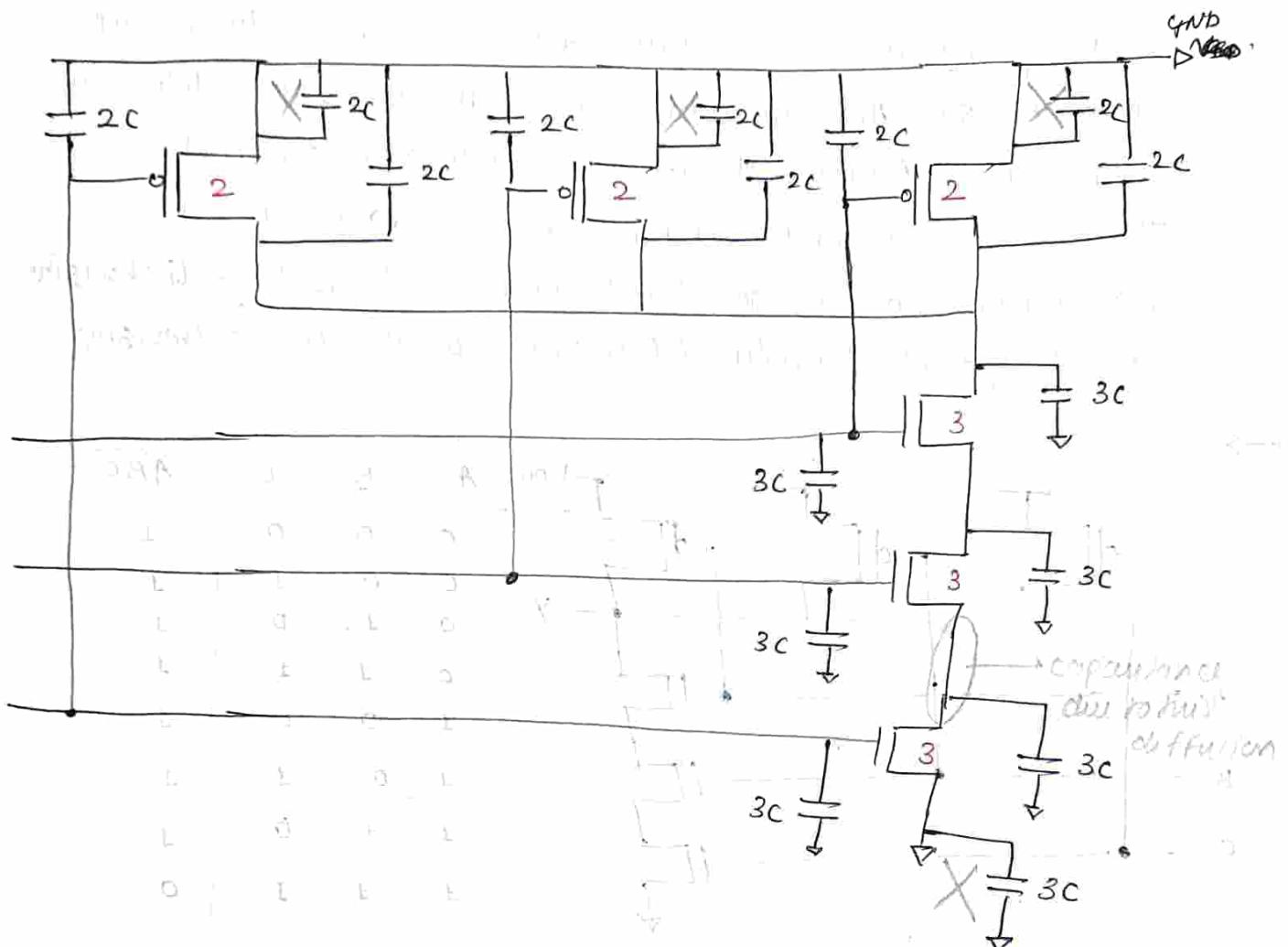
$$\therefore K_p = 2$$

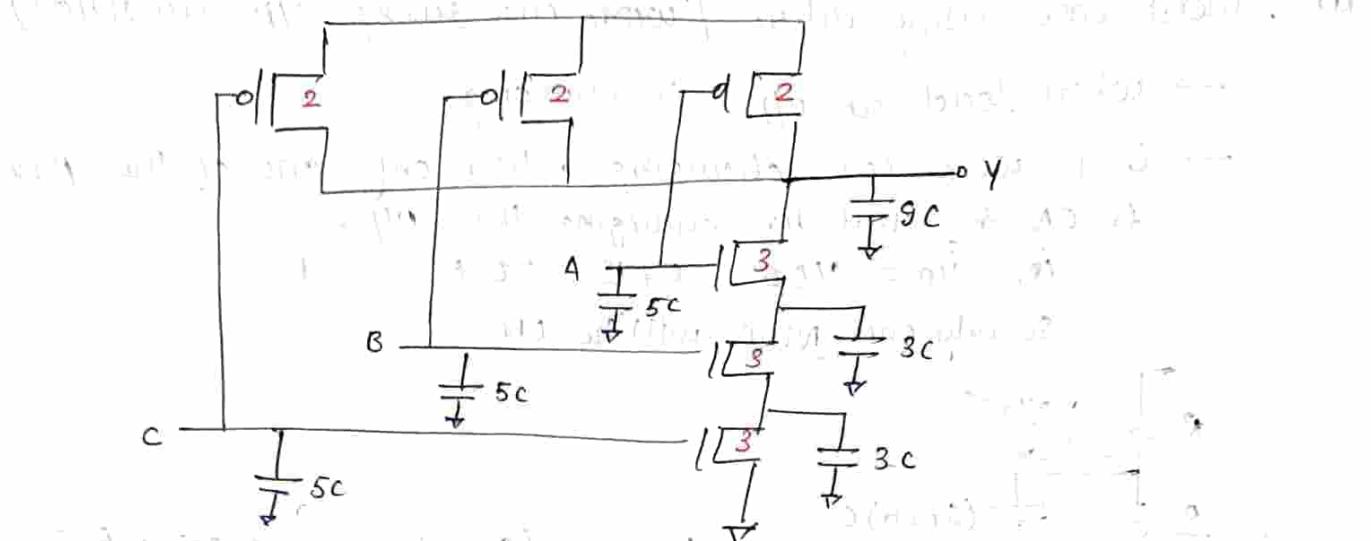
- (ii) If we want to discharge a load at o/p 'Y' then we can discharge it through discharging path, for max resistance of N1 should be $R/2$ & for N2 $\frac{R}{2}$

$$\therefore \text{Resistor path} = \frac{R}{2} + \frac{R}{2} = R \quad [\text{as they will sh. series}]$$



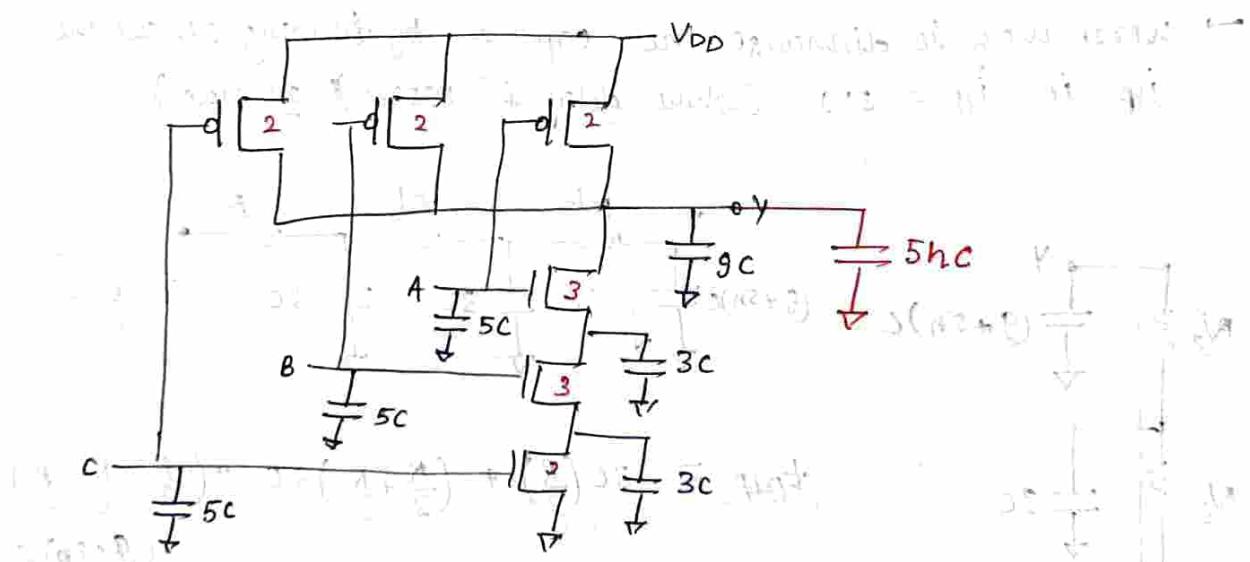
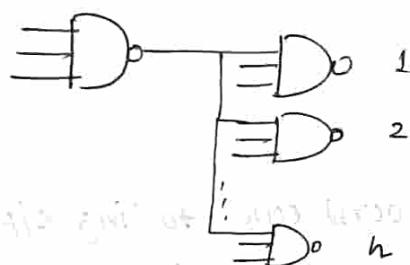
- Annotate the 3-input NAND gate with gain & diffusion (drain & source) capacitance.





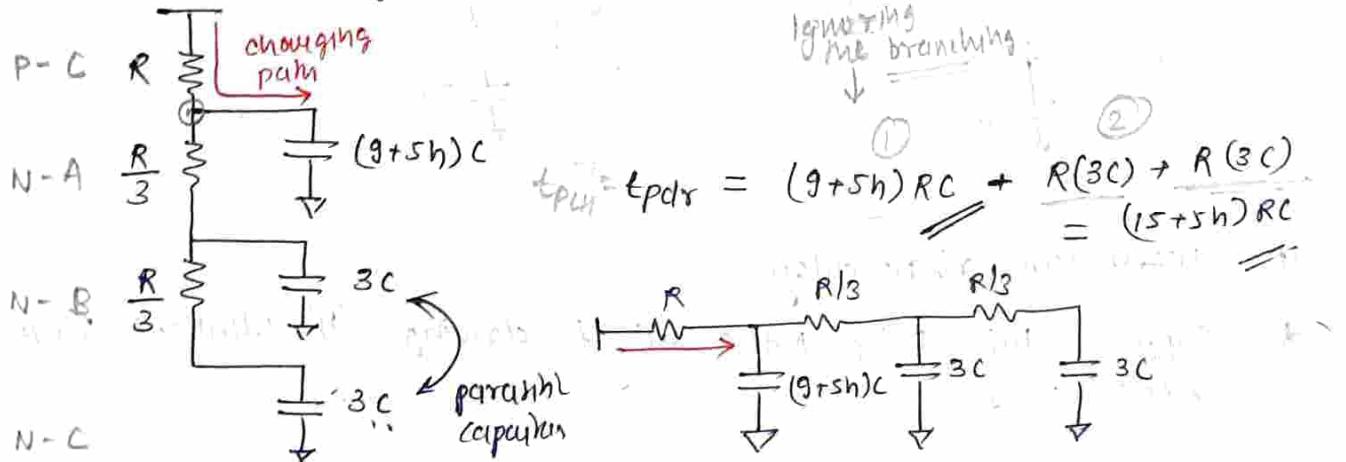
(ii) worst case rising delay

- * Suppose n 3-i/p NAND gate is driving n^2 identical 3-i/p NAND gates ($f_{output} = f_n$)

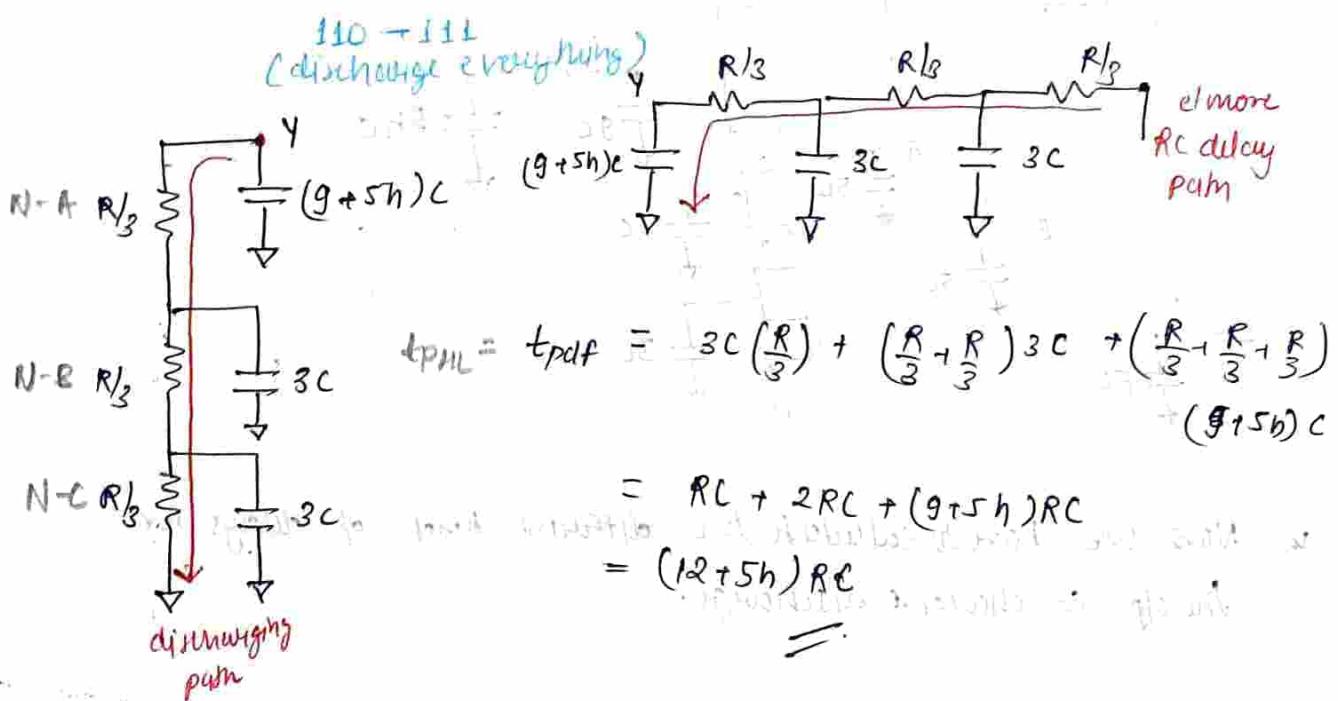


- Now we have to calculate the different kind of delays for the op to charge & discharge.

- (i) Worst-case rising delay [worst case rising o/p transition]
- when load at o/p is changing.
 - it is worst case changing when only one of the PMOS is ON & used for charging the O/P.
ie, $i_{lp} = 110 \dots$ ~~extra path~~
so only one NMOS will be off



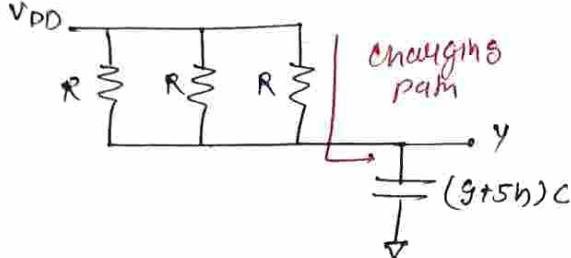
- (ii) Worst case falling delay [worst case falling o/p transition]
- when load at o/p is discharging through NMOS
 - worst way to discharge the cap is by turning on all the NMOS ie, $i_{lp} = 111$ (whole delay is worst or max)



(iii) Bust case rising delay

→ the best case will be when all the PMOS together charge the cap at o/p.
ie, $i_{IP} = 000$

*when only o/p cap. will be charged
④ least capacitor needs to be charged*



$$\begin{aligned} & \text{III} \rightarrow 000 \\ & \text{sub discharge} \quad \text{charge o/p } \text{A} \\ & \text{no sig} \quad \text{charging} \\ t_{CDR} &= \frac{R}{3} \times (g+sh)C \\ &= \left(\frac{g+sh}{3}\right) R C \end{aligned}$$

$$\left\{ \begin{array}{l} t_{CDR} < t_{PDRI} \\ t_{CAF} < t_{PDF} \end{array} \right\}$$

(iv) Bust case falling delay

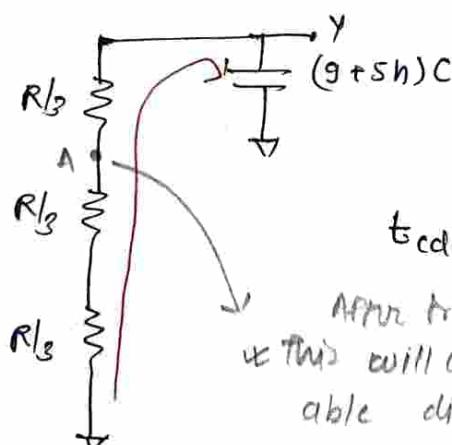
→ the best case for discharging will be when

$$i_{IP} = \overset{ABC}{0001} \rightarrow 011 \rightarrow III$$

→ when we are discharging the 3C of N3 (N-C) to GND then we discharge the 3C of N2 (N-B) to GND then we discharge the $(g+sh)C$ of N1 (N-A) to GND.

$$\begin{cases} \text{for } i_{IP} = 001 : \text{ delay} = RC \\ \text{for } i_{IP} = 011 : \text{ delay} = 2RC \\ \text{for } i_{IP} = III : (g+sh)RC = t_{CDF} \end{cases} \quad \left. \begin{array}{l} \text{when } i_{IP}s \text{ are} \\ \text{given one by one.} \end{array} \right\}$$

* The bust case is that the bottom two NMOS transistors are already ON when top one (N-1/NA) turns ON.
(N-B, N-C) [i_{IP} transition 011 → III]



*011 → III
↓
other cap. are
discharged
except for o/p
capacitance
now only
o/p cap
discharging*

$$t_{CDF} = (g+sh)CR$$

*After transition of 011 → III
it this will cut off GND so last 3C will not be
able to discharge complete $(g+sh)C$
so for 100C at 'A' then we can say discharged
 $(g+sh)C$ will discharge*

$$\therefore t_{CAF} = (g+sh) \frac{C}{R}$$

2. Logical Efforts

• Delay

• No. of stages

No. of stages = 2^{n-1}

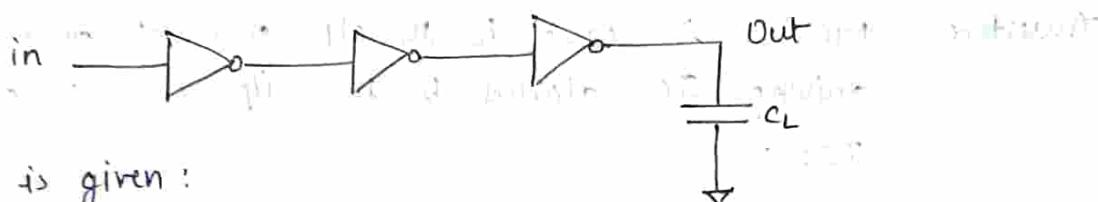
Point 2

Previously we have calculated the delay when size of the MOSFETs are given.

Now, we should be able to size the transistors for the desired delay.

* The concept will help us to size the devices & chains of ~~device~~ gates.

Transistor chain: The signal passes through a series of



If CL is given:

Q: How many stages are needed to minimize the delay?

Q: How to size the transistors?

→ May need some additional constraints

Delay in logic gates (CMOS Inverter)

First of all we have to express delays in process-independent unit. [Delay should be same for different technologies like 45nm, 90nm, 180nm]

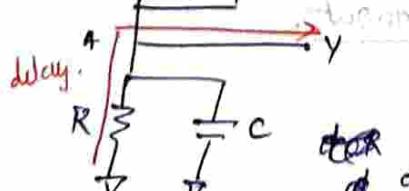
$$d = \frac{d_{abs}}{\tau} \quad \begin{array}{l} \text{absolute delay} \\ \text{delay} \end{array}$$

delay of CMOS

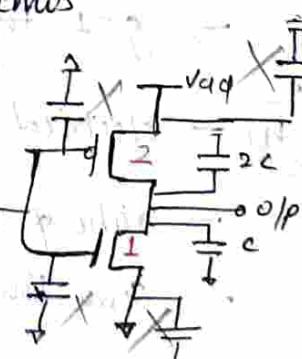
where,

$$\tau = 3RC$$

for $i/p = 1$,



$$\text{delay} = RC + R(2C) \\ = 3RC$$

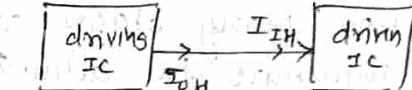


$T = 3RC$ → delay of cmos inverter w/o any fanout.
 $\approx 3\text{psec}$ in 0.5nm process
 $\approx 6\text{ops}$ in 0.6um process

$T \rightarrow$ if the inverter is ideal & has no self loading
 $= \frac{I_{OL}}{I_{IH}}$ then the delay will be one fanout.

fanout \rightarrow it is the number of cmos logic ilp that can be driven by one cmos logic olp.

Therefore, fanout is equal to the olp current of the driving IC divided by the ilp current of driven ICs:

$$\text{fanout} = \frac{I_{OH}}{I_{IH}} \quad \text{or} \quad \frac{I_{OL}}{I_{IL}}$$


Delay in logic paths

- Delay has two components
 $d = f + p \Rightarrow d = gh + p$
 - $f \rightarrow$ effort delay = gh (stage effort)
 - $g \rightarrow$ logical effort (relative to inverter)
 - measured relative ability of gate to deliver current.
 - $g=1$ for inverter.
 - $p \rightarrow$ electrical effort = $\frac{C_{out}}{C_{in}}$
 - ratio of olp to ilp capacitance
 - Some times called fanout.

• P → parasitic delay

→ Represents delay of gate driving, no load.
↳ delay of gate w/o any load)

$\therefore P = 3RC$ for minimum of propagation

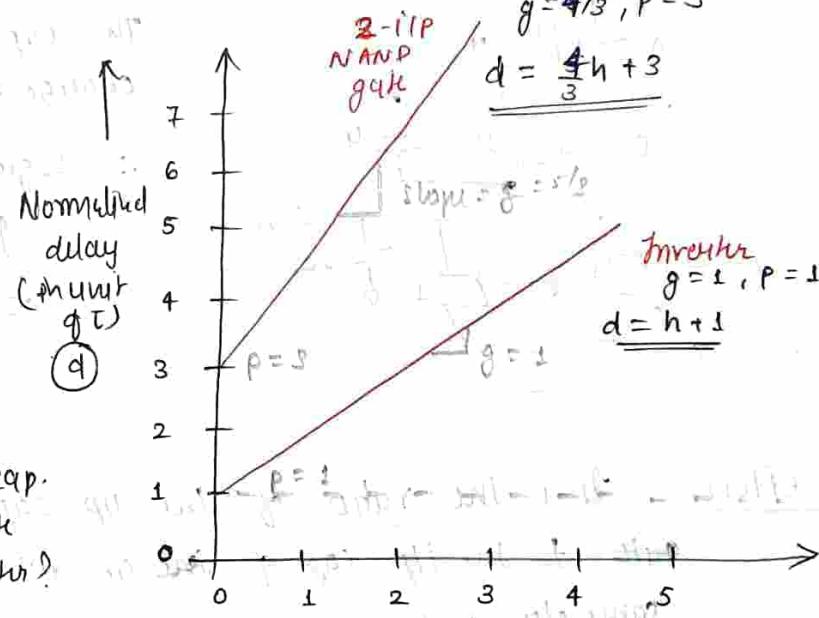
$$\textcircled{R} \quad d_{ans} = 3RC \quad] \quad d = \underline{\underline{1}} \quad \text{for an inverter.}$$

$$[\quad C = 3RC \quad]$$

"1 unit delay"

Delay plot:

- Q: Now do I force V_{DD} to a 2-inp NAND gate for example? So that it gives O/P that we desire.
- By changing its i/p capacitance.

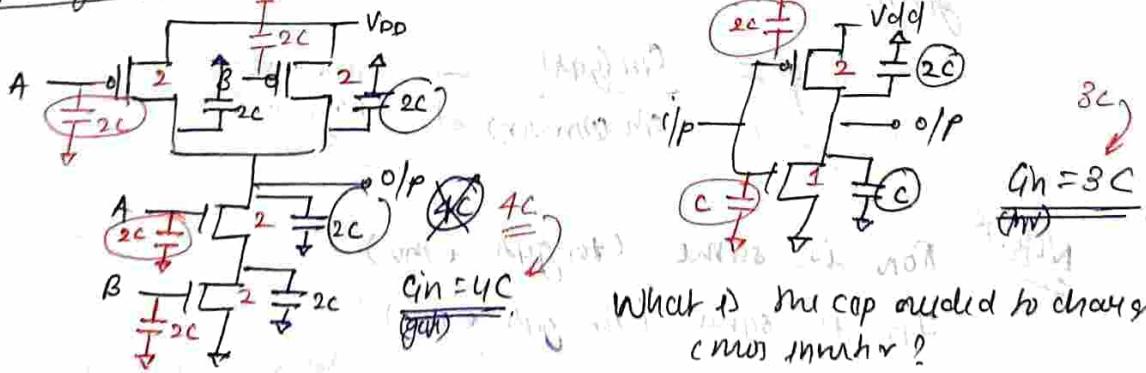


- Q: Can you see the relation b/w i/p cap. of two i/p NAND gate & i/p cap. of inverter?

→ Yes.

* By changing the i/p cap. means, when apply the i/p signal it will change the gate-body capacitance (i.e., gate capacitance) ~~or~~ i/p ~~capacitance~~.

Computing logical effort



What is the cap needed to charge CMOS inverter?

What is the cap. need to charge for this 2-i/p NAND gate?

$$\rightarrow \underline{\underline{4C}}$$

$$\rightarrow \underline{\underline{Gh = 3C}}$$

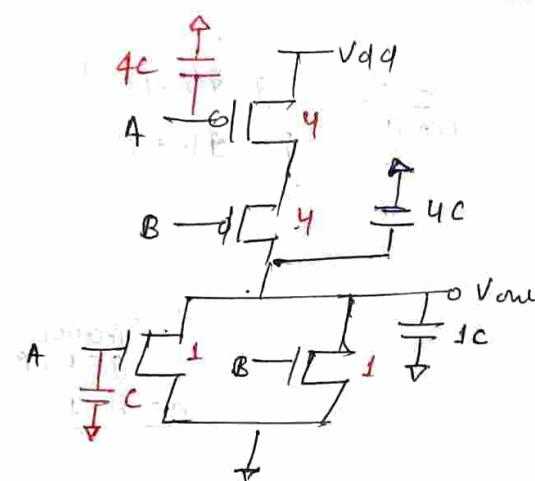
$$g = \underline{\underline{\frac{4C}{3C}}} = \underline{\underline{\frac{4}{3}}}$$

therefore,

logical effort \rightarrow how much more effort do I need to put to charge ~~any~~ cap with any logic as ^{with same value} compared to the inverter.

$$\therefore g = \frac{C_{\text{logic}}}{C_{\text{inverter}}} =$$

Q: logical effort of NOR-2?



$$4C + C = 5C$$

$$C_{\text{in}} = 5C$$

The cap. that we need to charge for NOR-2 is $5C$.

\therefore logical effort

$$g = \frac{5C}{3C} = \frac{5}{3}$$

logical effort: It is the ratio of the i/p capacitance of the gate to the i/p cap. of an inverter delivering the same o/p current. (Drive same power also)

(Q) how much more effort do I have to put to charge the i/p cap. of gate as compared to inverter to get the same response as the inverter would have.

given.

$$g = \frac{C_{\text{in}}(\text{gate})}{C_{\text{in}}(\text{inverter})} \leftarrow \begin{array}{l} \text{gate & inverter} \\ \text{have total} \\ \text{size} \end{array}$$

Note: R_{on} is same (for gate & inverter)

I_{D} is same (for gate & inverter)

\Rightarrow Both will drive same power

$$\text{as } V_{\text{DS}} = I_{\text{D}} \cdot R_{\text{on}}$$

(sizing of mos)

***^{**}** These sizings are done to ensure that, we'll have same equivalent resistance both for sink & fall time.

And not for the delay we have achieved. (wrong)

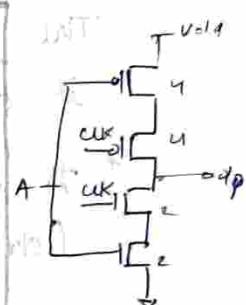
It lowers the logical effort, lowers the delay

because $\downarrow d = \underline{g_{le} f}$.

Logical Effort for common gates

Q: Does size of gates (w_p, w_n) affect the logical effort or not?

Gate Type	Number of Inputs	1	2	3	4	n
Inverter	1	1	1	1	1	1
NAND		$4/3$	$5/3$	$6/3$	$(n+2)/3$	
NOR		$5/3$	$7/3$	$9/3$	$(2n+1)/3$	
Tristate/mux	2	2	2	2	2	2
XOR, XNOR		$4, 4$	$6, 12, 6$	$8, 16, 16, 8$		



logical Efforts

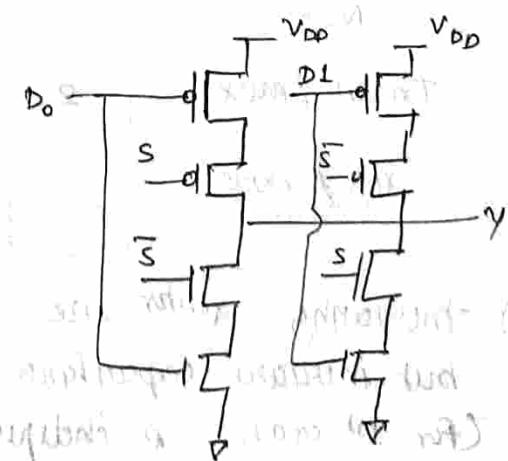
- For complex gates like XOR, XNOR it can have different logical effort for different i/p combination

* Logical effort is strong funⁿt of no. of i/p's.

- Multiplex built from ganged tristates, has logical effort of 2, independent of no. of i/p's.

- ~~size~~ But parasitic delay with multiplexer size.

∴ Therefore it is recommended to use of 1-i/p MUX, which is faster. (ie, least delay)



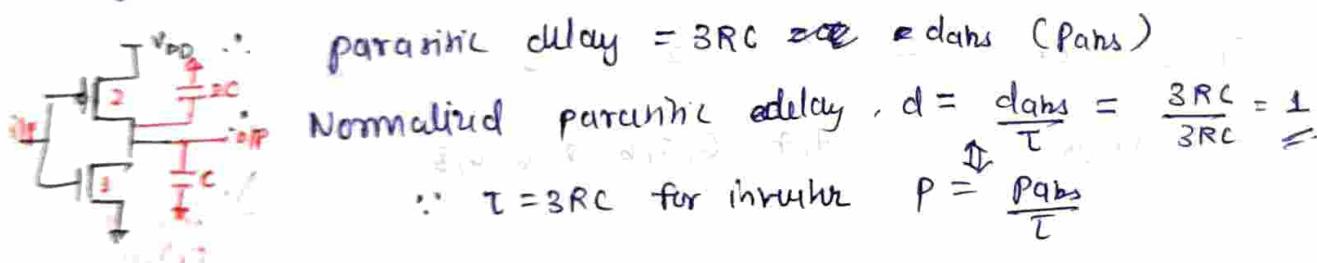
Q. Now can you "benefit" from that some gates have higher logical effort and some has less for XOR & XNOR gates?

- we can put the ILPs which don't switch much ~~area~~ at higher logic efforts. (It will save power)
- The gate ILP which has ~~less~~ ^{less} logical effort you will apply ~~more~~ signal which creates late.

Parasitic delay of common gates :

- The parasitic delay of a ^{logic} gate, is the delay of the gate when it drives zero load. (Delay due to self load only)
- It can be calculated by counting only diffusion capacitance (drain + source) on the o/p node.

Eg: Inverter has 3-ILP of diffusion capacitance on the o/p.



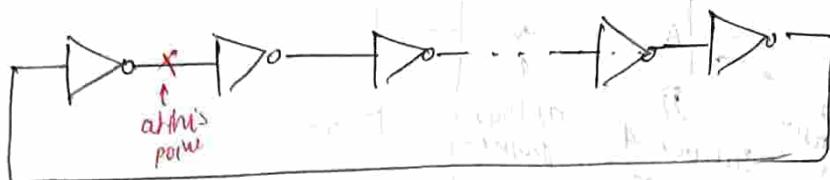
Gate type	Number of Inputs				
	1	2	3	4	n
Inverter	1				
NAND	2	3	4	n	
NOR	2	3	4	n	
Tri-state / mux	2	4	6	8	$2n$
XOR / XNOR	4	6	8		

* Increasing ~~gate~~ ^{W/L} size reduces resistance but increases capacitance, so parasitic delay (for 1st order) is independent of ^{W/L} gate size.

Normalized
parasitic
delay (P)

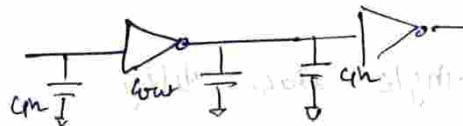
Example : Calculate the frequency of N -stage ring oscillator

(where $N \rightarrow$ odd integer)



→ for single stage,

(i) logical effort, $g = 1$.



(ii) Electrical effort (how many gates are connected to one inverter) = 1
(i.e., fanout)

$$h = 1$$

④ $h = \frac{C_{out}}{C_{in}} = 1$ (as both have same capacitance)

(iii) parasitic delay, for inverter $= 1$

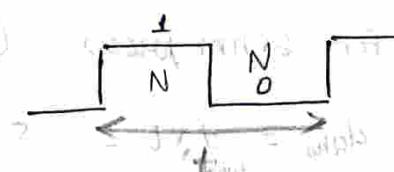
* for any inverter at any stage, how much effort does the previous stage needs to apply in compensation to an inverter so that its sum o/p as an inverter would give.

(iv) stage delay, $d = gh + p = 1 + 1 = 2$

$$\underline{\underline{d = 2}}$$

for N -stage

(v) Frequency of N -stage



waveform at the o/p of oscillator.

$$\begin{aligned} \text{total delay} &= 2N \times d \rightarrow \text{we have to go through all } N\text{-stages twice + each stage has delay } d. \\ &= 2N \times 2 \end{aligned}$$

$$\therefore \text{freq.} = \frac{1}{4N}$$

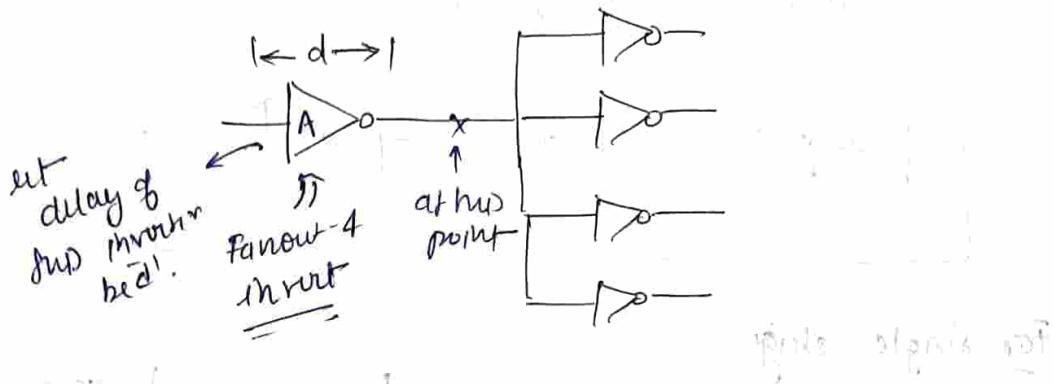
$$2\pi f = \frac{1}{2 \times 31 \times D} \Rightarrow D = \frac{1}{2 \times 31 \times 60 \times 31}$$

* $N = 31$ stages of ring oscillator in $0.64\mu m$ technology have freq. of $\approx 200\text{ MHz}$. [$d_{abs} = \frac{\text{total delay}}{T}$] for $T = 60\text{ ps}$

$$\Rightarrow d_{abs} = 31 \times 4N \times T = 4 \times 31 \times 60\text{ ps} = 7.2\text{ nsec}$$

$$f_{abs} = \frac{1}{d_{abs}} = 138.9\text{ MHz}$$

Example : Estimate the delay of a fanout-of-4 (FO4) inverter.



→ Single stage delay : $t_{single} = 1.0 \text{ units sample}$

(i) Logical effort, $g+h=1$ but $g=0.5$ for 4 fanout

(ii) Electrical effort, $h=4$ (fanout)

(iii) Parasitic delay, $p=1$

(iv) Stage delay, $d = gh + p = 4 + 1$

* The FO4 delay is about (i) 300ps in 0.6μm process

$\rightarrow [t = 60\text{ps}] = 3RC \rightarrow 1 \text{ unit delay}$

for 0.6μm process, $[t = 60\text{ps}] = 3RC$

$$(i) d_{ah} = d \times t = 5 \times 60\text{ps}$$

$$= 300 \text{ ps}$$

(ii) for 65nm process ($t = 3\text{ps}$) $= 3RC$

$$d_{ah} = d \times t = 5 \times 3 \text{ ps}$$

$$= 15 \text{ ps}$$

$$= 15 \text{ ps}$$

Note: We will only look at driver A when we are calculating parasitic delay.

- Because parasitic delay does not depends on load.

- It only depends on ref-load of driver logic.

Multistage logic N/W

- logical effort generalizes to multistage N/W

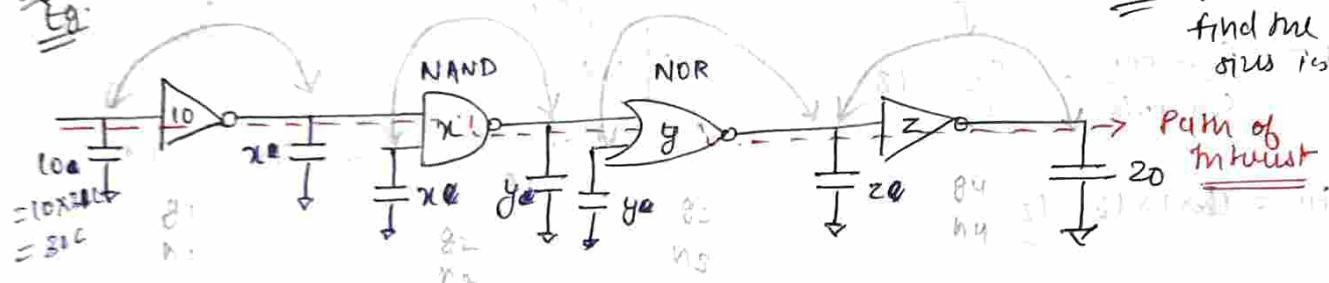
(i) Path logical effort $G = \prod g_i$

(ii) Path Electrical effort

$$H = \frac{\text{Cout-path}}{C_{in} - path}$$

Multistage logic Effort (iii) Path Effort $F = \prod f_i = \prod g_i h_i$

Eg:



Q: Now do we find the value of size x_1, y_1, z_1 ?

- All the capacitances drawn are i/p capacitances of logic gates.

$\therefore g_1 = \frac{C_{in}(\text{inverter}) - \text{given}}{C_{in}(\text{inverter}) - \text{given of same size}} = \frac{8e}{5e} = \frac{10 \times 5RC}{10 \times 3RC} = 1$

for calculating logical effort of any gate size n , we compare it with given of size n .

NAND, $g_2 = \frac{C_{in}(\text{NAND}) - \text{given}}{C_{in}(\text{inverter}) - \text{given of same size}} = \frac{8e}{5e} = \frac{x \times 4e}{x \times 3e} = \frac{4}{3}$

$$G = g_1 g_2 g_3 g_4$$

$$= 1 \times \frac{4}{3} \times \frac{5}{8} \times 1$$

$$= \frac{20}{27}$$

* logical effort is independent of size. (How?)

* Electrical effort depends on size

$\therefore h_1 = \frac{\text{Cout}_1}{C_{in} + \text{Cout}_1} = \frac{x e}{10e} = \frac{x}{10}$

$$H = \prod h_i$$

$$H = h_1 h_2 h_3 h_4$$

$\therefore h_2 = \frac{\text{Cout}_2}{C_{in} + \text{Cout}_2} = \frac{y e}{8e} = \frac{y}{8}$

$$= \frac{x y z \times 20}{10 x y z} = \frac{20}{10} = 2$$

$\therefore h_3 = \frac{\text{Cout}_3}{C_{in} + \text{Cout}_3} = \frac{z e}{8e} = \frac{z}{8}$

$\therefore H = \frac{\text{Cout-path}}{C_{in} - path} = \frac{20e}{10e} = 2$

Q: Can we write $F = GH$?

→ NO (how?) → due to path branching.

$$F \neq GH \text{ but } F = GBH$$

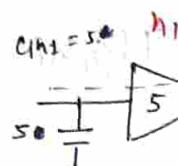
we are not considering the self-loop here.

Paths that Branch

- We have to find things for the below path only.

$$(i) G = g_1 g_2$$

$$= 1 \cdot 1 = 1$$



$$H = \frac{\text{Out-path}}{\text{In-path}} = \frac{90}{5} = 18$$

$$GH = 1 \times 18 = 18$$

$$(ii) h_1 = \frac{\text{Out}_1}{\text{In}_2} = \frac{(15 + 15)}{5} = 6$$

$$h_2 = \frac{\text{Out}_2}{\text{In}_2} = \frac{90}{15} = 6$$

$$\therefore F_{\text{path effort}} = g_1 h_1 g_2 h_2 =$$

$$= 1 \times 6 \times 1 \times 6 = 36$$

$$= 2 \times 18$$

$$= 2GH$$

Q: How does this '2' come? →

$h_i = \text{fanout} = 2$ only when the size of all the branches are same.

$$\textcircled{a} = \frac{15 \times 3 \times RC + 15 \times 3 \times RC}{5 \times 3 \times RC} = \frac{15 + 15}{5} = 6$$

Branching Effort [like Elmore delay] as branches.

- It accounts for branching b/w stages

$$b = \frac{\text{Con-path} + \text{Off-path}}{\text{Con-path}}$$

in path

Now we are also considering the delay of branches which are not in our own path as the Elmore delay.

$$B = TTB_i$$

$$TTh_i = BN$$

∴ path effort:

$$F = GBH$$

MP is only valid when all the branches reach same time.

Multistage Delay

(i) Path Effort Delay : $D_F = \sum f_i$ $d = gh + P$

(ii) Path Parasitic Delay : $P = \sum p_i$

(iii) Path Delay : $D = \sum d_i = \underline{D_F + P}$
 $= \sum (f_i + p_i)$

* The product of stage effort ($f = gh$) is path stage effort (F).

$$\textcircled{3} F = \prod f_i$$

* The sum of stage effort ($f_i = g_i h_i$) is path effort delay (D_F).

$$\textcircled{4} D_F = \sum f_i$$

Designing Fast Circuit \rightarrow Minimum path Delay

- The sum of a set of numbers whose product is constant is minimized by choosing all the no.s to be equal.
 - the path delay is minimum when each stage bears the same effort.
 - If path has 'N' stages & each bears the same effort, then efforts must be
- The min. path delay (D) of a N-stage path, with path effort (F) path parasitic delay (P) is

$$\hat{f} = g_i h_i = F^{1/N}$$

$$\therefore D = D_F + P$$

$$= \sum f_i + P$$

$$= \sum g_i h_i + P = \sum \hat{f} + P$$

$$D = N \cdot \hat{f} + P$$

Nb. of stages
stage effort
each stage.

- * The key result of logical effort is that the minimum delay of the path can be estimated knowing only
 - no. of stages
 - path effort
 - parasitic delay

w/o the need to assign transistor sizes.

- # Above method is better than simulation, because in simulation delay depends on sizes. And you will never know the size for min. delay. (?)

Gate Sizes

- How wide should the gates be for shortest delay?

→ ∵ optimal stage effort,

$$F^{IN} = \hat{f} = g_i h = g_i \cdot \frac{C_{out}}{g_i h} \quad \text{for } i^{\text{th}} \text{ stage}$$

$$\Rightarrow g_i C_{in,i} = \frac{g_i C_{out,i}}{\hat{f}} \quad \text{when OLP, capacitance of the path is given. effort}$$

and logical effort being constant & appropriate stage delay is also known.

- Starting with the load at the end of the path, work backwards applying the capacitance transformation to determine the size of each stage.

- Check the calculation by verifying that the size of the initial stage matches the specification.

Example: What will be the size of each gate for least delay?



$$g = 1$$

$$h = a$$

$$g = 5/b$$

$$h = b/a$$

$$g = 5/c$$

$$h = c/b$$

$$g = 1$$

$$h = 5/c$$

- (i) Effective fanout, $N = h_1 h_2 h_3 h_4 = a \times \frac{b}{a} \times \frac{c}{b} \times \frac{d}{c} = 5$
- (ii) $G = g_1 g_2 g_3 g_4 = 1 \times \frac{5}{3} \times \frac{5}{3} \times 1 = \frac{25}{9}$
- (iii) $F = BGH = 1 \times \frac{25}{9} \times \frac{25}{9} = \frac{625}{81} \rightarrow (B=1)$ since only one branch
- (iv) ~~\hat{f}~~ $\hat{f} = F^{1/N} = \left(\frac{625}{81}\right)^{1/4} = 1.9305 \rightarrow$ optimal stage effort

Therefore,

$$\textcircled{a} \quad \hat{f} = g \cdot \frac{C_{out}}{C_{in}}$$

$$C_{in} = \frac{g \cdot C_{out}}{\hat{f}} \quad \textcircled{b}$$

$$\textcircled{c} \quad a =$$

(i) for NOR NOT gate,

$$C_{in} = \frac{g_4 \cdot C_{out}}{\hat{f}} = \frac{1 \times \frac{5}{3}}{1.93} = 2.59$$

$$\therefore C = 2.59$$

(ii) for NOR-2

$$C_{in} = \frac{g_3 \times C_{out}(\text{NOR})}{\hat{f}} = \frac{\frac{5}{3} \times 2.59}{1.93} = \frac{\frac{5}{3} \times 2.59}{1.93}$$

$$= 2.236$$

$$\therefore b = 2.23$$

(iii) for NAND-3

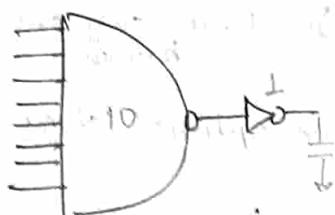
$$C_{in} = \frac{g_2 \times C_{out}(\text{NAND})}{\hat{f}} = \frac{\frac{5}{3} \times b}{1.93} = \frac{\frac{5}{3} \times 2.23}{1.93}$$

$$C_{in} = 1.9257$$

$$\therefore a = 1.93$$

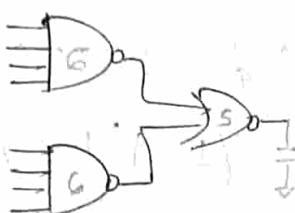
* The back hand calculation (logical effort) is just an estimation & not that accurate as we are not considering parasitic delays. & simulation will tell you the exact delay but we can't simulate everything like complex signals to cells hence we do the logical effort.

Example: 2 Which of the following 8-input AND gate has least delay & also calculate the area for each cases.



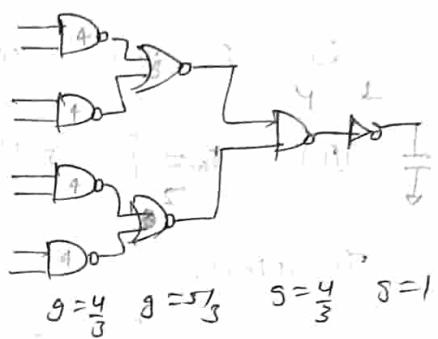
$$g = \frac{10}{3}$$

(a)



$$g = 2 + 4 = 6$$

(b)



$$g = 4 + 4 = 8$$

(c)

No. of IP = Parallel

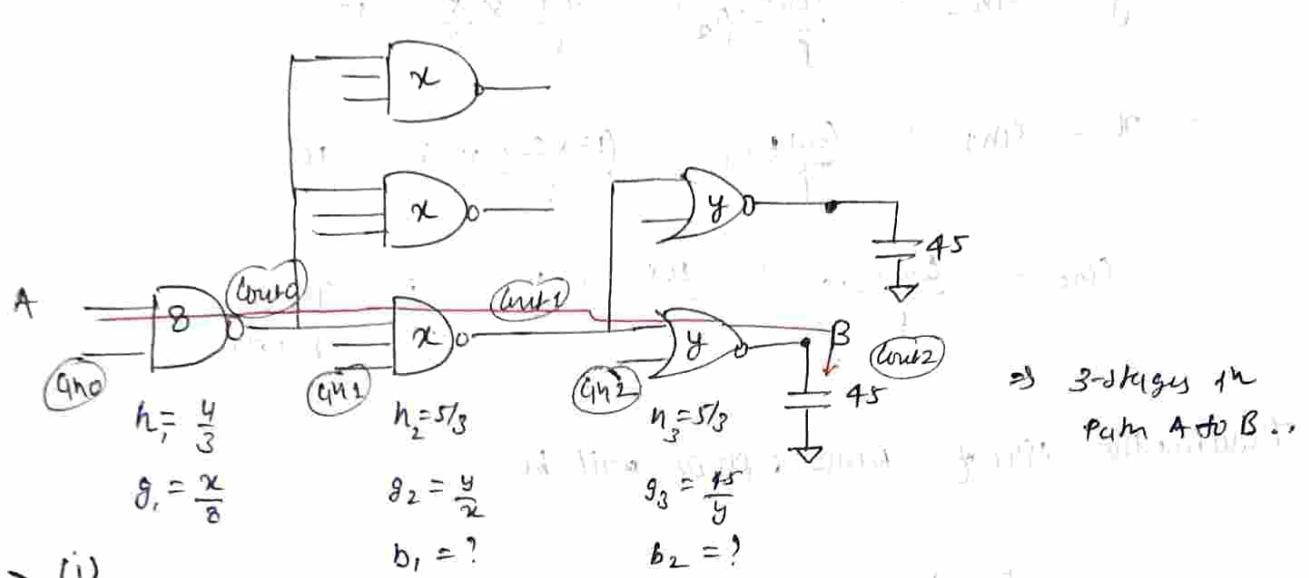
delay = $\frac{1}{f_{max}} = \frac{1}{1000} = 1\mu s$

area = $T^2 f_{max} = 1^2 \times 1000 = 1000 \mu m^2$

area = $1000 \mu m^2$

which one has less delay? \rightarrow parallel delay, parallel area will be less than serial delay. So serial delay will be more than parallel delay. If we consider serial delay will be more than parallel delay. So parallel delay will be less than serial delay. So parallel delay will be less than serial delay.

Example : 3 (i) Calculate the least delay from A to B.
 (ii) Select sizes x & y for atleast delay from A to B.



→ (i) Logical Effort: $G = \frac{4}{3} \cdot \frac{5}{3} \cdot \frac{5}{3} = \frac{100}{27}$

(ii) Electrical Effort: $H = \frac{\text{Conc-path}}{\text{Con-path}} = \frac{45}{8}$

(iii) Branching effort: $B = \prod b_i$

$$b_i = \frac{\text{Conc-path} + \text{Coff-path}}{\text{Con-path}}$$

$$b_1 = \frac{2L + 2K}{x} = 3, \quad b_2 = \frac{4+4}{y} = 2$$

$$\therefore B = b_1 b_2 = 3 \times 2 = 6$$

(iv) Path Effort: $F = GBH = \frac{100}{27} \times 6 \times \frac{45}{8} = 125$

(v) Best stage effort: $\hat{F} = (F)^{1/N} = (F)^{1/3} = (125)^{1/3} = 5$

(vi) Parasitic delay: $P = 2+3+2 = 7$

(vii) Delay ~~D~~: $D = N\hat{F} + P = 3 \times 5 + 7 = 22$

$$D = \frac{22}{5} = 4.4 \quad (?)$$

Delay $\Rightarrow \frac{D}{\hat{F}} = \frac{22}{5} = 4.4$ m sec and since $\rightarrow N = 3 \times (1.2) \approx 4$

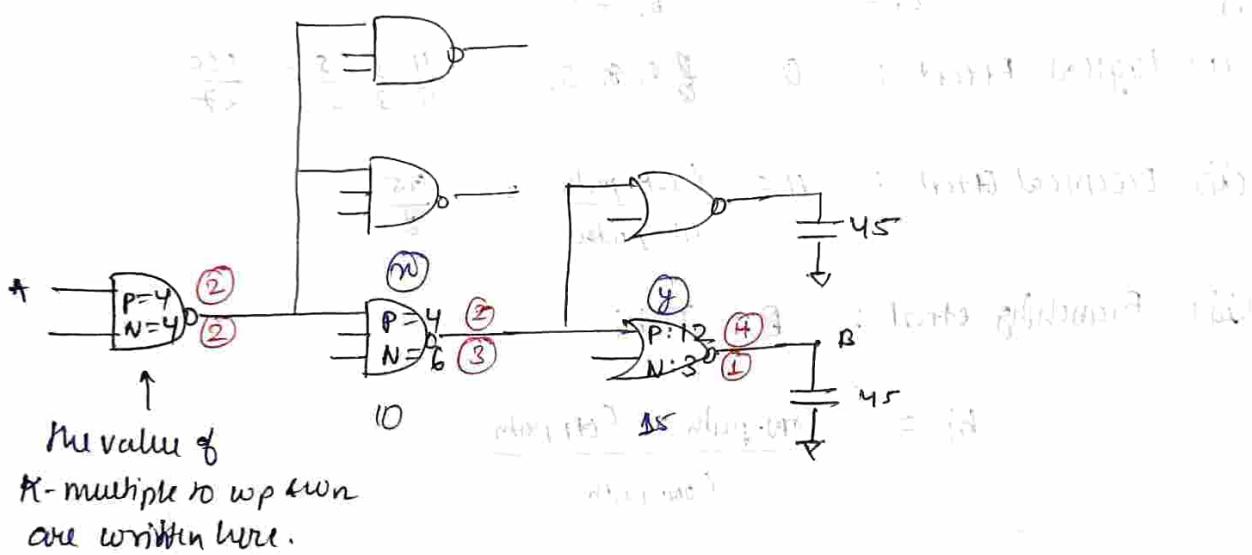
(ii) Now work backward for sizing :-

$$\hat{f} = g \cdot h$$

$$\hat{f} = g \cdot \frac{C_{out}}{g_h}$$

- $y = C_{in2} = \frac{C_{out2}}{\hat{f}} \cdot g_3 = \frac{45}{5} \times \frac{5}{3} = 15$
- $x = C_{in1} = \frac{C_{out1}}{\hat{f}} \cdot g_2 = \frac{(15 \times 2)}{5} \times \frac{5}{3} = 10$
- $C_{in0} = \frac{C_{out0}}{\hat{f}} \cdot g_1 = \frac{10 \times 3}{5} \times \frac{4}{3} = 8$ Hence power

Therefore the size of NMOS & PMOS will be



Conclusion → Using logical effort you have got the

- delay estimation
- area estimation

Best Number of Stages → cost delay + area both

- How many stages should a path use?
 - Increasing the no. of stages \downarrow total delay \downarrow .
 - Increasing the no. of stages \uparrow area always \uparrow .

$$D = NF^{1/N} + P$$

$$= N(F^{1/N}) + N$$

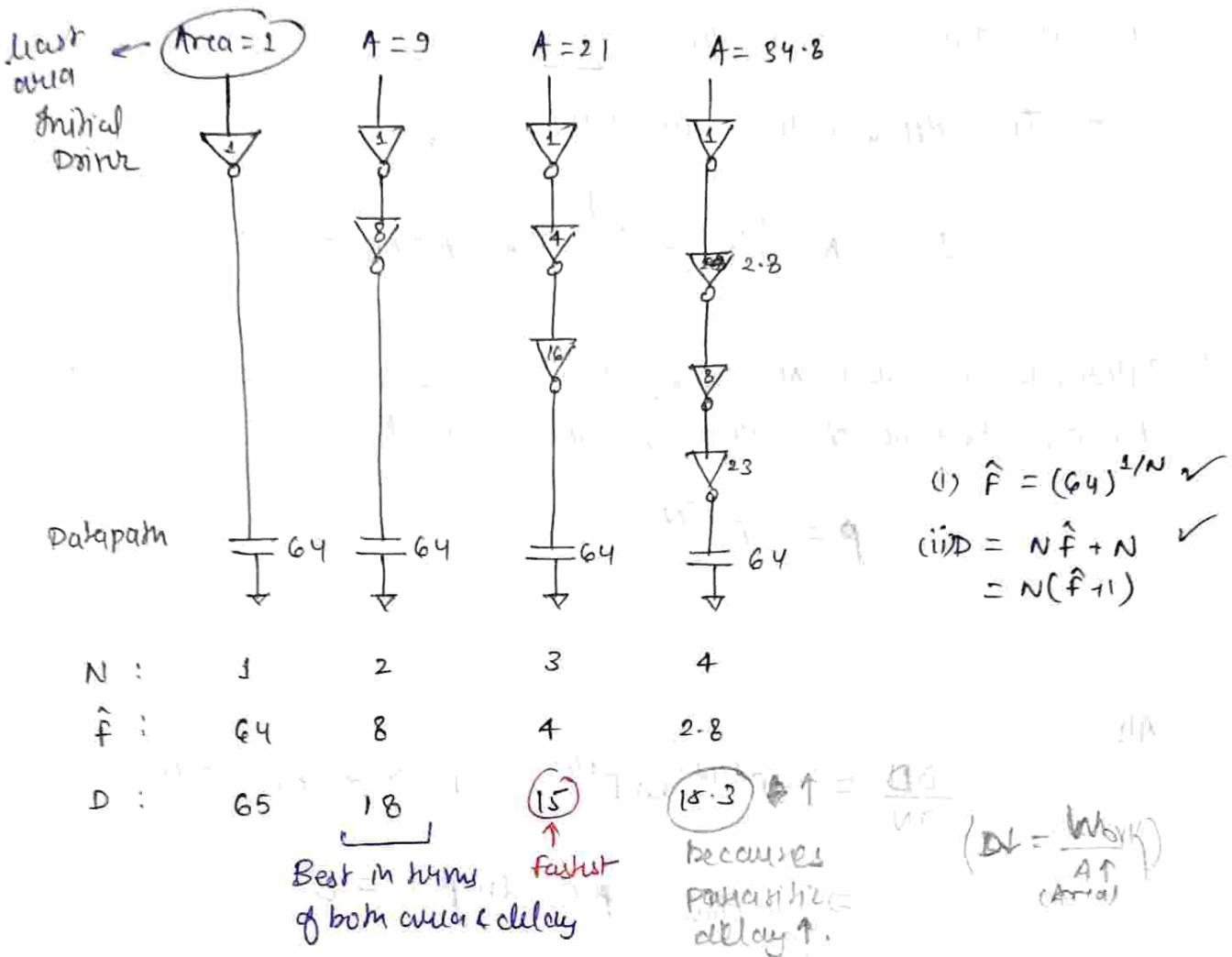
$$= N(G_4)^{1/N} + N \rightarrow \text{since we are only using the } \cancel{\text{inverter}}$$

$F = GBH$

$$= 4 \times 1 \times \frac{64}{4} = 64$$

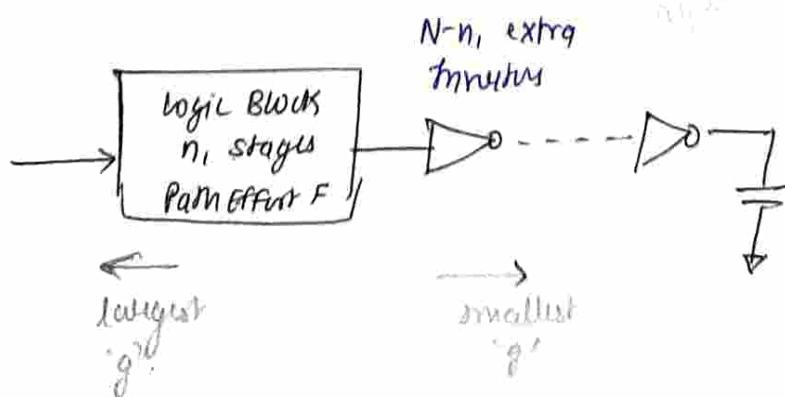
∴ No. of stages = ~~total~~ parallel delay.

Eg: Drive 64-bit datapath with unit inverters



doubt

- Note: (i) Gaths with fewer ips are better than gaths with many ips.
- (ii) It is believed that the fastest ckt are those that use the fewer stages of logic. (which is not true always)
 - often, the gate delay actually depends on the Electrical effort, so sometimes using fewer stage results in more delay.
- (iii) In general you can add extra inverters to the end of the path for decreasing the effort w/o changing the func' of original combinational block. (why?)



- The extra inverters don't change the path logical effort but do add parasitic delay
 - The delay of the new path is

$$D = NF^{1/N} + \sum_{i=1}^{n_L} p_i + (N-n_L) P_{inv}$$

- Differentiating w.r.t N & equating it to 0 allows us to solve for the best no. of stages, which we call \hat{N} .

$p = F^{1/N}$

Also,

$$\frac{\partial D}{\partial N} = -F^{1/N} \ln F^{1/N} + F^{1/N} + P_{inv} = 0$$

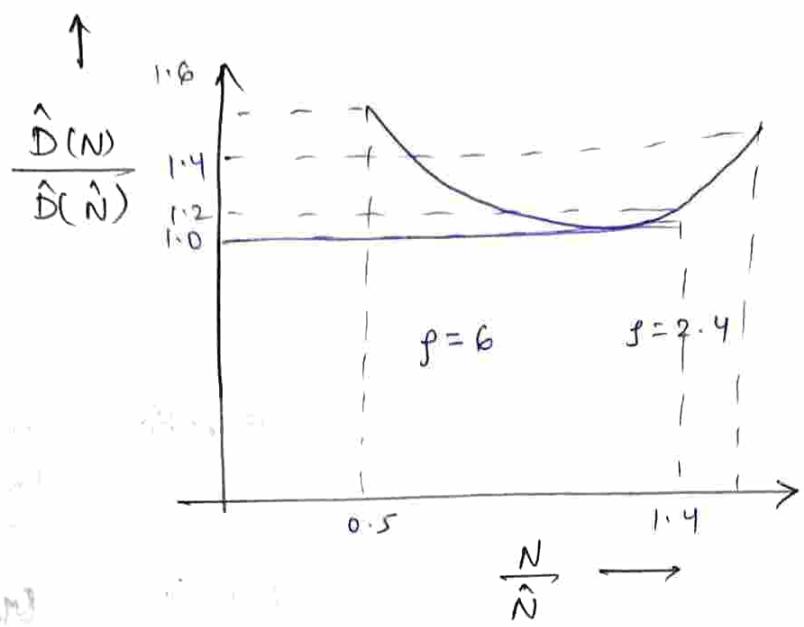
$$\Rightarrow P_{inv} + p(1 - \ln(p)) = 0$$

it has no closed form solution

- * Parasitic delay means each inverter is somewhat more costly to add.
- * As a result it is better to use fewer stages, or equivalently a higher stage effort than e.
- * solving numerically, when $P_{inv} = 1$, we find $p = 3.53$

* the path achieves best delay by using,

$$\hat{N} = \underline{\log_p F \text{ stages}}$$



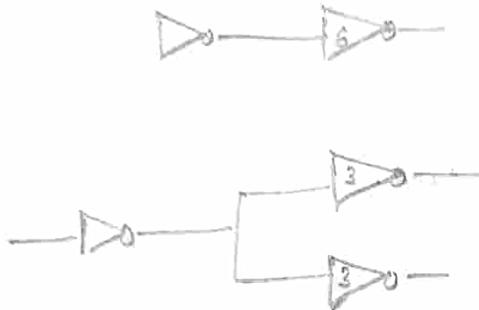
$$\frac{1}{\sqrt{f}} = \frac{1}{\sqrt{2(2n+1)}} \times T_b$$

$$10 \text{ MHz} = \frac{1}{2(2n+1)T_b}$$

$$T_b(2n+1) = 10^{-7}$$

$$T_b \cdot N = 10^{-7}$$

if $T_b = 2$



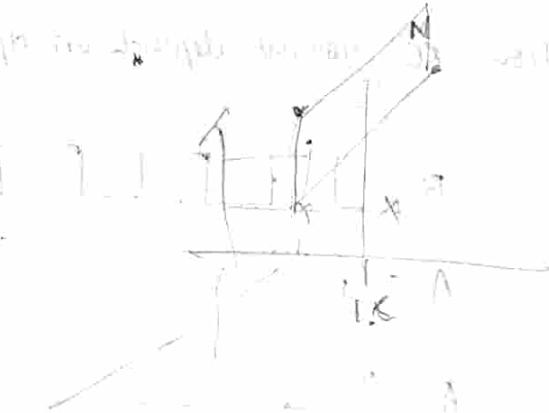
① INV - NOR4 - INV

Delay = ?

Stages = 3

② NOR4 - INV - INV

③ NAND4 - INV - INV - INV



$\Rightarrow V_T$ & L_{eff} are independent for Large channel device.
 V_T & L_{eff} are dependent for small channel device.

$\Rightarrow \pi \rightarrow FF$

use trim which changes both
NMOS & PMOS.

* Uncorrelated
Variations

we get $I_p \neq I_n$

• For fast corners:

push all corners, voltages, temperature
(V_{DD})
together.

• For calculating delay & ~~min max~~ freq.
we have calculate for SS.

$$\text{B-ratio} : \frac{I_{ON-PMN}}{I_{OFF}}$$

$$= \frac{I_{DP}}{I_{DN}}$$

• Functional corners (FFF)

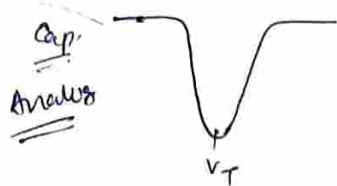
- Latches (most noise in our case)

• Characterization corners (FF)

$$\text{B-ratio} = \frac{W_p C_{ox} p}{W_n C_{ox}}$$

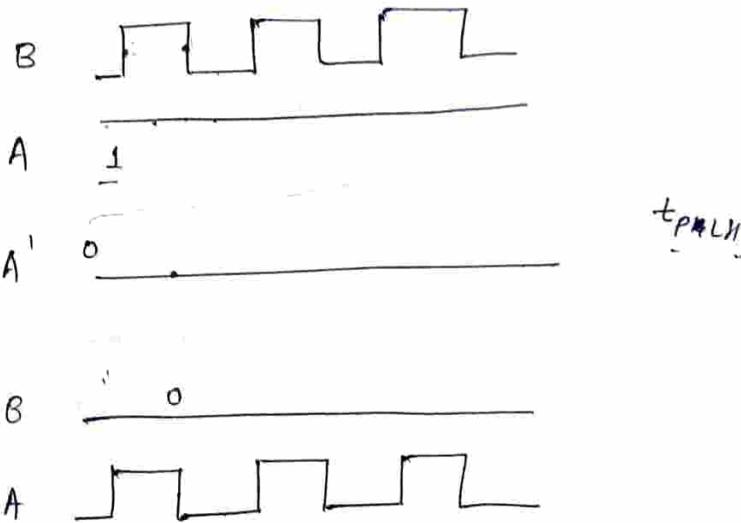
Q. What is V_T ?

→ The value of V_{in} at which $I_D = 14A$.



Q. How R_B current depends on applied voltage?

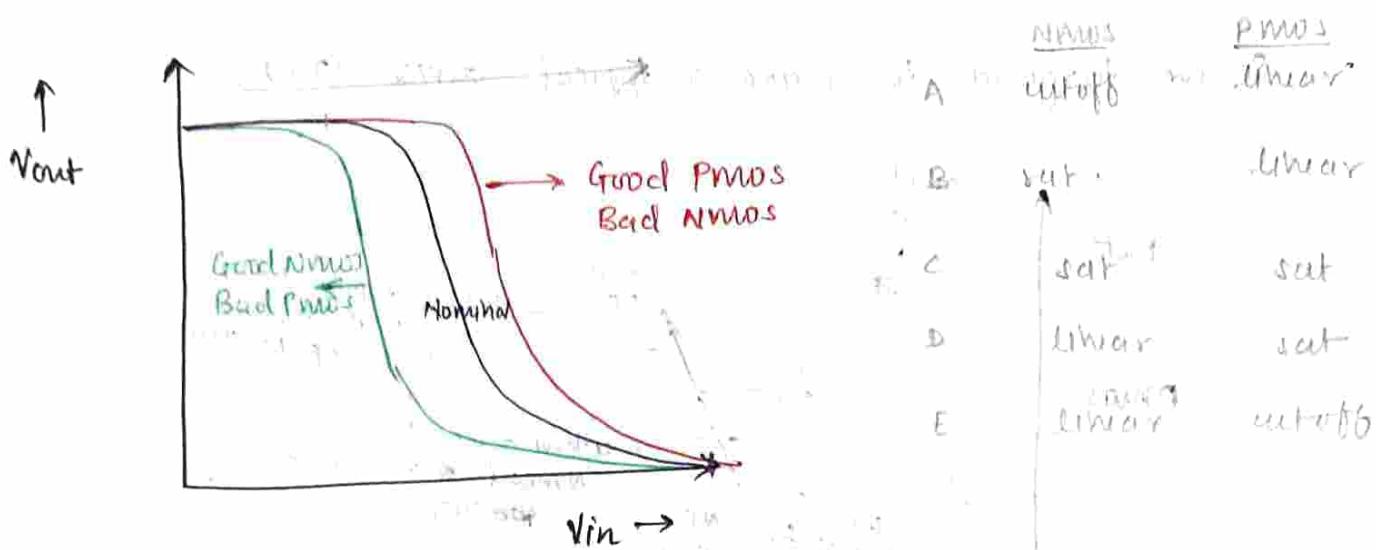
DFM:-



Voltage

3. Process Variation & Temperature (PVT)

Impact of process variation:



- When curve \Rightarrow shifting towards right then it means that PMOS is driving more current compared to NMOS. i.e., PMOS is fast (Good) & NMOS is slow (Bad).

— since PMOS is driving more current then capacitor will be charging for more time & curve shifts right.

Q: How do we have a good PMOS & Bad NMOS?

- Good PMOS \rightarrow Fast PMOS \rightarrow more current will be obtained
(It can be due to ~~a~~ bigger size of PMOS or higher doping)

$$* \beta\text{-ratio of model} \Rightarrow \frac{B_P}{B_n} = \frac{\mu_p C_{ox}}{\mu_n C_{ox}} = \frac{k_p}{k_n}$$

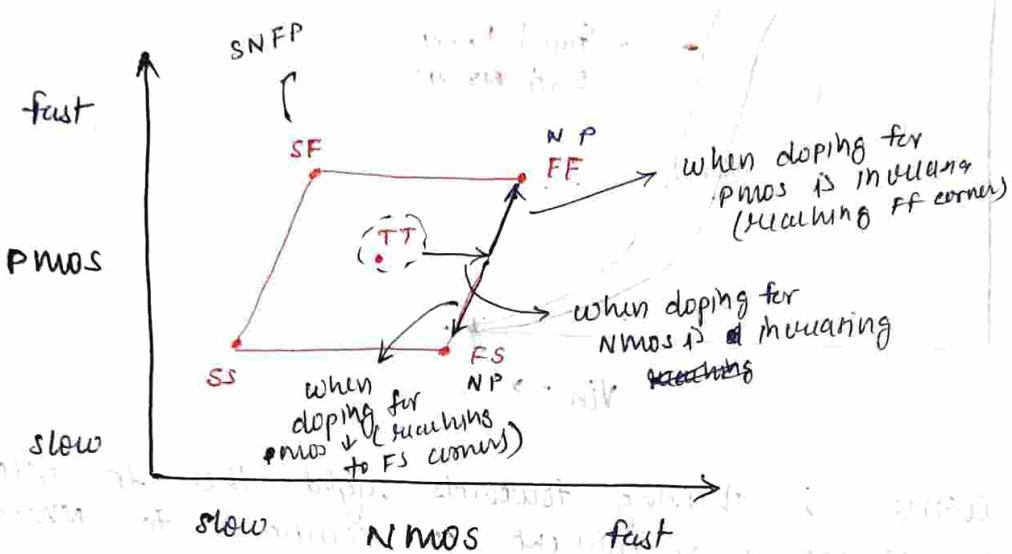
$$* \beta\text{-ratio} \Rightarrow \frac{I_{ON(PMOS)}}{I_{OFF(PMOS)}}$$

if $\beta\text{-ratio} > 1 \Rightarrow$ It is ON for more time & off for less time.
of MOSFET

- * As we can see that this PMOS & NMOS can have different driving strength depending on how they have been manufactured.
- It is possible that due to manufacturing differences any one of either PMOS & NMOS can be slow or fast. And this case called process slots.

Process Variations (PV)

- Variation in process is due to the following things:
 - (i) L_{eff}
 - (ii) V_t
 - (iii) Box of pmos & NMOS
- So we targeted to make a typical device (TT).

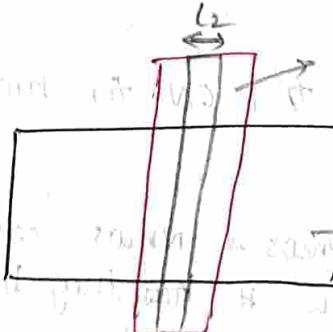


→ We started from TT.

* (why?) If we say that etching of poly (when we are etching the gates) that also happens for few more milliseconds more

so both NMOS + PMOS become little faster now ab. note: Compacting is going on the surface of the poly (etching time)

— when poly is etched more than length \rightarrow shorter \rightarrow $I_D \propto \frac{1}{L_{eff}}$ \Rightarrow more current will flow + becomes fast.



diffusion.

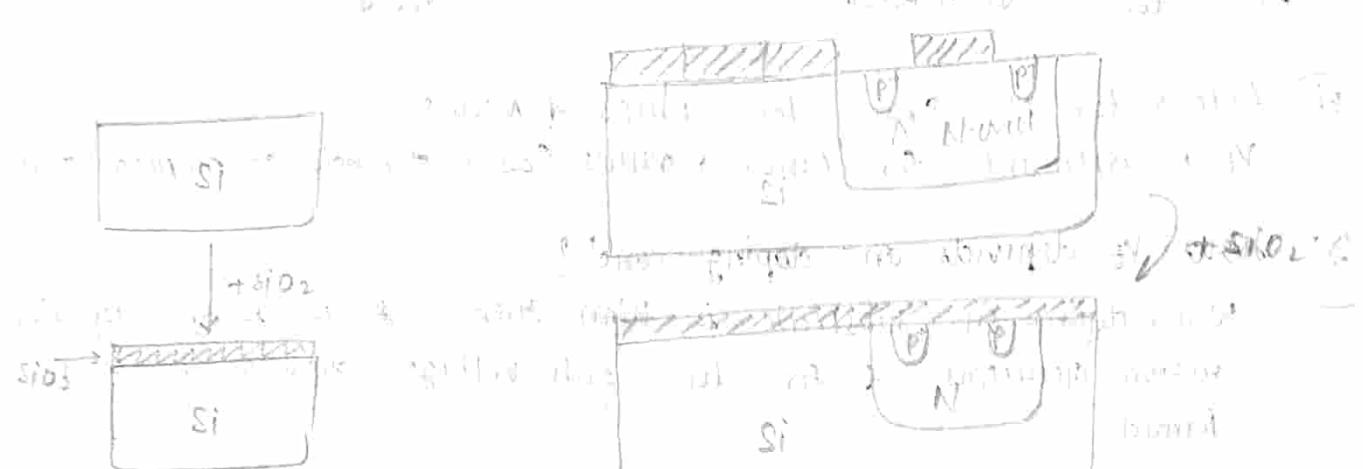
self-alignment mask with patterned photoresist and no photoresist intermixes with poly dopant to form with nitride, and bottom most layer mask "Self Aligned Gate" which has poly on

Q. What is the process of manufacturing?

- (i) First we will have the substrate & then we will have the N-well.
- (ii) Substrate → SiO_2 → etching area for n-well → etched
- N-well → SiO_2 → etching area for poly
→ SiO_2 → poly (gate is formed) → etched area for diffusion area.
→ diffusion → etched.

[Visualize all this]

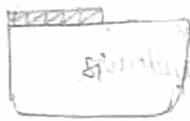
& then source & drain will always be in contact with each other.



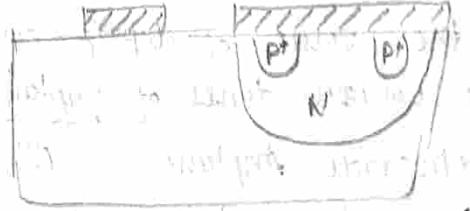
Adding P+ to Si



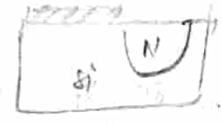
etching for n-diffusion



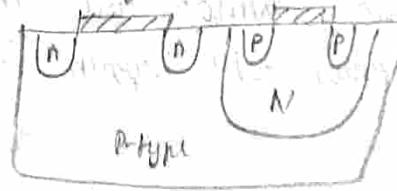
etching



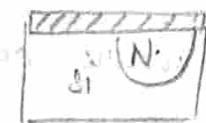
Adding n-diffusion



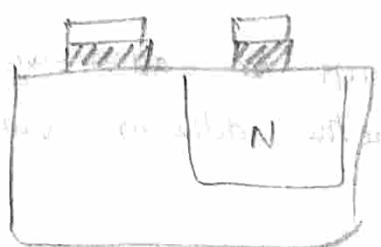
$\downarrow +\text{O}_2$ (etching)



P-type



etching



It is not a self-aligned poly.

- In the previous method of fabrication there is a possibility of poly-mismatch portion.
- But when poly is added ^{before} the n+ & p+ diffusion then poly is self aligned & no mismatch is observed

For a fast device :

- $L_{eff} \downarrow$ (short) $I_D = \frac{4}{3} \mu_0 C_{ox} \frac{W}{L} [(V_{GS} - V_t) V_{DS} - \frac{V_D^2}{2}]$
- $V_t \downarrow$ (low)
- $t_{ox} \downarrow$ (high) $\propto t_{ox} = \frac{\varepsilon A W}{t_{ox}}$

it L_{eff} & t_{ox} is same for PMOS & NMOS
 V_t is different for PMOS & NMOS (as it depends on doping concn)

Q: Now V_t depends on doping concn?

- When doping in substrate is high then width of depletion region decreases & for low gate voltage depletion region is formed
- ⇒ Doping \uparrow , width of depletion region \downarrow , $t_{ox}(eff)$ for capacitance \downarrow , $\propto t_{ox} = \frac{\varepsilon A W}{t_{ox}}$, $\downarrow V \propto \frac{1}{C} \uparrow$

it Changing the value of doping concn is called implant:
 There are various types of implant:

- | | |
|-------------------------|---|
| (i) retrograde implant | (ii) halo implant |
| (iii) superhalo implant | (iv) does not only depend on amount of doping, it depends on depth of doping & many other factors |

* Doping of PMOS & NMOS will independent step, so V_t shift of NMOS & PMOS will happen in different direction of process corners.

* On wafer a large no. of dies are present where some have FF, some have FS, some have SS & so on.

→ * Do we check ~~for~~ our design for all the corners & compare the delay for various process corners.

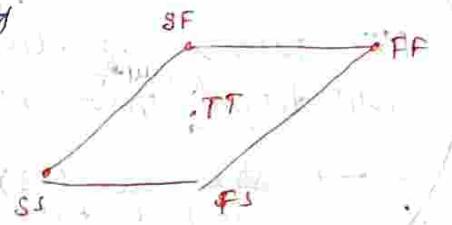
- In actual the device will lie anywhere in b/w four process corners.

Q: On what other factors the delay will change?

- (i) process corners (due to V_T , t_{ox} , L_{eff})
- (ii) Temperature
- (iii) Voltage of operation (V_{dd})
- (iv) Input slope
- (v) width of the mosfets (NMOS + PMOS)

* Manufactured IC will not have as bad as you designed it for.
* It will not be as good as you designed it for.

- As a designer we can't check at every point inside the process corners
∴ we only check at the corners & run we generate an analysis.



Environment Variations:

- mobility (μ_n) variation wrt temperature.

$$\mu_n \propto \frac{T^{3/2}}{N_D}$$

fast (at lower temp.)

$$\mu_n \propto \frac{F}{T^{3/2}}$$

slow (at higher temp.)

due to ionized impurity scattering due to phonon scattering

- At low temp & high voltage we get fast device
- At high temp & low Vdd voltage we get slow device

- V_{DD} & Temp. also vary ^{wrt time & space}.

Corner	Voltage	Temp
F	1.98	0°C
S	1.8	70°C
S	1.62	125°C

Ans: At high temp. I_{DN} \downarrow & due to low voltage less current carriers are generated. Hence slow device.

— At low temp. $I_{DN} \uparrow$ & due to high voltage more carriers are generated. Hence fast device.

Process corners:

- Process corners describe worst case variations
 - If design works in all corners, it will probably work for any variation.

Ques. Describe corner with four letters (T, F, S)!

(a) ~~process speed~~ → that can be described by more 4 variations

PVT means process, voltage + Temperature for which you will design your chip.

→ ~~process speed~~ (SS), low temp + high voltage
Fast - Fast (FF)

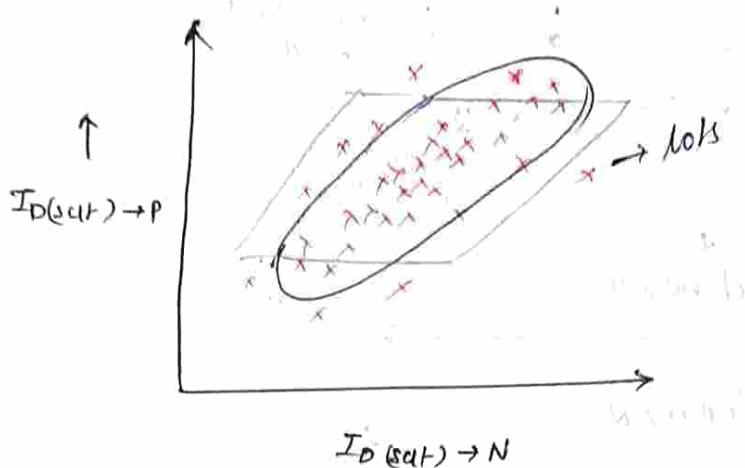
- NMOS speed
- PMOS speed
- Voltage
- Temperature

new?

- Contamination delay will be ⁽¹⁾ ~~worst~~ ^{max} in FF, low temp & high voltage
- Propagation delay will be worst in SS, high temp & low voltage
- Some critical simulation corner includes: (for worst case)

Purpose	NMOS	PMOS	VDD	Temp
Cycle time	S	S	S	S
Power	F	F	F	F
Subthreshold leakage	F	F	F	S

Model is a fitting exercise:



- (1) A large no. of wafers a manufacturer simultaneously & each wafer is called a lot.
- (2) When we characterize a particular lot then if it is FF, FS or SF then it is called lot variation.
- (3) And suppose we pick 100 ~~x~~ points in a lot, and when we characterize them then it is called device variation.

- In ~~modeling~~ equality we don't get ~~rectangle~~ we get ellipse. In ellipse you define your process corners.
- It is possible that in the lots outside the ~~rectangle~~ you will get some extra yield loss.
- These small lots are for different process corners at which our design or device may operate.

Conclusion:- (i) Our chip can be fabricated at any of this (x) i.e. lots. (called lot variation)

(ii) When we take a single ~~wafer~~ which is ~~operated~~ is fabricated at a particular lot then we do one "monte carlo simulation" to ~~be~~ validate a design with device variation like SS, SF, FF, FS.

Q. What is Monte Carlo Analysis?

- It is the analysis in which you will vary different parameters which are known to exhibit variation.

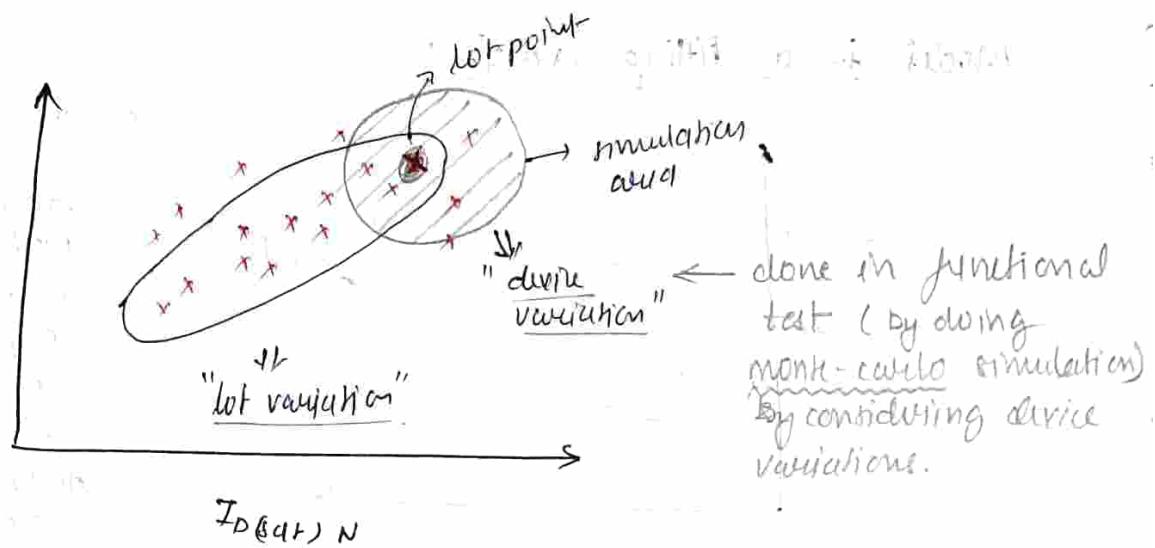
- Parameters that exhibit variations are (Monte Carlo parameters):

• Doping (V_t)

• L_{eff}

• t_{ox}

This variation can be put by putting variation corners SS, SF, FS, FF.



- When you do the monte-carlo simulation around the selected lot then you get the big shaded area around the lot point
 - N is the sample of x itors that could be manufactured on the particular lot.

- When we measure current in all the x itors inside the lot, and we list down all such current values & say my current distribution in this lot & this distribution is feed into simulator & we run monte-carlo simulation.

$$X_1 : \text{PMOS} \rightarrow V_{tss} + \Delta V_{t1}$$

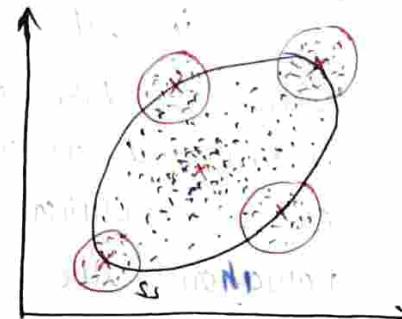
$$X_2 : \text{PMOS} \rightarrow V_{tss} + \Delta V_{t2}$$

$$X_3 : \text{NMOS} \rightarrow V_{tss} + \Delta V_{t3}$$

$$X_4 : \text{NMOS} \rightarrow V_{tss} + \Delta V_{t4}$$

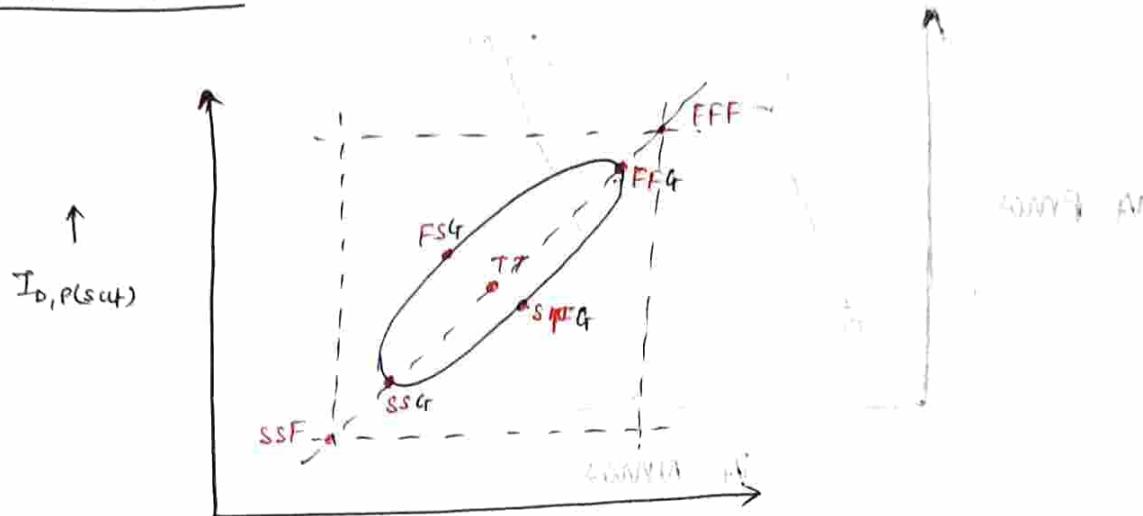
After selecting V_{tss} , we add ΔV_t thus is called monte-carlo simulation

but is for
monte-carlo variation
characterization



- characterization \rightarrow You have a design (on silicon or ~~anywhere~~) & you measure current, leakage, delay \rightarrow called characterization
- Verification \rightarrow You put the parametric variation like SS, SF, OTT & you put monte-carlo variation & then you verify your device will work fine or not
- Validation \rightarrow ~~studying~~ studying the ~~variations~~ variations in silicon & you validate whether it is meeting your ~~specification~~ specification or not.

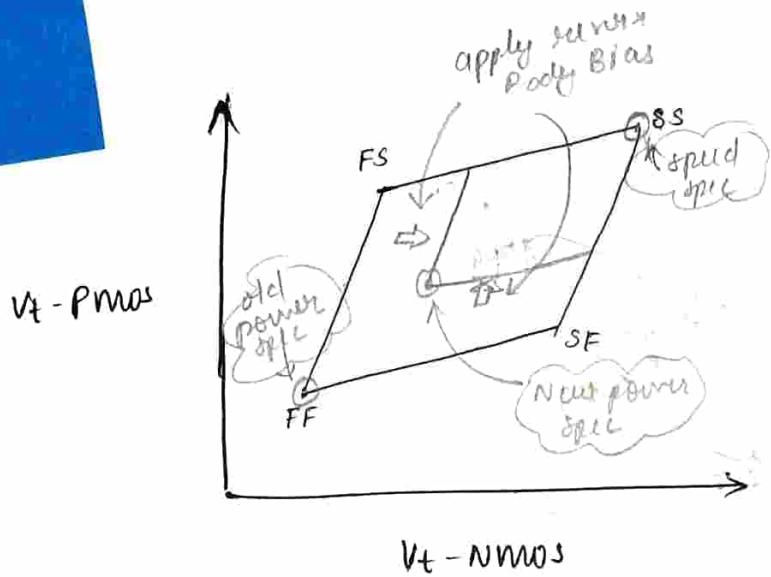
Functional models:



- Some foundries provide extra models called SSF & FFF for verifying the device.
- ⇒ mc command for more accurate simulation in ednlib

Process Compensation: [Area of research]

- Suppose a device is designed fast but more leakage is also very high & delays are much slower
- Eg. You designed something to operate at 2GHz, but after manufacturing it is operating at 3GHz. But we can't operate it at 3GHz because it can have resonance issue, inductance effect and so on.
- So how to operate the device at 2GHz which is fast & leakage is also very high?
- We can try to reduce the leakage. For that we can do multiple things:
 - Apply Body bias (reverse body bias), so device moves towards slow & Vt stays to inactivate (if leakage ↓) (doubt)
Adaptive Body Bias
 - Lower the Vdd from 1.2V to 1.15V.
At 1.2V VDD, operating freq is 3GHz
At 1.15V VDD, operating freq decreases
Adaptive Voltage scaling



- There are constraints on how much you can do reverse body bias, because tunneling can occur.
- * In FDSOI (fully depleted silicon on insulator), there is no junction in source to drain region, there is a buried oxide
 - Therefore there is no concept of junction leakage.

Power Dissipation

Q. Where does the power go in CMOS?

- (i) Dynamic power consumption
 - during charging & discharging of capacitor
- (ii) short ckt currents
 - Short ckt path b/w supply rails during switching
 - when short ckt path exists b/w VDD & ground b/w the supply rails during switching operations.
- (iii) Leakage
 - Leaking diode & transistors

⇒ Power is drawn from a voltage source attached to the VDD pins of a chip.

(i) Instantaneous Power :

$$P(t) = I(t)V(t)$$

- Designer must be interested in instantaneous power as it gives us the estimate of what kind of current is flowing & with that we can design our power grid. The current should not be that high that my power grid collapses.

(ii) Energy:

$$E = \int_0^T P(t) dt$$

- As a designer we must know the energy consumed for one calculation, so we can calculate the total energy required by a calculator to do 1 billion calculation and we place the battery according to that.

(iii) Average power :

$$P_{avg} = \frac{E}{T} = \frac{1}{T} \int_0^T P(t) dt$$

- When you quote that your bulb consumes 10watt then it is avg. power so this data is available for the user to understand & so the designer must define it.

* As a designer you are interested in all three of them.

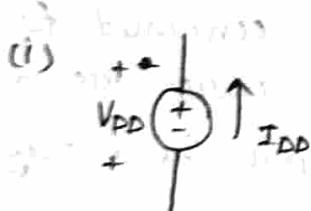
- Q: On what factor / parameter I (current) depends on? for CMOS
-
- Resistance of wire
 - V_{DD}
 - Resistance of channel or $\frac{W}{L}$ or U_n or V_T
 - slim → that will tell for how long short circuit current will flow
 - load
 - charge = $\int i(t) dt$
 - Frequency → No. of times the CMOS inverter will toggle decides the current through inverter at a particular instant.

$$T = \frac{dQ}{dt} \rightarrow \text{charge consumed per unit time}$$

* Now we are talking about no. of times the toggling is happening, so freq has very imp role to play & associated with it is called activity factor.

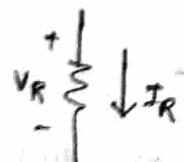
* Instantaneous power depends on freq. of operation, if you want to limit your instantaneous power you have to reduce the freq. of operation

Power in circuit elements



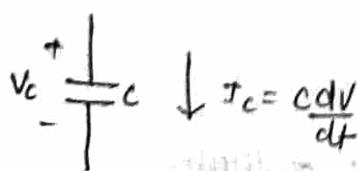
$$P_{VDD}(t) = I_{DD}(t) \cdot V_{DD}$$

(ii)



$$P_R(t) = \frac{V_R^2(t)}{R} = \underline{\underline{I_R^2(t) \cdot R}}$$

(iii)



$$E_C = \int I(t)V(t) dt = \int C \frac{dV}{dt} \cdot V(t) dt$$

$$E_C = C \int V(t) dV = \frac{1}{2} CV^2 - \textcircled{1}$$

Sources of Power Dissipation

- $P_{\text{total}} = P_{\text{dynamic}} + P_{\text{static}}$

- (i) Dynamic power: $P_{\text{dynamic}} = P_{\text{switching}} + P_{\text{short circuit}}$

- switching local capacitance

- short circuit current (C when driver gets on & load connects to GND).

- (ii) Static power: $P_{\text{static}} = (I_{\text{sub}} + I_{\text{gak}} + I_{\text{junc}}) V_{\text{DD}}$

$$P_{\text{static}} = I_{\text{sub}} + I_{\text{gak}} + I_{\text{junc}}$$

↓ ↓ ↓
 subthreshold Gate Junction
 leakage leakage leakage

While going through the dynamic & static power, we will consider a chip having 1 Billion transistors:

For 50M logic transistors

- Avg width = 12λ

- Activity factor = 0.1

• 950M memory transistors

- Avg width = 4λ (small)

- Activity factor = 0.02

• 1.0V, 65nm process

• $C = \text{EFF}/\mu\text{m} (\text{gate}) + 0.8\text{FF}/\mu\text{m} (\text{diffusion})$

- Above all are assumptions

$w_p \text{ or } w_n = 4\lambda$

$\lambda \rightarrow$ minimum width of transistor

1. Energy delivered by power supply:

$$E_{\text{del}} = \int_0^{\infty} i(t) \cdot V_{\text{supply}}(t) dt = \int_0^{\infty} C_L \frac{dV_o}{dt} \cdot V_{\text{DD}} dt = C_L V_{\text{DD}}^2$$

2. Energy consumed by capacitor:

$$E_C = \int_0^{\infty} V_o(t) \cdot i(t) dt = \int_0^{\infty} V_o \cdot \frac{C_L dV_o}{dt} dt = C_L \int_0^{\infty} V_o dV_o$$

$$\frac{1}{2} C_L V_{\text{DD}}^2$$

1. Dynamic Power : \rightarrow Power because of charging & discharging

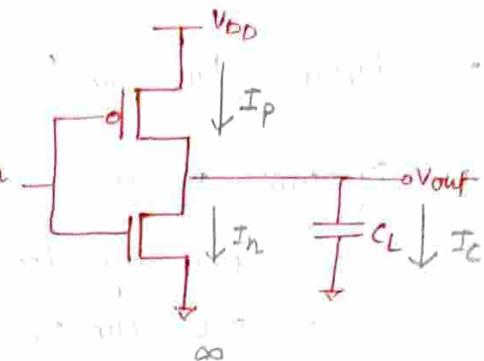
(i) When the gate o/p rises:

- Energy stored in capacitor is

$$E_C = \frac{1}{2} C_L V_{DD}^2 \quad (\text{from } ①)$$

- Energy drawn from the supply is

$$\begin{aligned} E_{\text{add}} &= \int_0^\infty I(t) V_{DD} dt = \int_0^\infty C_L \frac{dV}{dt} V_{DD} dt = \int_0^\infty C_L V_{DD} dV \quad \text{not a} \\ &= C_L V_{DD} \int_0^\infty dV = C_L V_{DD}^2 \end{aligned}$$



① Energy dissipated is function of Resistance (R_f) of NMOS & PMOS hence it is not dependent on design parameters such as width of itself.

- Half the energy from V_{DD} is dissipated in the PMOS transistor as heat, and other half is stored in capacitor.

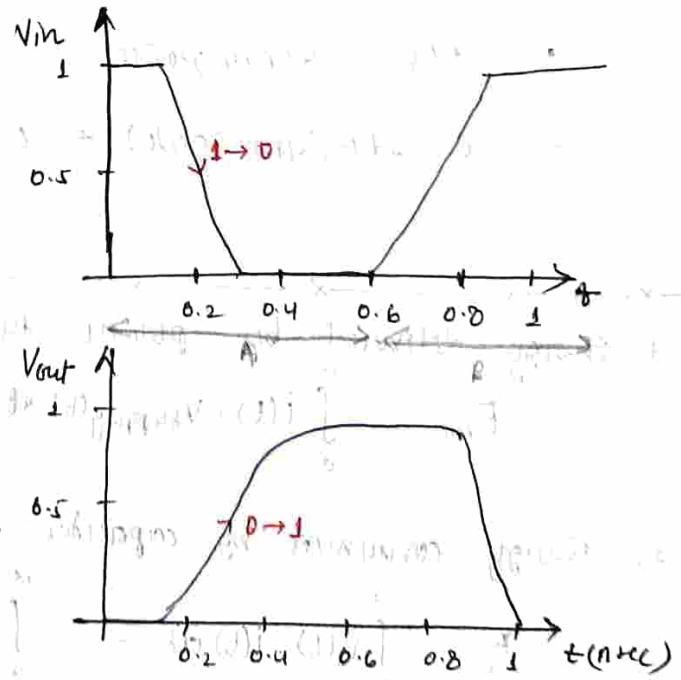
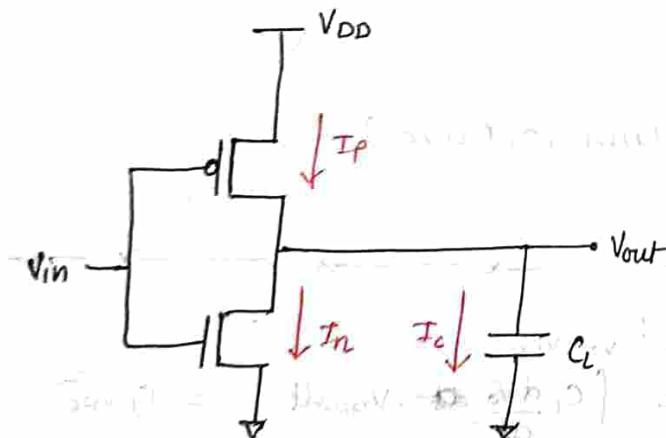
② But depends on W because of its own parasitic capacitance as well.

(ii) When gate o/p falls - dissipated as heat in NMOS

- Energy $\frac{1}{2} C_L V_{DD}^2$ gets dissipated as heat in NMOS

$$\text{Energy in } \Rightarrow E_C = \int_{V_{DD}} V_o \cdot i(t) dt = \int_{V_{DD}} V_o \cdot C_L dV_o = \frac{1}{2} C_L V_{DD}^2$$

Switching waveforms:



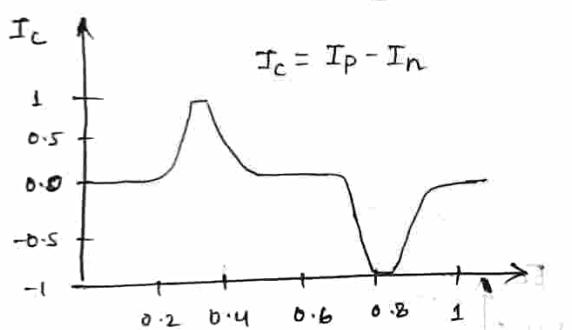
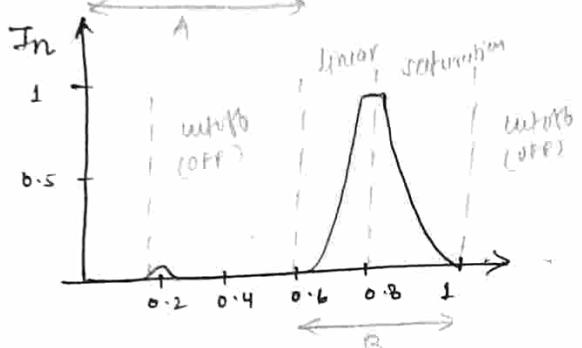
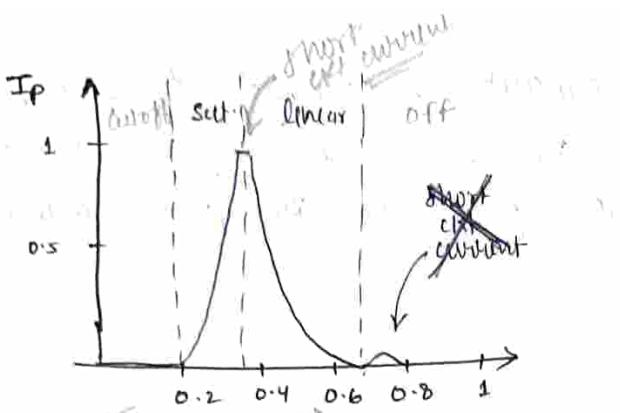
- When the o/p switches from 1 to 0, the NMOS transistor turns ON & charges the load to V_{DD} .

- Energy stored in capacitor

$$E_C = \frac{1}{2} C_L V_{DD}^2$$

$$V_{out} = \frac{1}{C} \int I_C dt$$

$$\therefore I_C(t) = C_L \frac{dV_{out}}{dt}$$



~~the~~ current will flow through PMOS but ~~the~~ value of current I_D starts decreasing. & after complete charging $I_D = 0A$.

Now when I_D transitions from 0 to + during the transition, NMOS is ON & the current through NMOS increases as it will be in linear region.

After some time NMOS enters into the saturation region & current starts decreasing & when cap. discharges completely then $I_N = 0A$

* A short ~~circuit~~ current will be flowing when both NMOS & PMOS are off & no current is supposed to flow.

- when V_{in} rises from '0' to '1', the PMOS starts to turn off. However, there is a small blip of current while the partially ON PMOS fights against the NMOS, this is called short-ckt current.

- But the power supply ~~current~~ $= CV_{DD}^2$

- Power dissipated in the PMOS transistor ~~is proportional to current squared~~ is

$$= \frac{1}{2} C_L V_{DD}^2$$

* The power dissipated depends only on the load capacitance, ~~is it not on the size of the transistor~~ ~~do we not have to consider the speed at which gate voltage changes?~~

- When I_D is transitioning from 1 to 0 then current is flowing through PMOS only & $I_N = 0A$ since capacitor is charging & $I_C = 0A$.
- But when transition stops & $I_D = 0$ then current is at peak value. Capacitor is charged completely.
- And when $I_D = 0$ for sometime & cap. is completely charged.

~~the~~ current will flow through PMOS but ~~the~~ value of current I_D starts decreasing. & after complete charging $I_D = 0A$.

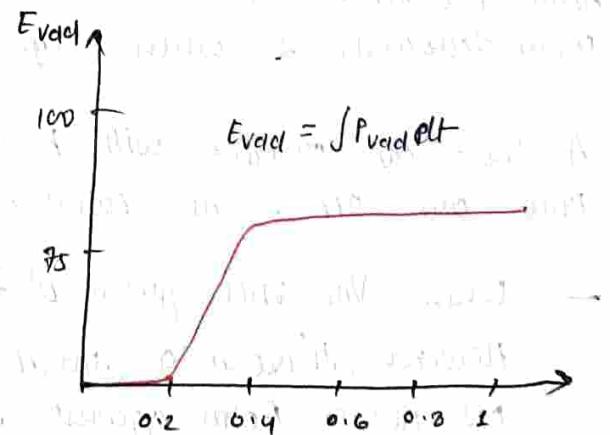
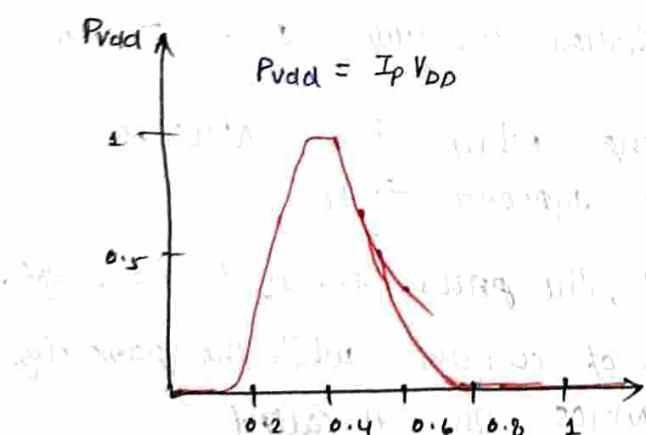
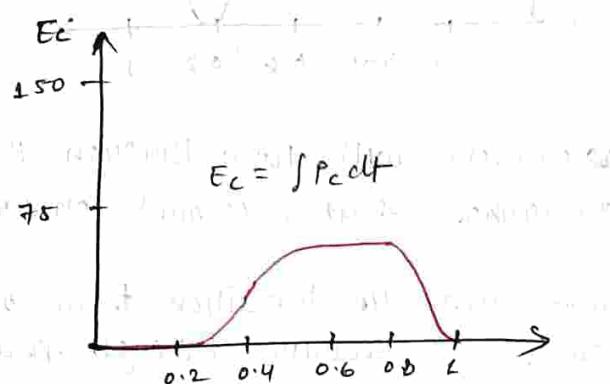
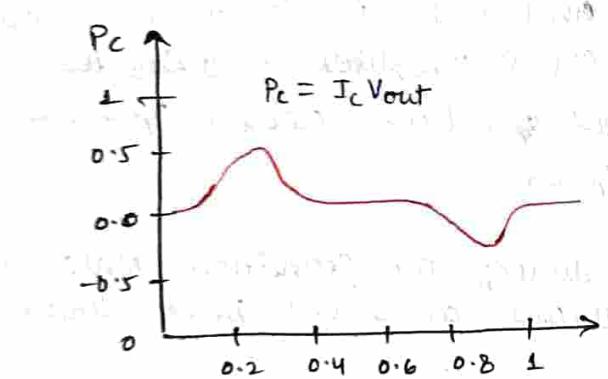
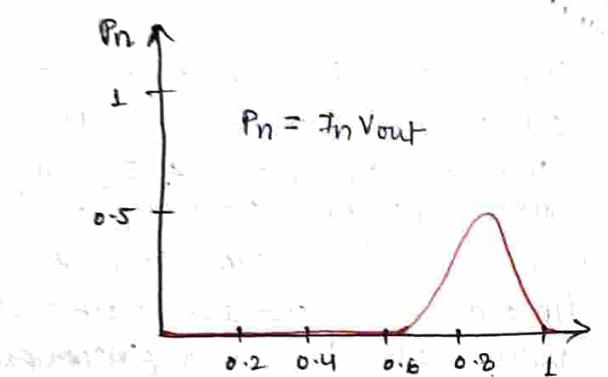
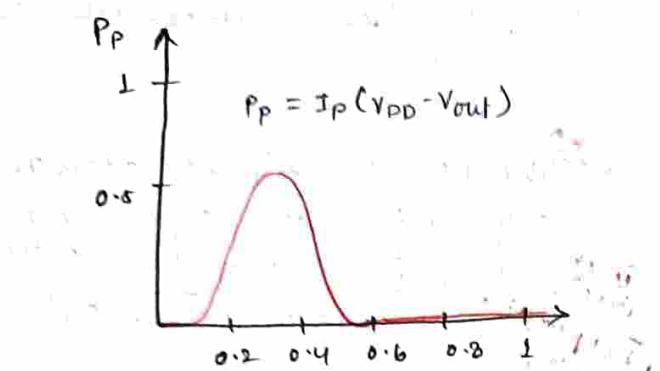
Now when I_D transitions from 0 to + during the transition, NMOS is ON & the current through NMOS increases as it will be in linear region.

After some time NMOS enters into the saturation region & current starts decreasing & when cap. discharges completely then $I_N = 0A$

* A short ~~circuit~~ current will be flowing when both NMOS & PMOS are off & no current is supposed to flow.

- when V_{in} rises from '0' to '1', the PMOS starts to turn off. However, there is a small blip of current while the partially ON PMOS fights against the NMOS, this is called short-ckt current.

- The energy stored on the capacitor is dumped to GND, the short-circuit current consumes an almost imperceptibly small 2.7 fJ of additional energy from V_{DD} during the transition.



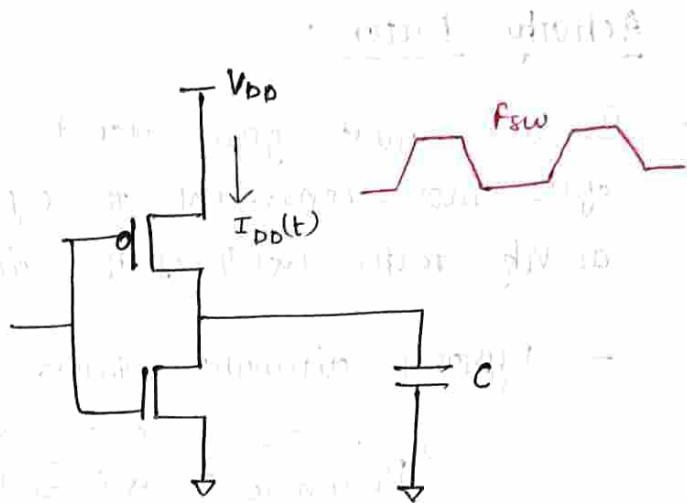
Switching Power:

$$P_{\text{switching}} = \frac{1}{T} \int_0^T I_{DD}(t) \cdot V_{DD} dt$$

$$= \frac{V_{DD}}{T} \int_0^T I_{DD}(t) dt$$

$$= \frac{V_{DD}}{T} [T f_{sw} (V_{DD})]$$

$$= C V_{DD}^2 f_{sw}$$



$$\textcircled{a} P_{\text{switching}} = \frac{E}{T} = \frac{T f_{sw} C V_{DD}^2}{T} = C V_{DD}^2 f_{sw}$$

~~when gate switches at some average freq. f_{sw} . Over time period T , min. load will be charged & discharged $T f_{sw}$ times.~~

* since, Energy drawn from the supply

$$E_{\text{dd}} = C_L V_{DD}^2$$

if gate switches at freq. f_{sw} over interval T :

$$\therefore E_{\text{dd}} = T f_{sw} C_L V_{DD}^2$$

\therefore switching power will be

$$\text{Power} = \frac{E_{\text{dd}}}{T} = C_L V_{DD}^2 f_{sw}$$

no. of times the
ckt switches or toggles
in a time period T

~~** $\frac{1}{2} C V_{DD}^2$ when ckt rises, $\frac{1}{2} C V_{DD}^2$ when ckt falls~~

$\Rightarrow [p.u \text{ cycles } C V_{DD}^2] \times \text{no. of times the ckt will toggle}$

$$\Rightarrow C_L V_{DD}^2 \times f_{sw}$$

Activity factor :

- Because most gates do not switch every clock cycle, it is often more convenient to express switching freq fno as an activity factor (α) times the clock freq. f.
- Dynamic dissipation power may be written as:

$$P_{\text{switching}} = \alpha C V_{DD}^2 f$$

clock freq.

where,

$$f_{sw} = \alpha f$$

α → activity factor

No. of times any node in the circuit is switching per CLK cycle.

- If signal is a clock, $\alpha = 1$.
- If the signal switches one per cycle of clock, $\alpha = 1/2$.
- If the signal switches twice per cycle of clock, $\alpha = 2/2 = 1$.
- If the signal switches three per cycle of clock, $\alpha = 3/2$.

CLK



$$\alpha = 1$$

(i) signals 1



$$\alpha = 1/2$$

(ii) signal 2



$$\alpha = 2/2 = 1$$

(iii) signal 3



$$\alpha = 3/2$$

(iv) signal 4



$$\alpha = 3/2$$

* for 1 glitch, $\alpha = 3/2$ ~~reduces~~
2 glitch, $\alpha = 5/2$

} th signal 4

it we don't want glitches in our system because

(a) It is a signal integrity issue

(b) Power issue (more power dissipation)

Short circuit current

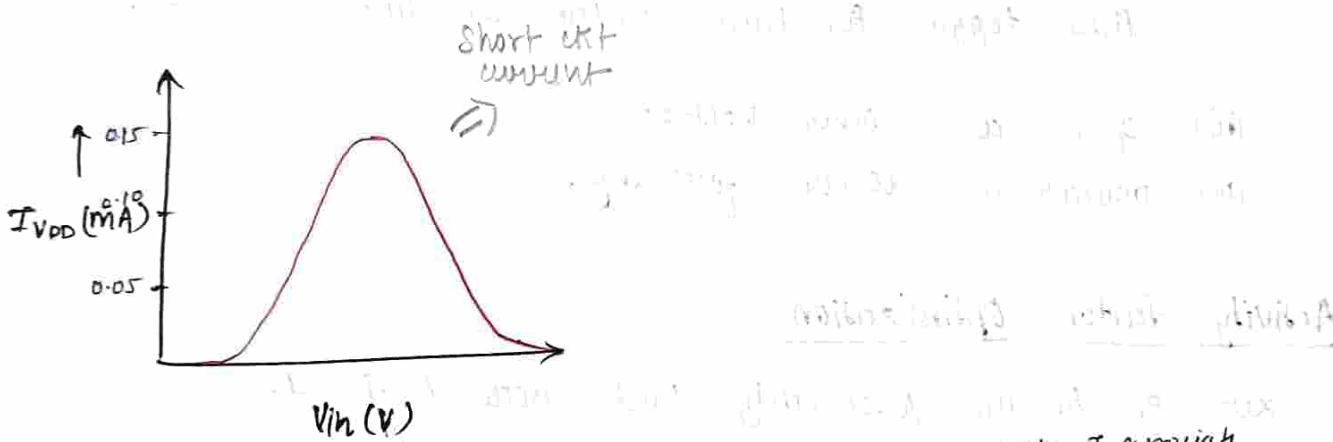
- When Xistor switch both NMOS & PMOS networks may be momentarily ON at once.

→ It leads to a blip of "short ckt" current

It because some current may flow from VDD to GND (very less) + some current may flow from VDD to Cload + Cload to GND

- < 10% of dynamic power if rise/fall time of IIP + OLP are comparable. [rise/fall time (1), short ckt current (2)]

- We generally ignore this component while doing the calculation but consider it while doing the simulation.



Dynamical power for 1B Xistor chip : $\frac{\text{how much width I can work}}{\text{width } \lambda, \lambda = 25 \text{ nm here}}$

$$(i) \text{C logic} = (50 \times 10^6) (12\lambda) \left(\frac{0.025 \mu\text{m}}{\lambda}\right) \left(\frac{1.8 \text{ fF}}{\mu\text{m}}\right) = 27 \text{ nF}$$

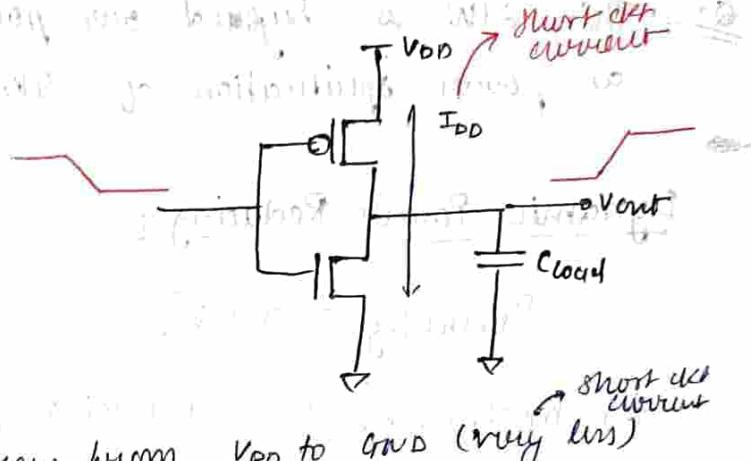
$$= 50 \times 10^6 \times 12 \times 0.025 \times 1.8 \text{ fF}$$

$$(ii) \text{C memory} = (950 \times 10^6) (4\lambda) \left(\frac{0.025 \mu\text{m}}{\lambda}\right) \left(\frac{1.8 \text{ fF}}{\mu\text{m}}\right) = 17 \text{ nF}$$

$$(iii) \text{Pdynamic} = (0.1 \text{ C logic} + 0.02 \text{ C memory}) (1.0)^2 (2 \text{ GHz}) = 6.1 \text{ W}$$

α for logic
 α for memory
(Activity factor)

While taking the value of ' α ' we have considered most this kind of glitches can be more.



Q: This 6.1W is beyond our power budget we can only operate at power specification of 4W. So what can be done?

Dynamic Power Reducing:

$$P_{\text{switching}} = \alpha C V_{DD}^2 f$$

- for minimizing the switching power:

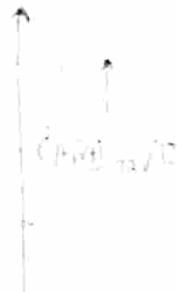
(i) Decrease Activity factor per ckt (how to do that?)

• - by gating the device, it will put NAND at some places & enables only the part of the ckt & not other parts of the ckt, so if it reduce the activity factor of the ckt, it reduces power.

(ii) Reduce the capacitance, at least the capacitances that toggle for that we can use smaller devices.

(iii) operate at lower voltage

(iv) operate at lower frequency.



Activity factor Optimization

- Let P_i be the probability that node i is 1.

$$\therefore P_i = 1 - \bar{P}_i$$

Therefore And,

$$\alpha = P_i * \bar{P}_i$$

- completely random data has $P = \frac{1}{2} = 0.5$ + $\alpha = 0.25$

$$\begin{aligned}\therefore \alpha &= P_i * \bar{P}_i \\ &= \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} = 0.25\end{aligned}$$

- Data is often not completely random

- E.g. Upper bits of 64-bit words representing bank account balance are usually 0.

- Data propagating through ANDs & ORs has a lower activity factor since each stage adds to the total time
- it depends on design, but typically, $\alpha \approx 0.1$.

Q: Why do AND & OR gate have lower activity factor?

* switching Probability:

- what is prob. of o/p of AND2 being 1?
- o/p of AND2 is 1 when both A & B is 1.

$$\therefore P = P_A \cdot P_B = \frac{1}{4}$$

A	B	AB
0	0	0
0	1	0
1	0	0
1	1	1

$$P = \frac{1}{4}$$

- similarly for NOR2, if o/p is 1 when any of input is 1 \Rightarrow 1 - when both inputs are 0.

$$\therefore P = 1 - \bar{P}_A \bar{P}_B = 1 - \frac{1}{4} = \frac{3}{4}$$

Gate	P_Y
AND2	$P_A P_B$
AND3	$P_A P_B P_C$
OR2	$1 - \bar{P}_A \bar{P}_B$
NAND2	$1 - P_A P_B$
NOR2	$\bar{P}_A \bar{P}_B$
XOR2	$P_A \bar{P}_B + \bar{P}_A P_B$

$$\alpha_{OR2} = P \times \bar{P} = \frac{3}{4} \times \frac{1}{4} = \frac{3}{16} \text{ (min)}$$

$$\alpha_{AND2} = P \times \bar{P} = \frac{1}{4} \times \frac{3}{4} = \frac{3}{16} \text{ (min)}$$

$$\alpha_{XOR2} = P \times \bar{P} = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4}$$

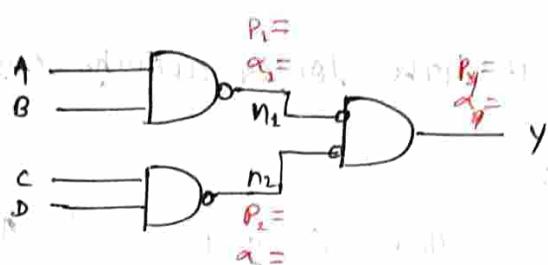
$$\alpha_{NAND2} = \frac{1}{4} \times \frac{3}{4} = \frac{3}{16}$$

$$\alpha_{NOR2} = \frac{3}{4} \times \frac{1}{4} = \frac{3}{16}$$

Example: A 4-input AND gate is built out of two levels of gates. Estimate the activity factor at each node if the input have $P = 0.5$.

* $P = 0.5$ means

Prob. of A being 1 is $\frac{1}{2}$.
Prob. of B being 1 is $\frac{1}{2}$.



P_1 = Prob. of o/p of NAND gate being 1

$$P_1 = \frac{3}{4}, \quad \bar{P}_1 = \frac{1}{4}$$

$$\alpha_1 = \frac{3}{4} \cdot \frac{1}{4} = \frac{3}{16}$$

AB	\bar{AB}
00	1
01	1
10	1
11	0

P_2 → Prob. of o/p of NAND gate being 1

$$P_2 = \frac{3}{4}, \quad \bar{P}_2 = \frac{1}{4}$$

$$\alpha_2 = \frac{3}{16}$$

P_y → Prob. of o/p of Gate being 1

$$P_y = \frac{1}{16}, \quad \bar{P}_y = \frac{15}{16}$$

$$\alpha_y = \frac{15}{256}$$

R S	$\bar{R} \bar{S}$
00	1
01	0
10	0
11	0

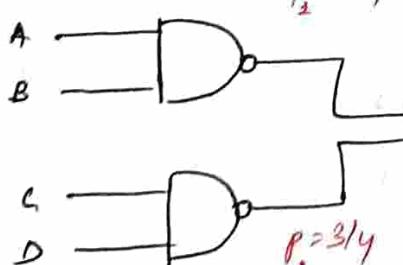
$$Y = \bar{A}\bar{B} + \bar{B}\bar{C}\bar{D}$$

$$= \bar{A}B + \bar{C}D$$

$$Y = AB\bar{C}D$$

$$Y = L \text{ for } ABCD = 1$$

OR

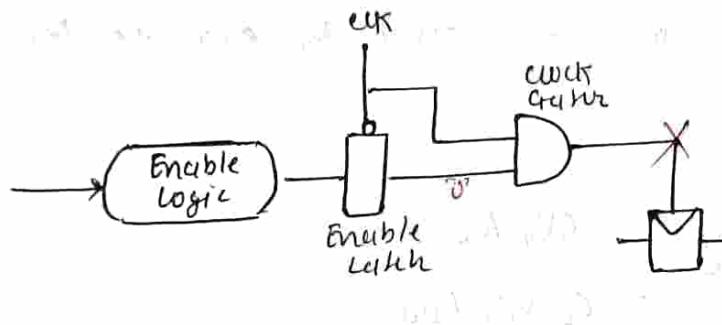


$$P_y = \frac{1}{4} \cdot \frac{1}{4} \text{ (for AND gate)}$$

$$\bar{A}B + \bar{B}\bar{C}\bar{D} = \bar{A}B + \bar{C}D$$

Reducing Dynamic power

1. Clock Gating: It reduces the activity factor.
 - The best way to reduce the activity factor is to turn off the clock to registers in unused blocks.
 - We can put an enable logic, such that when enable logic is 0 then flop will not toggle, and if DIP does not toggle, we save power.
 - Saves clock activity ($\alpha = 1$)
 - Eliminates all switching activity in the block
 - Requires determining if block will be used.



Gate capacitance

- fewer stages of logic
- small gate sizes

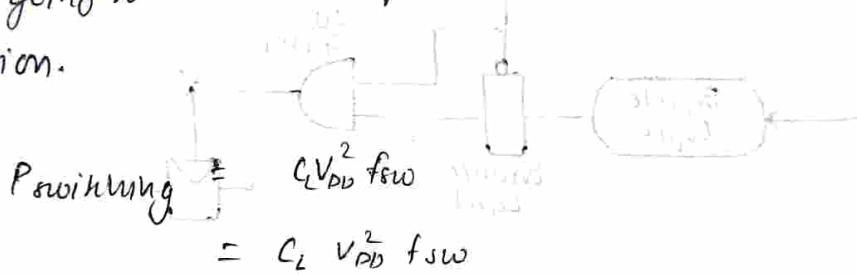
Wire capacitance

- Good floor planning to keep communication blocks close to each other
- Drive long wires with drivers or buffers rather than complex gates
- Things that we have said to reduce logical effort can be used to reduce the ~~dynamic~~ dynamic power also.

Dynamic Voltage & Frequency scaling

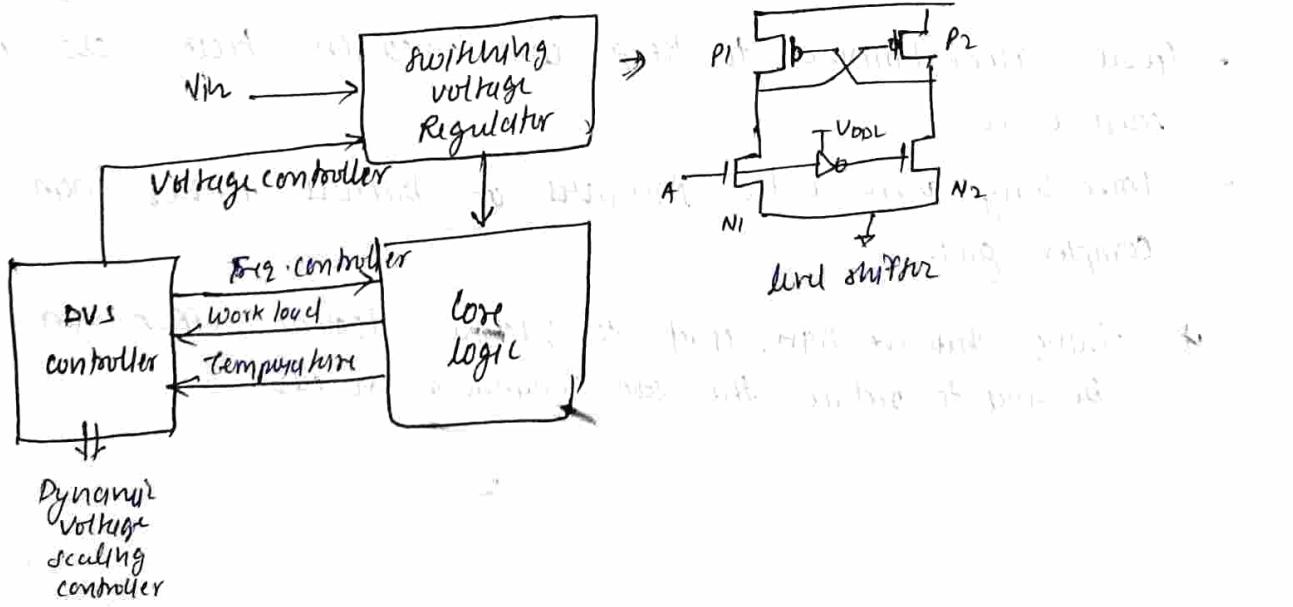
10.09.2019, 10:00 AM

- while using the cell-phone and just browsing the internet mobile will start operating at 200MHz instead of 2GHz i.e. thus ~~so~~ dynamic power consumption is reduced. called Dynamic Frequency scaling.
- If we are operating at 200MHz for part of a minute then CPU will also lower the voltage to reduce the power consumption. So we use 0.7V supply instead of 1.2V because 200MHz can be easily ~~be~~ operated by ~~at~~ 0.7V (as an ex.) [Dynamic Voltage scaling]
- At ~~200~~ low voltage & 200MHz my delay will increase but at this time we don't need that much speed so we are going to save the power by reducing the ~~the~~ voltage of operation.



so we reduced power consumption ~~by~~ cubically.

- Now how to reduce the voltage?
- by switching voltage regulator with the help of DVS controller



- * First reduce the freq. down voltage
 - * First increase the voltage down freq.
 - means at whatever freq you are operating, the voltage ~~is~~ must be able to do at more general freq.

Quick Refresh

- # Dynamic Power Consumption
 - # Reducing switching power
 - Switching power = $\alpha C V_{DD}^2 f_{clk}$
 - short clk power
 - Activity factor
 - fully Random Data, AND/NAND gate
 - # short - CLK power
 - Clock Transition slope
 - Demand glitches

max delay : SS 0.7V 125°C (slowest) Dynamic voltage & freq.
 min delay : FF 1.82V 25°C (fast)

[max buruge :
153244 - 153250
153251 - 153257
153258 - 153264
153265 - 153271
153272 - 153278
153279 - 153285]

to obtain a quiet place
to do some writing
and a chance to
recharge with a
few hours of playtime
among people at a level
they can relate to.
A simple thing like
an hour with a billion

Thermal stability with quinones and transition metals

agosto 1911 Blaupunkt

Barack Obama
Barack Obama

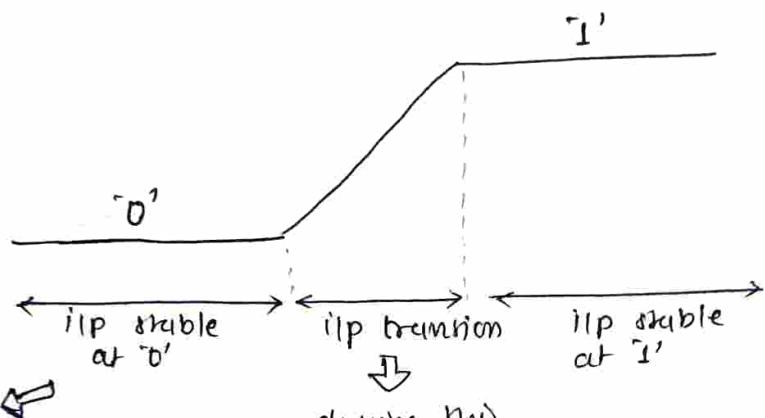
THE NIGHT WILL DISAPPEAR IN
THE NIGHT AND VICE VERSA

variable : rapid shift

skip April 1: annual report

2. Static Power :

- When there is no change in input means everything is stable (no transition is occurring) even then some power is getting consumed that is called static power (because some current is still flowing).
- In CMOS inverter, in a stable state there is no current neither from the NMOS nor from the PMOS because only one of the NMOS or PMOS will be ON, not both will be ON.



when i_{lp} is stable at '0' then PMOS is ON & NMOS is OFF hence at this time capacitor is completely charged & there is no path from V_{DD} to ground. But still some current flows & called static power

during this transition both NMOS & PMOS flow ~~some~~ some current & the power calculated \Rightarrow dynamic power

Revisit 4B transistor chip for static power

• Subthreshold leakage

- Normal V_t : $100nA/mm$
- High V_t : $10nA/mm$

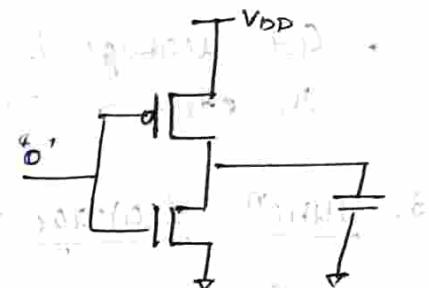
* All memories use High V_t
95% of logic gates use High V_t

- Gate leakage : $5nA/mm$

- Tunnel leakage : Negligible

Note: When $i_{lp} = 0$, and NMOS is off then still some current flows in NMOS i.e., NMOS is leaking means leakage current will flow through NMOS.

- * In static condn, the PMOS will not govern the amount of current flowing through it, the NMOS will govern the amount of current flowing through PMOS also.
- The overall current that can flow from the PMOS is limited by the current that will flow from the NMOS (why?)
- PMOS can't supply more current than the NMOS can sink.
- The capacitor is charged at V_{DD} + cap. can't source current from itself. So if NMOS wants to sink any current, & PMOS is ready to give so much of current will the current come from the capacitor or PMOS?
- from PMOS



Static power calculation for 1B Xistor

- $W_{\text{nominal}-V_t} = (50 \times 10^6) (12A) \left(\frac{0.025 \mu m}{\lambda} \right) (0.05) = 6.25 \times 10^6 \mu m$
- $W_{\text{high}-V_t (\text{logic})} = (50 \times 10^6) (12A) \left(\frac{0.025 \mu m}{\lambda} \right) (0.95)$
- $W_{\text{high}-V_t (\text{memory})} = (950 \times 10^6) (4A) \left(\frac{0.025 \mu m}{\lambda} \right) (1)$
- $\therefore W_{\text{high}-V_t} = W_{\text{high}-V_t (\text{logic})} + W_{\text{high}-V_t (\text{memory})}$
 $= 105.25 \times 10^6 \mu m$
- $I_{\text{sub}} = \frac{[W_{\text{nominal}-V_t} \times \frac{10mA}{\mu m} + W_{\text{high}-V_t} \times \frac{10mA}{\mu m}]}{2}$ → divide by 2 because at half time either of PMOS or NMOS will be ON. → Half of half will occur during any of them.
- $I_{\text{high}} = [W_{\text{nominal}-V_t} + W_{\text{high}-V_t}] \times \frac{5mA}{\mu m} = 275mA$
- $P_{\text{static}} = (584mA + 275mA)(1V) = 859mW$

Sources of Leakage

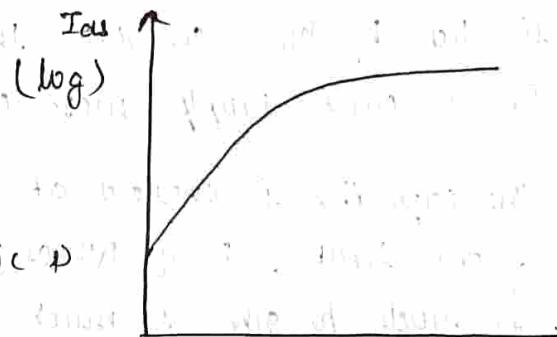
1. subthreshold conduction:

- Transistor also conduct current when $V_{GS} \leq V_{TH}$ & this conduction of current is called subthreshold cond'.

* Transistor can't abruptly turn ON or OFF.

2. Gate leakage:

- Funneling through ultra-thin gate dielectric



- Physical thickness of dielectric is of order 1nm, 2nm etc.

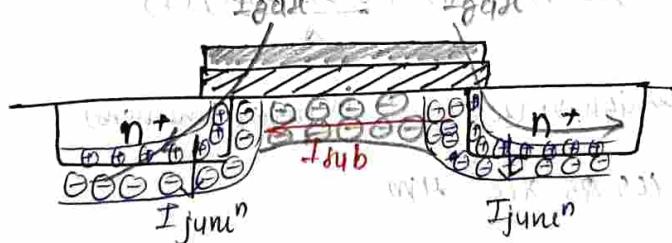
- Gate leakage is reduced by reducing V_{GS} & using high k -dielectric material & using high k -dielectric. Dominant source in contemporary transistors

3. Juncⁿ leakage:

- RB pn-juncⁿ diode current b/w source-body + drain-body.

$$I_{JFF} = \frac{W}{L} I_0 e^{\frac{(V_{GS}-V_T)}{nV_T}} \quad (\text{Volts})$$

$$(I_{JFF})_{\text{forward}} = (I_{JFF})_{\text{reverse}} = (I_{JFF})_{\text{saturation}}$$

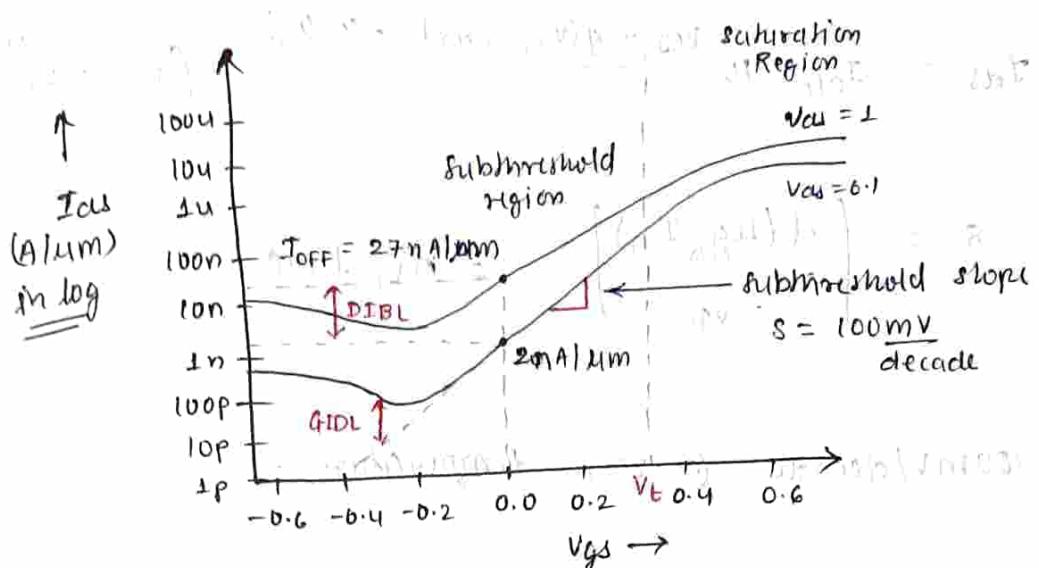


$$\text{P-type substrate}$$

- * There is some subthreshold cond' when current flows pbelow from drain to source where m-o-s has a depletion region & np-juncⁿ also has depletion region.

$$I_{JFF} = (V_S)(A_{eff} + A_{eff}^2)$$

① Subthreshold Leakage



- * In a good transistor, the subthreshold slope is ~~s' < 60mV/decade~~.
- Q: why lower value of 's' is better?
→ * when we say subthreshold slope, $s = 65\text{mV/decade}$. ~~Because it~~
we get $\frac{1}{10}$ in the current below when we go below V_t by 65mV .
- * so if we get $\frac{1}{10}$ in current at lower value of voltage below V_t then it is good.
- * subthreshold leakage means source to drain path. for e^- flow.
- * thus, it is also possible that w/o depletion region current can also flow.
- * subthreshold depends on V_t (if $V_t \downarrow$, subthreshold current \uparrow)
- * subthreshold current also depends on body bias voltage i.e., the relationship b/w source + body.
- * subthreshold depends on temperature, by the virtue of the component $\frac{kT}{q}$.
- * subthreshold leakage exponential with V_{gs} .

$$I_{ds} = I_{d0} \cdot e^{\left(\frac{(V_{gs} - V_{to} + nV_{ds} - K_V V_{sb})}{nV_T} \right)} \left(1 - e^{-\frac{V_{ds}}{V_I}} \right)$$

$n \rightarrow$ process dependent
typically : $1.3 - 1.7$

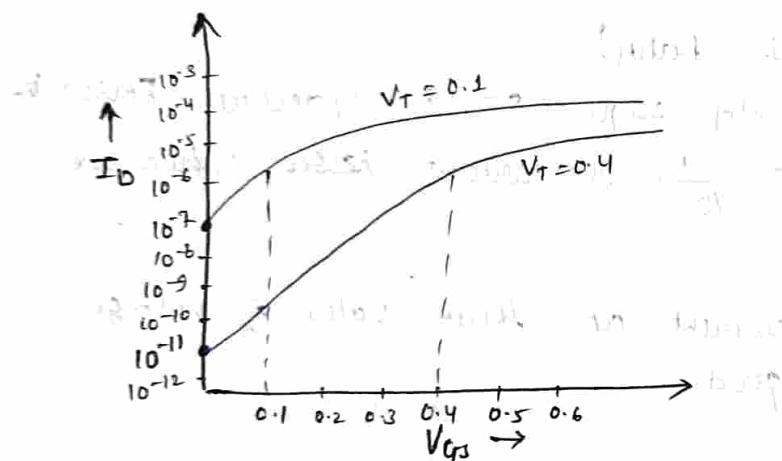
- Rewriting relation to I_{off} on log-scale

$$I_{ds} = I_{off} \cdot 10^{\frac{V_{gs} + \eta(V_{ds} - V_{dd}) - KV_{sp}}{s}} \cdot (1 - e^{-\frac{V_{ds}}{V_T}})$$

$$s = \left[\frac{d(\log_{10} I_{ds})}{dV_{gs}} \right]^{-1} = nV_T \ln 10$$

$s \approx 100 \text{ mV/decade}$ @ room temperature

Subthreshold leakage wrt V_t



Typical value in 0.5nm tech.

for $\eta = 0.1$, $K_V = 0.1 \text{ V/decade}$
and $s = 100 \text{ mV/decade}$

- (i) $I_{off} = 100 \text{ nA/lum} @ V_t = 0.3 \text{ V}$
- (ii) $I_{off} = 1 \text{ nA/lum} @ V_t = 0.4 \text{ V}$
- (iii) $I_{off} = 1 \text{ nA/lum} @ V_t = 0.5 \text{ V}$

* As we decrease V_t (threshold voltage), subthreshold leakage current increases.

1. Drain-Induced Barrier Lowering (DIBL)

- EF from drain affects the channel
- More pronounced in small transistors where the drain is closer to the channel.
- Drain voltage also affects V_t

$$V_t' = V_t - nV_{ds}$$

- High drain voltage causes current to increase.
 - more effective for small channel device.

* We can also get the subthreshold leakage w/o applying any gate voltage (or $-V_G$) i.e. $V_G < V_t$ for NMOS or any drain voltage (V_D) i.e. $V_D < 0$ for NMOS in the M-O-S region.

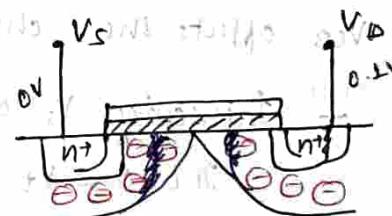
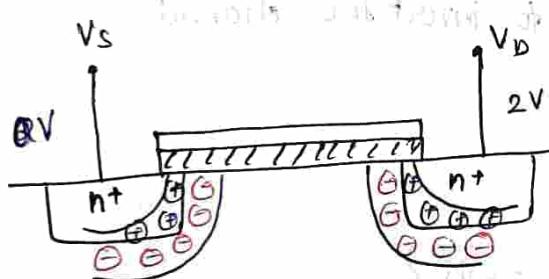
— By applying some drain voltage.

Subthreshold leakage \Rightarrow when device is off or $V_G < V_t$ for NMOS or no depletion region is present even though some current flows from source to drain, when we apply V_{DS} .

(*) Both ON current & off current depends on V_{DS} .

Short-channel effect: As technology is scaling we are using shorter length devices. & this devices starts to behave differently as compared to standard set of eqn when it is called short-channel effect.

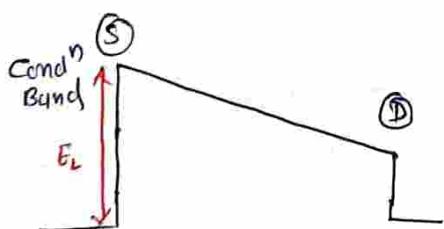
a) why this happens in short-channel device?



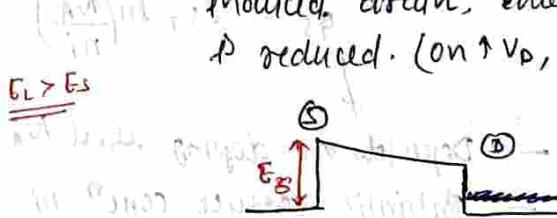
long channel

short channel

for short channel as we reduced drain energy barrier is reduced. ($V_D \uparrow$, $E_{barrier} \downarrow$)



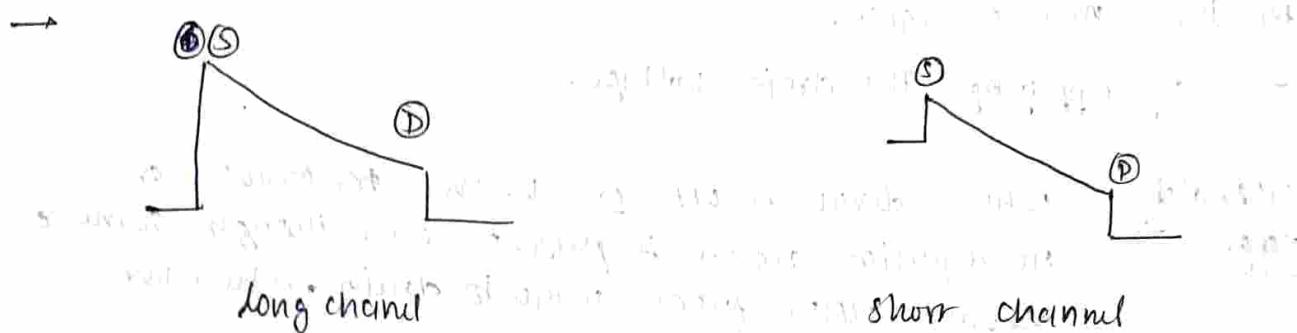
e⁻ require high energy to move from source to drain



Electrons require less energy to move from drain to source

(*) $E_S < E_D$ because at less V_D we can obtain the subthreshold current.

Q2 for most channel, how source is moving upwards as we apply V_d ?

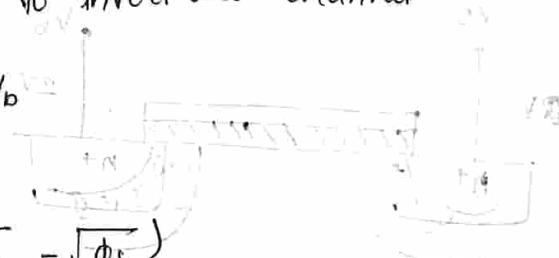


Because of the short channel, the depletion region below drain & body and drain & body starts to overlap & hence we see such kind of effect in energy difference.

2. Body effect : (doubt)

- As you change Body voltage, V_t will change, causing the change in leakage.
- Body is the fourth terminal of transistor.
- V_{SB} affects the v_t change required to invert the channel
increasing V_S or decreasing V_B increases V_t

$$\text{as, } V_t = V_{to} + \gamma (\sqrt{\phi_s + V_{sb}} - \sqrt{\phi_s})$$



$\phi_s \rightarrow$ surface potential at threshold point.

$$\text{from equations above } \phi_s = 2V_T \ln\left(\frac{N_A}{n_i}\right)$$

- Depends on doping level N_A
- Intrinsic carrier conc' n_i



- Body effect coefficient

$$\gamma = \frac{t_{ox}}{E_{ox}} \sqrt{2qE_{si}N_A} = \frac{\sqrt{q\epsilon_{si}N_A}}{8x}$$

Plots with V_B and V_S vs V_D show a shift in threshold voltage due to body effect.

- For small source-to-body voltage, neglect as linear

$$V_t = V_{to} + K_V V_{SB}$$

where, $K_V = \frac{\gamma}{2\sqrt{\phi_s}} = \frac{2 \cdot \epsilon_s N_A T}{C_V \ln(N_A/n)}$

$$2 \cdot C_{ox}$$

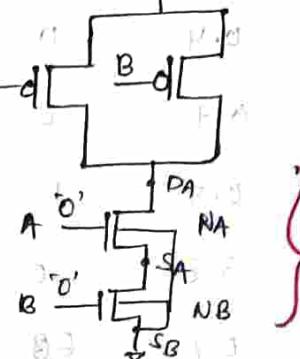
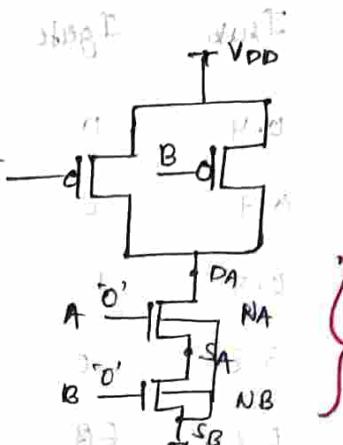
Stack effect:

- At $A=1, B=0 \Rightarrow S_A = V_{DD}$
- At $A=0, B=0 \Rightarrow S_A = \text{Any floating value}$
- At $A=0, B=0$ means both N_A & N_B are off and circuit of N_A (D_A) is cut. V_{DD} , so the stack of N_A & N_B will leak.

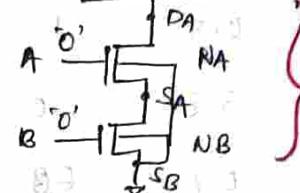
However, S_A settles to a point, that current.

- If voltage at S_A is small then N_A will see much smaller DIBL effect & will leak less
- As voltage at S_A rises, V_{GS} for N_A becomes negative, reducing its leakage
- Hence we would expect that the series resistor leak less.
- At $A=0, B=0$, voltage at S_A is floating hence due to body effect $V_{SB} \neq 0$ hence $V_t \neq 0 \Rightarrow$ leakage.

- Q: If we connect source to body directly, then we will never observe body effect. Then why don't we connect body to source directly.
- See through the physical structure of the device & layout.



leakage through
a-stack reduces
10 times



each transistor has the same

NAND3 Leakage Example

- 100nm process

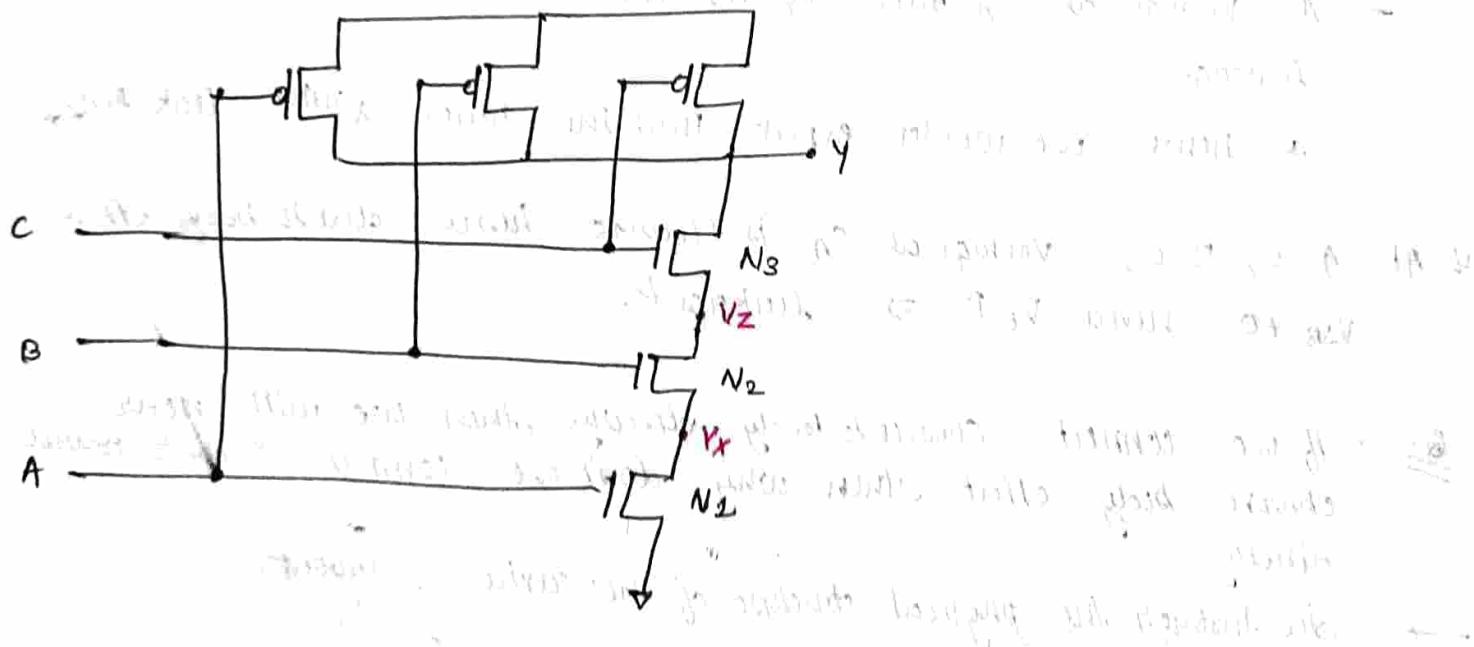
$$I_{gn} = 6.8 \text{nA}$$

$$I_{gp} = 0 \text{A}$$

$$I_{offn} = 5.63 \text{nA}$$

$$I_{offp} = 9.8 \text{nA}$$

Input state (ABC)	I_{sub}	I_{gate}	I_{total}	V_x	V_z
000	0.4	0	0.4	Stack effect	Stack effect
001	0.7	0	0.7	Stack effect	$V_{DD} - V_t$
010	0.7	4.3	2.0	Intermediate (floating value)	Intermediate
011	3.8	0	3.8	$V_{DD} + V_t$	$V_{DD} + V_t$
100	0.7	6.0	7.0	0	Stack effect
101	3.8	6.3	10.1	0	$V_{DD} - V_t$
110	5.6	12.6	18.2	0	0
111	28	18.9	46.9	0	0



Subthreshold Leakage control

- Leakage & Delay trade off:
 - Aim for low leakage in sleep & low delay in active mode
- To Reduce leakage:
 - (i) Increase V_t : multiple V_t
 - Use low V_t only in critical circuits
 - (ii) Increase V_s : stack effect
 - I/P vector control in sleep mode
 - (iii) Decrease V_b
 - Reverse Body Bias in sleep mode
 - Forward Body bias in active mode

Q. Difference b/w sleep mode & clock gating.
→ Sleep mode involves much more than just enable signal on the clock gating path.

② Gate Leakage:

- Carrier tunnel through very thin gate oxides
- Exponentially sensitive to t_{ox} & V_{DD}

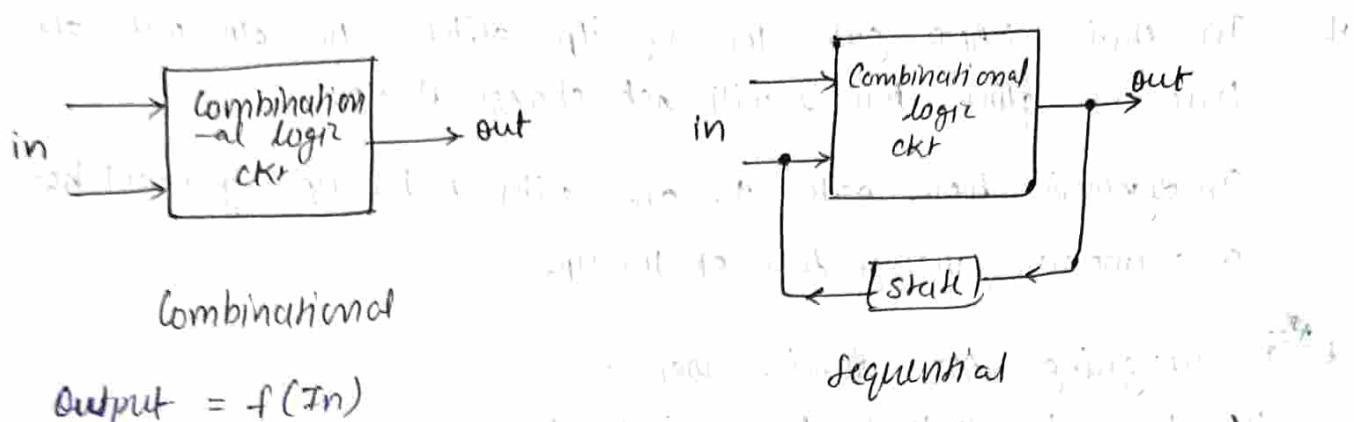
$$I_{gate} = WA \left(\frac{V_{DD}}{t_{ox}} \right)^2 e^{-\frac{B t_{ox}}{V_{DD}}}$$

① $V_{DD} \downarrow$, $I_{gate} \uparrow$ & ② $t_{ox} \downarrow$, $I_{gate} \uparrow$

- A & B are technology constants
- Greater for e's
- * So NMOS gate漏 more.
- Negligible for older processes ($t_{ox} > 20 \text{ \AA}$)
- Critically important at 65nm & below 10.54 t_{ox} .

Combinational Circuit Design

Combinational vs sequential logic



Static CMOS vs Dynamic

Static CMOS

- At every point in time (except during switching transients) one op is connected to either V_{DD} or V_{SS} via a low resistance path.
- Outputs of the gates are (driven) at all times - to the value of Boolean func. implemented by the ckt

Eg. In static memory, you will have transistors in memory cell & in dynamic memory, you will have capacitors in which you store charge.

Dynamic logic

- Rely on temporary storage of signal values on the capacitance of high impedance node
- Advantage → Area gets reduced
- Disadvantage → we get some race around cond? & cascading conc?
- Resembles clock

* So, in dynamic logic for some combinational state where the op is not driven by anything else but is dependent on some charge stored in capacitor.

- Now, if due to crosstalk or leakage that charge gets dissipated then there is no way to recover that charge during the cycle & data is lost.

- In static NAND gate, for any ilp glitch, the o/p will also have some glitch, but it will not change its state.
In dynamic logic gate, for any glitch in the ilp you will have a permanent voltage loss at the o/p.

* For designing AOI domino logic →

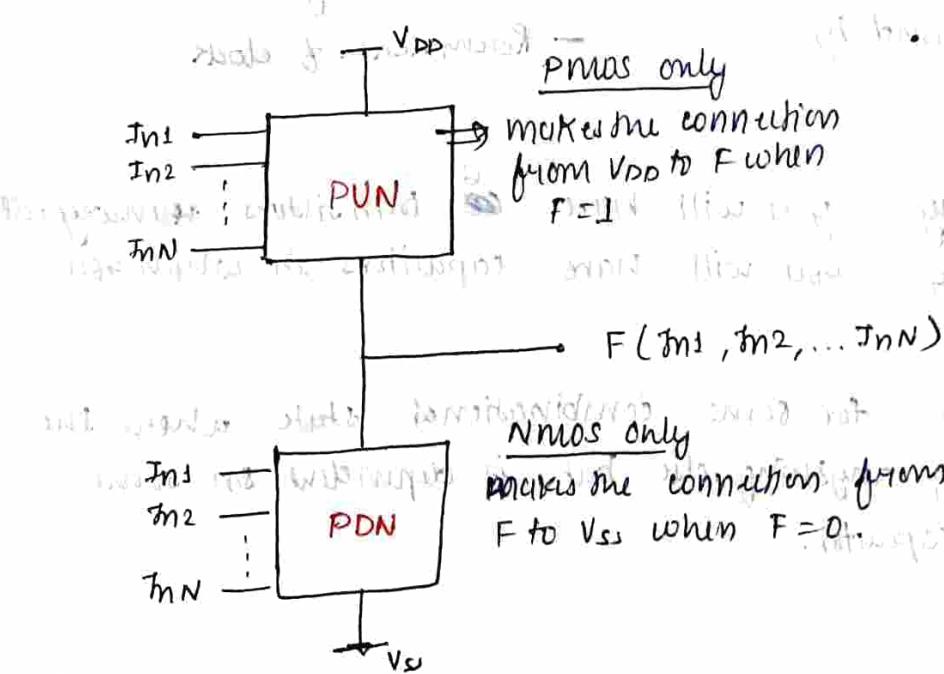
- You have to look at the various failure modes of dynamic logic
- You have to take care of the capacitances
- You have to check the noise immunity of your gate

Q: DRAM cell is not a dynamic logic then why dynamic is used?

→ It is called dynamic RAM cell & not dynamic logic

* dynamic because it is storing something on a capacitor & you need to refresh it you need to ensure that some leakage or whatever has happened it is kind of recovered on regular basis.

1. Static Complementary MOS (Static CMOS)



* The function of PUN o/p is to provide a common connection b/w the o/p & V_{DD} anytime the o/p of logic gate is meant to be 1 (based on ilps).

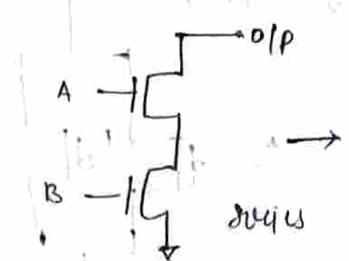
* PUN & PDN are dual

and regular static inverters could not replace each other

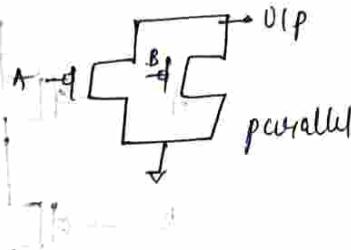
because the connection from F to V_{SS} when $F=0$.

Eg. For NAND gate

if PDN is,



then PUN is,



When both inputs A & B are high, then output is low

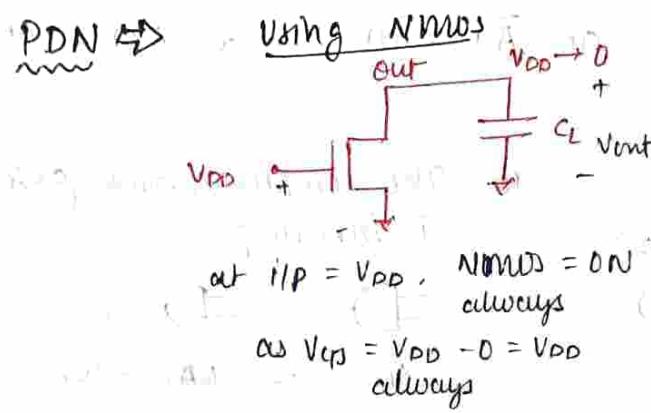


Threshold Drop

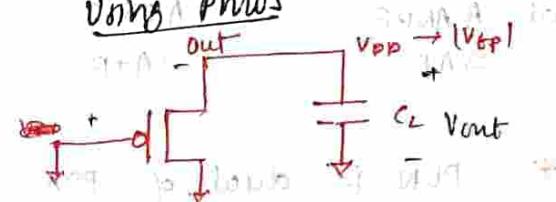
Q. Why PDN is constructed using NMOS & PUN using PMOS?

PDN \Rightarrow

Using NMOS



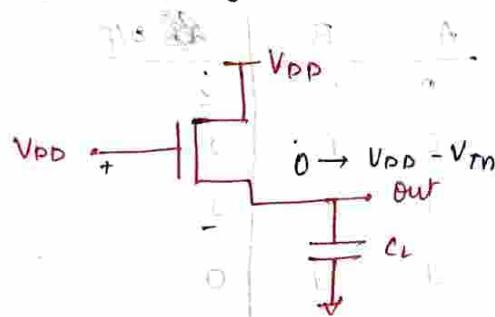
Using PMOS



at $I_{IP} = 0$, $V_{GS} = -V_{DD}$
as $V_{GS} = -V_{DD}$ as $V_{GS} = -V_{DD}$
then PMOS \rightarrow OFF & C_L will
not completely discharge.
as $|V_{GS}| \leq |V_{TP}|$

PUN \Rightarrow

Using NMOS



$$V_{GS} = V_{DD} - V_{DD}$$

\therefore NMOS is ON till $V_{GS} > V_{TN}$

$$\text{at } V_{DD} = V_{DD} - V_{TN}$$

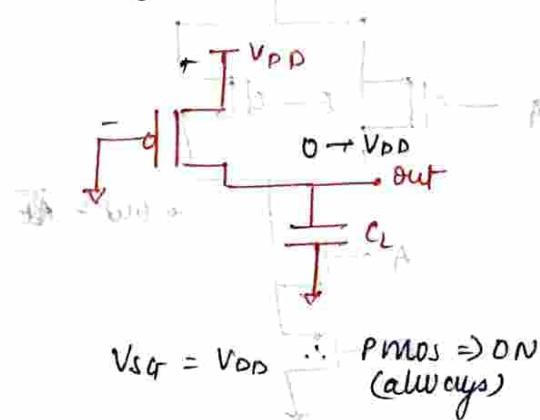
$$V_{GS} = V_{DD} - (V_{DD} - V_{TN})$$

$$V_{GS} = V_{TN}$$

\therefore NMOS becomes OFF

& No charging will occur.

Using PMOS

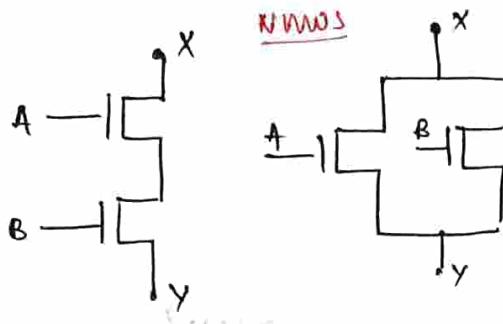


$$V_{GS} = V_{DD} - V_{DD} \quad \text{PMOS} \Rightarrow \text{ON (always)}$$

$\therefore C_L$ charge to V_{DD}

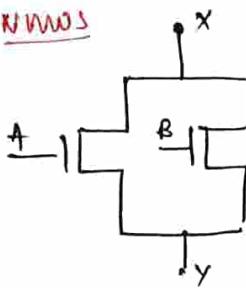
- * NMOS transistor pass a "strong" $\rightarrow 0$ but "weak" $\rightarrow 1$
- * PMOS transistor pass a "weak" $\rightarrow 0$ but "strong" $\rightarrow 1$

Transistor in Series & Parallel



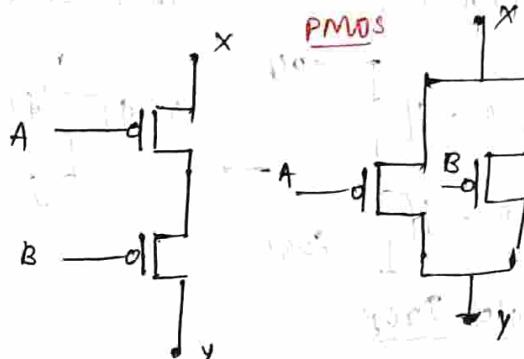
$$\textcircled{1} \quad Y = X$$

cut, A AND B
⇒ AB



$$\textcircled{2} \quad Y = X$$

out A OR B
⇒ A+B



$$\textcircled{3} \quad Y = X$$

cut $\bar{A} \cdot \bar{B} = \bar{A} + \bar{B}$ (INV) or $\bar{A} + \bar{B} = \bar{A} \cdot \bar{B}$

$$\textcircled{4} \quad Y = X$$

out $\bar{A} + \bar{B} = \bar{A} + \bar{B}$ (INV)

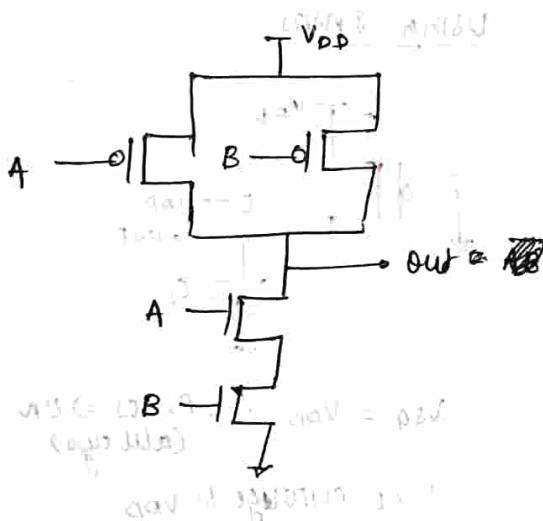
* PUN is dual of PDN

1+4 it, if $AB \xrightarrow{\text{dual}} \bar{A} \cdot \bar{B}$ (NAND)

2+3 if $A+B \xrightarrow{\text{dual}} \bar{A} \cdot \bar{B} + B$ (NOR)

1+2 \Rightarrow AND + INVERT must
• Function of inverter

Example : NAND Gate



NAND		O/P
A	B	O/P
0	0	1
0	1	1
1	0	1
1	1	0

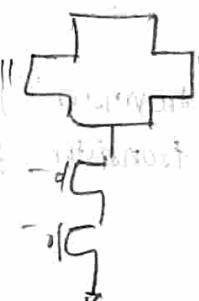
(OR)

OR		O/P
A	B	O/P
0	0	0
0	1	1
1	0	1
1	1	1

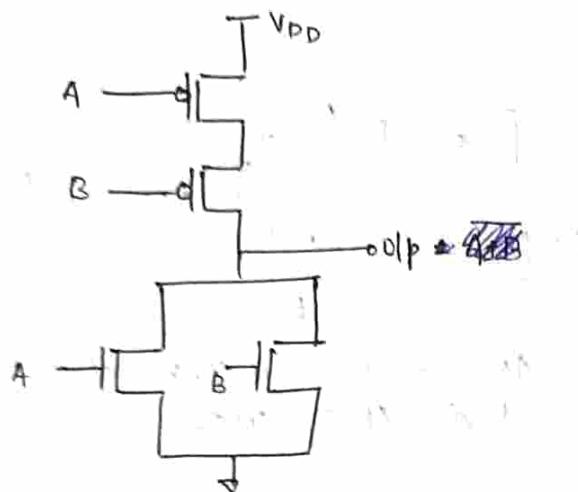
PDN : $G = AB \rightarrow$ condⁿ to GND

PUN : $F = \bar{A} + \bar{B} \rightarrow$ condⁿ to VDD

$$\overline{G(in_1, in_2, \dots, in_N)} = F(\overline{in_1}, \overline{in_2}, \dots, \overline{in_N})$$



Example : NOR Gate (Pg. No. 323) \Rightarrow Output of NOR gate is $A + B$



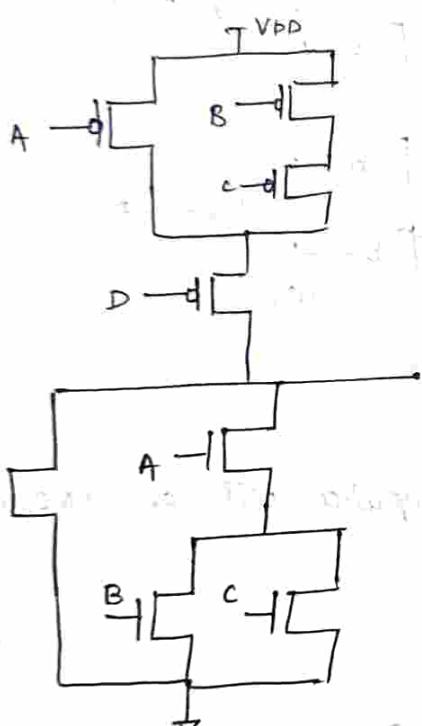
A	B	O/P
0	0	1
0	1	0
1	0	0
1	1	0

$$\text{Input of P/NPN} \Rightarrow A + B$$

$$\therefore O/P = \overline{A+B}$$

Example : Complex CMOS Gate

True

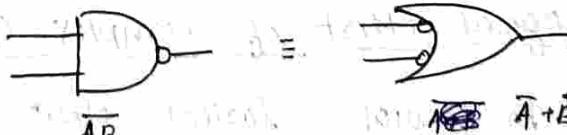


- * while constructing the complex CMOS gate from given expression
 - Draw PPN for expression provide complement
 - make the dual of PPN in PUN using PMOS.

Bubble Pushing : (Wash Pg. 323)

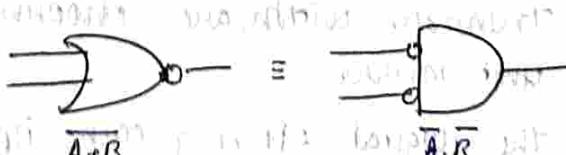
- A NAND is equivalent to an OR of inverted inputs.

$$A \cdot B \equiv \overline{\overline{A} + \overline{B}}$$



- A NOR is equivalent to an AND of inverted inputs.

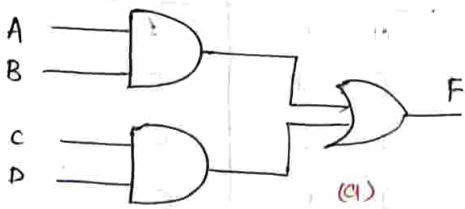
$$\overline{A + B} \equiv \overline{\overline{A} \cdot \overline{B}}$$



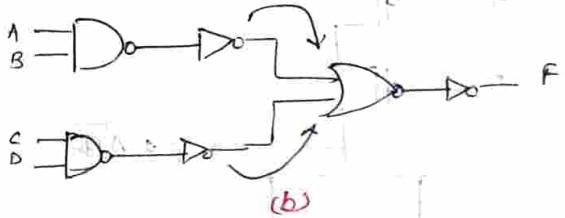
Ques - Design a ckt to compute $F = AB + CD$ using NORs.

NORs

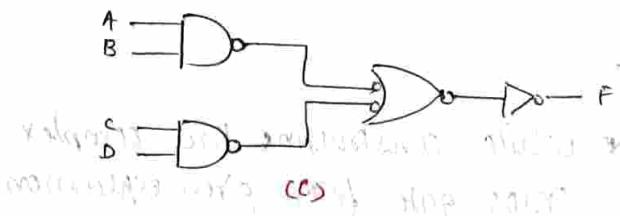
→



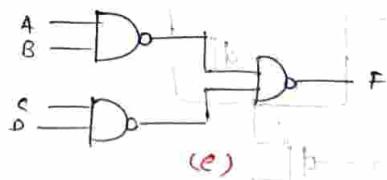
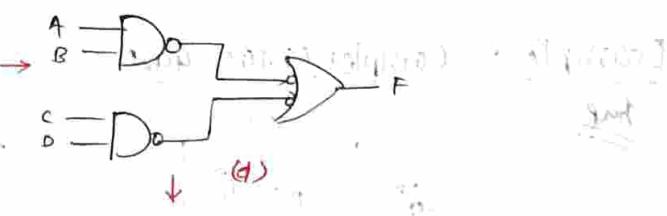
Ckt consists of two AND
+ one OR gate



ANDs + ORs are converted to
bun's emos stages

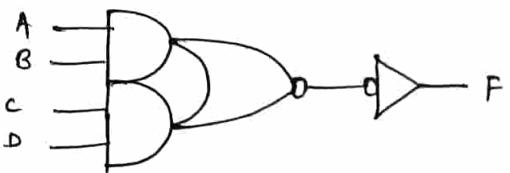


Bubble pushing is
used to simplify the
logic to three NORs



Complex Gates (WUIC Pg. 329)

- the funcⁿ $F = AB + CD$ can be computed with an AND-OR-INVERT -22 (AOI22) gate and an inverter.



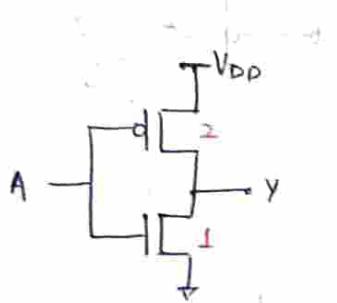
Logical Effort of Complex Gates

- In general, logical effort of compound gates can be different for different inputs
- transistor widths are chosen to give the same driving unit Miller
- the logical effort of each ilp is the ratio of the ilp cap of next ilp to the ilp cap of inverter.

- The parasitic delay is usually estimated from the total diffusion capacitance on the olp node by summing the size of the transistors attached to the olp.

Unit inverter

$$y = \bar{A}$$

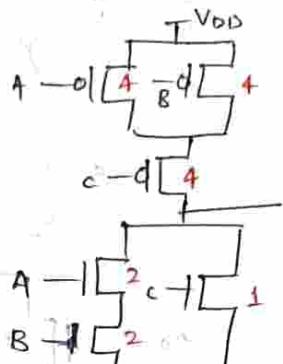


$$g_A = 3/3$$

$$P = 3/3$$

AOI 21

$$y = \overline{AB + C}$$



$$g_A = 6/3$$

$$g_B = 6/3$$

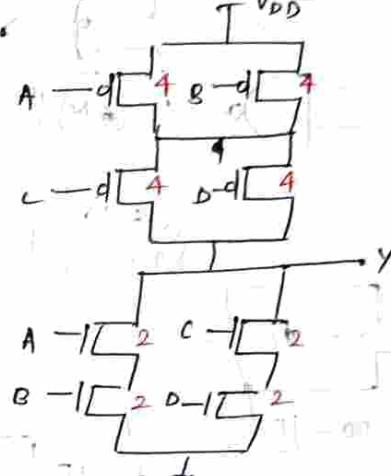
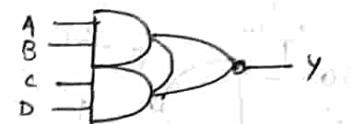
$$g_C = 5/3$$

$$P = 7/3$$

w/o
stack
Diagram

AOI 22

$$y = \overline{A \cdot B + C \cdot D}$$



$$g_A = 6/3, g_B = 9/3$$

$$g_C = 6/3, g_D = 6/3$$

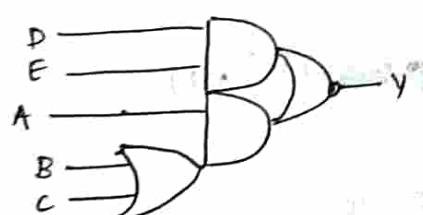
~~$$P = 12/3$$~~

while Designing
 $R_{PUN} = R_{PDN}$

& Area should be less.

Complex AOI

$$y = \overline{A(B+C)+DE}$$



$$g_A = 5/3$$

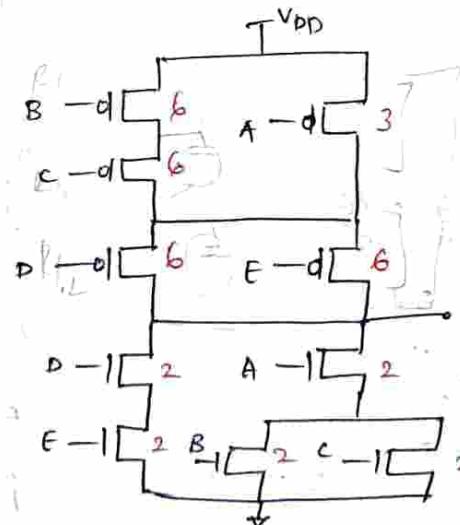
$$g_B = 8/3$$

$$g_C = 8/3$$

$$g_D = 8/3$$

$$g_E = 8/3$$

$$P = 16/3$$



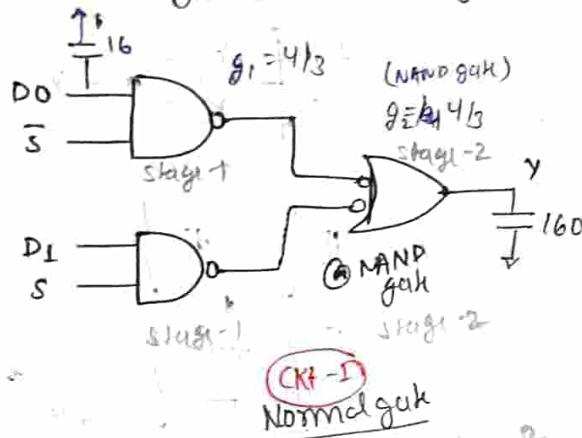
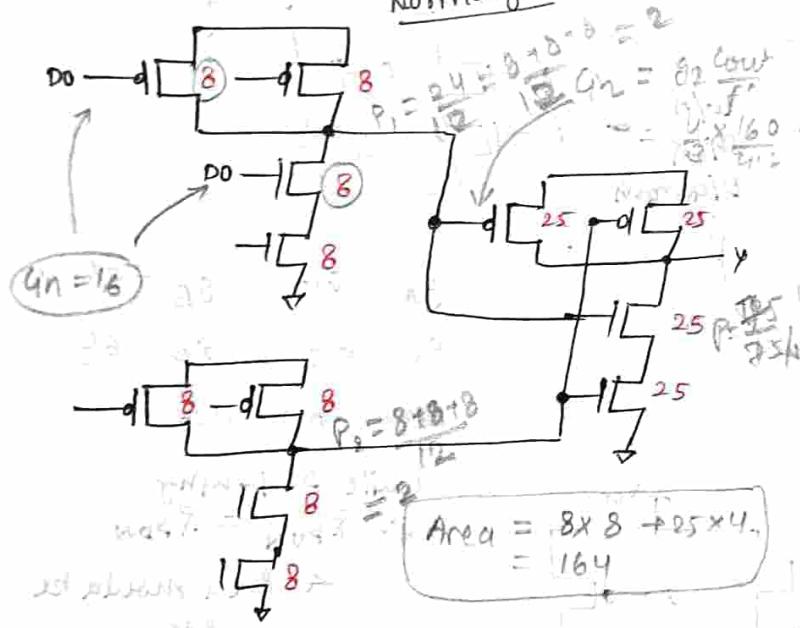
* For AOI 21 grid, the logical effort is lower for OR dynunical (C) than for the AND dynunical (A,B).

Q4

2. Delay & Gate size of complex gates

- The multiplexer has a maximum iip capacitance of 16 units on each iip. It must drive a load of 160 units. Estimate the delay of two designs

SOP ABCD

(CKT - I)
Normal gate

$$G = g_1 g_2 = \frac{4}{3} \times \frac{4}{3} = \frac{16}{9}$$

$$B = 1 \quad (\text{Fanout} = 1)$$

$$n = \frac{160}{16} = 10$$

$$F = GBH = \frac{160}{9}$$

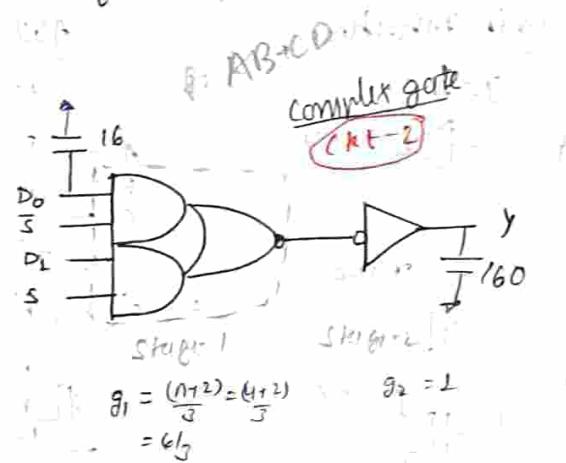
$$\hat{F} = (F)^{1/2} = 4.2$$

$$P = 2+2 = 4$$

$$D = NF + P = 2(4.2) + 4$$

$$D = 12.4$$

$$\Rightarrow D_{\text{obs}} = 12.4T$$

for iip D_0

$$G = \cancel{\frac{16}{9}} \left(\frac{6}{3}\right) \times (1)$$

$$B = 1$$

$$n = \frac{160}{16} = 10$$

$$F = GBH = \frac{60}{3} = 20$$

$$\hat{F} = (F)^{1/2} = 4.5$$

$$P = 4+1 = 5$$

$$D = NF + P$$

$$D = 2 \times (4.5) + 5$$

$$D = 14$$

$$D_{\text{obs}} = 14T$$

- * The delay for complex gates is more than the delay of normal gate. But area of complex gate is less as compared to normal gate.
- * OR lesser the i/p capacitance at the clock during the setup of the CKT. [Compare delay & Area of both the CKTs] since the area is lesser, the parasitic will also be lesser & hence in the practical you will see the delay of complex gate is less.
- Q: If my CKT has diff. g values for different i/p then for calculation of total delay which i/p should I consider?
 - For calculation of prop. delay, take worst g. (i.e. high g)
 - for calculation of cont. delay, take best g. (i.e. low g)

- Input Ordering Delay Effect:
- Complex gates like AOI21, are asymmetrical, that one i/p sees less capacitance than another.
 - NAND & NOR gates symmetrical but actually have slightly different logical effort & parasitic delay for the different inputs.
 - Consider a CKT of 2-i/p NAND gate annotated with diffusion parasitic C.

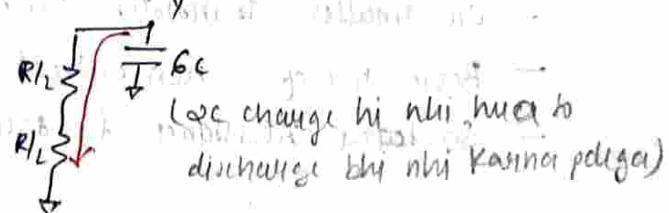
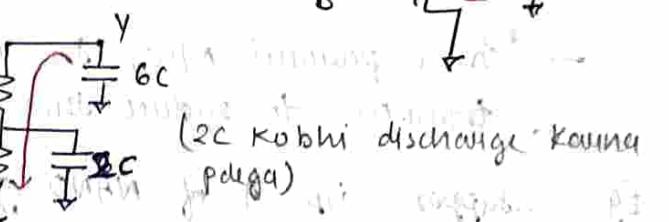
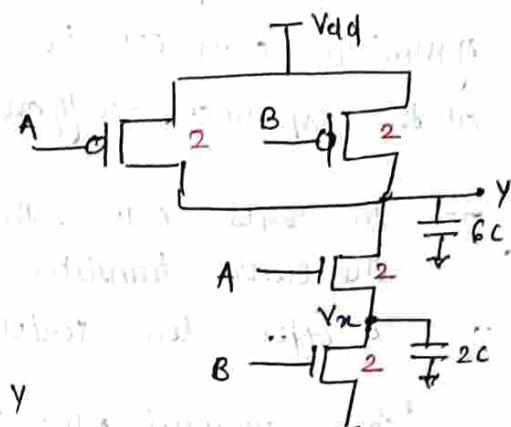
— Consider the falling o/p transition occurring when one i/p is held at '1' & other varies from '0' to '1'.

(i) Let $A = 1$, $B = 0$, $V_{in} = V_{DD} - V_t$
then $A = 1$, $B = 1$

$$\begin{aligned}\text{Elmore delay} &= 6C(R) + 2C\left(\frac{R}{2}\right) \\ (\text{Worst case delay}) &= 7RC \\ &= \frac{4}{3}(3RC) = 4.33T\end{aligned}$$

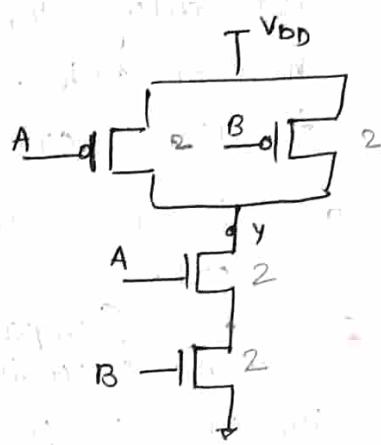
(ii) Let $A = 0$, $B = 1$, $V_{in} = 0$
then $A = 1$, $B = 1$.

$$\begin{aligned}\text{Elmore delay} &= 6C(R) = 6RC \\ (\text{Best case delay}) &= \frac{1}{2}(3RC) = 1.5T\end{aligned}$$



- * Outer ilp \rightarrow ilp closer to supply rail (eg. B)
- * Inner ilp \rightarrow ilp closer to the o/p (eg. A)

* Parasitic delay is smallest when inner ilp switches last because the intermediate nodes have already been discharged.



- Signal arriving at the inner gate ilp must occur late, if any one ilp is coming late

* The inner ilp has a slightly higher logical effort because the intermediate node x tends to rise & cause -ve f/b when inner ilp turns ON.

Six

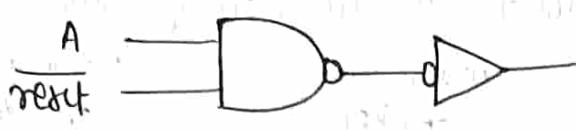
- * Which ever ilp provides shortest delay to the o/p should be designed to be arriving last in your path.
 - here delay is smallest for ilp A ~~because it will be already~~ means A should arrive last.

Asymmetric Gates: (White Pg 1332)

- when one ilp is less critical than another, even nominally symmetric gates can be made asymmetric to favour the less critical ilp at the expense of early one.
- In series N/W, the shunt connecting the early ilp to the output transistor & making the transistor wider so that it offers less series resistance when the critical ilp arrives
- In a parallel N/W, the early ilp is connected to a narrower transistor to reduce the parasitic capacitance.

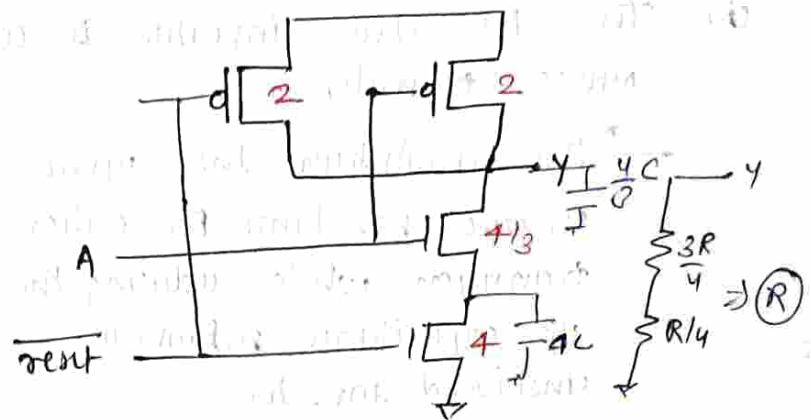
Eg. Suppose ilp A of NAND gate is most critical

- use smaller transistor on A (less capacitance)
- Boost size of non-critical path
- so total resistance is same



$$g_A = \frac{4+2}{3} = \frac{16}{9} = 1.78$$

$$g_B = \frac{4+2}{3} = 2$$



$$g_{\text{total}} = g_A + g_B = \frac{28}{9}$$

(*) Reduces the delay but

area increases

it but total logical effort \uparrow .

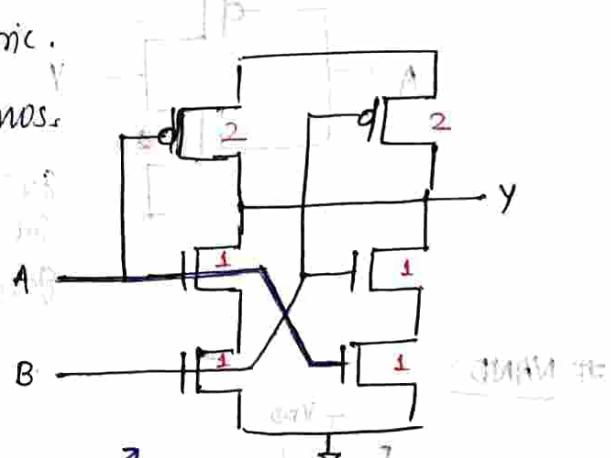
it logical effort for critical path ($A \rightarrow y$) \downarrow .

- the improvement in logical effort of i/p A comes at the cost of much higher effort on the least i/p.

(*) We will only use asymmetric gates when we have complete understanding of critical path.

Symmetric gates

- Input can be made perfectly symmetric.
- In symmetric gates you have two NMOS, each of same size but your both the i/p's are coming randomly at the same time. So I get equal delay from A to Y & B to Y.
 - If A is closer to o/p
 - If B is closer to o/p.



(*) This is used when your i/p data A & B is completely random.

Skewed Gates → This will be used in pulse width in sequential CKT.

- Hi-skew gates to favour the rising o/p transition.
- Lo-skew gates to favour the falling o/p transition.
 - This favouring can be done by decreasing the size of the non-critical transistor.

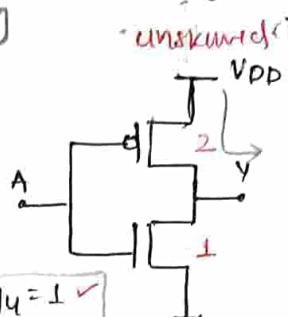
$g_u \rightarrow$ logical effort for rising (up) transition of o/p.

$g_d \rightarrow$ logical effort for falling (down) transition of o/p.

- (*) And then we use the ratio of the i/p capacitance of the skewed gate to the i/p capacitance of an unskewed inverter with equal drive for the transition.

(i) The HI-skew inverter is constructed by downsizing the NMOS transistor ($I_{DLP} \Rightarrow 0 \rightarrow 1$)

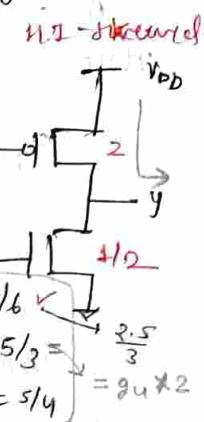
* This means keeps the same effective resistance for critical transition while reducing the input capacitance resulting to unskewed inverter.



$$g_u = 1 \checkmark$$

$$g_d = 1$$

$$g_{avg} = 1$$

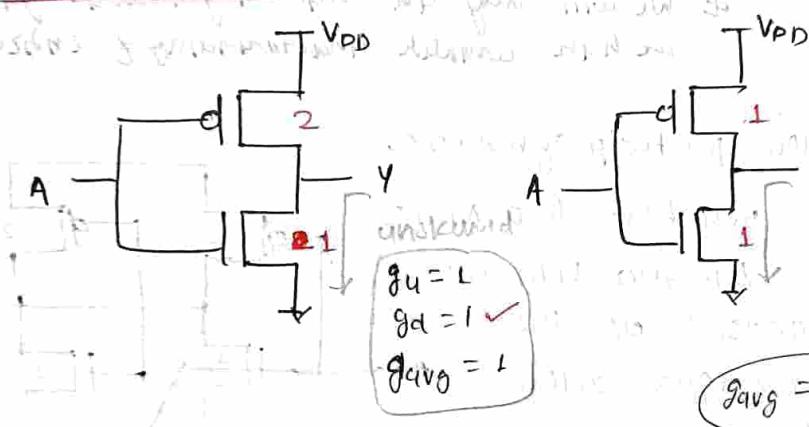


$$g_u = 5/6 \checkmark$$

$$g_d = 5/3 = \frac{5}{3}$$

$$g_{avg} = 5/4 = g_u * 2$$

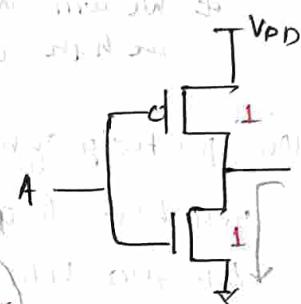
(ii) The LO-skew inverter is constructed by down-sizing the PMOS transistor & keeping the NMOS of same size so that effective resistance for critical path remains same & reduces the input capacitance hence logical effort decreases.



$$g_u = L$$

$$g_d = 1 \checkmark$$

$$g_{avg} = L$$

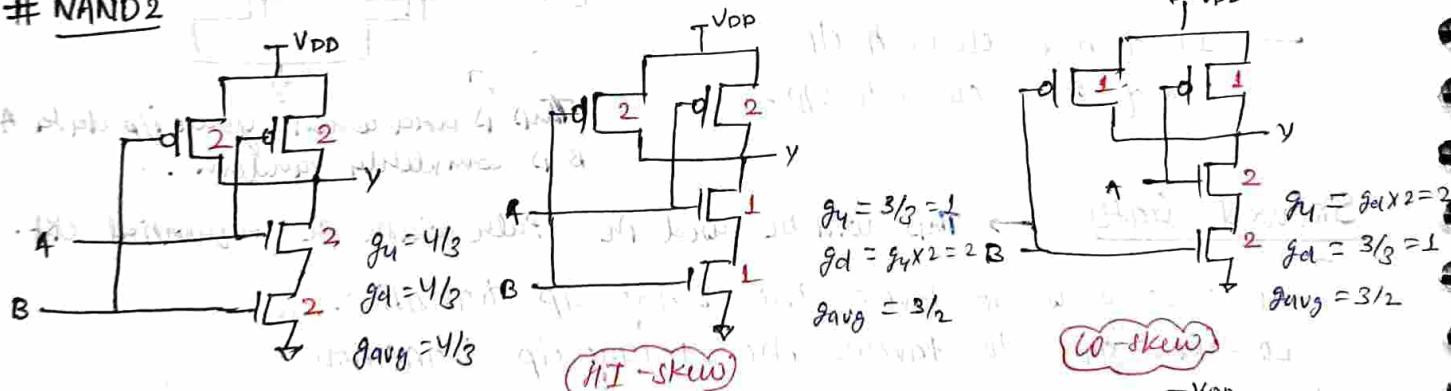


$$g_u = 4/3 = g_d * 2$$

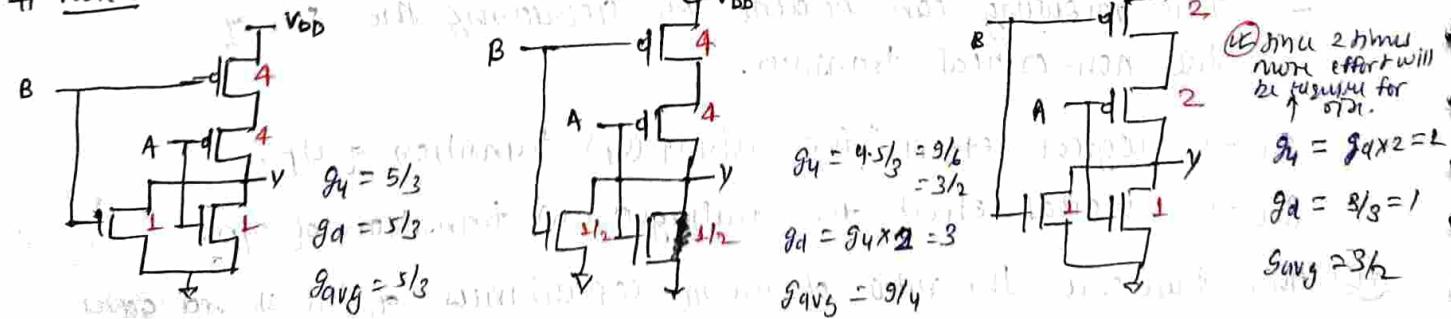
$$g_d = 2/3 \checkmark$$

$$g_{avg} = \frac{g_u + g_d}{2} = \frac{6}{3} = 2$$

NAND2



NOR2



* Skew gains are and only when one transition is imp.

* Skewed gains work particularly well with dynamic CMOS.

(*) since 2 times more effort will be required for it.

$$g_u = g_d * 2 = L$$

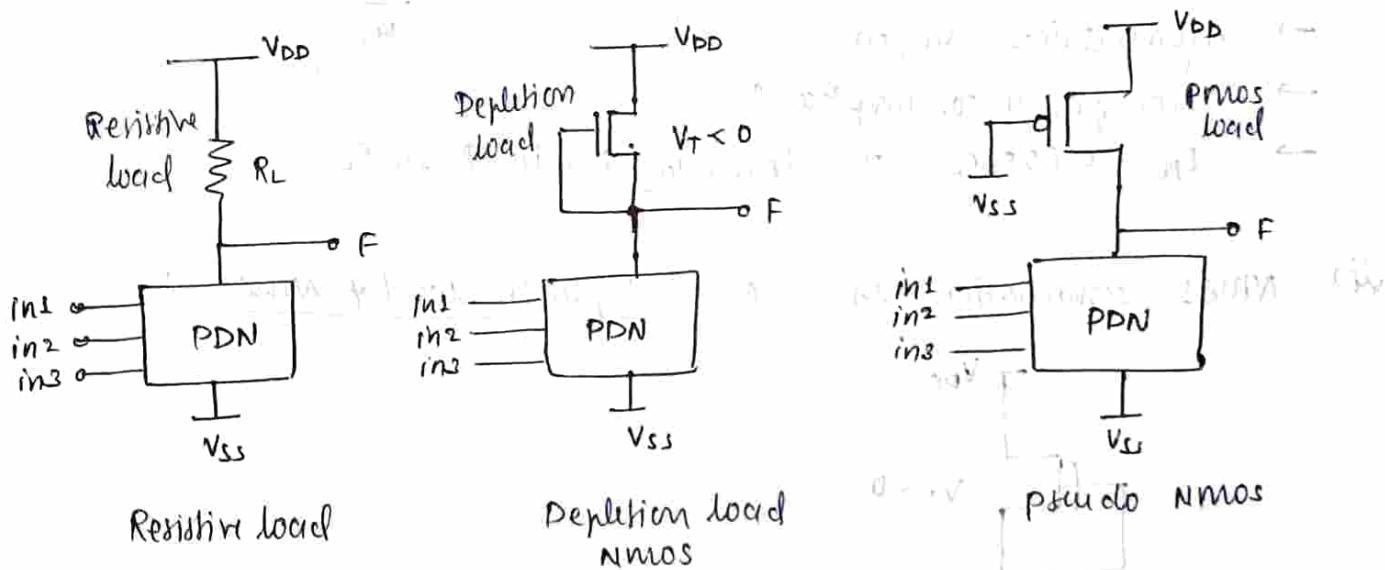
$$g_d = 2/3 = 1$$

$$g_{avg} = 3/2$$

Ratioed logic: (Rabacy Pg. 249)

- In static CMOS logic design, for N -no. of inputs we have to use atleast $2N$ no. of devices and on that N devices are PMOS which is larger in size.

Goal: To reduce the number of devices over complementary CMOS
 → for that we can use Resistive load, depletion load NMOS and Grounded PMOS (pseudo NMOS) in place of PUN.



- But the noble wins this kind of load if
 - noise margin \downarrow (bad)
 - propagation delay \uparrow (good) (flow!)
 - power dissipation \uparrow (leakage \uparrow) (bad)
 - output will never go to exact zero
 - a continuous static current from V_{DD} to o/p.

- Q. What if increase the resistance of our load such that we get clear '0'?
- If we do so then
 - current through load will be very less & charging time for o/p increases.
 - (?) High resistive load increases the area also.

(i) NMOS in PDN + susceptive load:

- $V_{ON} = V_{DD}$

↳ min. High o/p. voltage

- $V_{OL} = \frac{R_{PN}}{R_{PN} + R_L}$ Jauge pull up

↳ max. low o/p. voltage

$V_{OL} \Rightarrow$ must go to V_{DD} for better switching.

(\therefore noise margin \uparrow)

→ Asymmetrical response

→ static power consumption \uparrow .

- $t_{PL} = 0.69 R_{L} C_L \Rightarrow t_{prop. delay} \uparrow$ with \uparrow in R_L .

↳ signal ringing

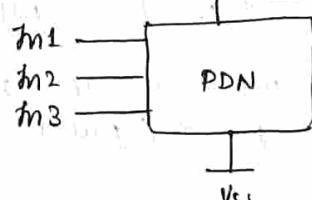
↳ what happens when $R_L \gg \infty$

↳ what happens when $R_L \downarrow$?

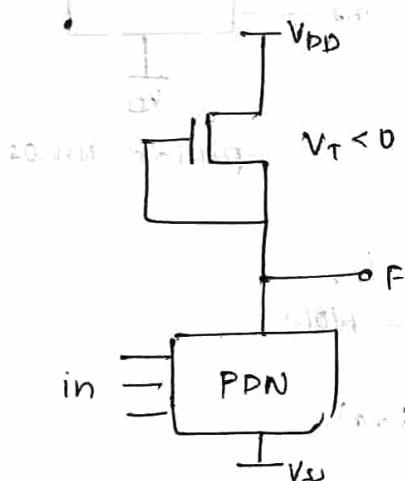
↳ output voltage

↳ PDN

↳ V_{SS}



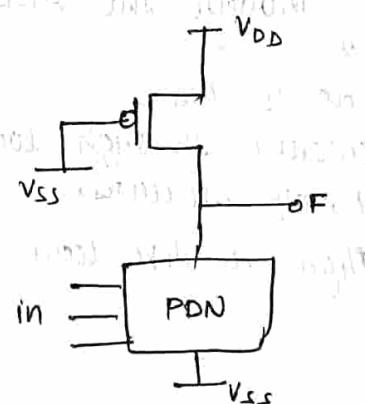
(ii) NMOS combination with PDN + Depletion load of NMOS



↳ load capacitance
↳ gain

(iii) NMOS comb. in PDN + pseudo-NMOS

- Computing the DC transfer characteristics of the pseudo-NMOS proceeds along paths similar to those used for its complementary CMOS counterpart.
- The value of V_{OL} is obtained by equating the current through driver + load driver for $V_{IN} = V_{DD}$.



- At this operating point, it is reasonable to assume that the NMOS device operates in linear mode (since V_{OL} should ideally be close to 0V), while PMOS load is saturated.

$$K_n \left[(V_{DD} - V_{TN}) V_{OL} - \frac{V_{OL}^2}{2} \right] = \frac{K_p}{2} (V_{DD} - V_{TP})^2$$

assuming, $V_T = V_{TN} = V_{TP}$

$$\therefore V_{OL} = (V_{DD} - V_T) \left[1 - \sqrt{1 - \frac{K_p}{K_n}} \right] \quad \text{"from slides"}$$

$$\textcircled{a} \quad V_{OL} \approx \frac{K_p (-V_{DD} - V_{TP}) \cdot V_{DSAT}}{K_n (V_{DD} - V_{TN})} \quad \text{"from Pabency"}$$

$$V_{OL} \approx \frac{W_P W_P}{4 n W_N} \cdot V_{DSAT}$$

- In order to make V_{OL} as small as possible, the PMOS device must be sized much smaller than the NMOS PDN devices.
[size of PMOS < size of NMOS] \Rightarrow take n as large as possible

- This will create a -ve impact on the propagation delay for changing up the output node since current provided by the PMOS device is limited. (Rising time degrades)

- * A major disadvantage of the pseudo-NMOS gate is the static power dissipated when the output is low through the direct current path that exists b/w V_{DD} & GND. (takes short cut current \uparrow)

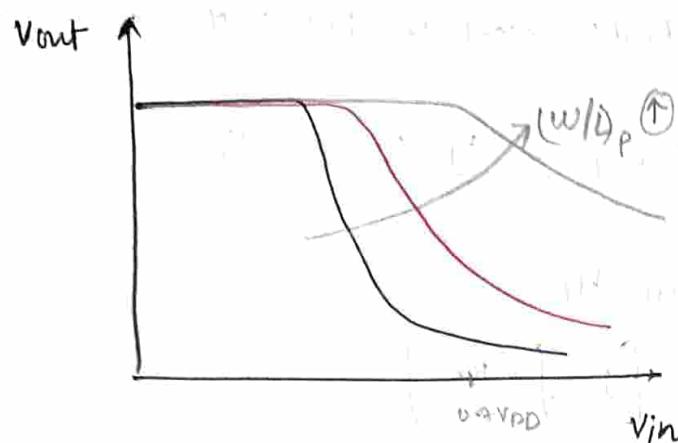
- ~~**~~ ^{True} A larger pull-up device $\uparrow (R_L \downarrow)$
- improves the performance (t_{PD})
 - reduces static power dissipation (P_{DD}) \propto size
 - lower noise margin (ie $V_{OL} \uparrow$) (R_{AD})
 \hookrightarrow Aache n' d' nji Ban rha hai.

\therefore To decrease V_{OL} , small PMOS is required.

But small pull-up device, $(R_L \uparrow)$

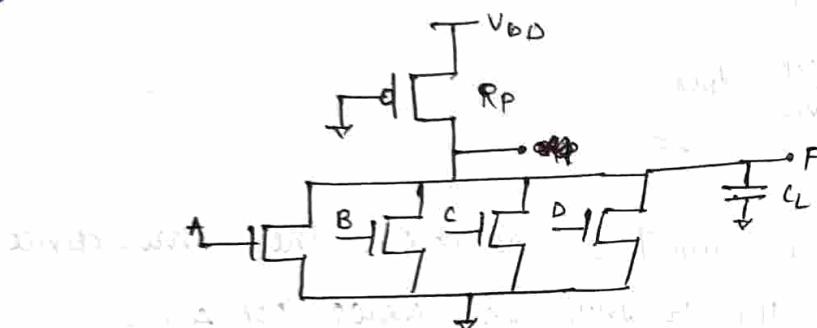
- increases rise time (t_{RD})
- decreases static power (G_{DD}) \propto size
- $V_{OL} \downarrow$. (G_{DD}) (Noise margin \uparrow)

Pseudo NMOS - Voltage Transfer Cc



PMOS size $\uparrow \Rightarrow \left(\frac{W}{L}\right)_P \uparrow$
 $\Rightarrow V_{OL} \uparrow \Rightarrow N_M \downarrow$

pseudo pmos :



on $\downarrow R_P$, $V_{OL} \uparrow$

$$(i) R_P = 10k\Omega, R_n = 2k\Omega$$

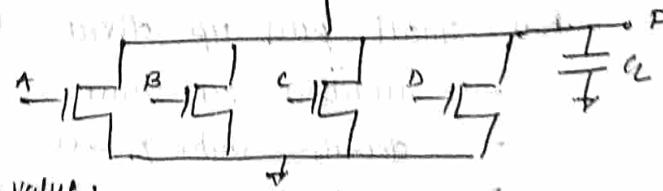
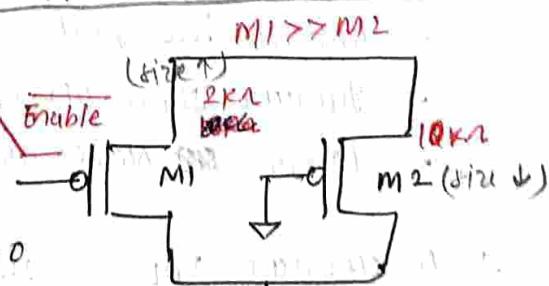
$$V_{OL} = \frac{R_P}{R_P + R_n} V_{DD} = \frac{1}{10+2} V_{DD} = 5\% V_{DD}$$

$$(ii) R_P = 2k\Omega, R_n = 2k\Omega$$

$$V_{OL} = \frac{0.5k\Omega}{2.5k\Omega} V_{DD} = 20\% V_{DD}$$

(iv) Improved load :

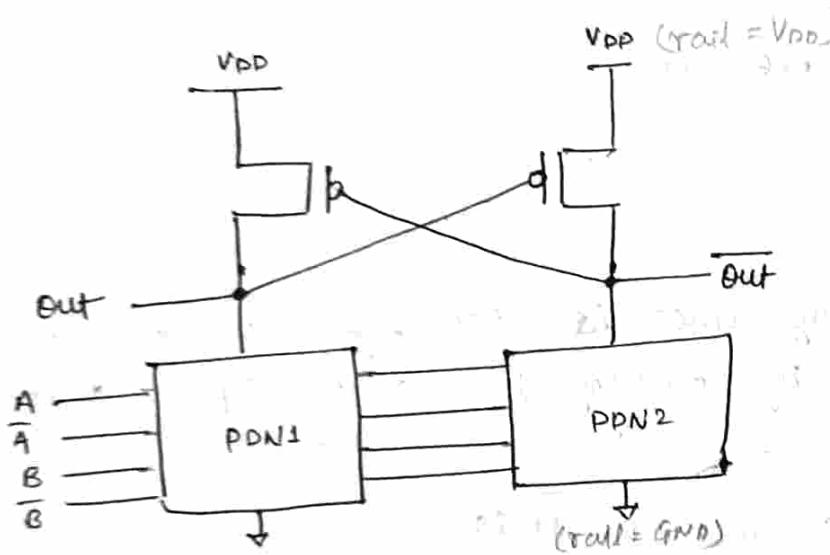
- Previously, on $\downarrow R_P$, $V_{OL} \uparrow$ & delay also \uparrow , so we have to overcome this
 - Add less ($2k\Omega$) resistance PMOS in parallel to the already present PMOS, so that overall resistance will be equal to the resistance of low sensitive PMOS. (static power dissipation \uparrow)
- During the evaluation only top (idle) PMOS will be ON & effective resistance is $9k\Omega$ so our CM becomes fast.
- And due to top enable signal the M_1 becomes ON at $enable = 0$ so o/p gets to V_{DD} from low sensitive path. (\therefore min delay \downarrow)
- And from this high sensitive M_2 the deviu becomes gets to $V_{OL} = 0$ value.



- Enable = 1, charging path = m_2 (Highly resistive) hence $V_{OL} \downarrow$.
 - Enable = 0, charging path = m_1 (low resistance) hence delay for rise \downarrow .
- ⇒ for this load all the things are achieved.

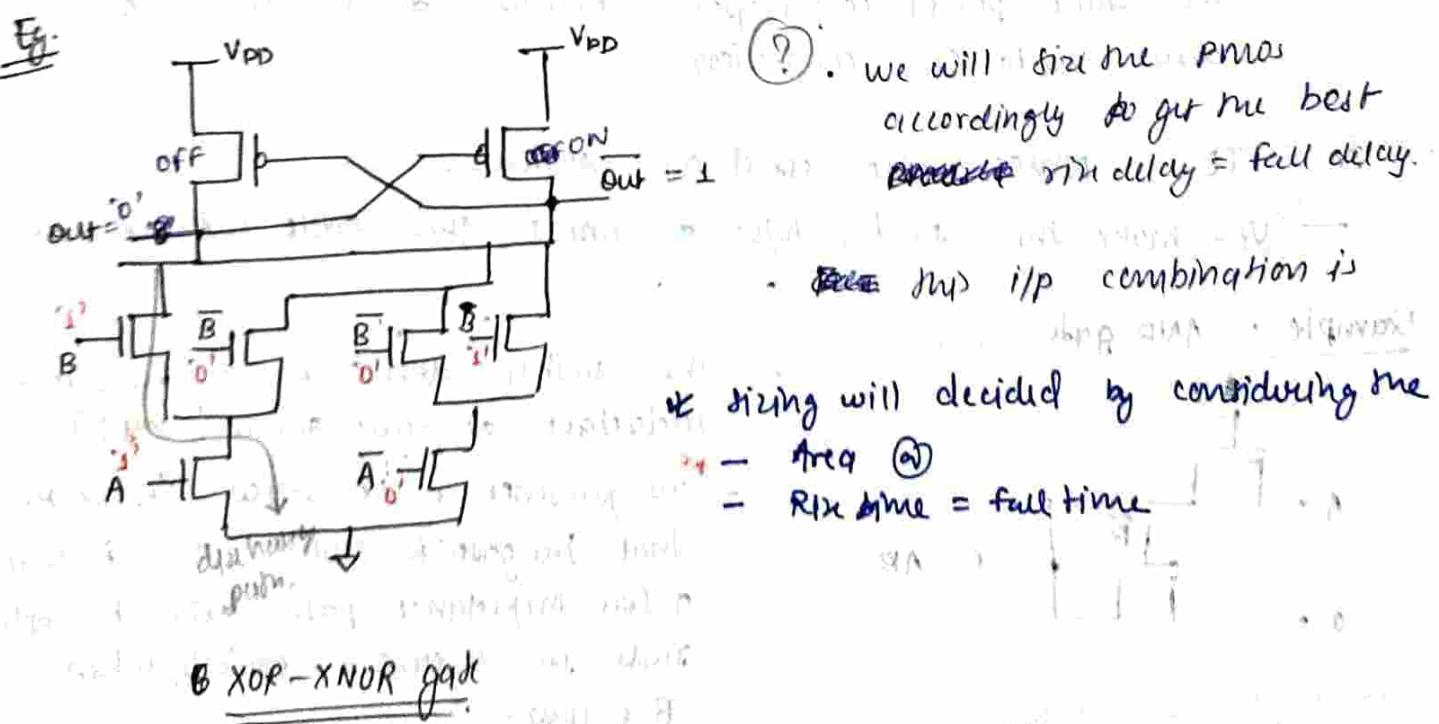
(v) DCVSL : "Differential cascode voltage switching logic"

- Using this logic style we completely eliminates static currents & provides rail-to-rail swing. such a gate combines
 - differential logic
 - $+ve f/B$ → it ensures that load devia is turned off when not needed. (static power dissipation \downarrow)

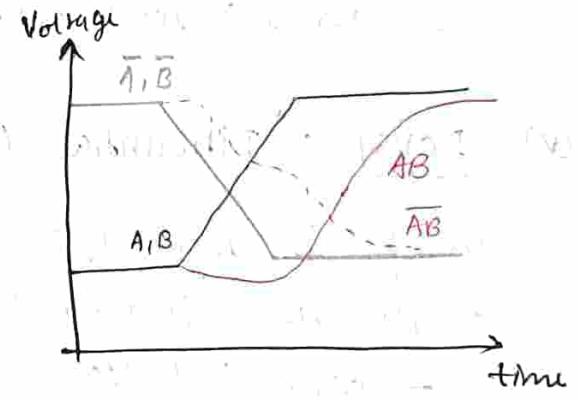
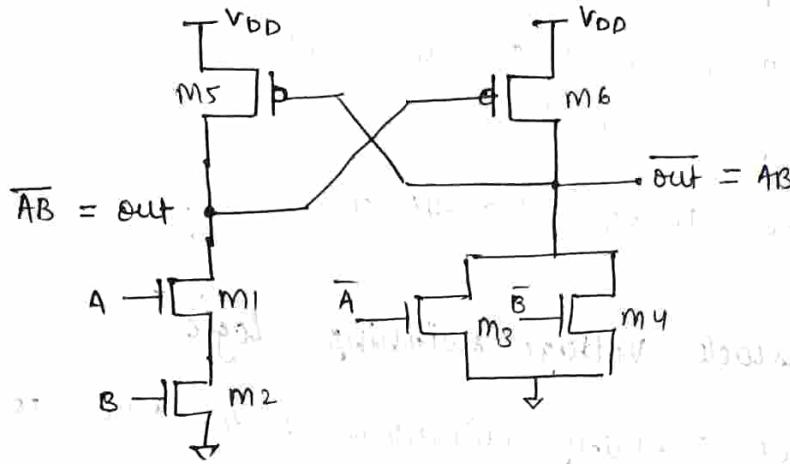


The pull down N/w
PDN1 + PDN2 are NMOS devices & are mutually exclusive

— when PDN1 conducts
PDN2 \rightarrow off & vice versa



Example : AND + NAND gate



- $A = 1, B = 1, \text{out} = 0, \overline{\text{out}} = 1$,
 $\downarrow \quad \downarrow$
 $m_6 \rightarrow \text{ON} \quad m_5 \rightarrow \text{OFF}$

- $A = 0, B = 1, \overline{\text{out}} = 0, \text{out} = 1$,
 $\downarrow \quad \downarrow$
 $m_5 \rightarrow \text{ON} \quad m_6 \rightarrow \text{OFF}$

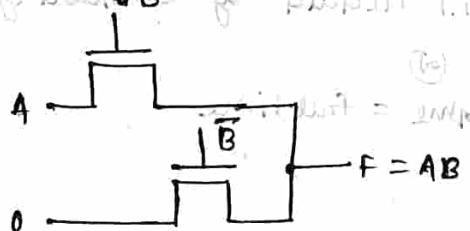
Pass Transistor logic : (Rubey pg. 259)

Alternative to complementary CMOS is PTL logic, which uses less number of transistors by allowing the primary inputs to drive gate, source, drain terminals.

- PTL will only NMOS and no PMOS
- No static power consumption because of no PMOS no direct V_{DD} to GND connection

- In PTL, NMOS are used as switches.
 — You make the switch N/W of NMOS run circle buffer to it.

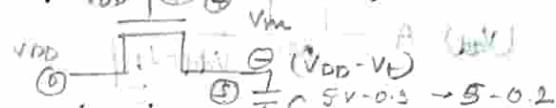
Example : AND gate



The switch driven by B seems to be redundant at first glance. (why?)
 Its presence is essential to ensure that the gate is static, i.e., it has a low impedance path to supply rails in all cases & especially when B is low.

- * Now we can implement AB (AND gate) with only 4 transistors (including the threshold to input B). ~~while~~ CMOS require 6 transistor to implement AND gate.

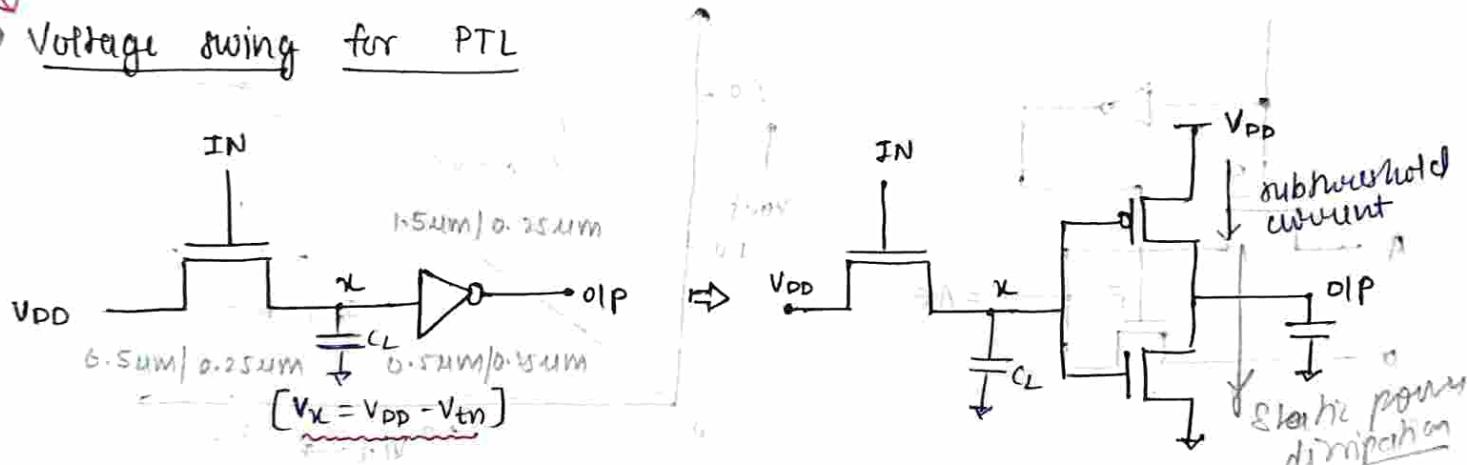
Note: An NMOS is effective at pulling '0' but poor at a node to V_{DD}. When the pass transistor pull a node high, the o/p only charges upto $V_{DD} - V_{tn}$.



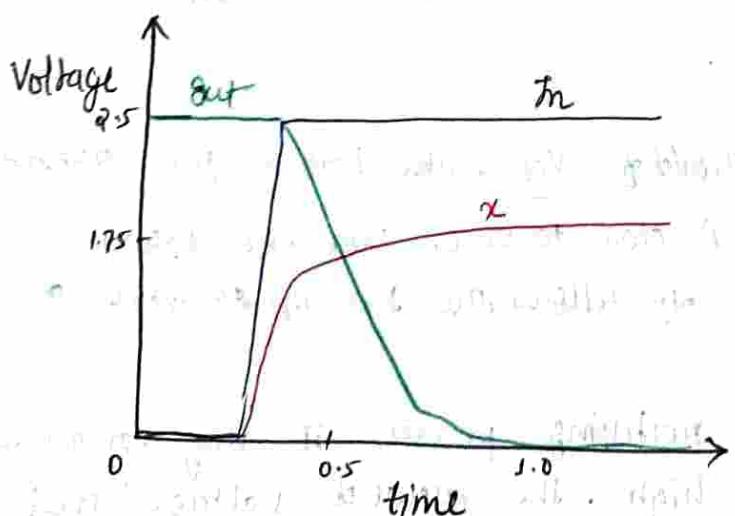
- * In PTL you can have AND functionality & NAND functionality also. So it is no longer limited to inverted o/p.

- * To improve the voltage swing at the o/p of the PTL, we add a resistor so that we can get proper V_{DD} .

Voltage swing for PTL



- The drive (gate to source voltage) reduces significantly as the o/p approaches $V_{DD} - V_{tn}$ & the current available to charge up node x reduces drastically.

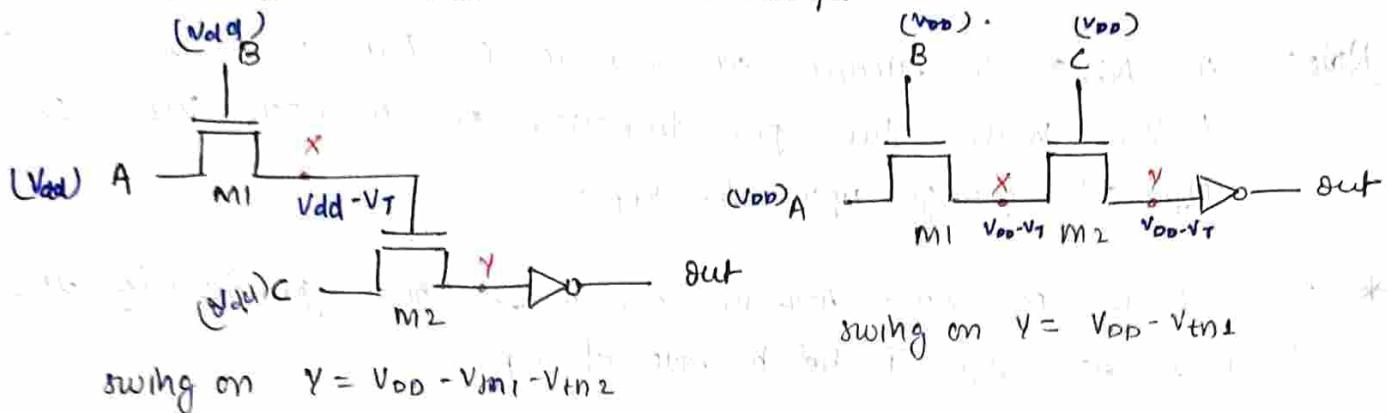


- Now the PMOS is not fully off so the prior current (subthreshold current) will not be zero.
- NMOS is completely on.
- It is in saturation.

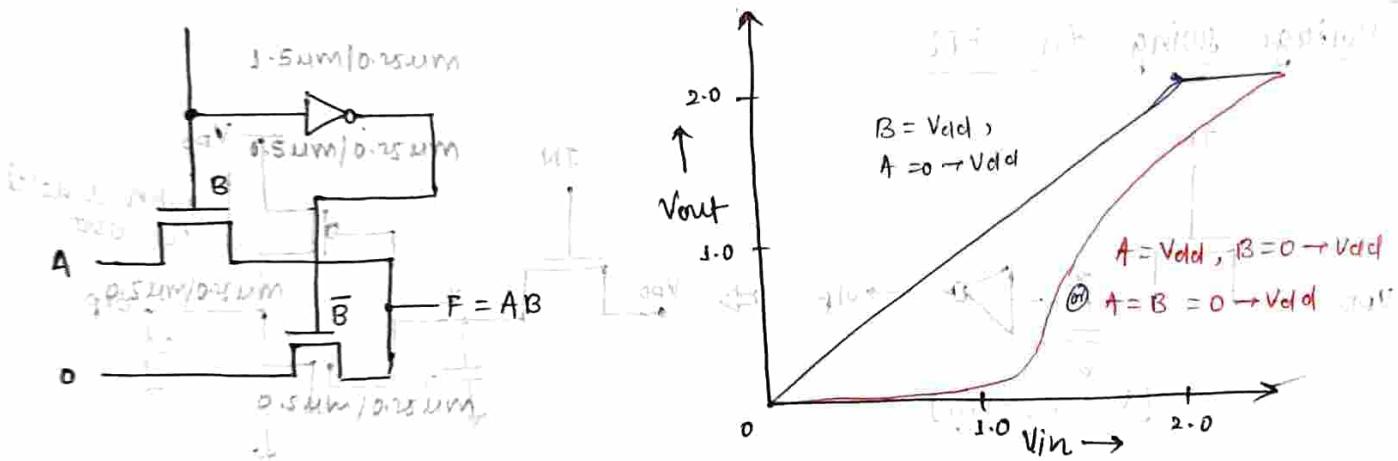
- It in PTL there should not be any static current, but here we are observing static current.

- * V_x does not pull up to 2.5V but 0.5V - V_{tn} static o/p - threshold voltage loss causes static power consumption.

- NMOS has higher threshold than PMOS due to body effect.
- PTD/P (Drain/Sources) terminal should not drive other gate terminal to avoid threshold drops.

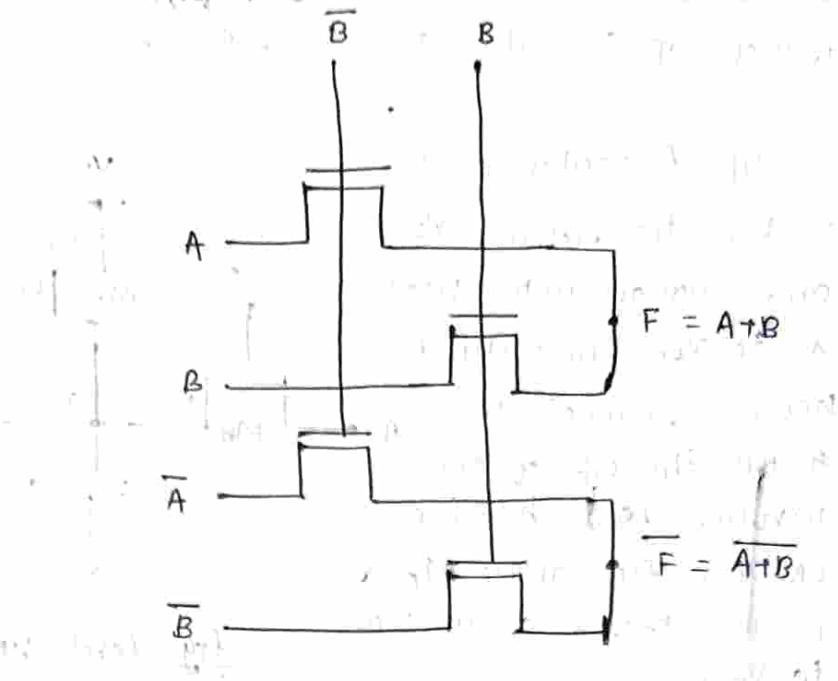
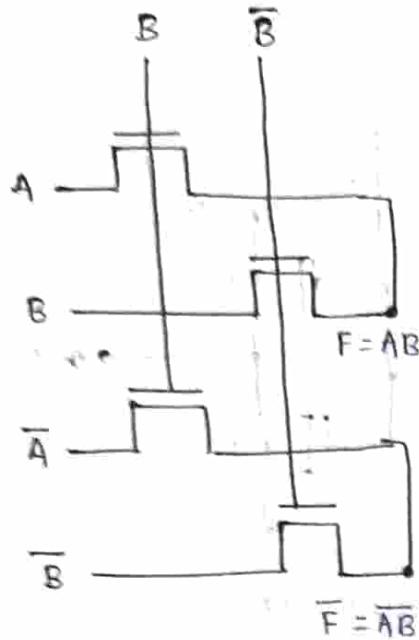


VTC of pass transistor AND gate



- (i) When $B = V_{dd}$, top pass transistor \rightarrow ON
bottom pass transistor \rightarrow OFF
the o/p just follows the i/p A until the i/p is high enough to turn off the top pass transistor (i.e., reaches $V_{dd} + V_{thn}$).
- (ii) When $A = V_{dd}$, $B = 0 \rightarrow 1$.
since the inverter has the threshold of $\frac{V_{dd}}{2}$, the bottom pass transistor is turned ON till the o/p is close to zero. Once the bottom pass transistor turns OFF, the o/p follows the i/p B without a threshold drop.
- * While the circuit exhibits low switching power, it may consumes static power when the o/p is high - the required voltage level may be insufficient to turn OFF the PMOS transistor of cross switch.

Complementary PTL & Differential PTL



AND/NAND

OR/NOR

- * CPL (complementary PTL) or DPL (Differential PTL) like single ended PTL suffers from static power dissipation & reduced noise margin, since the high ilp to the signal-restoring inverter only changes upto $V_{DD} - V_{TN}$. — Robust & Efficient PTL design are given below.

Robust & Efficient PTL

- Unfortunately, differential PTL, like single ended PTL suffers from static power dissipation & reduced noise margin, since the high ilp to the signal restoring inverter only changes upto $V_{DD} - V_{TN}$.

1. NMOS only logic: Level Restoring Transistor

- the common solution to the voltage drop problem is the use of level restorer, which is a single PMOS configured in a feedback path.

↳ it solves the problem of

- static power dissipation (it removes it)
- it ~~improves~~ improves the Rail-to-Rail swing (i.e. noise margin)

- Assume that node X is at $0V$ (out is at V_{DD} and the M_F is turned OFF) with, $B = V_{DD}$ & $A = 0$.

(i) If input A makes a 0 to V_{DD} transition, M_N only charges upto node X to $V_{DD} - V_{thN}$. This is however, enough to switch the o/p of the inverter low, turning ON the f/b device M_F & pulling node X all the way to V_{DD} .

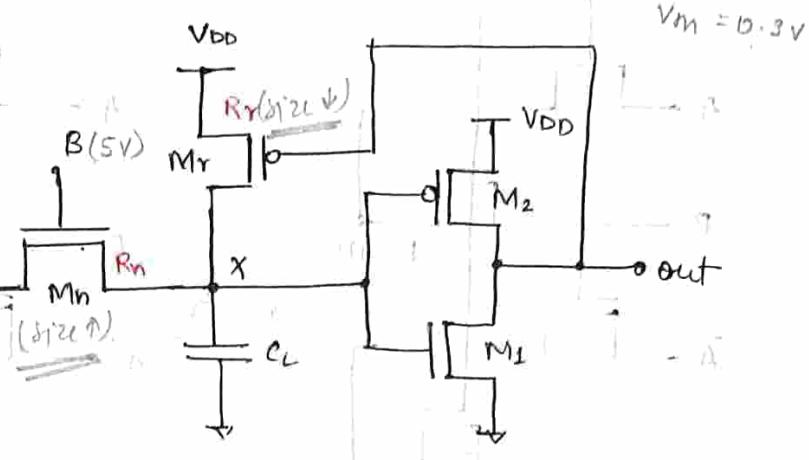


fig. Level restoring ckt

- This eliminates any static power dissipation in the inverter. Furthermore, no static current path can exist through the level restorer + pass-transistor, since all the restorers are active when A is high.

Problem with Restorer \rightarrow

- (ii) The problem occurs during the transition of node X from high to low. The pass transistor N/W attempts to pull down node X , while the level restorer pulls node X to V_{DD} . Therefore, the pull-down device must be stronger than the pull-up device in order to switch node X and the output.

from (1 to 0) (so have resistance)

We have to take the size of transistors in such a way that the voltage at node X drops below the threshold of the inverter, i.e., $V_M = f(R_1, R_2)$.

* \hookrightarrow this condn is sufficient to guarantee a switching of the o/p voltage V_{out} to V_{DD} & a turning off the level-restoring transistor.

\Rightarrow size \uparrow , $R \downarrow$, $V \downarrow$,

(M_F) decreasing the size of M_F , $R_r \uparrow$, hence $V_{(CL)} \downarrow$.

* Imp. point is that sizing of M_F is critical for PC functionality, not just performance.

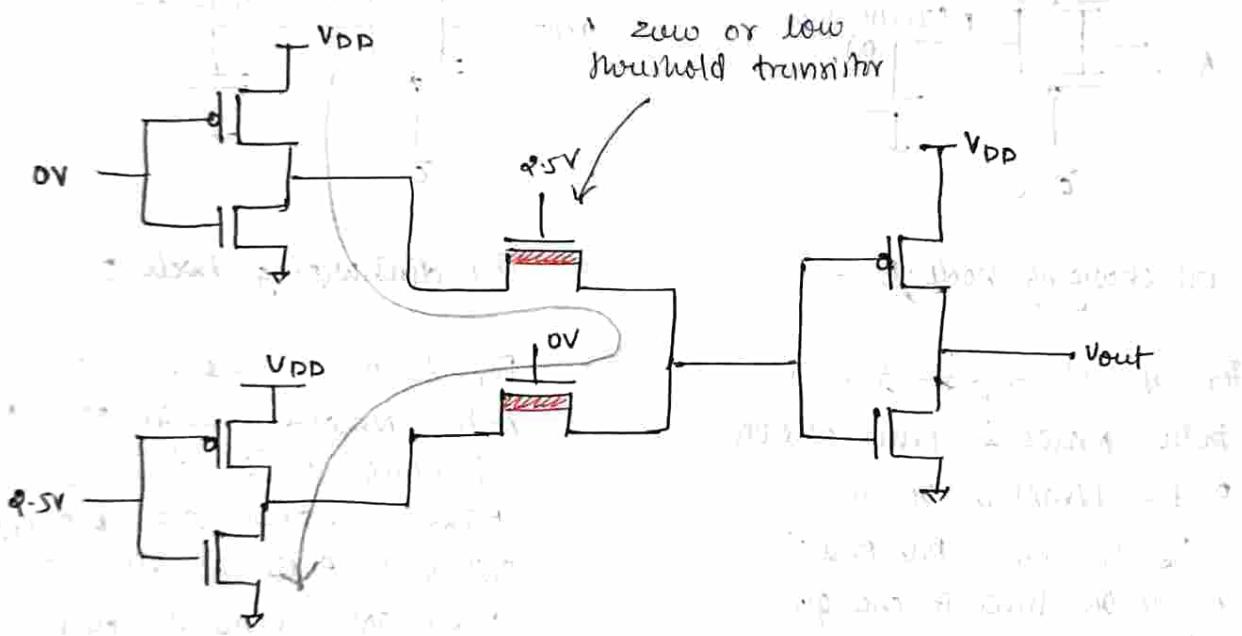
(M_N) size of M_N should be high.

P2

Another problem is on adding restoring devices in parallel in capacitance at the internal node x , slowing the gate down.

- the rise time of the gate is - very affected, since, the level restoring transistor M_2 fights the change in voltage at node x before being switched off.
- On the other hand, the level restorer enables the fall time, since the PMOS transistor, once turned ON, speeds the pull up action.

(ii) Multiple threshold transistor

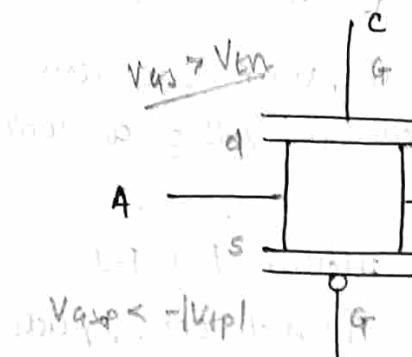


Using "zero threshold" device for the NMOS pass-transistor eliminates most of threshold drop & passes a signal close to V_{DD} .

the use of zero threshold transistor, can be dangerous due to subthreshold current that can flow through one pass transistor, even if V_{GS} is slightly below V_T .

— while sub leakage paths are not critical when the device is switching, they do pose a significant energy overhead when the device is in the idle state.

(iii) Transmission Gate Logic



A

B

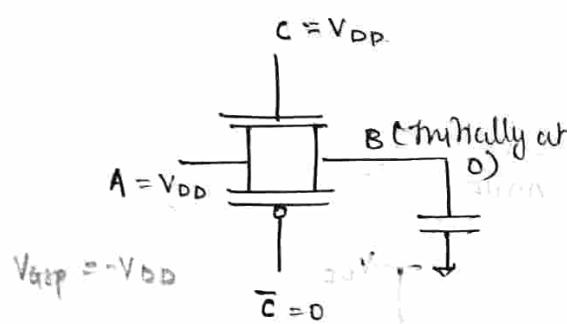
A

B

C

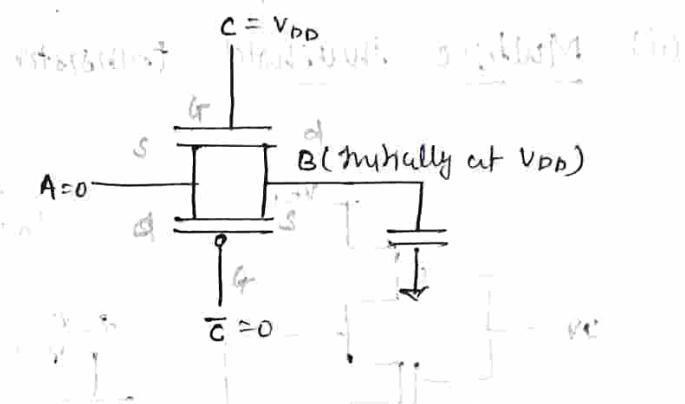
\bar{C}

"circuit"



(a) charging node B

- For $A=1 \wedge C=1 \wedge \bar{C}=0$
both NMOS & PMOS are ON
 $\wedge B=1$. NMOS is ON still
 $V_B = V_{DD} - V_{tn}$ but PMOS
 \rightarrow still ON hence B charges
to V_{DD} .



(b) discharging node B

- For $A=0 \wedge C=1 \wedge \bar{C}=0$
both NMOS & PMOS are ON
discharges to '0'. initially
PMOS pulls down the B to V_{tp}
after that PMOS gets off but NMOS
 \rightarrow still ON & hence it discharges
B to zero.

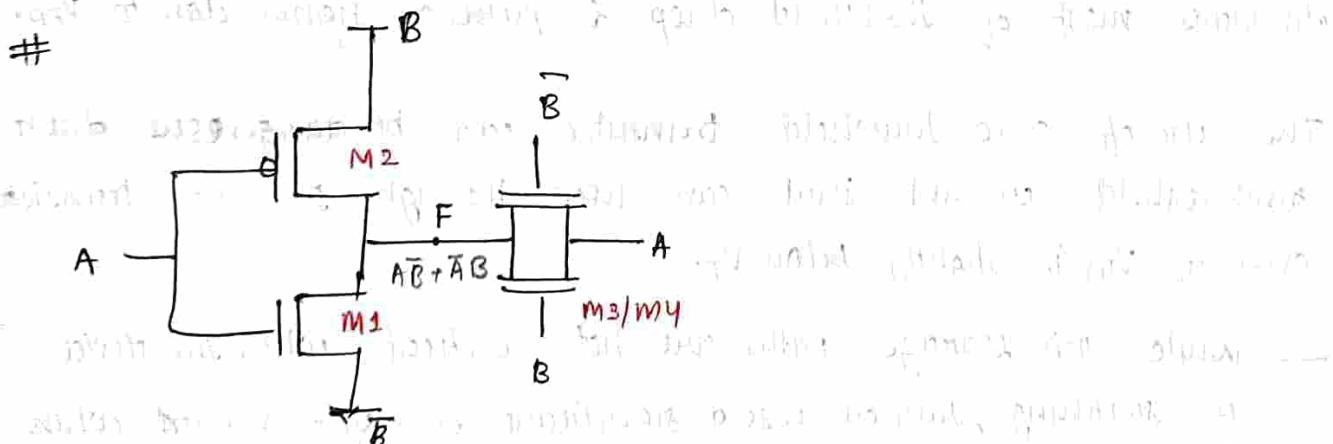


fig. Transmission gate
XOR

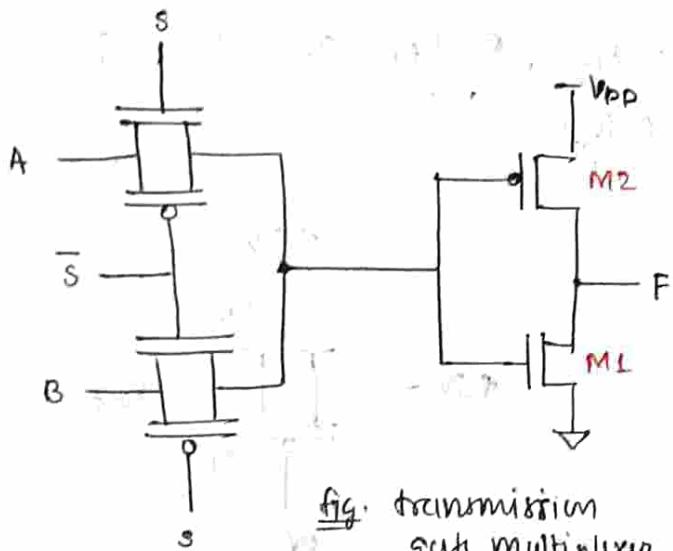


fig: transmission gate multiplexer

Delay in Transmission gate

- there is a series resistance associated with the transmission gate.
- the effective resistance of the switch D is modelled as a parallel connection of the resistance R_n & R_p of the NMOS & PMOS devices.
- For low-to-high transition of V_{out} , for low value of V_{out} - the NMOS is in saturation & resistance is high.

$$R_n = \frac{V_{DD} - V_{out}}{I_N} = \frac{V_{DD} - V_{out}}{K_n \left(\frac{W}{L}\right)_N \left[(V_{DD} - V_{out} - V_{TN}) V_{DSAT} - \frac{V_{DSAT}^2}{2} \right]}$$

$$R_n \approx \frac{V_{DD} - V_{out}}{K_n (V_{DD} - V_{out} - V_{TN}) V_{DSAT}}$$

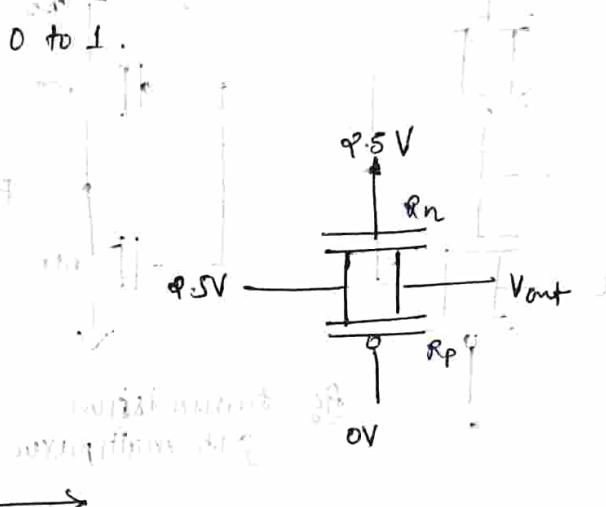
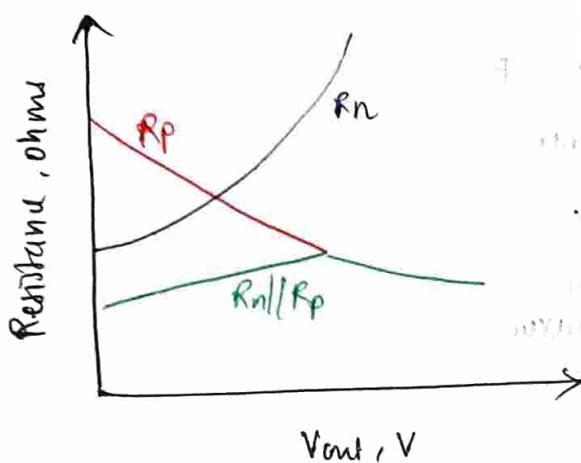
* The resistance goes up as V_{out} increases. R_n approaches ∞ as V_{out} reaches $V_{DD} - V_{TN}$ when NMOS turns off.

- Similarly, for PMOS transistor, when V_{out} is small, the PMOS is saturated but it enters the linear mode of operation when V_{out} goes to V_{DD} then R_p is high.

$$R_p = \frac{V_{DD} - V_{out}}{I_P} = \frac{V_{DD} - V_{out}}{K_p \left((V_{DD} - V_{TP}) (V_{out} - V_{DD}) - \frac{(V_{out} - V_{DD})^2}{2} \right)}$$

$$R_p \approx \frac{1}{K_p (V_{DD} - V_{TP})}$$

- * The simulated value of $R_{eq} = R_n \parallel R_p$ as function of V_{out} when V_{out} moves from 0 to 1.

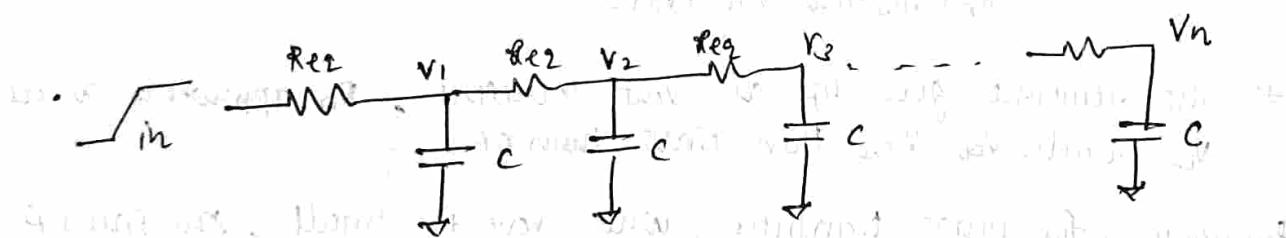


Eq. consideration of model

- The N/W of transmission gate are replaced by their equivalent resistance $R_{eq} = R_n \parallel R_p$ as function of V_{out}
- the delay of a n/w 'n' transmission gate in sequence can be obtained using elmore delay

$$t_p = 0.69 \sum_{k=0}^n (R_{eq} \cdot R) = 0.69 C R_{eq} \frac{n(n+1)}{2}$$

- * $t_p \propto n^2 \rightarrow$ as no. of switches (ie, n) \uparrow
prob. delay \uparrow



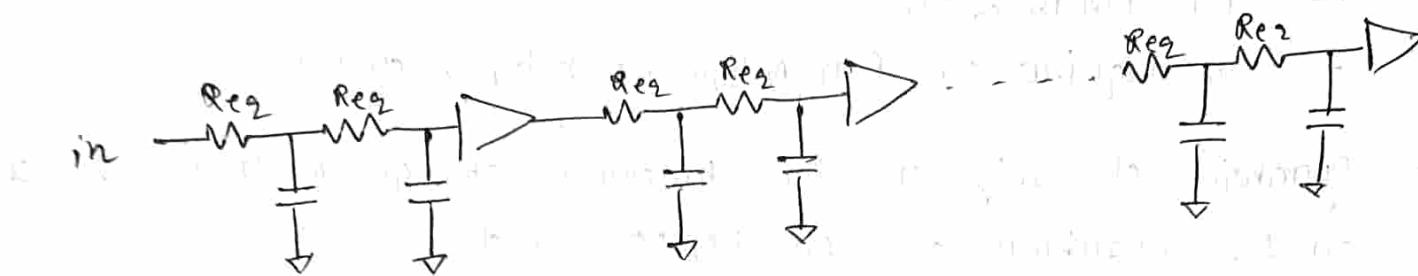
- the most common approach for dealing with the long delay is to break the chain & by inserting buffer every 'm' switches. Then delay becomes

$$t_p = 0.69 \left[\frac{n}{m} C R_{eq} \frac{m(m+1)}{2} \right] + \left(\frac{n}{m} - 1 \right) t_{buf}$$

prob. delay
of each buffer.

* the optimal no. of switches b/w buffers can be found by setting the derivative $\frac{\partial t_{p}}{\partial m} = 0$ to 0. $m_{\text{optimal}} = 1.7 \sqrt{\frac{t_{p, \text{BUR}}}{C_{Req}}}$

$$\therefore m_{\text{optimal}} = 1.7 \sqrt{\frac{t_{p, \text{BUR}}}{C_{Req}}}$$

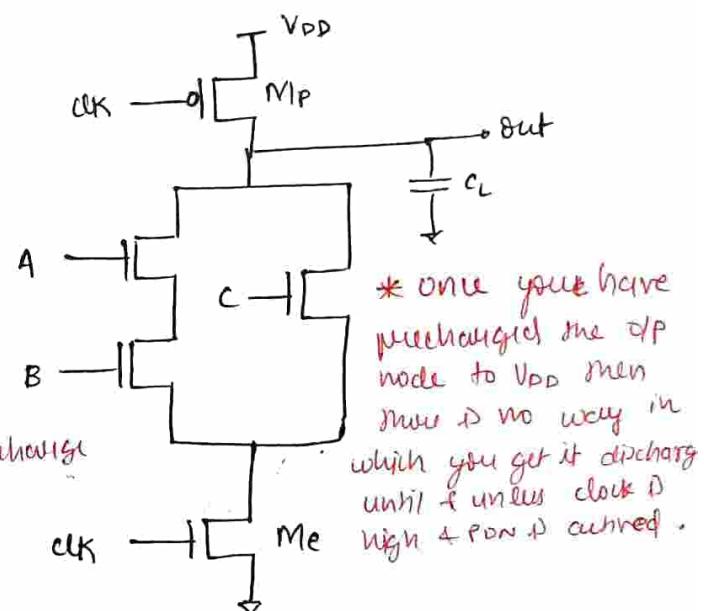
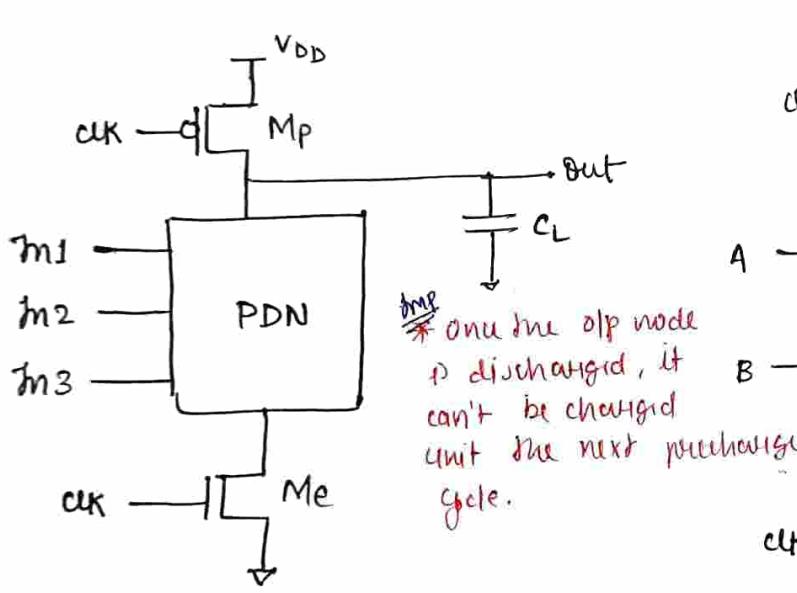


* $T_p = \sum_{i=1}^m t_{p,i}$ (Time taken by bus to transfer data)
 $t_{p,i} = \sum_{j=1}^{n_i} t_{p,j}$ (Time taken by i-th module to transfer data)
 $t_{p,j} = \sum_{k=1}^{m_j} t_{p,k}$ (Time taken by k-th stage to transfer data)
 $t_{p,k} = \sum_{l=1}^{d_k} t_{p,l}$ (Time taken by l-th component to transfer data)



2. Dynamic logic

- In static CKT at every point in time (except when switching) the o/p is connected to either GND or VDD via a low resistance path.
 - For fan-in of n
 - n requires $2n$ (n N-type + n P-type) devices
- Dynamic CKT rely on the temporary storage of signal values on the capacitance of high impedance nodes.
 - n requires $n+2$ ($n+1$ N-type + 1 P-type) devices.



Two phase operation:

(i) Precharge ($\bar{clk} = 0$)

$$\begin{aligned} \therefore M_p &\rightarrow \text{ON} \\ M_n &\rightarrow \text{OFF} \end{aligned} \quad \left. \begin{array}{l} \text{Out} = 1 \\ \end{array} \right\}$$

$$\begin{aligned} \text{Out} &= \bar{A} \bar{B} \bar{C} + (\bar{A} \bar{B} \bar{C}) \bar{clk} \\ &= \bar{clk} + (\bar{X} \bar{X} \dots \bar{X}) \bar{clk} \end{aligned}$$

"charging"

↓ Boolean expression

(ii) Evaluate ($\bar{clk} = 1$)

$$\begin{aligned} \therefore M_p &\rightarrow \text{OFF} \\ M_n &\rightarrow \text{ON} \end{aligned} \quad \left. \begin{array}{l} \text{Out} = 0 \\ \end{array} \right\}$$

"Discharging"

$$\wedge \bar{A} \bar{B} \bar{C} = 0$$

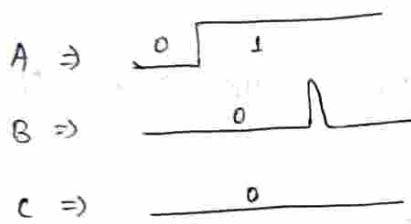
concept

* when during evaluation, when $\text{clk} = 1$ in the past } worst case
 'worst case' but $\overline{AB+C} \neq 0$ }
 (i) let $A=1, B=0, C=0$ }
 $\therefore \text{op} \text{ is kept floating or it will}$
 $\text{be in high impedance state}$

it therefore it is called dynamic logic.

- so there is some change showing (charge leakage) happening b/w ~~one~~ node but & A & we can't recover that change.
- there is no path to ground

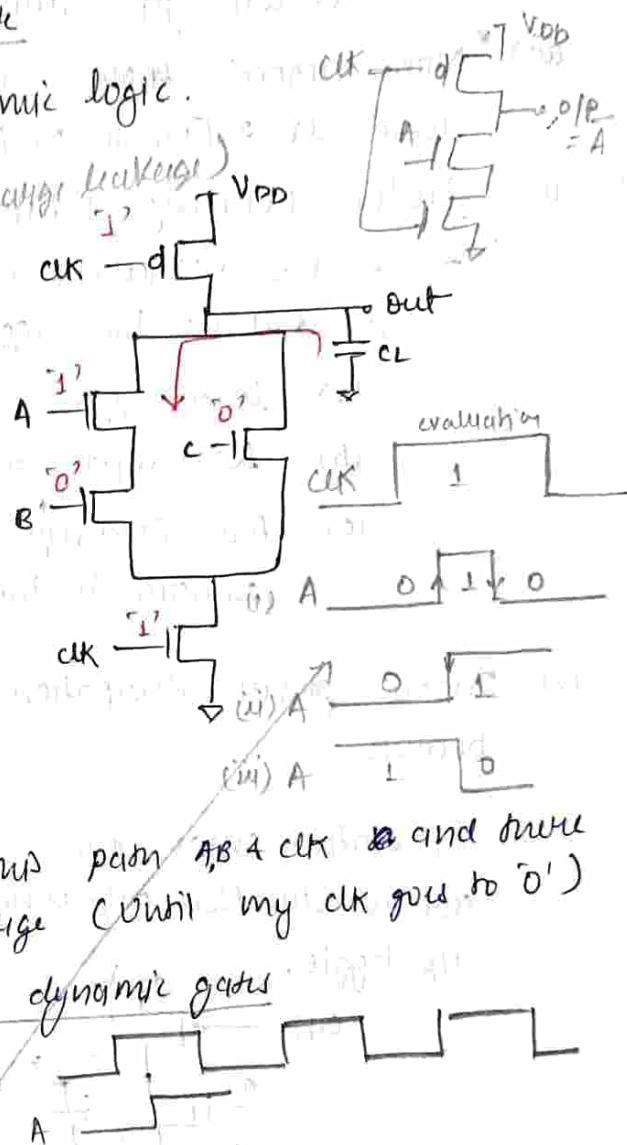
- (ii) $A=1, B=0, C=0$
 and 'B' has a glitch



- then out will discharge through this path AB & CLK & and there is no way to recover this charge (until my CLK goes to '0')
- the glitches are not acceptable in dynamic gates

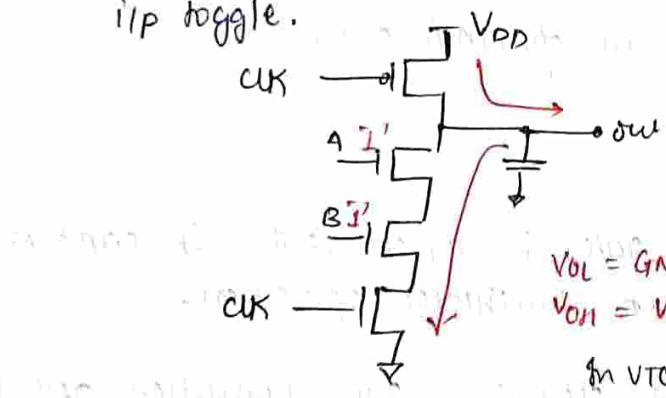
Conditions on Output

- Once the op of a dynamic gate is discharged, it can't be changed again until the next "charge" operation.
- ~~Input to the gate can make at most one transition during evaluation. and that must be a rising transition (0 to 1).~~
- ~~Output can be in high impedance state during + after the evaluation (PDN off), state is stored on C.~~
 - we can use level restorer to recover that change.
- ~~Output can only go from '1' to '0'.~~



Properties of Dynamic Gates

- (i) logic funcn is implemented by PDN only.
- (ii) It has full swing outputs ($V_{OL} = GND$ & $V_{OH} = V_{DD}$)
 - we will get full '0' & full '1' \rightarrow because of clock
- (iii) * Non-ratioed logic - sizing of the devices does not affect the logic levels. [But in static CMOS sizing decides the functionality]
- (iv) Faster switching speeds
 - As we reduces the no. of transistors that an i/p goes to hence we reduces the logical effort & path delay reduces.
 - (a) load cap. \downarrow as i/p cap. (R_{in}) \downarrow
 - (b) load cap \downarrow as smaller o/p loading (C_{out})
 - (c) No short-ckt current (I_{sc}), so all the current provided by the PDN goes into discharging C_L .
- (v) Overall power dissipation usually higher than static CMOS. because
 - In static CMOS, the current will flow only when i/p toggles.
 - But in dynamic logic, the current will flow even when i/p does not toggle + CLK is changing.



* Hence frequent changing & discharging is done due to CLK.

switching. more power dissipation.

$$V_M = V_{IH} = V_{IL} = V_{th}$$

In VTC, above which the value is read as '1' & below it the value is read as '0'.

- No static current path ever exist b/w V_{DD} & GND .
 $(\because$ No static power loss) including no short-ckt current.
- But large dynamic power is there
- No glitching is allowed
- higher transition probabilities
- extra load on CLK should be more

Q: Why no static power dissipation?

- No static power means no charge leakage will occur.
- because will never be at const. '0' or '1' as o/p resp. may be switching in charge to the PDN NMOS transistors.

(vi) PDN starts to work as soon as the i/p signal exceeds V_{TH}
so V_M , V_{IH} , V_{IL} equal to V_{TH} .

→ low noise margin

(v) Dynamic logic gets required precharge/evaluate clocks.

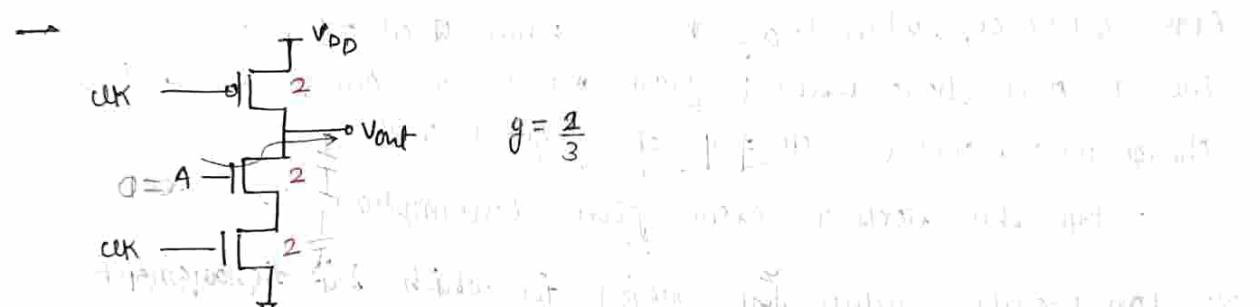
Q: Can we put CLK before the PDN & PDN after CLK (i.e. CLK near the OLP)? What will be the failure?

→ It is not usually done, but we can do it

- It occurs soon as CLK toggles we see the noise in OLP.

- As soon as one CLK comes the other things have already been evaluated

Q: What will be the logical effort of dynamic logic?



** In static CMOS, there is no extra load on CLK.
CLK only goes to FF.

In dynamic logic, CLK goes to every combination block.

∴ load on CLK is very high
for clock, $\alpha = 1$ (activity factor)

that is also the reason that the dynamic power of dynamic logic system is very high

* Once we know you do

In evaluation cycle, the i/p A, B & C will start switching only when the i/p voltage is large than the gate voltage of the individual transistors.

→ so you have to wait till the i/p voltage do the PDN crosses the threshold voltage.

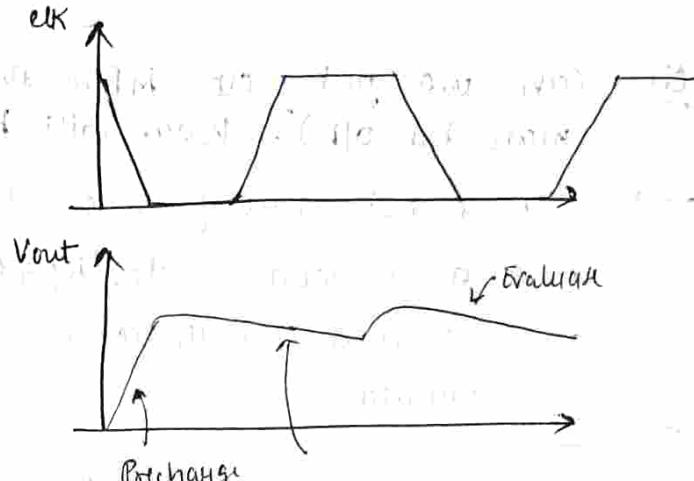
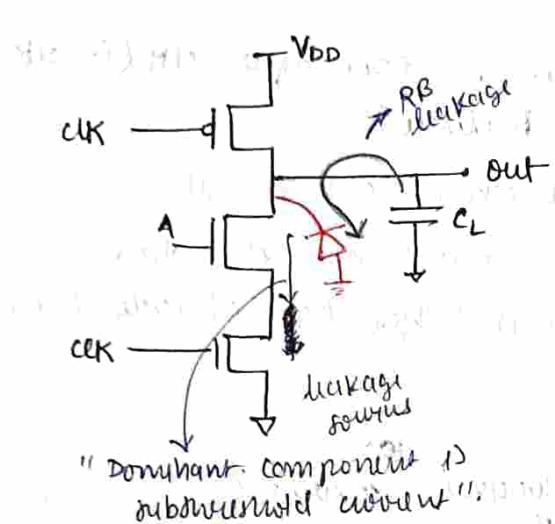
$$\left\{ \begin{array}{l} \text{switching} = V_M = V_{TH} \\ \text{threshold} \end{array} \right.$$

then NML rather than going to GND
it goes to V_{TH} (\because less NML) it does not mean V_{TH}

Issues in Dynamic Design

1. Charge Leakage

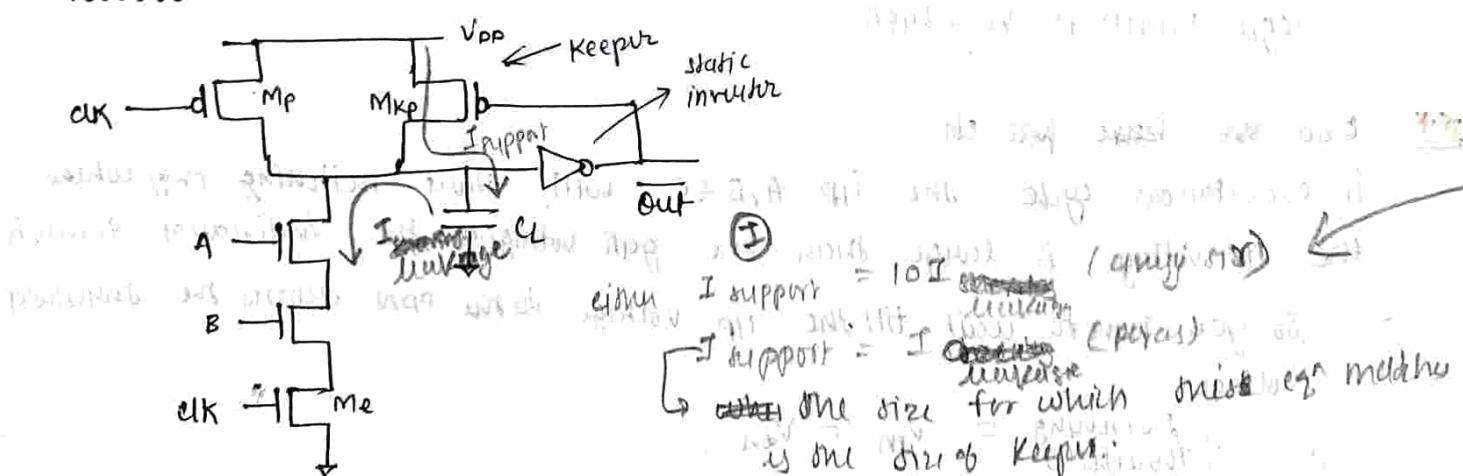
- Due to floating p/p must \Rightarrow some charge leakage.



- After precharge, when $V_{out} = V_{DD}$, when $clk = 1 \wedge A = 0$ there is some charge leakage from 'out' to A. And to sustain that charge we are giving the $clk = 1$ at frequent rate.
- But this leads to extra power consumption.

- Q. Can we also reduce the period for which this discharging is happening?
- A period of CLK in combination of logic is determined by the overall evaluation time of that combinational logic. So you can't change the time period of CLK.
- * ~~evaluation~~. Thus will be determined by the critical path.

Solution: Same approach as used earlier for PTL.



- Q. Sizing of keeper gate must be dependent on the input & the load capacitance at out.
- Our main work is to get the output from O/P of inverter only.
 - Why? We have to use very small keeper, it will restore the charge at C_L even if it has very high resistance.
 - so that keeper is fully ON & PMOS is in saturation
 - saturation current of PMOS v/s off current of NMOS
- * In project, we use small sized keeper & when $W = 0.135 \mu\text{m}$ &/or $L > 0.6 \mu\text{m}$
- Now keeper gate is a F/B gate, where we can use ratioed logic.
 - what size should we take of keeper?
 - It depends on of your device
- or PVT \Rightarrow f_S 125°C at high voltage
-
- since PMOS is slow so the PMOS keeper should still be able to compensate for the charge loss through the NMOS which are fast. Temp. is also very high, so leakage is also high.
- Q. Why don't I prefer FF 125°C 1.32V over FS 225°C 1.32V?
- Because H_D will also make the PMOS keeper very fast, it will anyway offer more current
- * keeper size should be tested at F-NMOS + S-PMOS at high voltage & high temperature, the leakage is max. & leakage strength of PMOS will also reduce.

- Q. Sizing of keeper gate is very important for design A02/A12.
How do we decide it?
- Humko aise point me keeper ka size nikalna hai, jab keeper PMOS jo hai vo slow ho, high temp & high voltage pe operate kar raha ho.

$$V_{DD} - V_{TH} = V_A$$

Part 3

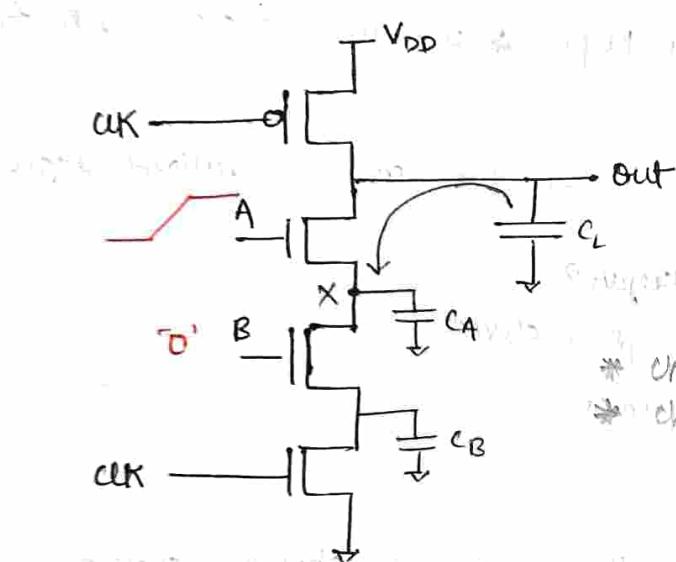
- * When there are gates cascaded one after the other, or if gate 1) cascaded after a dynamic gate then that gate has to come after the inverter, it can't come before the inverter.
(To avoid backgate coupling).

- * There is no time period for signal (ilp) A & B.

- Thus \Rightarrow timing constraint b/w clock, A & B is A & B should be settled before the clock arrives (similar to setup)

2. Charge sharing

- Initially both $A=B=0$ (no charging)



• When ilp A goes from '0' to '1'
& $B=0$ then charge gets
redistributed from C_L to C_A
leading to reduced robustness.

* Charge leakage when $A=0$ (OFF)
* Charge sharing when $A=1$ (ON)

(i) If $\Delta V_{out} < V_{tn}$ \rightarrow Final value of V_x equals $V_{DD} - V_{tn}(V_x)$

From charge conservation,

$$C_L V_{DD} = C_L V_{out}(t) + C_A [V_{DD} - V_{tn}(V_x)]$$

(1)

$$C_L V_{DD} = C_L V_{out}(t) + C_A V_{DD} - C_A V_{tn}(V_x)$$

$$- [C_A (V_{DD} - V_{tn}(V_x))] = C_L (V_{out}(t) - V_{DD})$$

$$V_{out}(t) - V_{DD} = - \frac{C_A}{C_L} [V_{DD} - V_{tn}(V_x)]$$

where,

$$\Delta V_{out} = V_{out}(t) - V_{DD}$$

$$\Delta V_{out} = - \frac{C_A}{C_L} [V_{DD} - V_{tn}(V_x)]$$

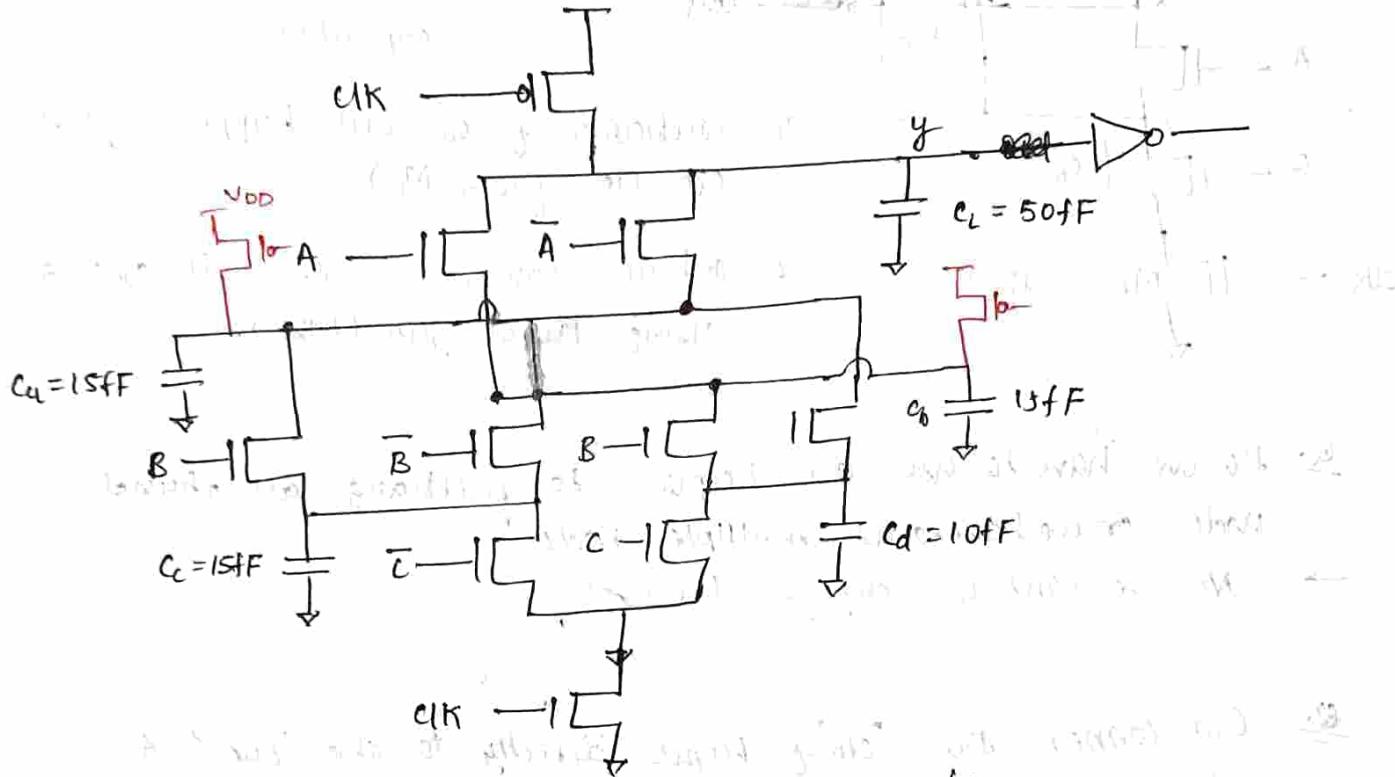
$$\therefore V_x = V_{DD} - V_{tn} =$$

(iii) If $\Delta V_{out} > V_{tn}$ \rightarrow V_{out} becomes equal to V_L

$$\therefore \Delta V_{out} = V_{DD} \left(\frac{C_A}{C_A + C_L} \right)$$

So output becomes V_L .

Example: Charge sharing example



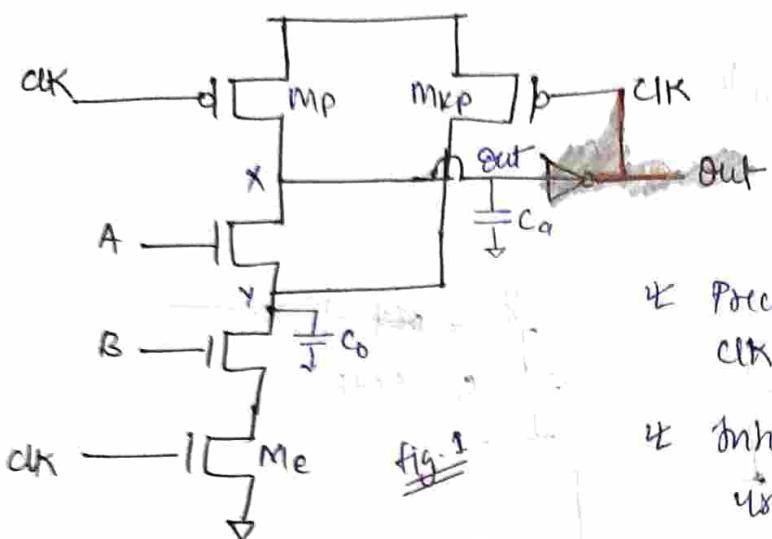
- New the total charge from C_L will distributed between the b/w all the capacitances C_A, C_B, C_C, C_D & C_L such that the total equivalent voltage across the output capacitance becomes nearly $\frac{V_{DD}}{2}$.

Q. What can be done to solve this charge sharing problem?

We can recharge all the capacitance (i.e. C_A, C_B, C_C, C_D) using the keeper gate.

→ we will not put the inverter at all the keeper but we only put inverter near the O/p.

Solution: * Precharge all the internal nodes using a clock-driven transistor (at the cost of inrush and extra power).



(i) At $\text{clk} = 0$, node $X = 1$ (i.e., M_p is ON).
(ii) If node $X = 1 \Rightarrow \text{out} = 0$.
M_{kp} = ON hence it will also charge the internal capacitance.

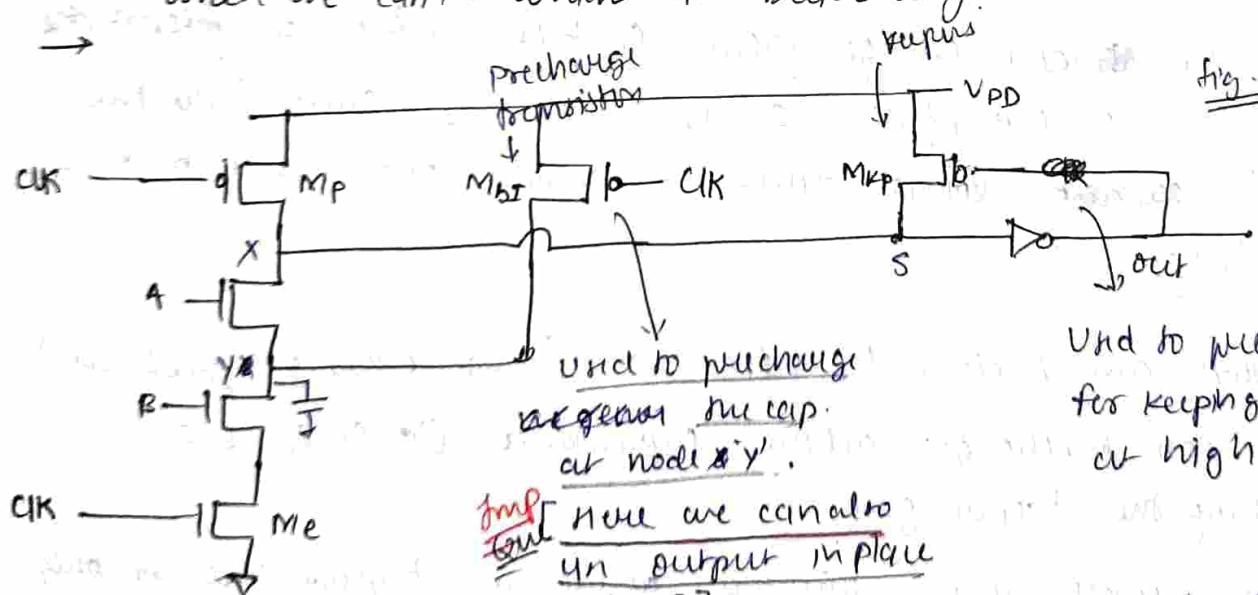
* Precharge of C_a will happen from clk (i.e., from M_p)

* Internal capacitance C_b will change using keeper gate (M_{kp}).

Q: Do we have to use one keeper to precharge all internal nodes or we have to use multiple nodes?

→ No we can't use only one keeper

Q: Can connect the 'clk' of keeper directly to the 'out' & when we can't. Which is better way?



Used to precharge or for keeping the nodes at high.

Now we can also un output in place of clk?

but we don't un M_p
because at A=0, B=1
a direct path from VDD to
ground is more.

keeper & Precharge transistors

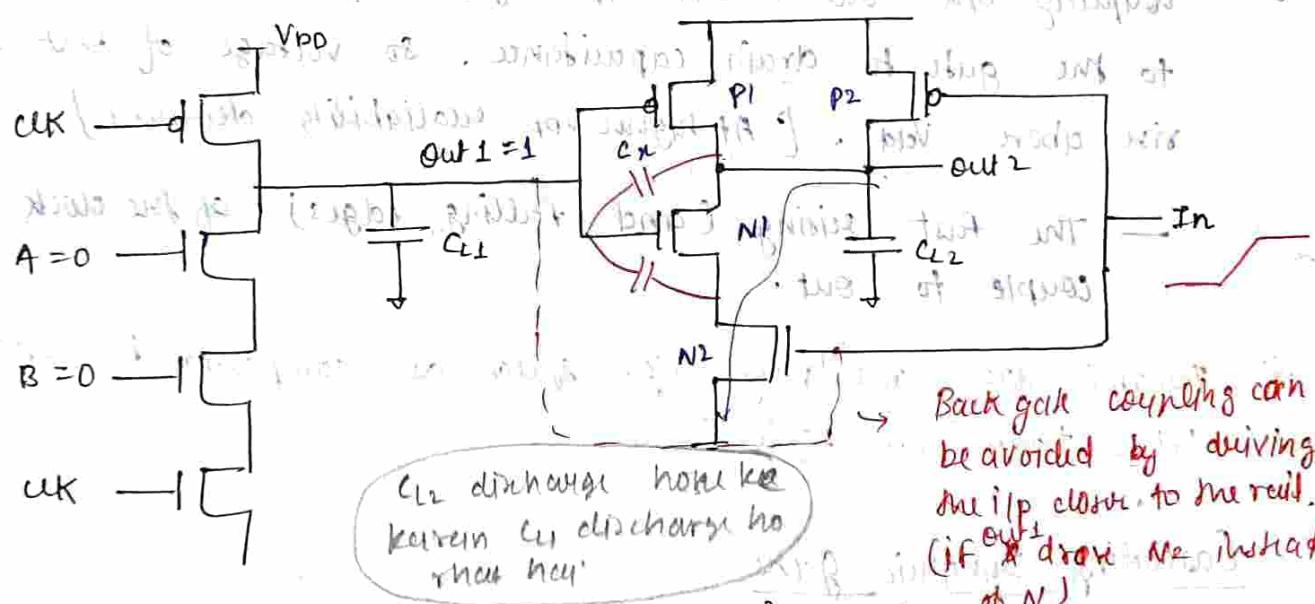
Q. Why we keep 'clk' on the bottom of PDN?

- As we know, slower signal must be closer to the o/p.
In evaluation phase, the o/p will discharge when A + B will go to '1'. So what is coming after 'clk' has gone to 1.
→ A + B. Therefore 'clk' needs to be at the bottom.

3. Backgate Coupling

To avoid this, we always add static CMOS at the o/p of dynamic CMOS.

- You could couple static gates directly at the o/p of the dynamic gate's w/o the inverter because static gate can take both kind of transitions ($0 \rightarrow 1$) & ($1 \rightarrow 0$) (so what kind of transition that dynamic gate can take?)
- But when you couple ~~static~~ ^{dynamic} gate right after the out then backgate coupling occurs.



$$f_n = (0 \rightarrow 1)_{out1} + (out1 = 1)_{[P1 OFF]}$$

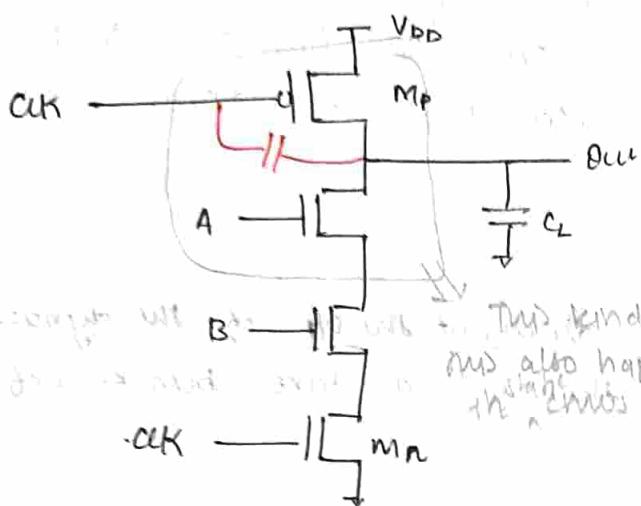
Then C_{12} at out_2 will discharge through N_1 & N_2 (P_2 off).

- When out_2 discharges, due to Miller coupling out_1 will also see some discharge happening called Back Gate Coupling. (why?)

* out_2 becomes 0, then ~~can't exist node~~ and coupling capacitance is b/w out_2 & out_1 . therefore bus ~~will exist~~ out_1 also discharges from bus Miller capacitance (C_b) path.

4. Clock Feedthrough

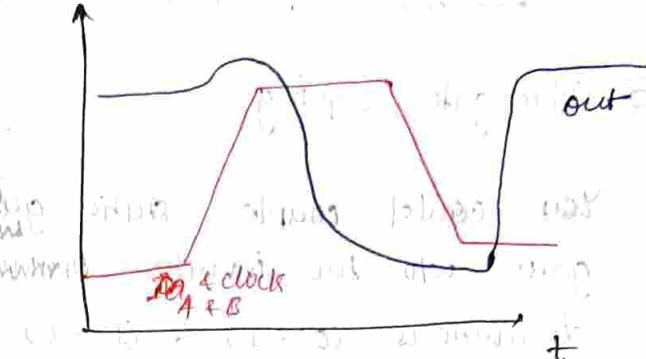
→ this can't be avoided.



This kind of
coupling also happens
in the static inverter

- when clock comes true
some feedthrough b/w out & ck.

Voltage



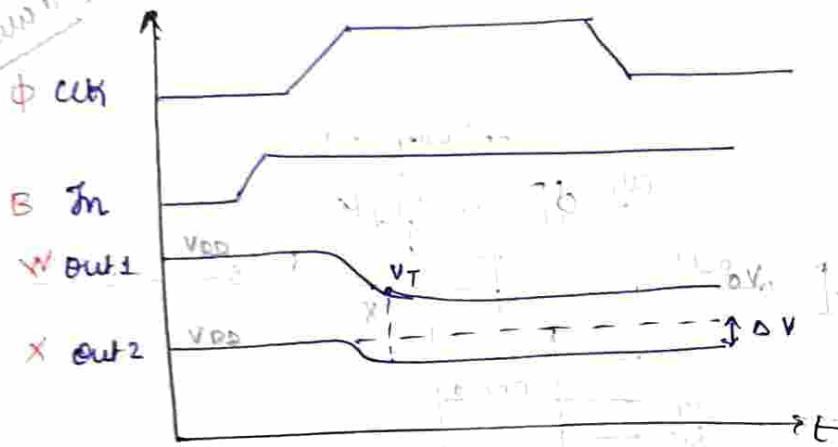
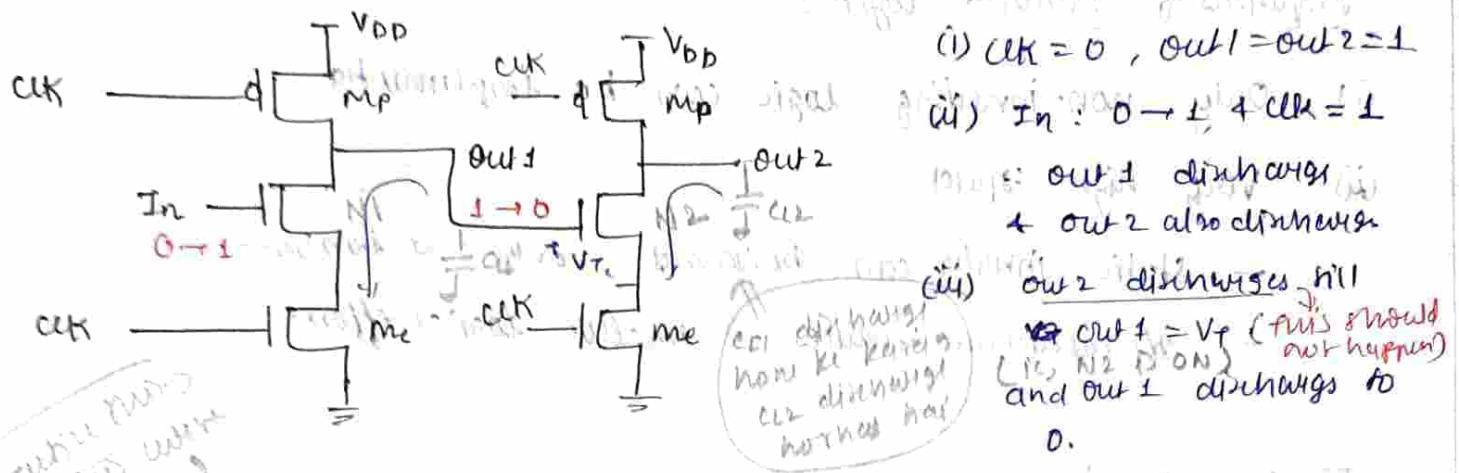
- coupling b/w out & ck if p of the precharge device due to the gate to drain capacitance. so voltage of out can rise above Vdd. [At higher Vdd reliability decreases]
- The fast rising (and falling edges) of the clock couple to out.

* Dynamic logic are to very high speed as compared to static CMOS.
PTL, transmission logic.

Cascading Dynamic Gates (Rabey 6.3.4 Pg- 283)

- During the precharge phase (i.e. CK=0) → the QP of both inverters are 1 (i.e. VDD). Inherent switches that show up in each inverter (bottom) open and both nodes gain enough potential until 3.8V.
- Assume that the primary I/P 'in' makes a $0 \rightarrow 1$ transition.

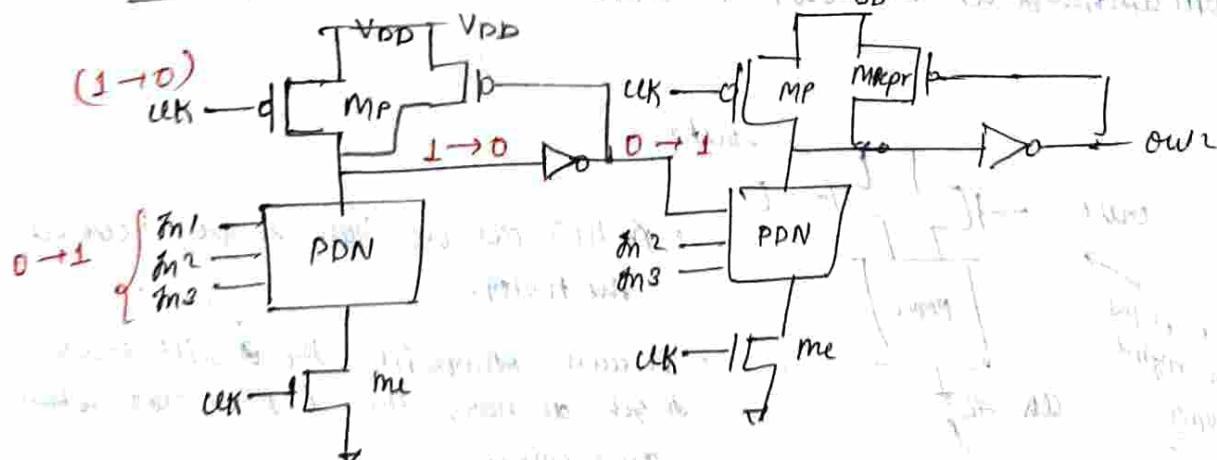
Switches have been closed but not yet turned on yet. A transient current flows through both nodes and inverts. Now it's easy with in inverter having (1.5) transistors will turn on turning triggering into a switch.



Therefore we use other ways to cascade dynamic gates.

① Domino logic :

- At 1 → 0 transition of dynamic logic, $0 \rightarrow 1$ transition is safe but $1 \rightarrow 0$ transition is not safe. (why?)



* This is called domino because one after the other (stage by stage).

* Cascading problem occurs because the outputs of each gate (is input of next stage are precharged to +)

Properties of Domino logic:

(i) Only non-inverting logic can be implemented

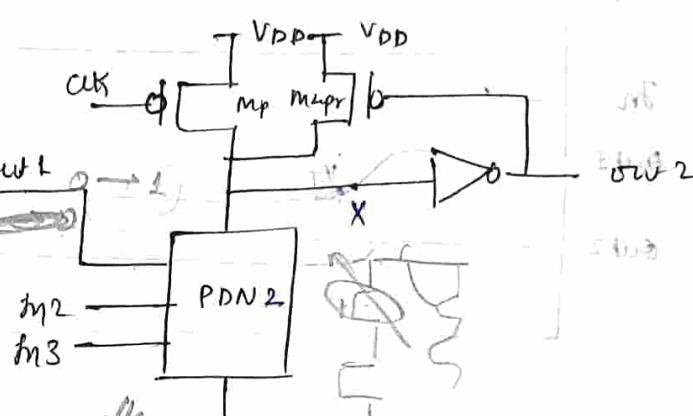
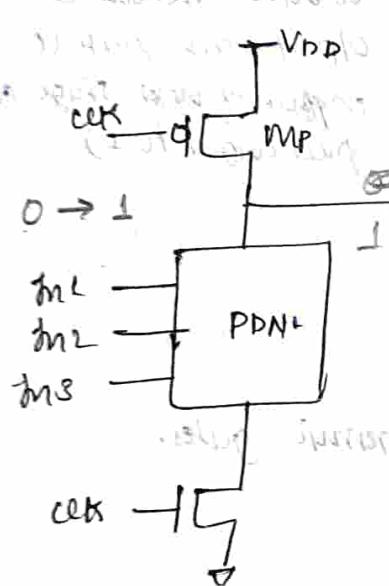
(ii) very high speed

- static charges can be stored, only $0 \rightarrow 1$ transitions

- H_p capacitance reduces \rightarrow smaller logical effort

iii) clocking + buffering

1.2 Footless Domino

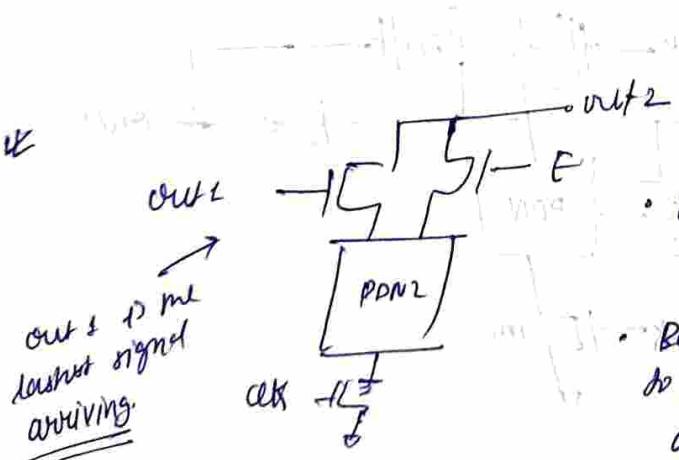


Q. What will happen when footless domino is used?

- out1 acts as a breaking point to next clk.

- and since it sufficient then only the node X will be discharged

- As $out1 = 1$ then only the node X will be discharged through PDNL2 (which does not require any clk). Node X will discharge as the $out1 = 1$ comes.



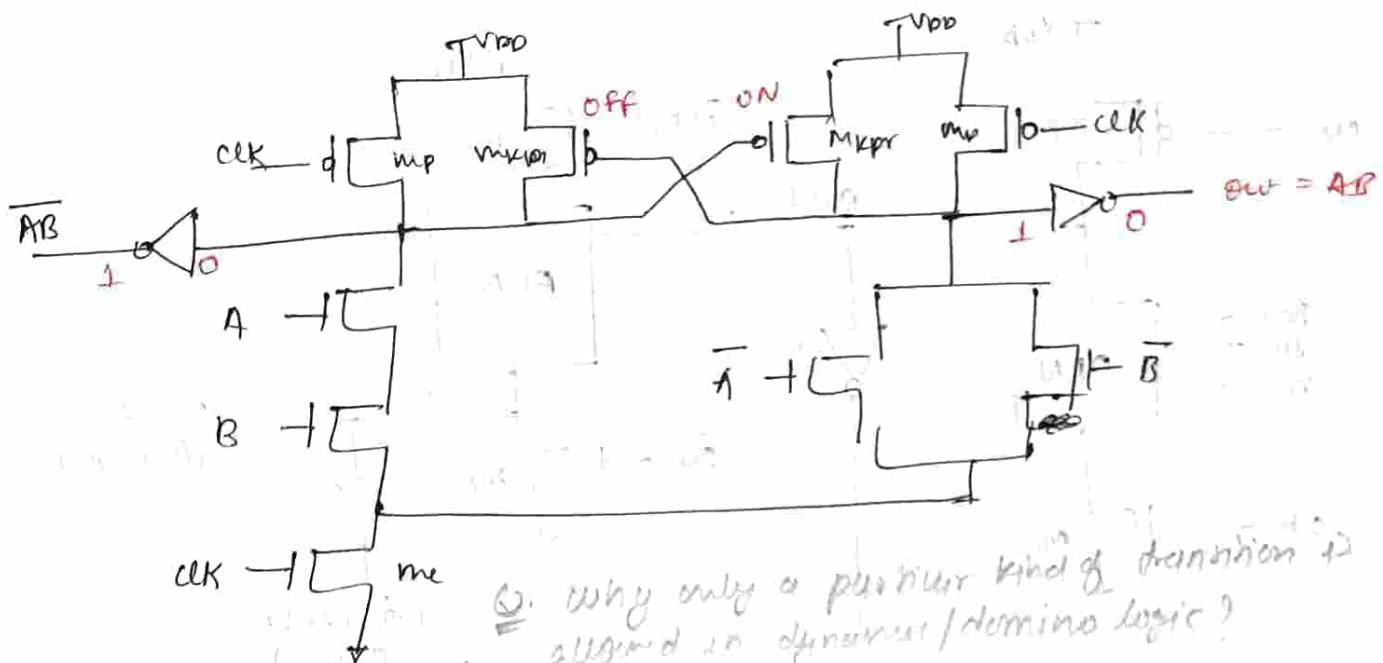
- In this case we have to use a clk at the footer.

- Because otherwise to get discharge the out2 comes.

Q. Problems with footless domino?

- only starting Nmos is allowed for parallel Nmos connection is problem

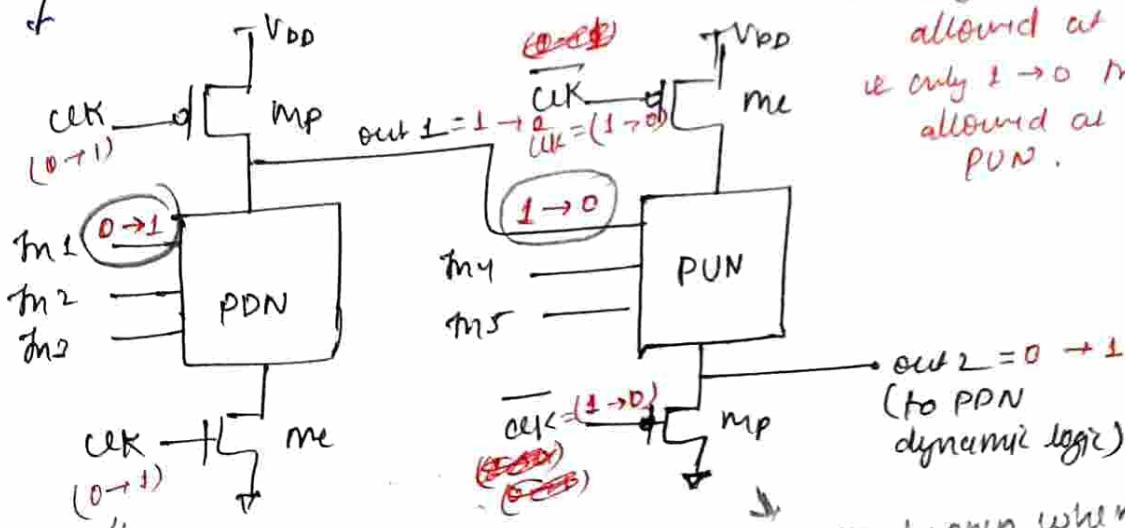
② Differential (Dual Rail) Domino



③ NP-CMOS

- when PDN is driving the PUN, then we can only allow $0 \rightarrow 1$ transition (why?)
as is allowed at the output and to get that transition one needs to invert in b/w the cascading.
- But when PUN is also connected to PDN, then transitions are only allowed so you don't need inversion in b/w the dynamic logic.

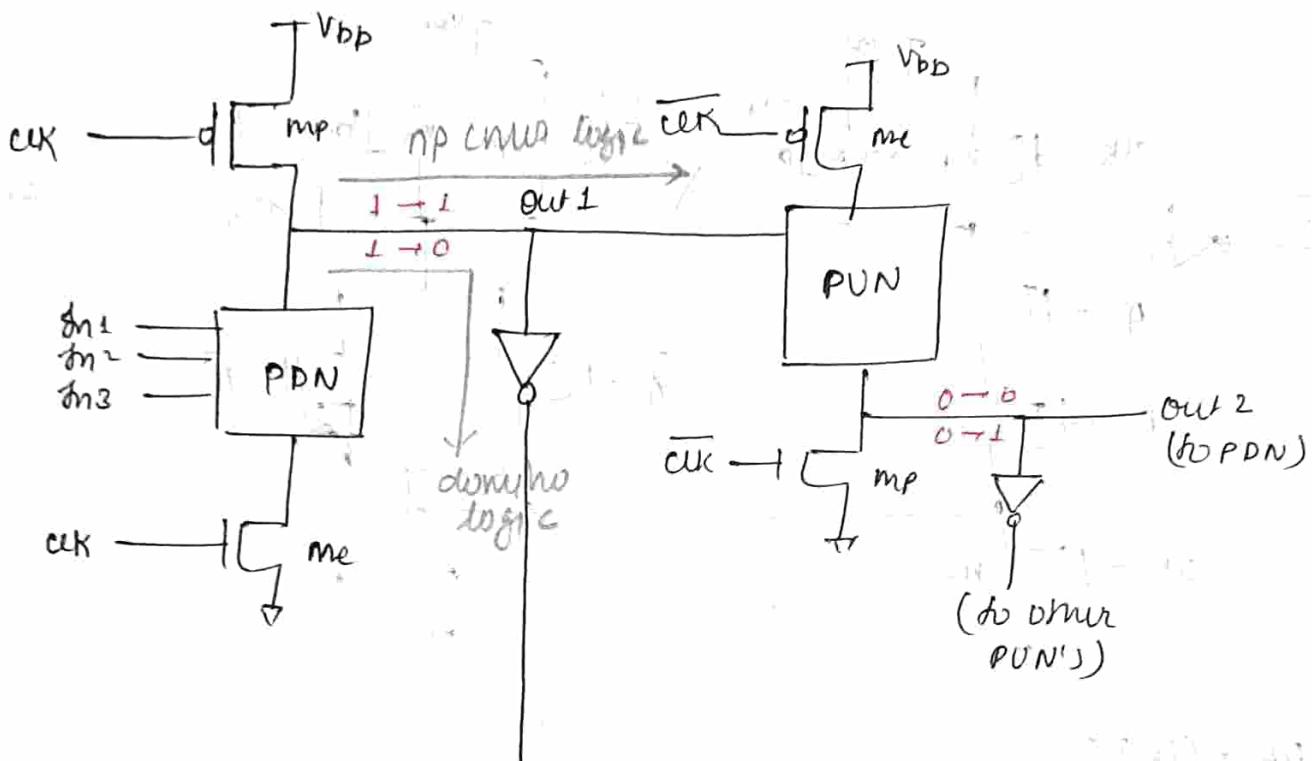
why?
Only $0 \rightarrow 1$ transition allowed at PDS & PUN
ie only $1 \rightarrow 0$ transition allowed at PPU & PUN.



Q. what will happen when m1 transition from $1 \rightarrow 0$?

Q. what will happen when transition from $0 \rightarrow 1$.
out 1

④ NORA logic:



Network 2 and single driver (to other PUN's) work well with switches in the middle.

- Using PUN increases area but it also decreases delay caused by me inverter b/w the dynamic logics



Sequential Circuit

Combinational ckt \rightarrow Output only depends on current input.

Sequential logic \rightarrow Output depends on current i/p & past o/p.
(or past i/p)

- As we need something to hold our previous state, present state or future state, for that we use state or tokens.

Eg. FSM, pipelines.

In FSM we use sequential elements to hold their state.

In pipelines we use sequential elements.

Q: What does exactly token means?

\rightarrow Token is something \Leftrightarrow when state of ~~FF~~ signal is zero
then only I will proceed forward.

* Only when my previous set of results are stored then
only I will start my second operation.

In Asynchronous ckt \rightarrow we call it tokens

In synchronous ckt \rightarrow we call it states

\hookrightarrow where in FFs when clock comes it
says okay let go to next stage.

It can have a special o/p token
give, which says that okay
next i/p.

Q: What is meant by pipeline?

\rightarrow Basically to make the o/p of ckt faster, we separate two
blocks of ckt & put flip-flops in b/w so that one o/p
of previous one is not able to interact with the i/p of the
next stage until I find one clock signal.

Q. do even if the delays ~~toff~~ till the transitions of i/p for both the ckt for stage 1 & stage 2 are similar, even then to make the ckt faster, ~~should~~ do we put FF or not.

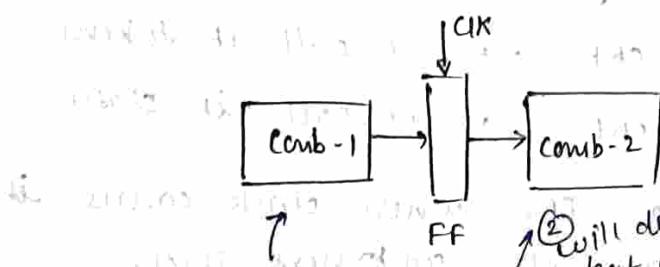
→ Yes, even then also we put the FFs. because functionality will not be affected even after removing the FF but by adding FF it makes the ckt faster. That's why we use it.

Q. There are two purposes of using such sequencing class

- 2 marks
- (i) To ensure sequencing, i.e. when all the i/p have come then only the o/p should be calculated.
→ earlier i/p should not get mixed with the present i/p.
 - (ii) To improve overall system throughput.
→ To make the system faster.

* We need to keep the sequencing elements to play safe for worst cases.

Note:



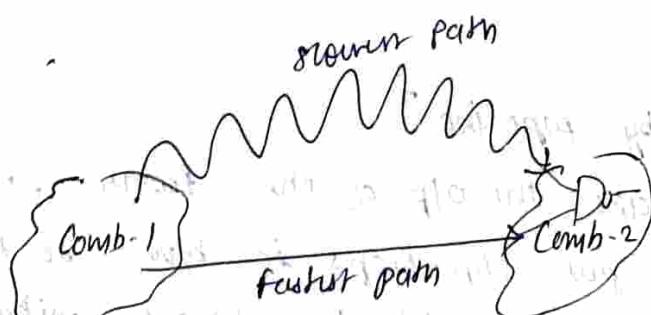
- Q. Why do pipelining?
→ ① To make ckt faster
② To improve throughput
③ To care about the slow coming signal
④ To tackle with spurious o/p

① What ever

happened in the
previous cycle

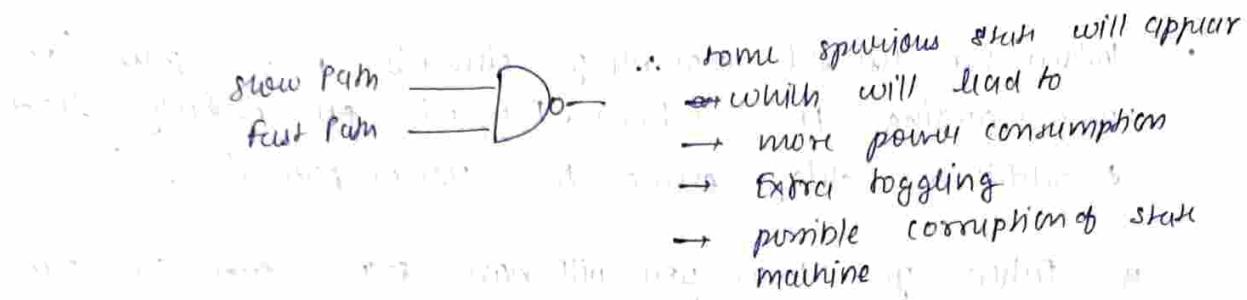
happened in the
previous cycle

~~last~~



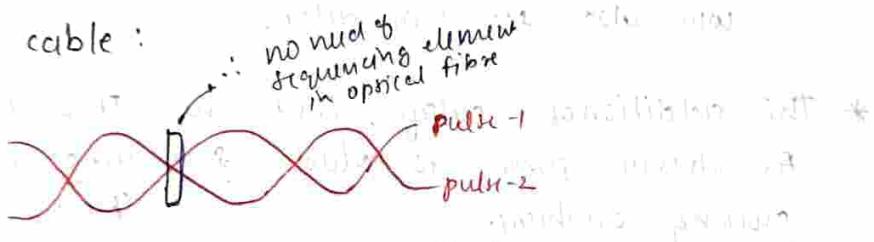
if faster path & slower path would never intersect
then we are safe.

⇒ there are two paths, one comes very fast & other comes slow. Then in intermediate state, when this path has reached but the other one has not reached I could get glitches. @ we get spurious transitions.



- * if we can't ensure that flow of information will always be synchronized we use what is called as tokens. @ FF.

Eg. fibre-optic cable:



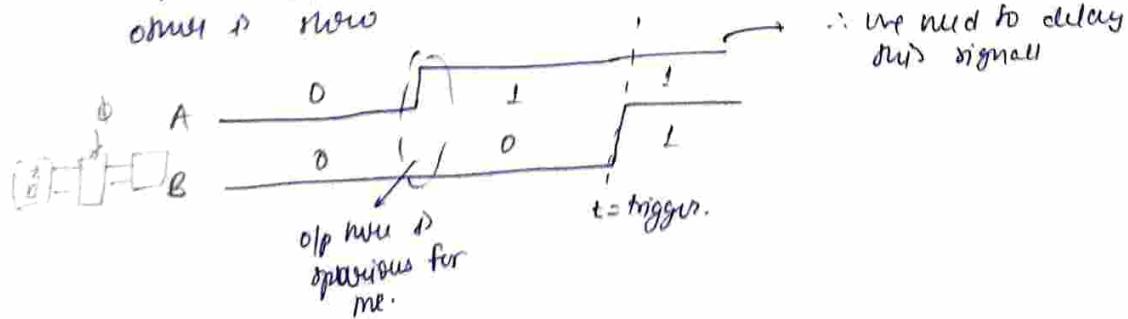
both will travel with speed of light

- * it is only the dispersion b/w the pulses that delays.
- what is the minimum time b/w two pulses?

- * But in terms of our hardware @ software we can't always ensure that the slow path would be able to catch up with fast path.
- ∴ We need FF or buffers to delay this fast signals

Q. What is meant by spurious OLP? Because whenever we get

- we get the OLP acc. to me functionality?
- suppose I want to add two bits & one is fast & other is slow



Note: If tokens move through pipeline at constant speed, then no sequencing element would be necessary.

Sequencing overhead :

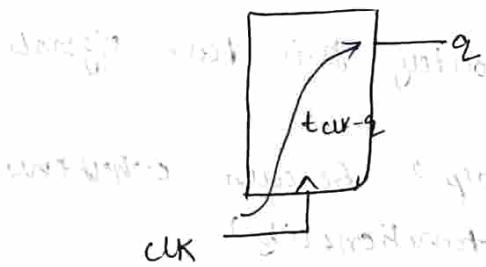
* Def When we use this sequencing elements in our path (ie, we are adding FF & latches) & they will ~~be~~ inevitably lead to additional delay even on the slower paths. (Why?)

* When you know you will have some slow path & fast path in your system, but you exactly don't know which path is fast & which is slow.

∴ We put FF on all the paths, then the slower path will also see some delay.

* This additional delay that this FF or latch will take for slower path is called sequencing overhead (better)

* This kind of overhead also occurs in token system. In token system (ie, in synchronous system) it is called as clocking overhead.



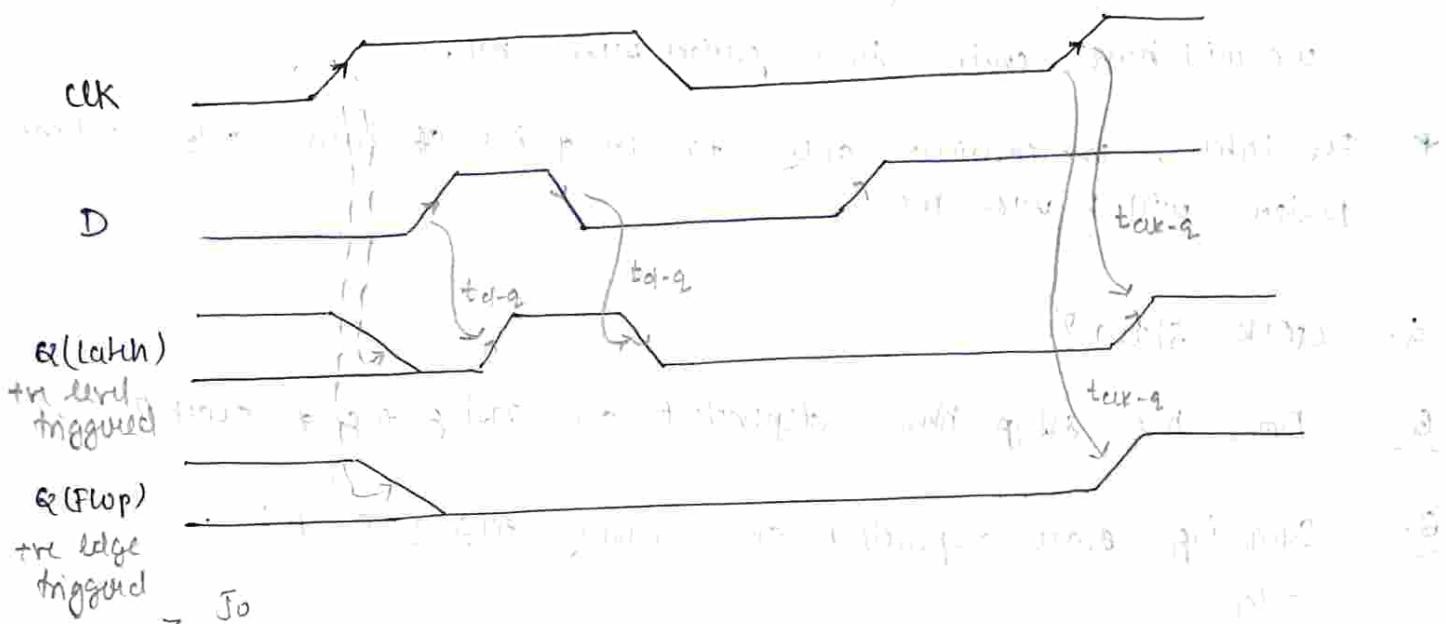
* the sequencing overhead will be the FF delay ie. clk to Q delay of the FF.

States	token
synchronous sequencing overhead	Asynchronous clocking overhead

→ because that is the extra delay that your logic has to travel.

Sequencing Elements

1. Latch → level sensitive (transparent Latch, D-Latch)
2. Flip-Flop → edge triggered (Master-Slave FF, D-FF, D-Register)



- * Latch can have both t_{d-q} delay & t_{ck-q} delay
 - if D comes late (i.e., after the CLK) then t_{d-q} delay will be ~~more~~
 - if CLK comes late (i.e., after the D) then t_{ck-q} delay will be ~~more~~
- * FF will only have t_{ck-q} delay.

- Note:
- Only on the CLK edge FF will be transparent and other than that it ~~is opaque~~ will hold its previous value.
 - Only on the ~~the~~ level of the clock latch will be transparent & other than that it ~~is opaque~~ will hold its previous value.

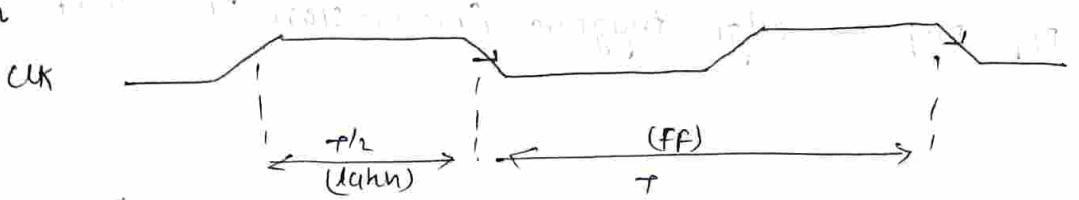
* FF will transmit the data whatever is given on the rising edge. (for edge triggered FF).

Ques: Why we prefer FF over the latch?

- ① Transparent window for FF is very low \therefore it will not transmit any glitch.
it is less likely that it will transmit any glitch.
But latch will also transmit the glitch.

Q. Why we prefer FF over latch when we study STA?

* (2) For FF, we consider sampling at the edge of the clock when



we will have entire time period with me.

* For latch, we consider only the end of the clk when only $T/2$ time period will be with me.

(3) To avoid race condition latch.

Q. Clock skew?

Q. Does setup time dependent on rising edge of clock?

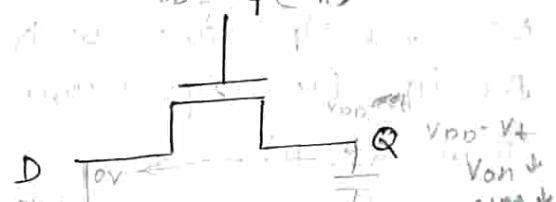
Q. Sampling time dependent on rising edge of clock?

→ Yes.

1. Latch Design:

$$V_{DD} = \phi(V_{DD})$$

(1) Pass Transistor Latch:



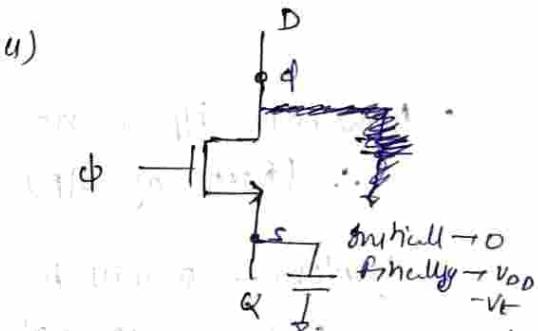
- Pros:
 - tiny
 - low clock load

~~will have floating voltage for all~~

Cons: (i) V_T drop at Q, the o/p cannot swing from 0 to V_{DD} .

Assume, D as (driven) & Q as (output)

o/p being able if $D = V_{DD}$, & $\phi = \text{logic-1}$
to drive IP when Q becomes $V_{DD} - V_T$

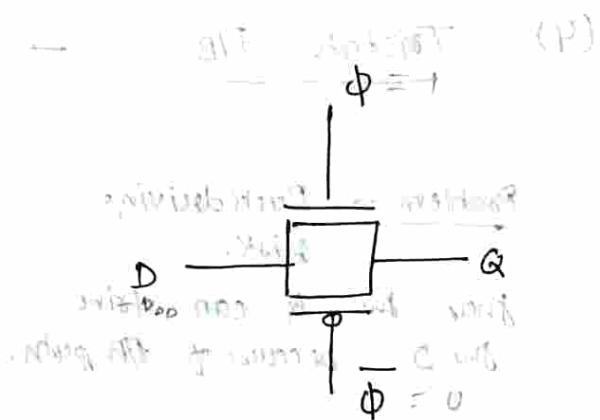


- Bulk driving: $V_{QI} = \phi - (V_{DD} - V_T) = V_{DD} - V_{DD} + V_T = V_T$
- Non-restoring
- Signal can flow in both directions.
- (ii) The o/p is dynamic @ the o/p floats when the latch is opaque.

- (iii) D drives the diffusion IP of a pass transistor directly, leading to potential noise issue & making delay harder to model with static timing analyzer.
- (iv) The state node is exposed, so noise on the o/p can corrupt the state

(2) transmission Gate Latch:

- Now it has no V_T drop
- But other problem remains same
- Extra problem is, it requires inverted clock also.

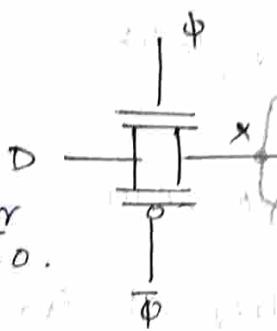


extra inverter stage
and inverter latches

(3) Cascading with Inverter (Inverting Buffer)

- Now the OLP can't back drive the IIP (i.e. signal can't flow from Q to D).

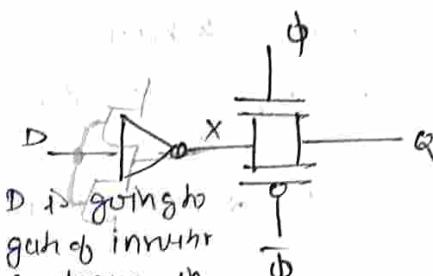
Problem → X will go on floating or dynamic in nature when both $\phi = 0$.



α is going to draw a lot of current ... it is not connected to diffusion IIP.

- Now the IIP is no longer driving the diffusion. \therefore (Restoring OLP)

Problem → Q will be floating or dynamic in nature when both $\phi = 0$.



D is going to go to Q so change in ϕ

α would not transmit onto D.

We can fix only either of two problems:-

(i) Output noise sensitizing or

(ii) Diffusion IIP is not able to drive

signals with different voltage value

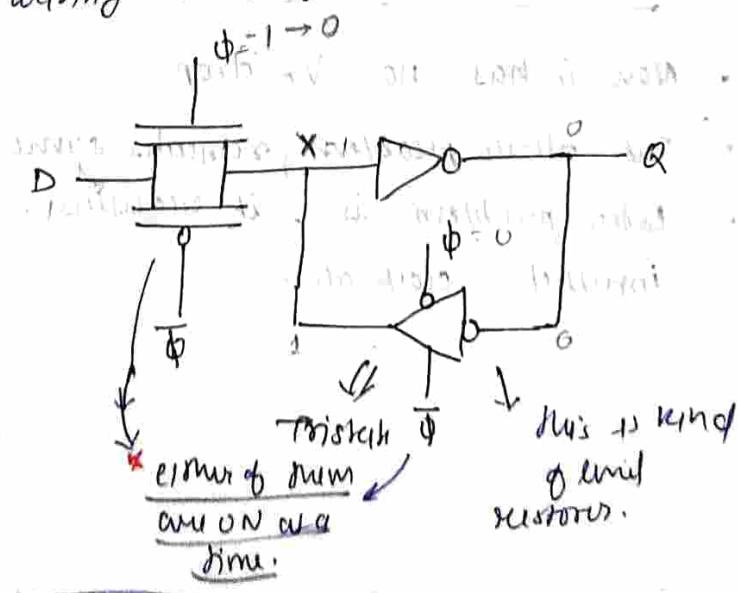
Q1 and Q2 will be in floating state \Rightarrow static noise problem

(4) To-STATE FIB \rightarrow this will solve the problem of floating value of X.

Problem → Backdriving risk.

that the α can drive the D because of FIB path.

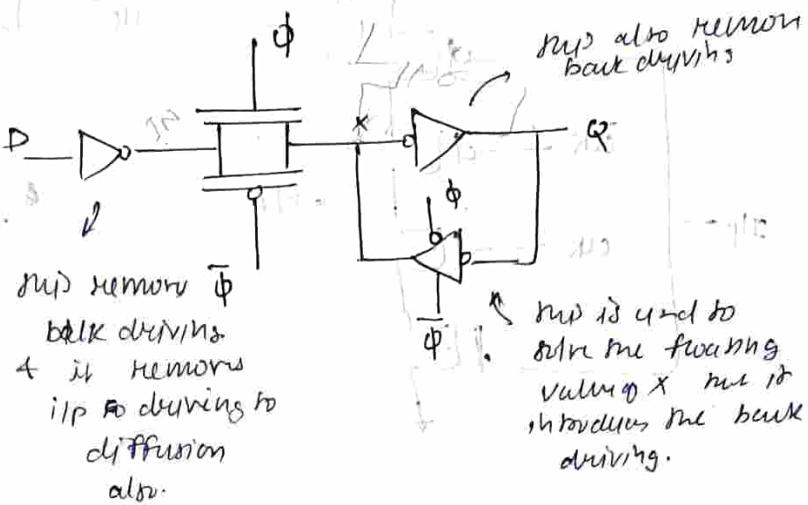
* static latencies are now essential because of leakage



This is kind of leaky restorer.

(5) Buffered i/p with f/b at o/p → To remove the risk of back driving.

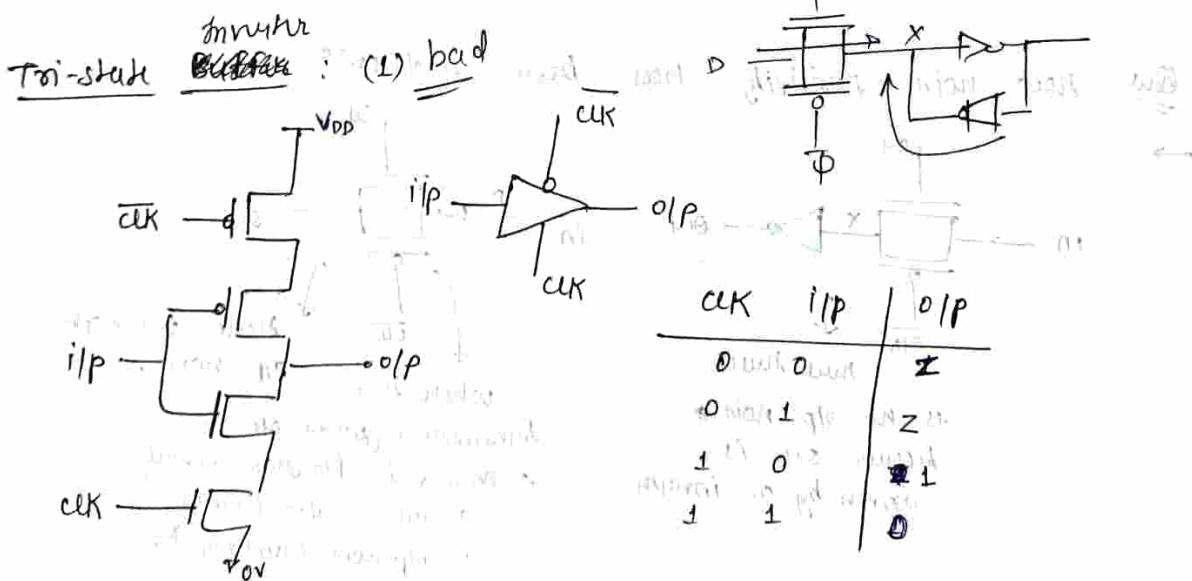
Problem: There can be some noise on Q that will impact X.



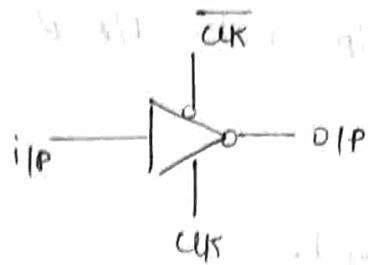
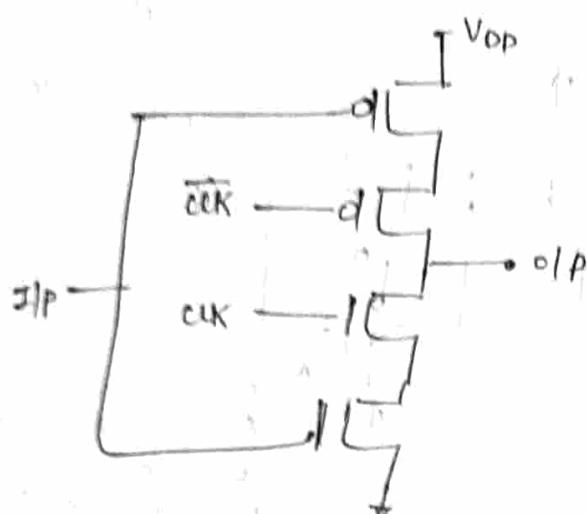
- We put tri-state dynamic logic as level restorer to keep X at 1.0 in most of the time.
- But now we put ~~inverter~~ for ~~restorer~~ as a level restorer so to keep X at both 1 & 0 at required times.

Q: why did we put tri-state & not a simple inverter?

- if tri-state = P, supplied by inverter then both will be driving X & inverts will also be driving X hence it would be difficult to study now.



Tristate buffer (e) good

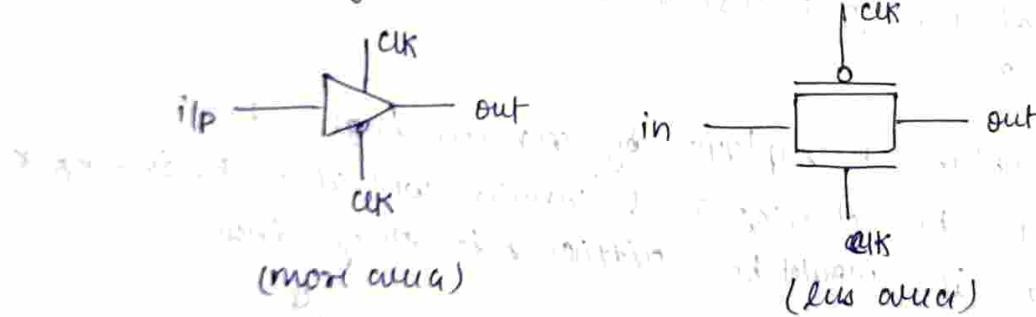


CLK	IOP	OIP
0	0	Z
0	1	Z
1	0	1
1	1	0

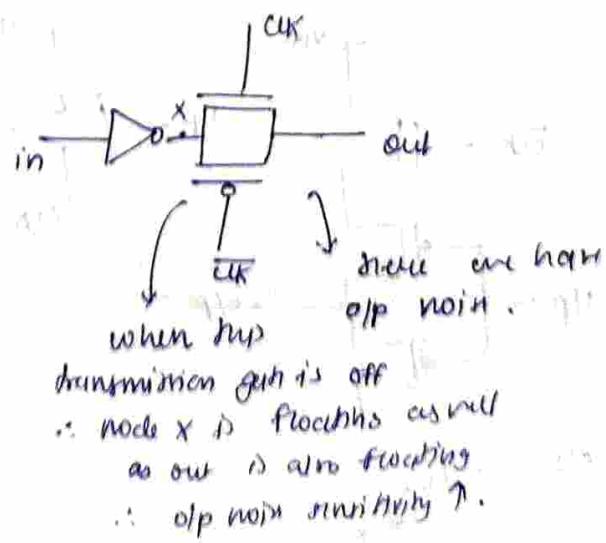
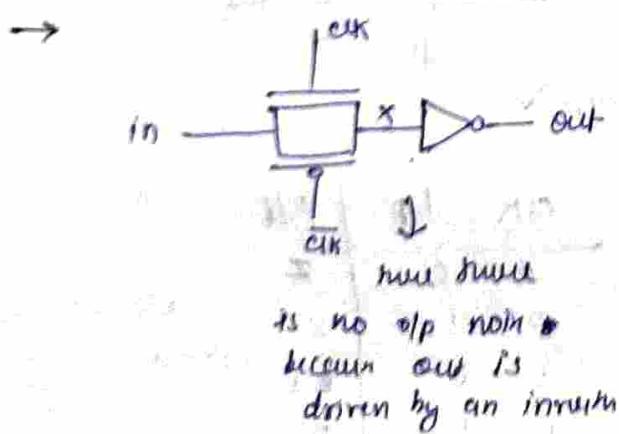
Q. What is the difference b/w the dip & the glitch?

→ Glitch means that there is something like or full, i.e., some transition on the signal is so large that the subsequent stage samples it as a change in input.

Note: Transmission gate → advanced version of tristate ~~buffer~~

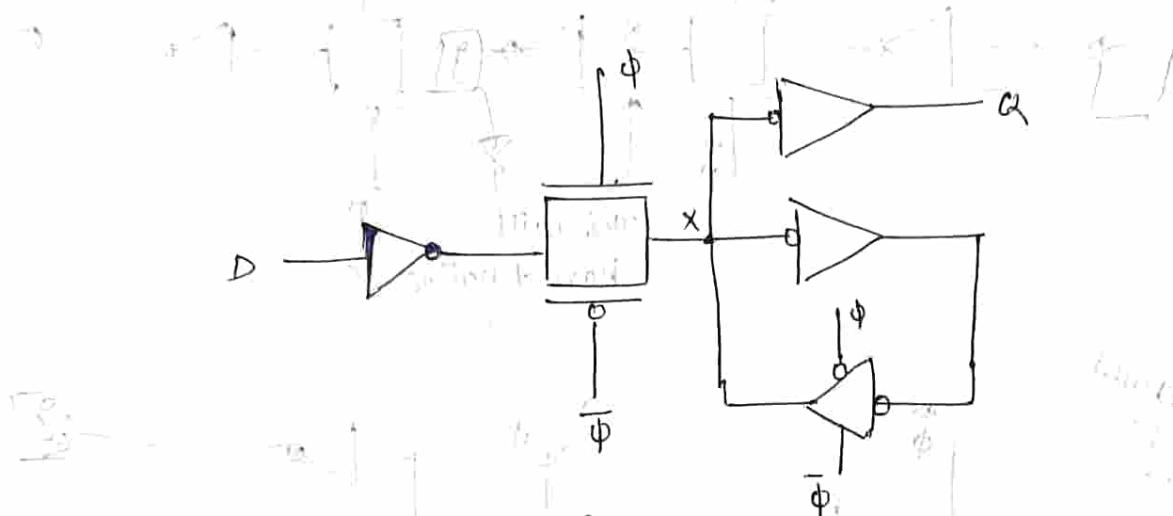


How noise sensitivity has been reduced?



(6) * Modified Buffered I/P Catch with f/B:

- In Buffered I/P Catch, noise on Q can impact the X.
- But in modified CRT noise on X can be reduced by,

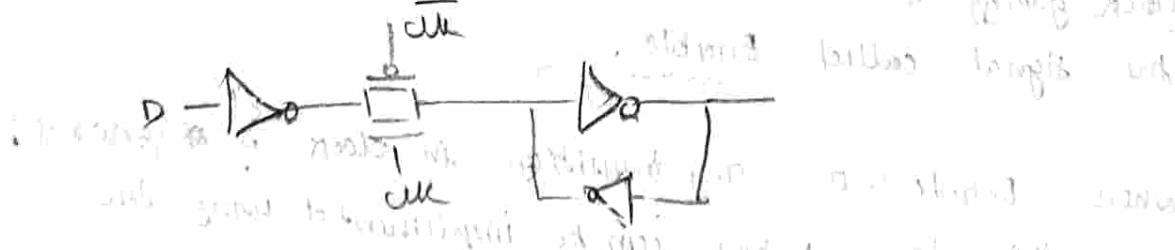


- It is widely used in standard cells.
- * Very robust. (most important)
 - Fanout large
 - Fanout slow (1.5 - 2 for delay)
 - High clock loading

Q: Are we taking $\text{clk}(\Phi)$ & $\text{clk}(\Phi\bar{\Phi})$ as independent signals or clk is o/p of inverter with i/p as clk .

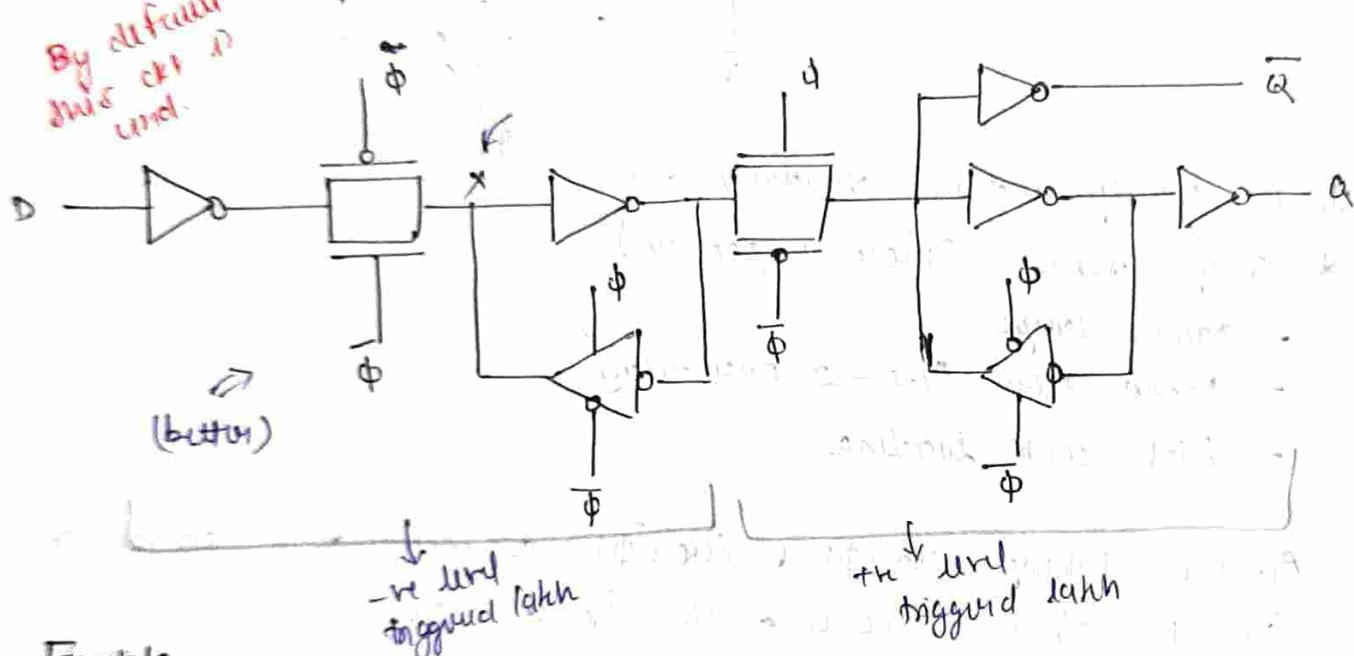
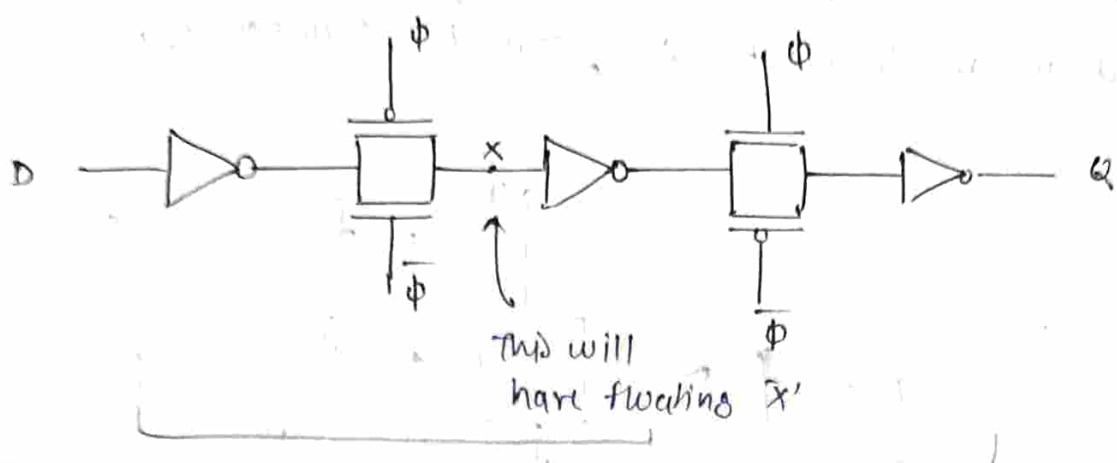
→ we will generate clk from PLL (clock generator) & pass signals clk through the inverter to get clk , where we will see some delay in clk w.r.t clk .

∴ there will always be delay b/w clk & clk



2. Flip-Flop Design : ~~remaining~~

- FF is built as pair of back-to-back latches.



- To reduce dynamic power, we will get some part of the ORT so that there is no clock toggling hence we save toggling power of the signal.
- clock gating is enabled in FF + latches by using the signal called Enable.
- When Enable = 0, any transition in clock is ignored; And this above system can be implemented using the
 - (i) MUX
 - (ii) AND & NAND gates

* When enable becomes zero then ~~data~~ out becomes

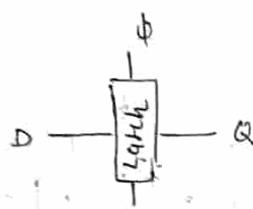
→ independent of CLK.

- MUX : inverter latch D-Q delay ($t_{d-q} \uparrow$)

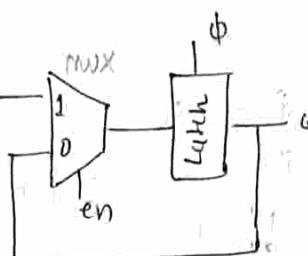
- Clock gating : inverter enables setup time, skew



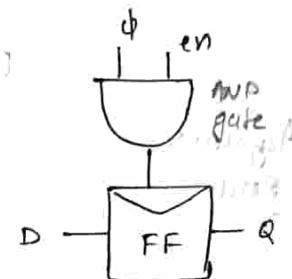
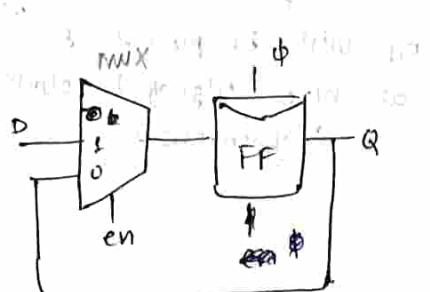
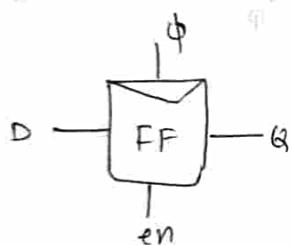
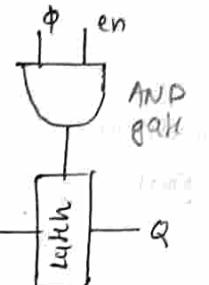
symbol



MUX design



Clock gating design



* Now signal & clock should come, thus to some constraints

(i) constraint b/w D & $\text{CLK}(\phi)$

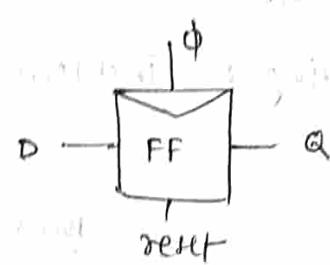
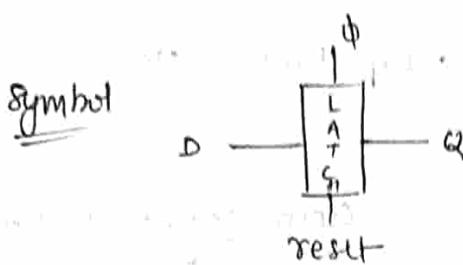
(ii) constraint b/w en & $\text{CLK}(\phi)$

(iii) No constraint b/w D & en.

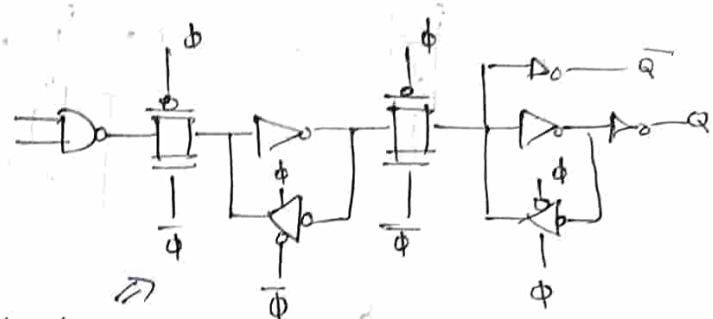
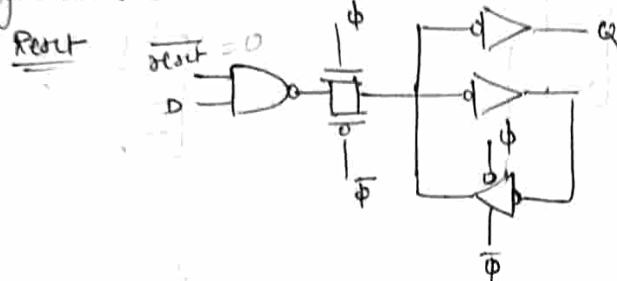
* when $\text{en} = 0$, Q will hold its value, independent of $\text{CLK}(\phi)$.
when $\text{en} = 1$, Q will hold D with some delay like $t_{\text{comb}} + t_{\text{clk-Q}}$.

Reset:

- * When next D arrives then Q is pushed to zero/low.



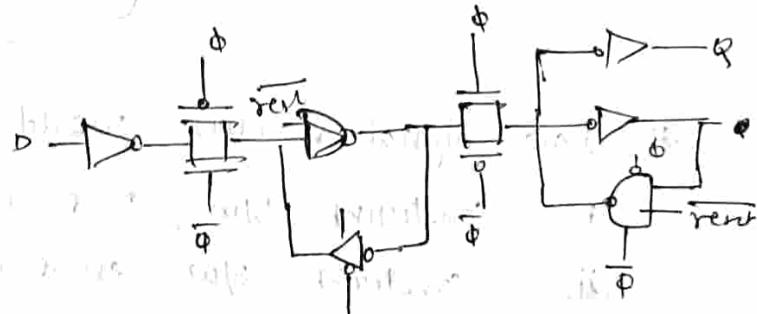
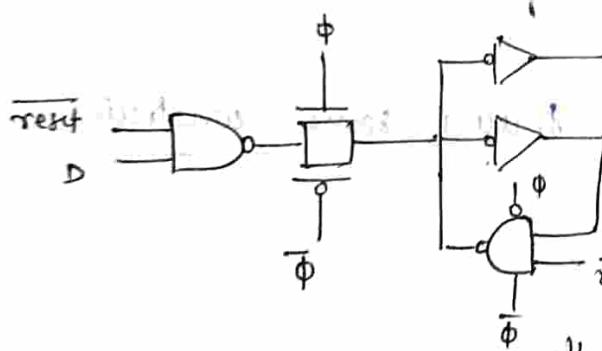
Synchronous



(Q will be pushed to zero
at next edge of the clock)
∴ synchronous

Asynchronous

Reset



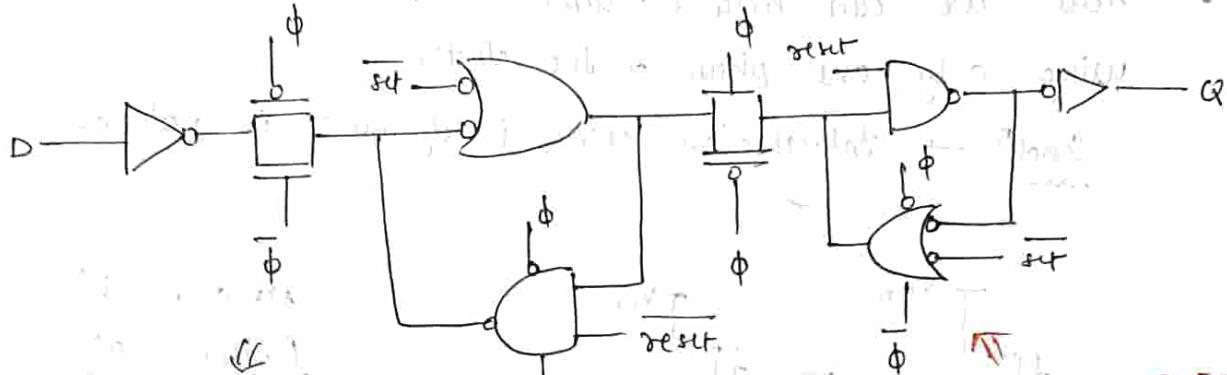
Asynchronous reset.

Q: how they are working & how they are going to get reset?

Working?

Set / Reset: ~~remains~~ ^{initial} state until signal $S/SET = 1$ & $R/RESET = 1$.

- (i) set forces the output to high value when Enable = 1.
- (ii) FF with asynchronous set & reset.



they are generally used for asynchronous machine.

Similar to SR FF in digital ckt.

Q. Why we need set & reset?

→ (i) When power cuts we want to set all the states to '0' ∵ reset is used.

(ii) For memory, some are expected to be at '1' & some are expected to be at '0' ∵ we need both set & reset.

→ As we have seen that, there is possibility of overlap b/w CLK & CK̄ which can lead to back gate driving, then we added Inverters in f/b to avoid back gate driving.

* Also we can use TSPC, we can avoid all the problems occurred in transmission gate latch & FF.

Q. How back driving occurs when, CLK & CK̄ are overlapping?



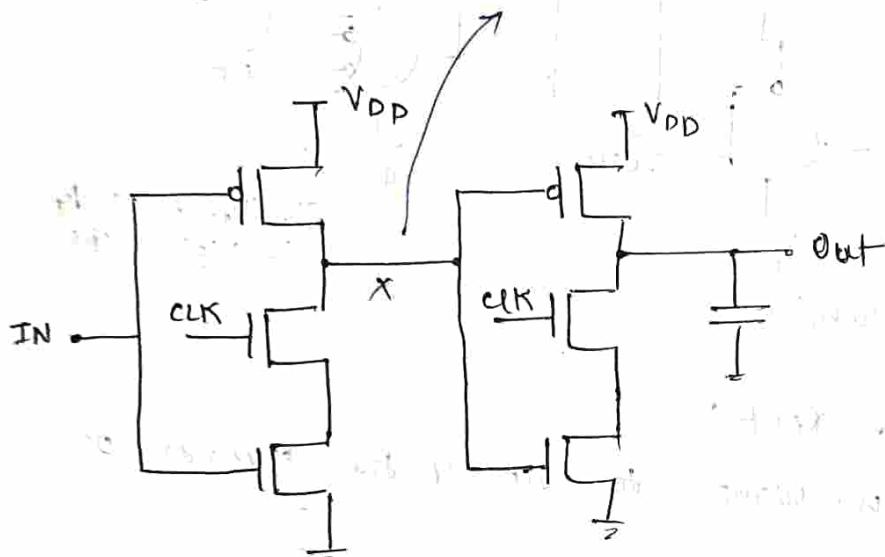
W ₁	X	W ₂	W ₃
0	1	0	0
1	0	1	0
0	0	0	1
1	1	1	1

3. TSPC [True single Phase clock]

~~Sum:~~ Using this FF design we can avoid all the problems like back ~~end~~ receiving o/p.

- Here we can hold the data or make it as a latch by using only one phase of the clock.

Condⁿ → Intermediate node is dynamic in nature.



(Answer for CLK)

Positive latch

(i.e transparent when $CLK = 1$)

for $CLK = 1$,

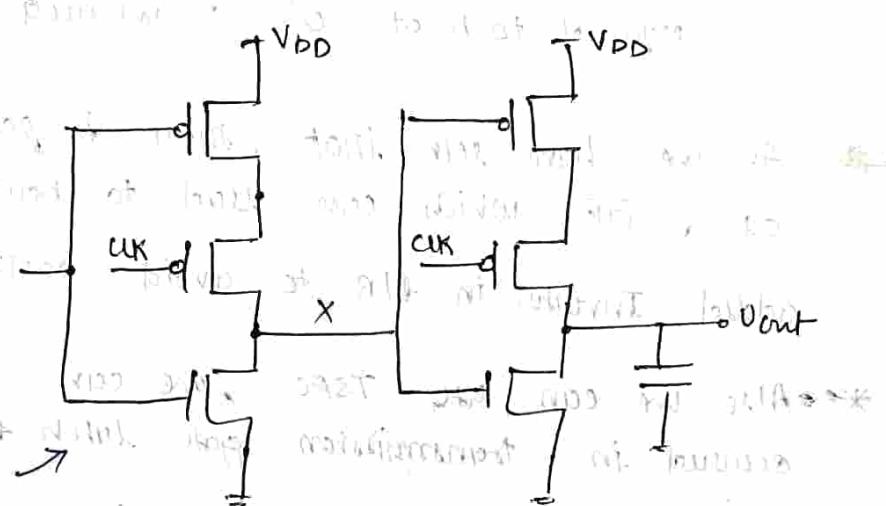
$$X = \overline{IN}$$

$$\therefore OUT = \overline{X} = \overline{IN}$$

~~Problem~~ Node X will be floating

for $CLK = 0$ & $IN = 1$

CLK	IN	X	Out
0	0	1	z
0	1	z	?
1	0	1	0
1	1	0	1



~~Problem~~ Node X will be floating for $CLK = 1$ & $IN = 0$.

Negative latch (Answer for CLK)

(i.e transparent when $CLK = 0$)

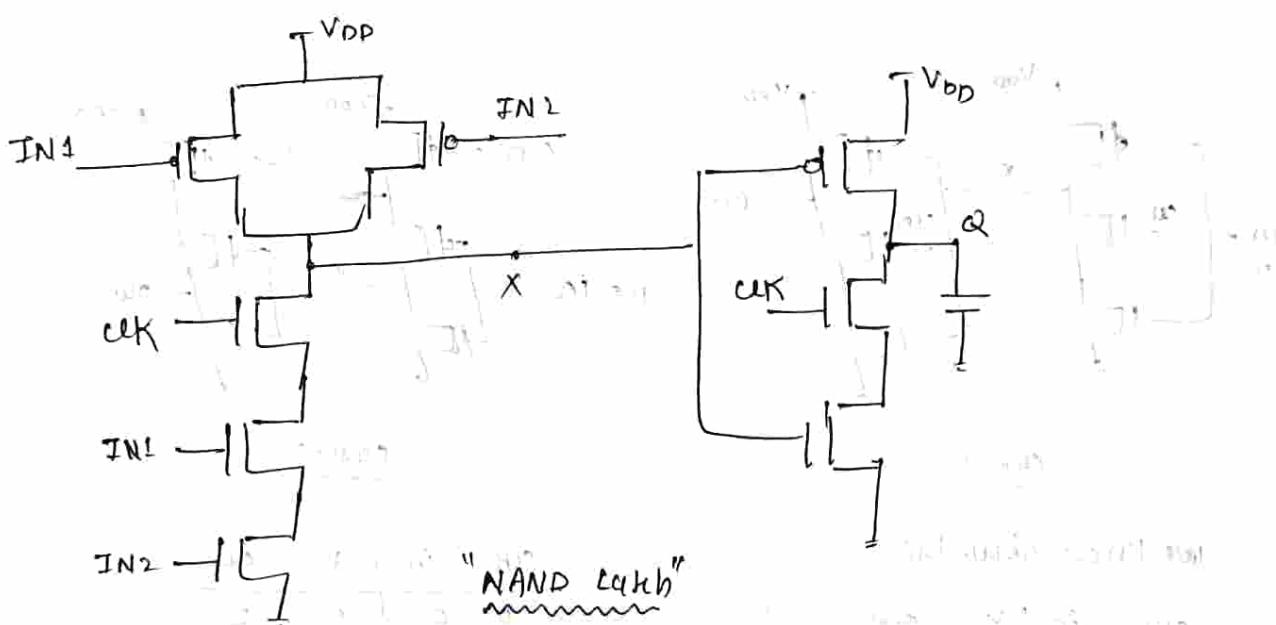
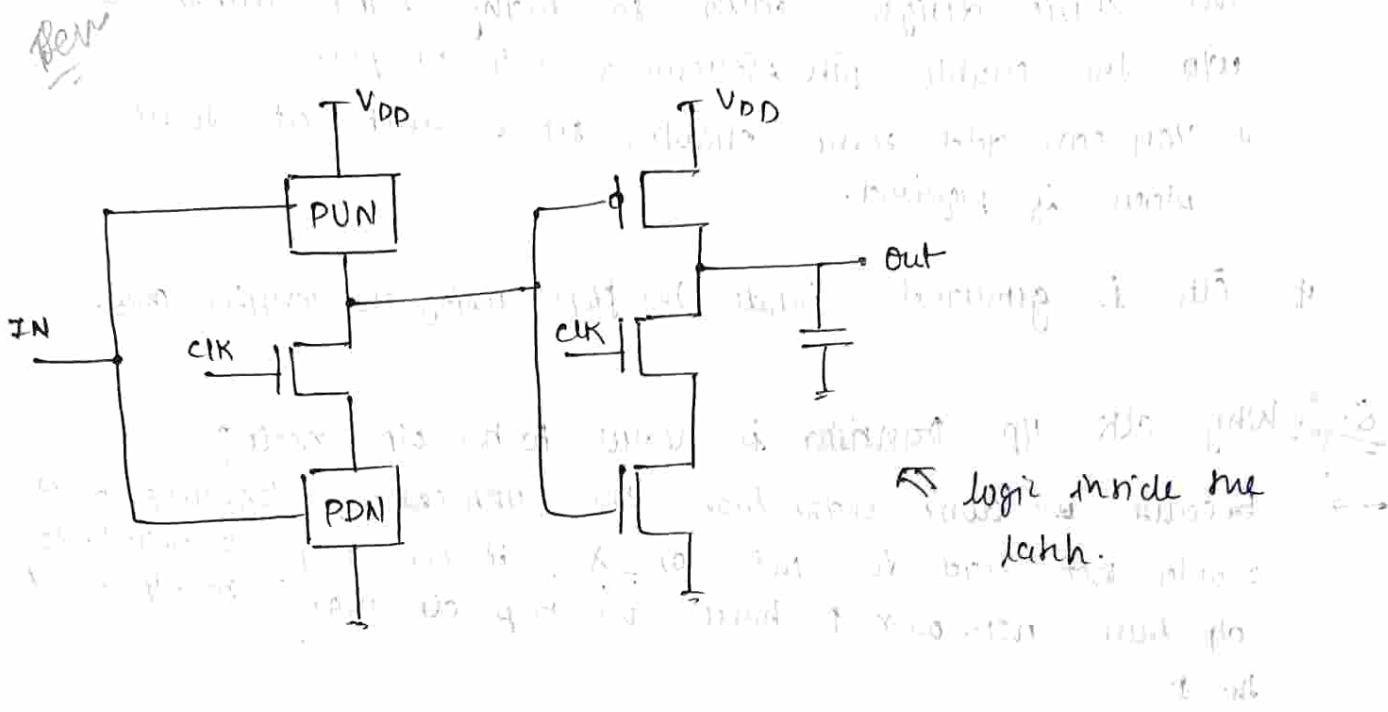
CLK	IN	X	Out
0	0	1	0
0	1	0	1
1	0	1	?
1	1	0	?

for $CLK = 0$,

$$X = \overline{IN}$$

$$\therefore OUT = \overline{X} = \overline{IN}$$

Including logic in TSPC

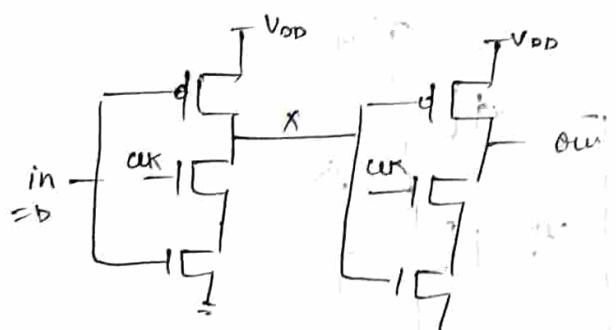


when, $IN_1 + IN_2 = 1$ & $CLK = 0$ then node 'X' will be floating.

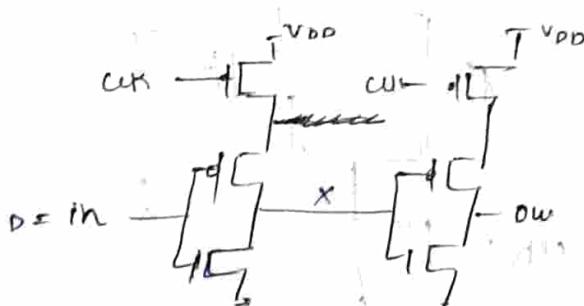
CLK	IN1	IN2	X	Out (output of NAND gate)
0	1	1	Z	?
0	1/0	0/1	L	0
0	0	0	1	0
1	0	1	1	0
1	1	0	1	0
1	1	1	0	1

- Q. By default what design do we use?
- The static design with so many extra transistors who see enable pin whenever & who set pin.
 - * You can add some enable, set & reset at some places if required.
 - * CLK is generated inside the flop using the inverter only.

- Q. Why CLK input transistor is nearer to the OLP node?
- Because we don't want that the unnecessary toggling of D should not lead to noise on X, if we keep D nearer to OLP then noise on X ↑ hence we keep CLK nearer to OLP & not the D.



eqn-1



eqn-2

We know about this:

clk	in	X	out
0	0	1	z
0	1	z	?
1	0	1	0
1	1	0	1

when CLK = 1,

out = in

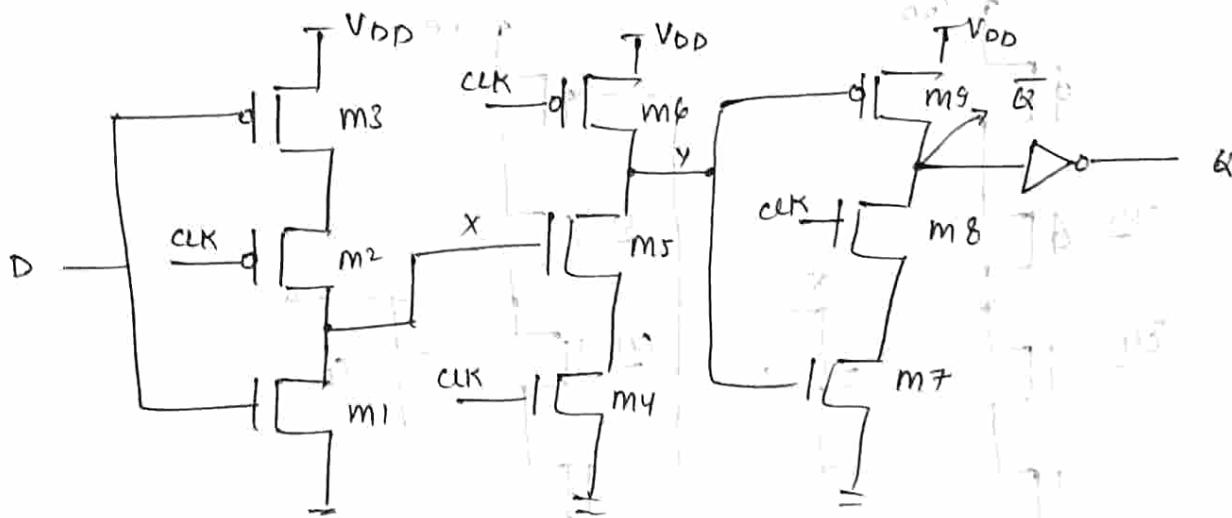
∴ latch with
and triggered
clk.

CLK	in	X	out
0	0	1	z
0	1	0	1
1	0	z	?
1	1	0	z

This is not working as
any transparent CLK.
(∴ this is wrong design)

* keeper can be added to make
circuit pseudo static.

TSPC - Flip Flop :

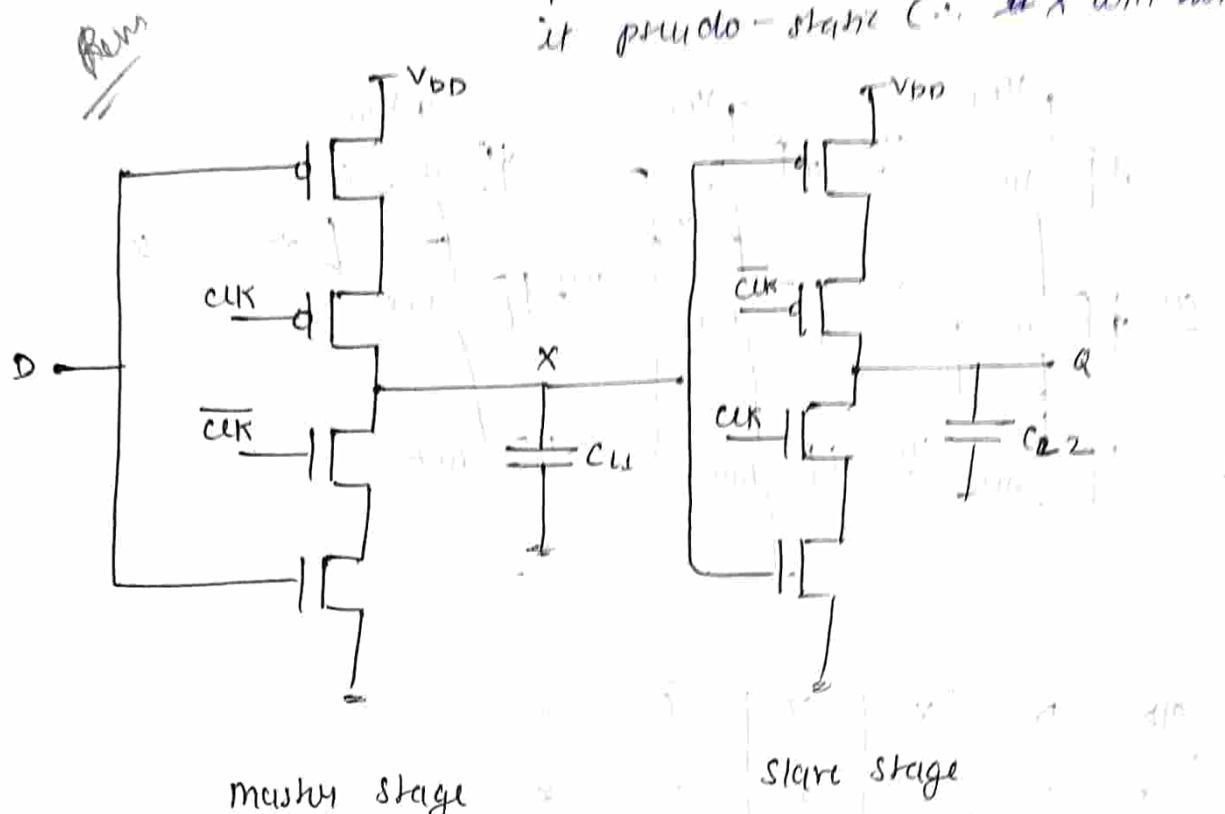


CLK	D	X	Y	\bar{Q}	Q
0	0	1	1	Z	?
0	1	0	1	Z	?
1	0	Z	?	?	?
1	1	0	Z		

Q: Can we put CLK upside? near V_{DD}. (clk on top, comn)
 since D is FF \Rightarrow edge triggered \therefore only one edge of CLK D will
 impact the node X.
 so if D is clear to the X then even when edge is
 not triggered then also it can impact the op. w/o triggering
 of the CLK.

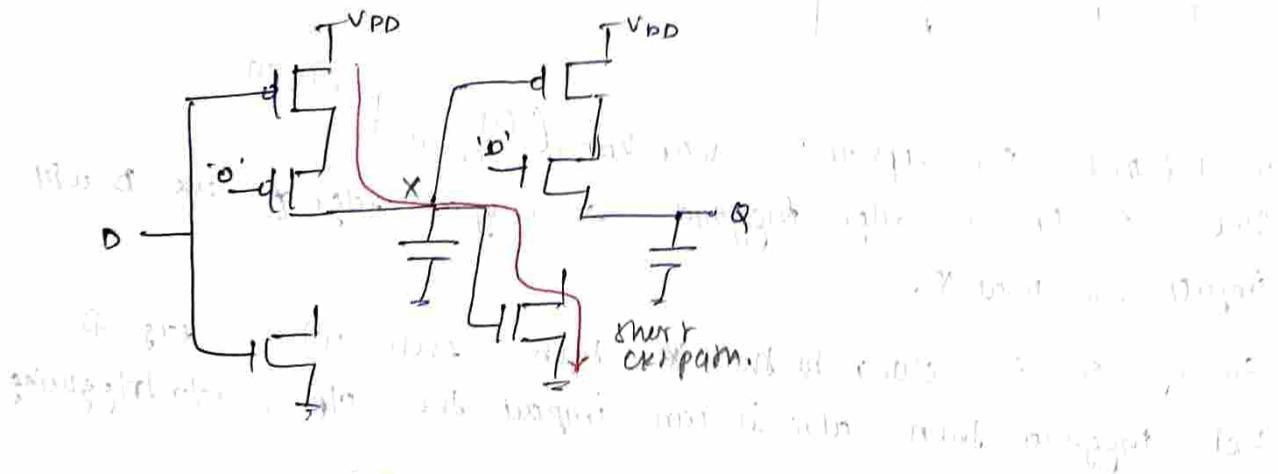
4. C²MOS

Keepers can be added to the CLK to make it pseudo-static (i.e. X will not float).

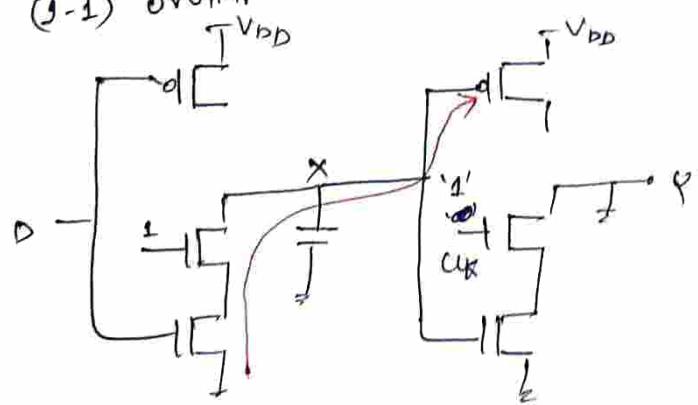


Transitivity of clock overlap:

- (i) (0-0) overlap



- (ii) (1-1) overlap



Sequential Circuit - Timing

Sampling Edge: The edge of the clock at which the data gets sampled is called the sampling edge.

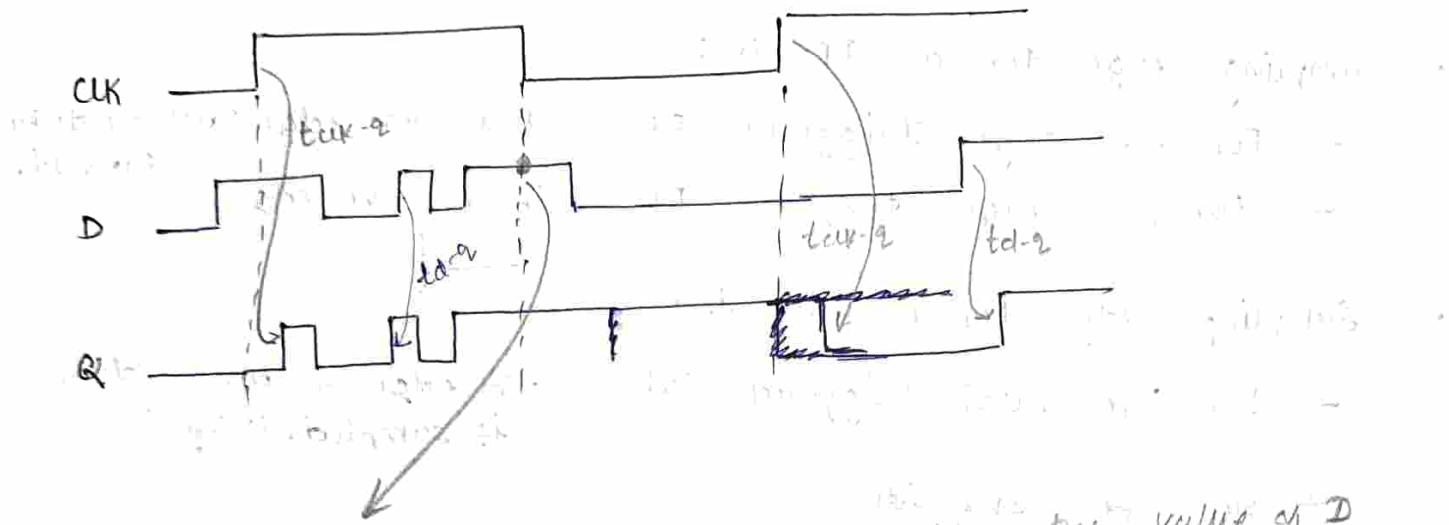
- Sampling edge for a FF is:
 - For +ve edge triggered FF, it is the edge where data is sampled.
 - For -ve edge triggered FF, it is the edge
- Sampling edge for a Latch is:
 - For +ve level triggered latch, -ve edge is where data is sampled. (why?)

 - because it starts to become opaque after the level.
 - For -ve level triggered latch, the edge is where data is sampled.
~~because it becomes opaque after the level~~

Note: As we already know that there are some delays b/w I/P & O/P called propagation or contamination delay.

- But FF & Latch have some additional feature called sampling edge & with respect to this there are delays associated to it.
- * Delays w.r.t. I/P (D) & O/P (Q) are:
 - i) propagation
 - ii) contamination
- * Timings associated with sampling edge are:
 - i) hold time
 - ii) setup time

- Q. Why the level latch sample at the edge of clock?
 → Because it became opaque after the level of the clock ie. at the edge of one clock period.



At opaque duration it is sampling the value of D which it has obtained during the edge of CLK .

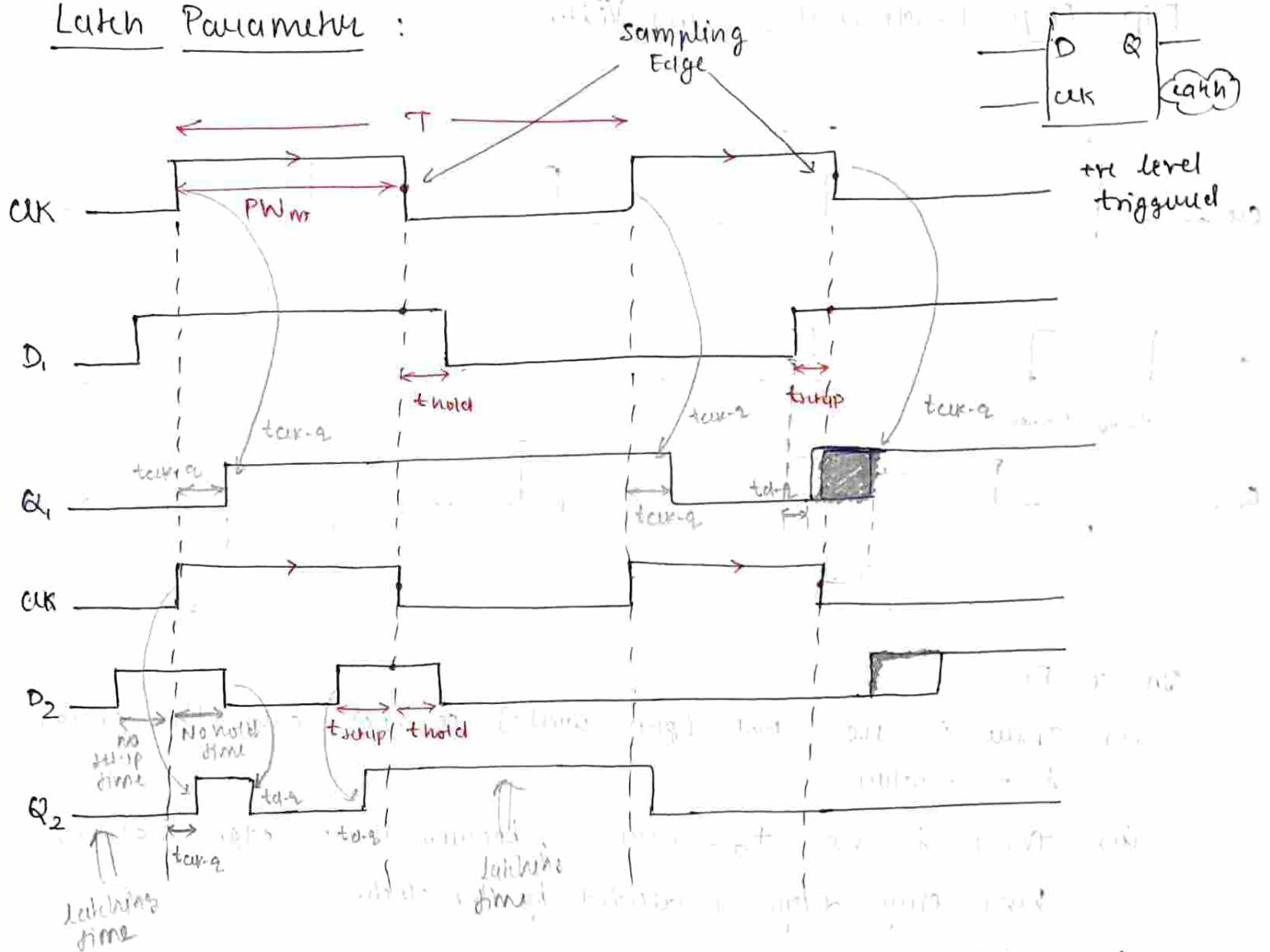
latching → fasten (a door or gate) with a latch.

- * We will only consider the setup & hold time at the falling edge of the clock for a level triggered latch.
- * we are considering that edge from where latching action starts that is the sampling edge.

Setup time :- It is the time before the sampling edge for which the input should be stable.
 (t_{setup}) [it is b/w two different sampling edge]

Hold time :- It is the time after the sampling edge for which the input should be stable.
 (t_{hold}) [it is b/w same sampling edge]

Latch Parameter :



* setup time (t_{setup}) & hold time (t_{hold}) are defined from the sampling edge of the clock.

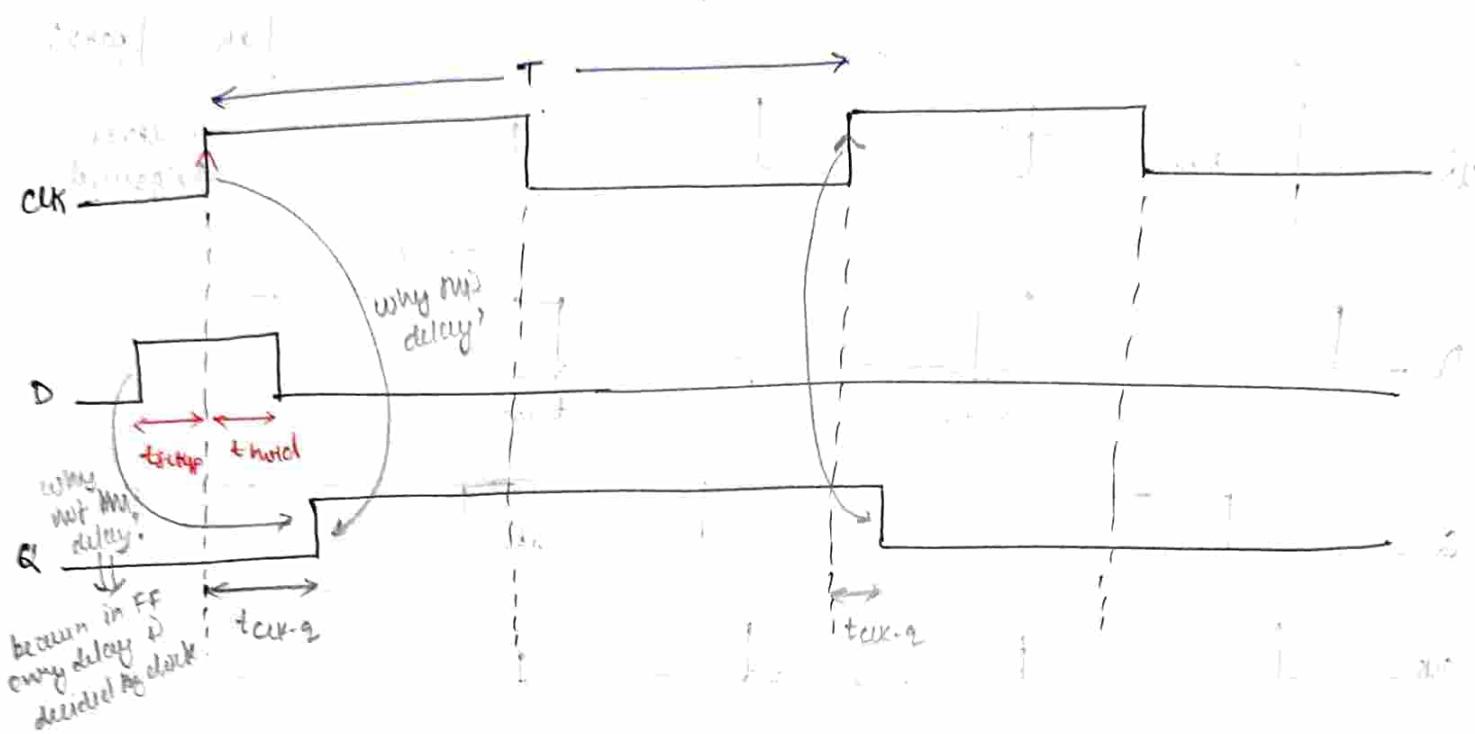
* We will only consider the falling edge of the the level triggered for calculating t_{hold} & t_{hold} .

Double edge triggered FF \rightarrow they are both edge of the clock for sampling so they are essentially doing half the path period for evaluation.

* Pulled Latch? \rightarrow A ~~square~~ pulse generator that behave as FF rather than a clock. It gets a pulse and holds it.

Note: * For -re level of the clock , the value of Q_1 will be ~~locked~~ to the value that it has sampled at the sampling edge.

Flip - Flop Parameters : "Imp. Video" (Video - 68)

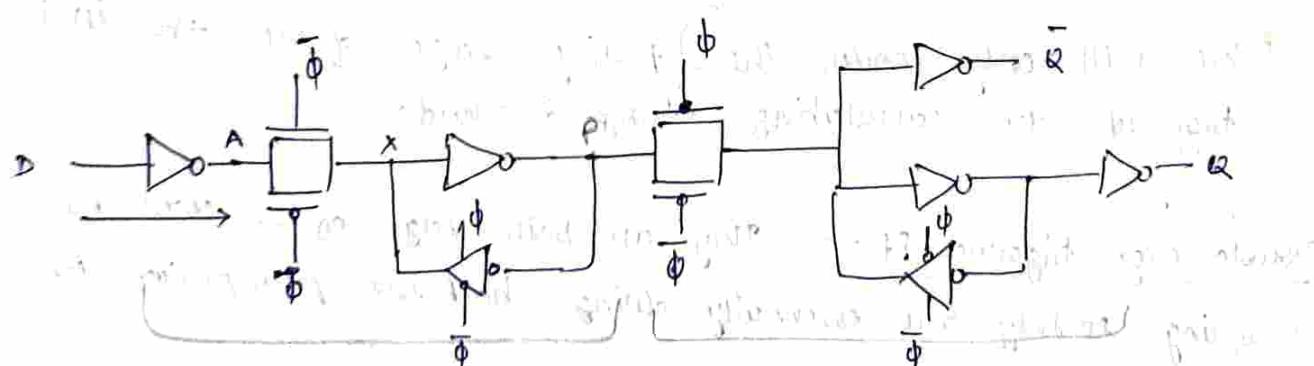


* In a FF :-

- there is no PW (pulse width) concept*, only time period that matters
- there is no t_{d-q} delay, because it is edge triggered, hence every delay is decided by the clock.

Q: Why we are taking only $clk-q$ delay & not $d-q$ delay?

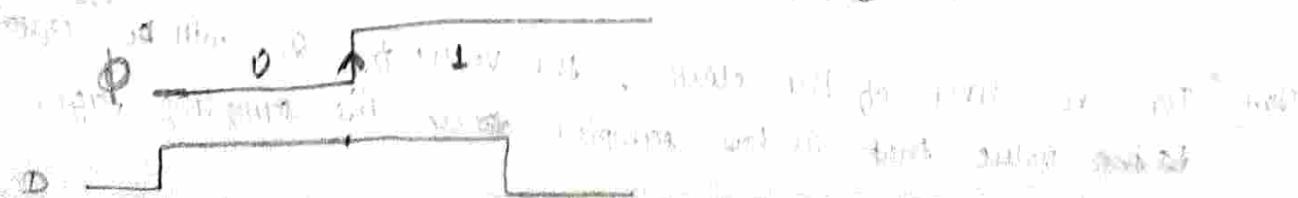
→ Let's take an example of FF,



when $\phi = 0$, min D can only go up to 4 X.P.

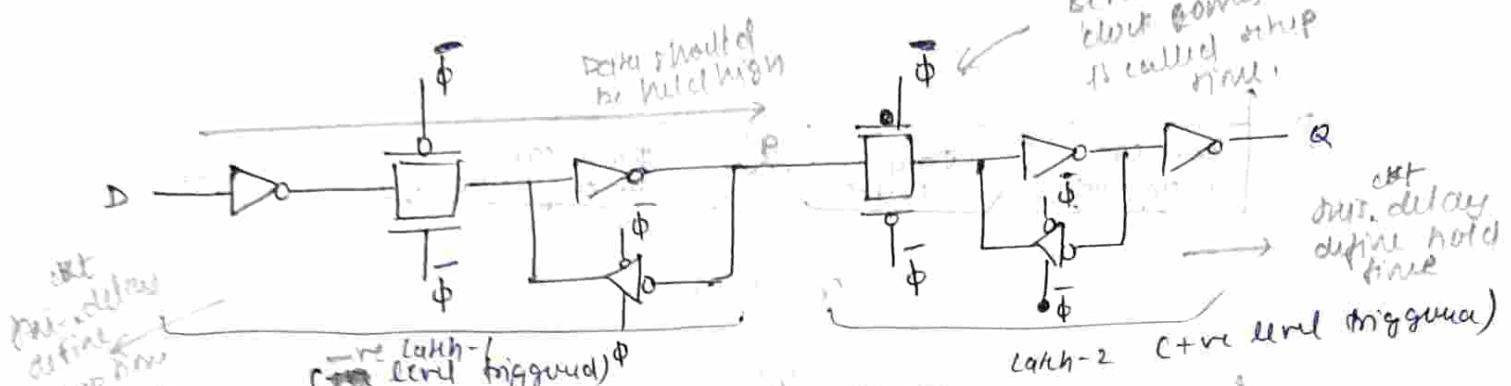
when $\phi = 1$, min ϕ can go to Q hence what is triggering one?

→ ϕ is clk is triggering one Q hence clk delay is considered.



Q Can your clk_2 delay be less than hold time?
 good question. can $\text{clk}_2 < \text{t}_{\text{hold}}$
 → It depends on the data path.

* also see how setup + hold time is defined.



* for the following edge of the clk_2 1st latch becomes

- For the falling edge of the clk_2 , the transmission gate becomes transparent and D is written into the 1st latch.

- when $\text{clk}_2 = 0$, latch-1 becomes opaque + latch-2 becomes transparent (so it can hold the sampled data from latch-1) so the data from latch-2 comes to Q .

* * D should be hold high until & unless my first point is totally opaque is called the setup time.

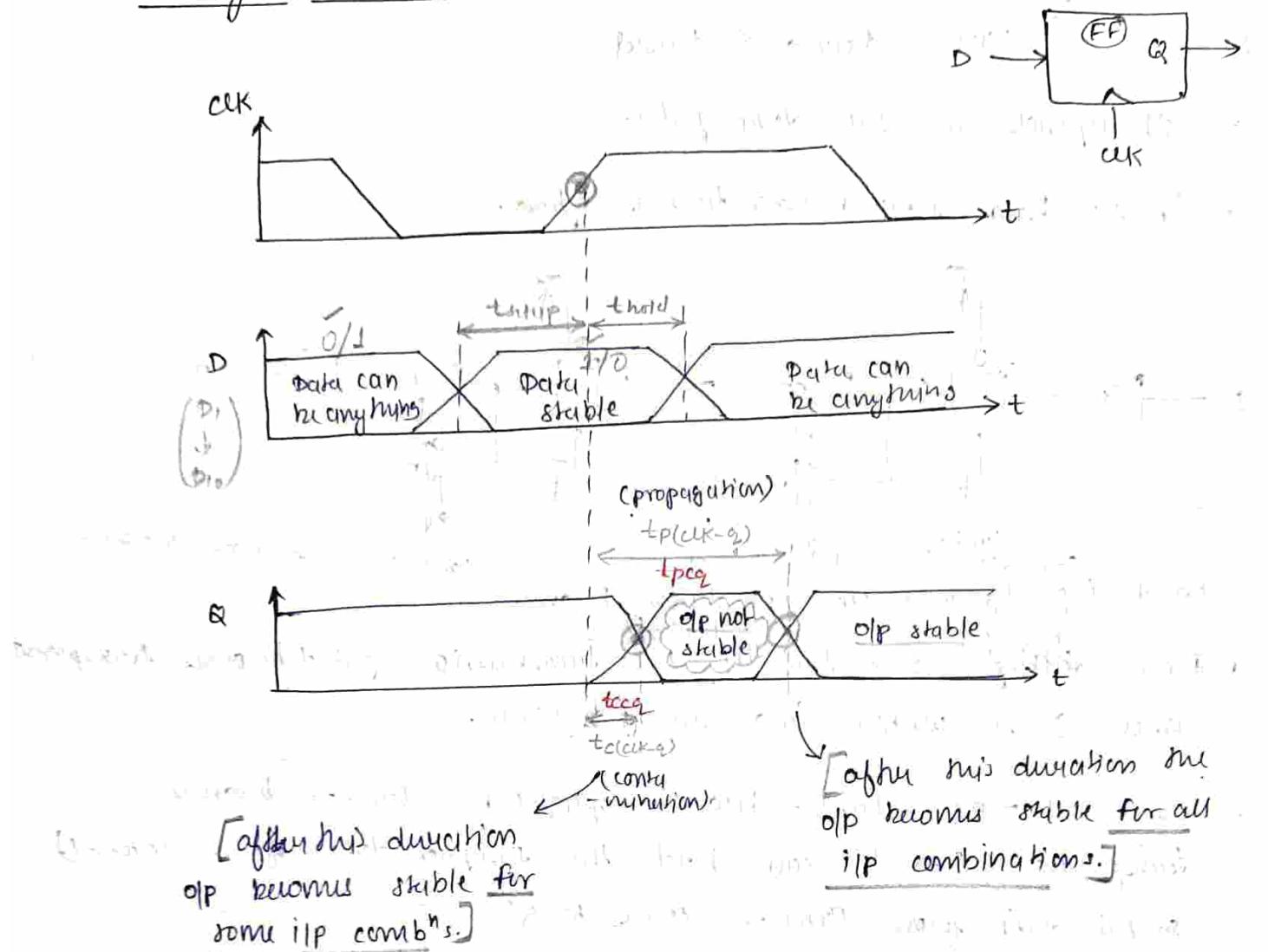
* the clk_2 delay will totally depend on the second part of FF.

* You define setup + hold time based on what is the impact of bringing D closer to clock edge on clk_2 delay.

Q. now setup + hold time are related & different from propagation delay?

* stick to the definition of setup + hold time & according to that look at the circuit now there can many delays that can be added to the clk setup + hold timings.

Timing Definitions:



Ex. For 8-bit Register,
 more are 8-FF, me 8-bits
 o/p will start to vary after
 time of $t_{\text{clk-Q}}(t_{\text{cq}})$ but only
 after t_{pcq} all the bits have
 been toggled and no further
 changes are observed.

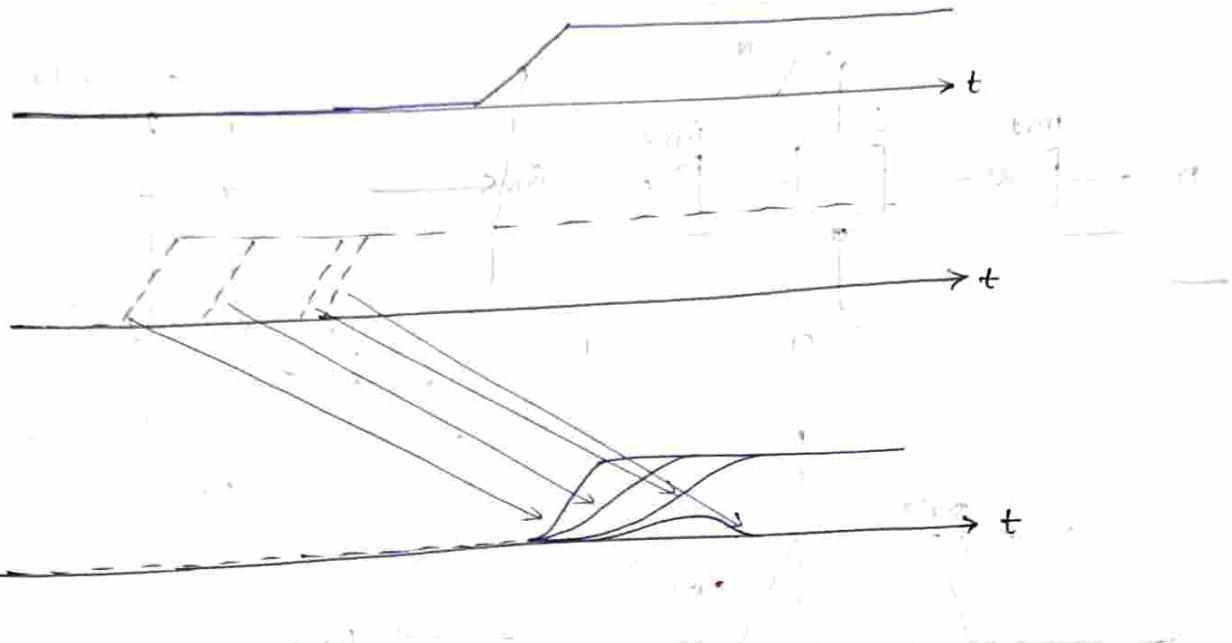
After continuation delay
 \rightarrow o/p starts to vary (change)
 After propagation delay
 \rightarrow o/p settles to whatever
 changed value it was
 running.

Characterizing Setup / Hold Time:

- As I bring my data close to the rising-edge then the o/p delay or clk-Q delay starts increasing because ↑ in the o/p slope.
- Q: Why does sup o/p slope increase? \rightarrow bringing D closer to clk edge?

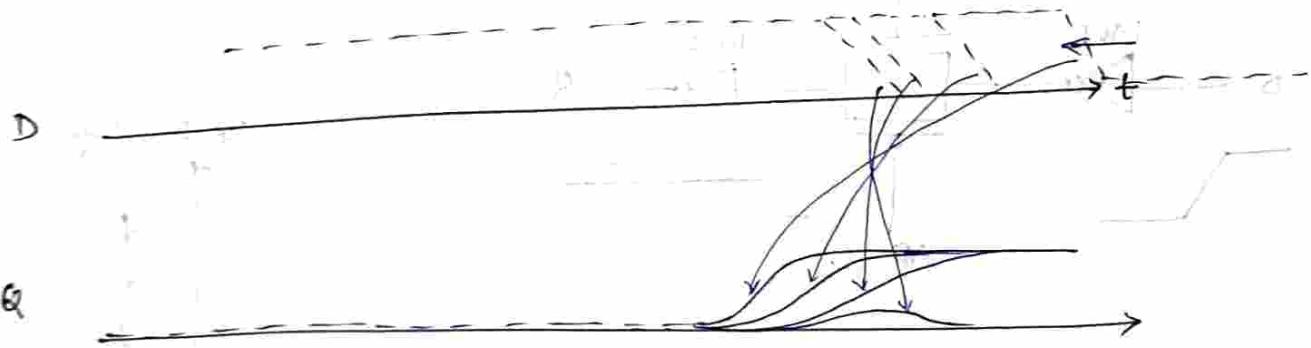
for setup time

CLK



for hold time

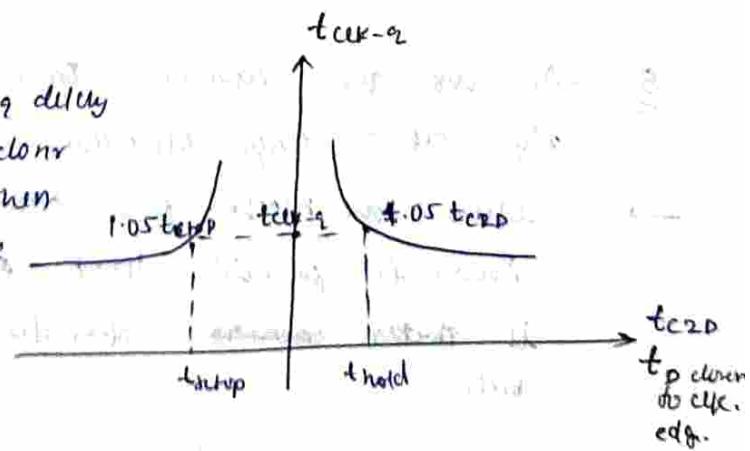
CLK



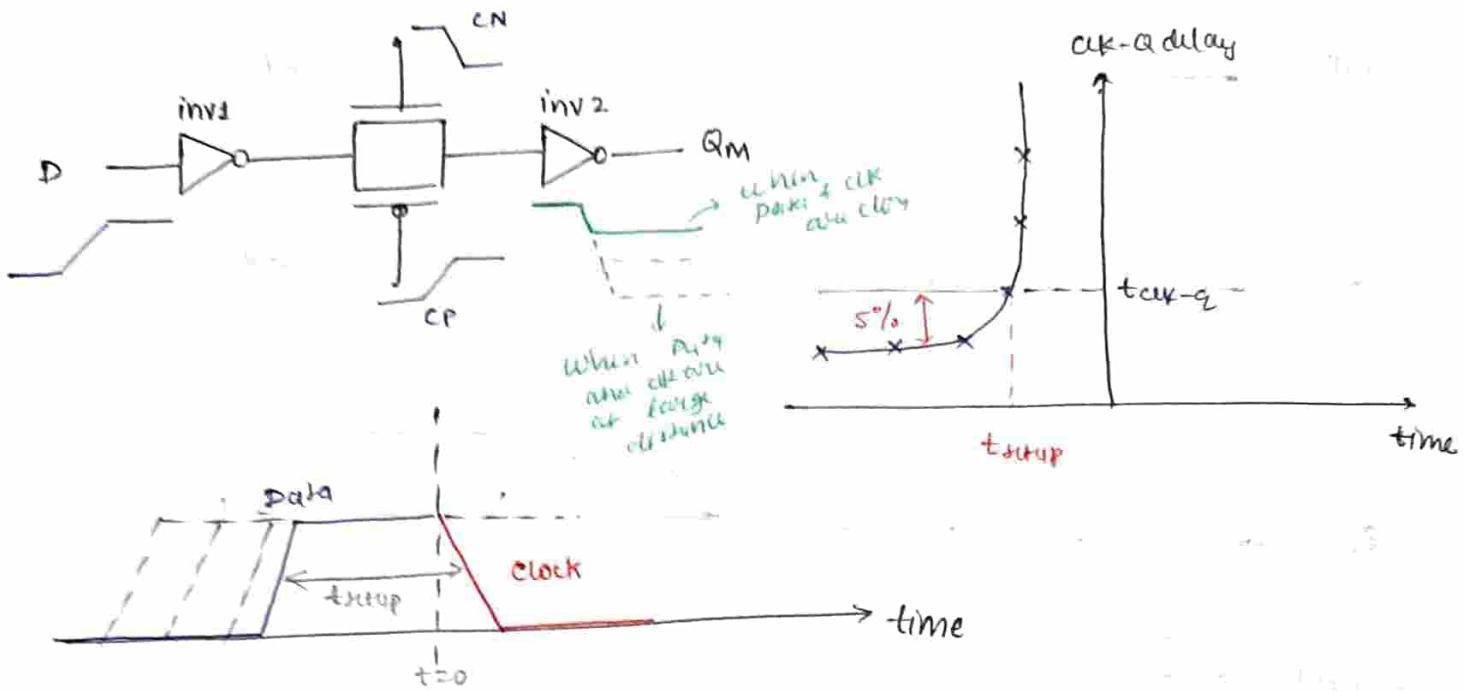
- * As the edge of data comes closer to the clk-edge the clk-q delay starts increasing:

— This is valid for both setup time & hold time.

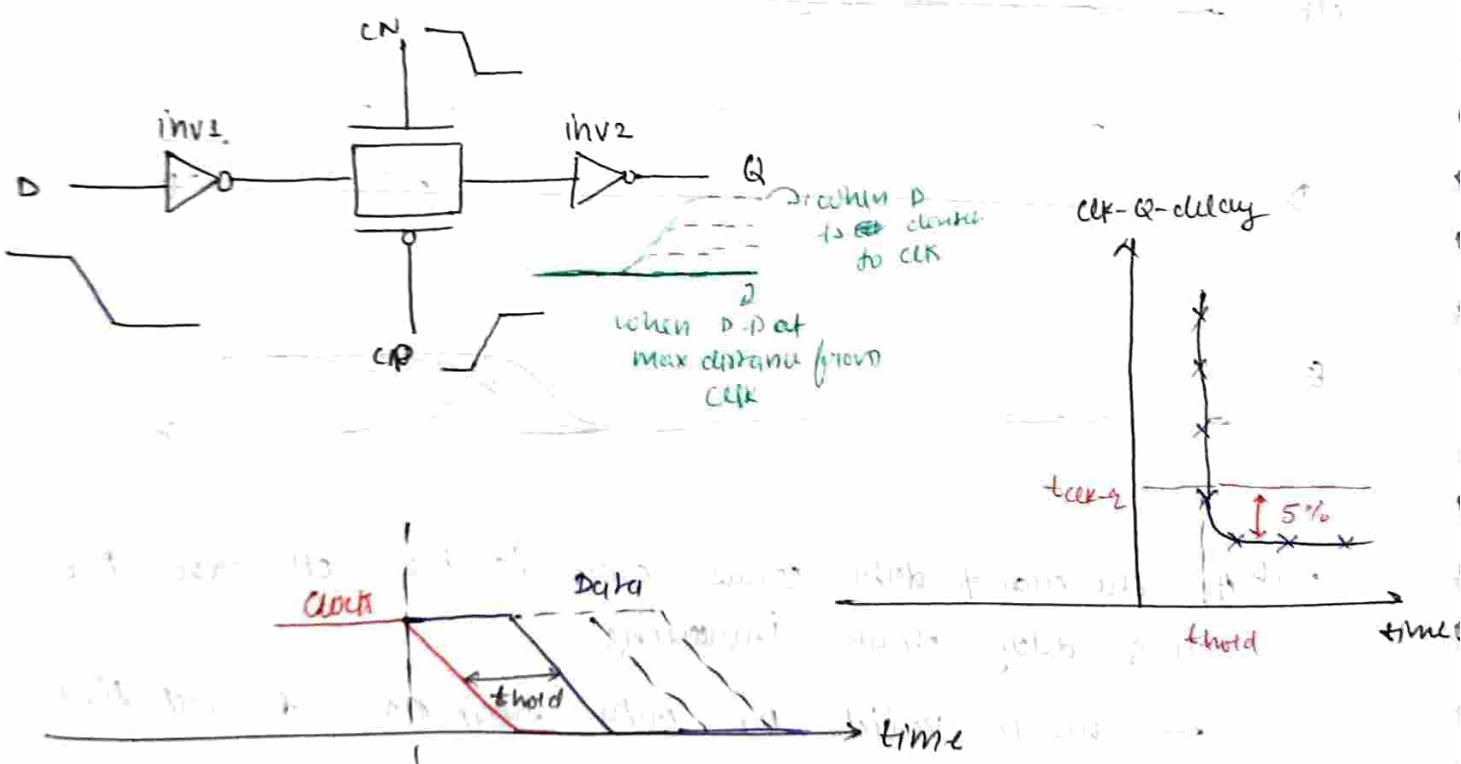
- * when edge of D is far away from edge of CLK then there is some t_{clk-q} delay which is less. Then on bringing D closer to edge of CLK we see ↑ in t_{clk-q} when t_{clk-q} becomes 1.05 times its max value then at that time diff. b/w edge of D to edge of CLK is called tramp.



Setup-time illustration



Hold-time illustration

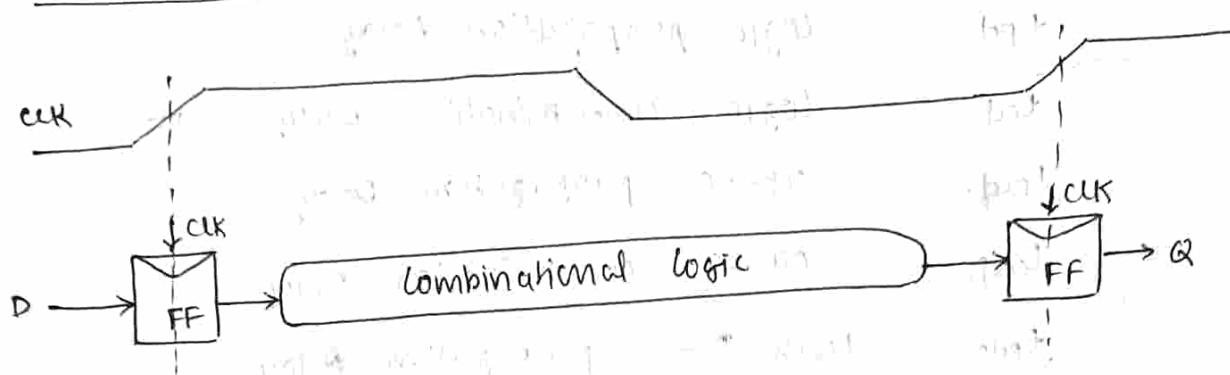


- Q. As we are bringing data edge closer to the clock edge then why $ck-Q$ delay increases?
- when we make our data to be stable for less time then there is possibility that some disturbance may go to & hence it ~~degrades~~ degrades the slope of the transition on second side.

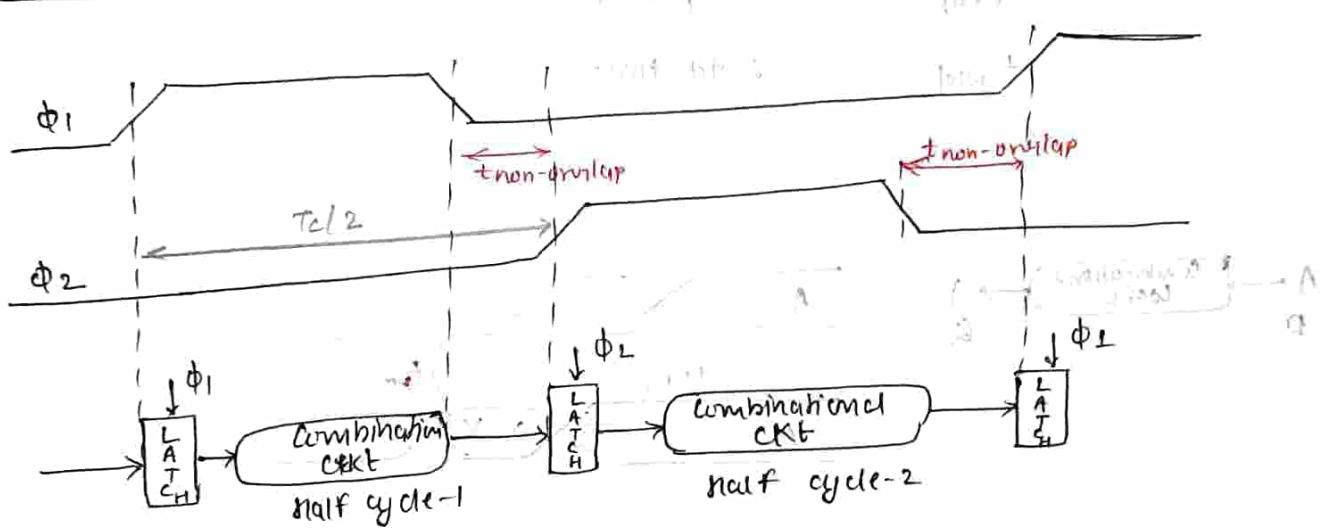
Sequencing Methods :

multiple output

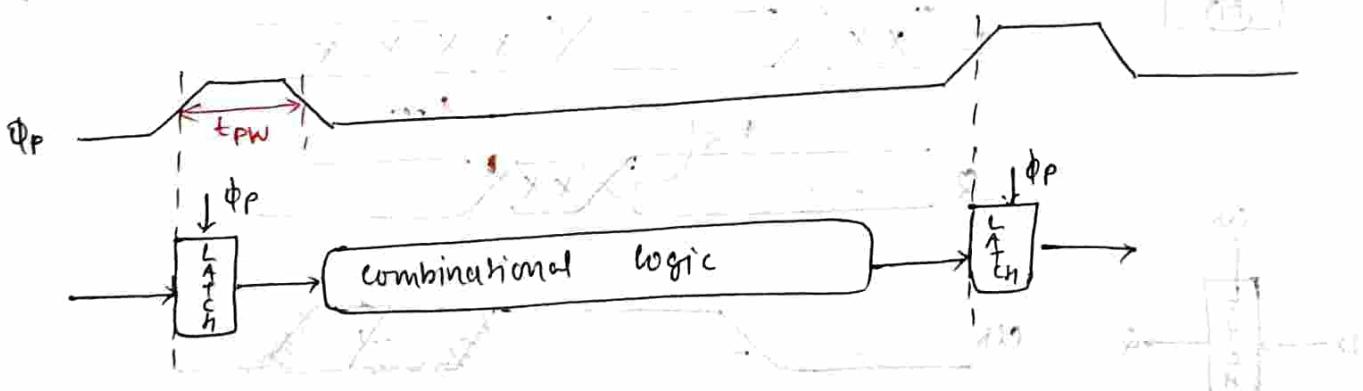
1. Flip - Flops



2. Two - phase Transparent Latch

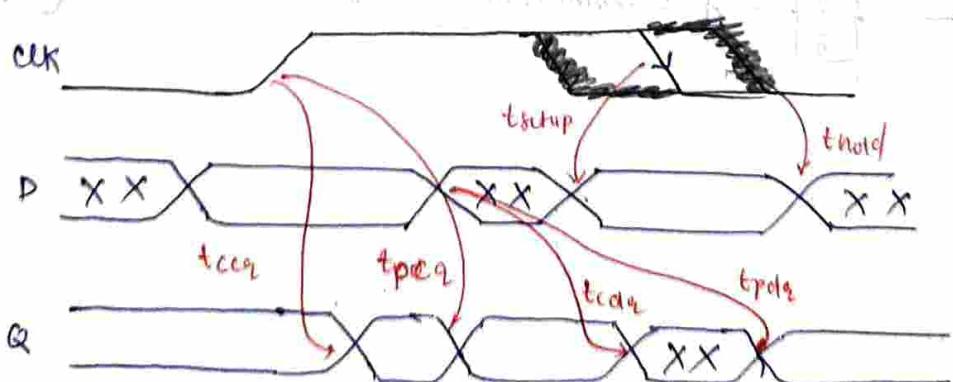
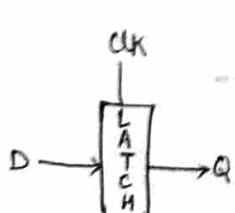
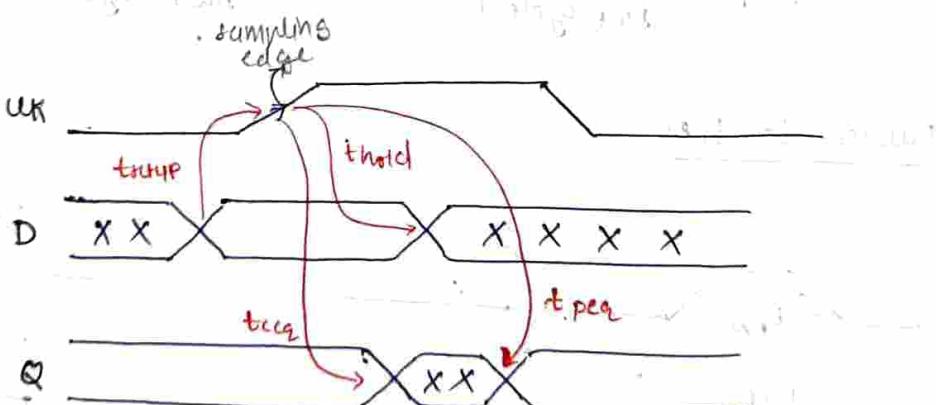
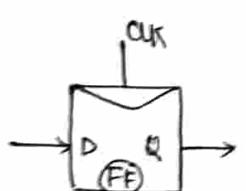
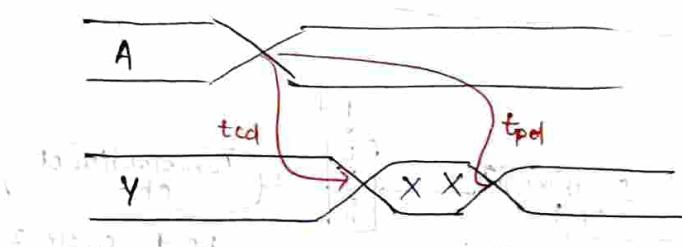
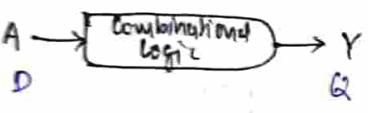


3. Pulsed Latches



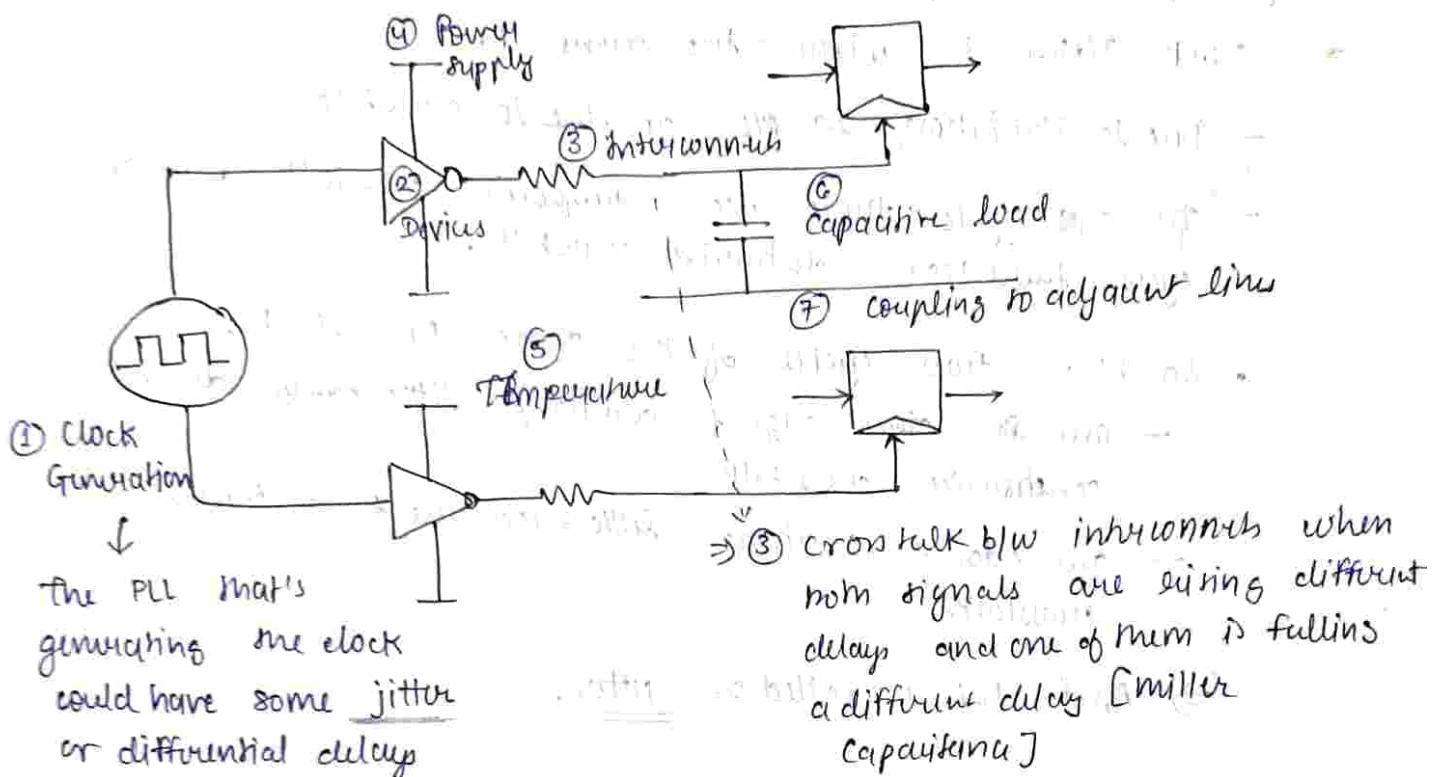
Timing Diagram

t_{pd}	logic propagation Delay
t_{cd}	logic contamination Delay
t_{pcd}	CLK \rightarrow Q propagation Delay
t_{ccd}	CLK \rightarrow Q contamination Delay
t_{pdq}	Latch D \rightarrow Q propagation Delay
t_{cdq}	Latch D \rightarrow Q contamination Delay
t_{trup}	Setup time
t_{hold}	Hold time



Timing Constraints

Clock Uncertainties: All of these are sources of uncertainty.



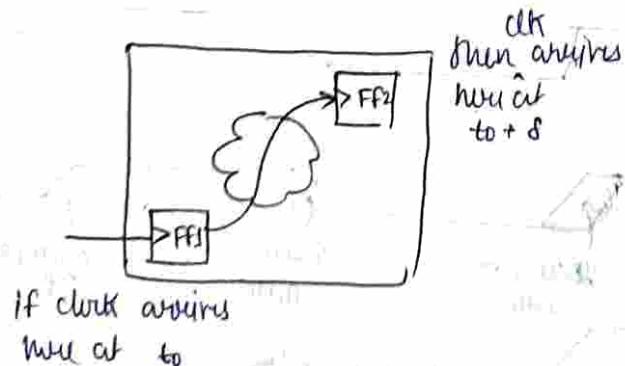
- * One clock may experience delay via one path & other clk experiencing delay through another path.
 - This is called clock uncertainty.

Clock Non-idealities:

- Skew is because of the routing for different clk paths.
- Jitter is because of temporal variations.

1. Clock Skew:

- Spatial variations in temporally equivalent clock edges;
[deterministic + random = tsk]



* clock skew is b/w 2 different ff's

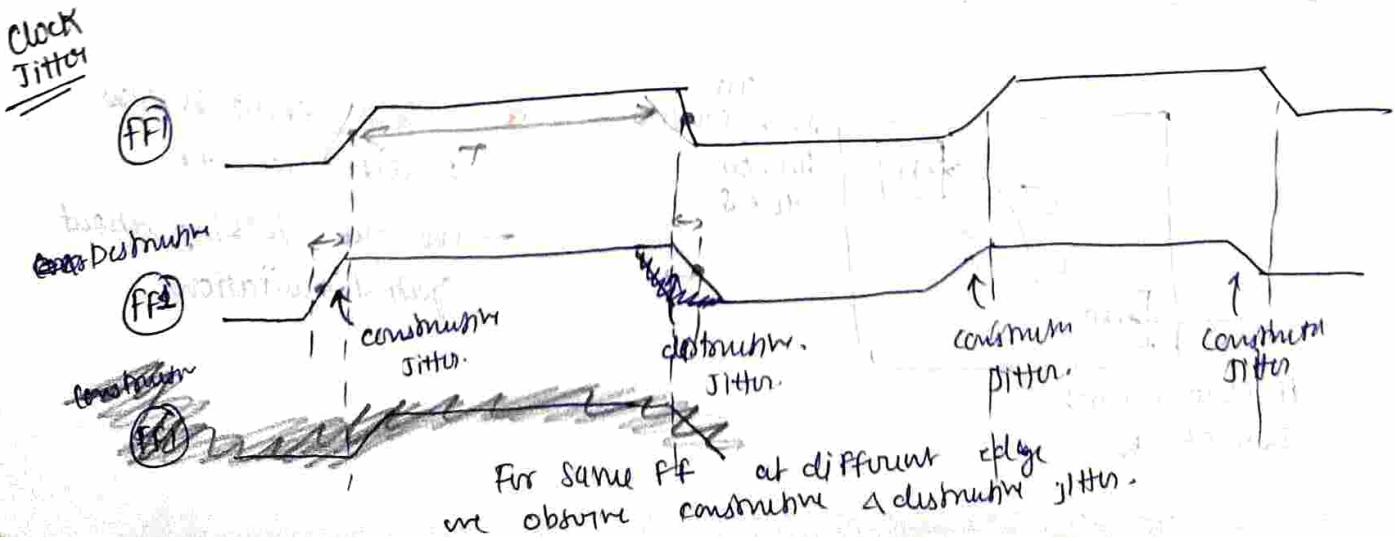
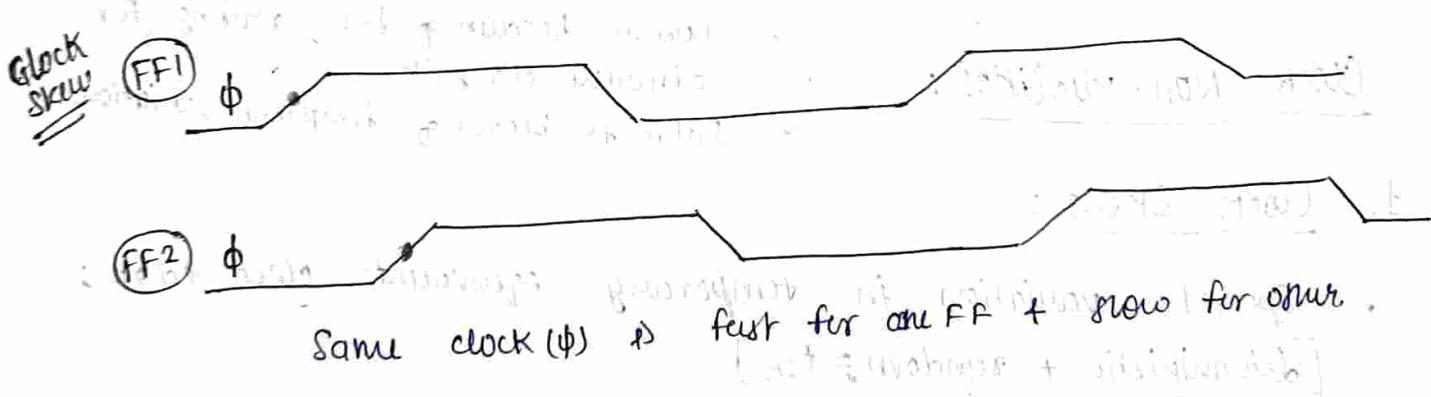
— we are talking about spatial variations

2. Clock Jitter:

- temporal variation in consecutive edges of the clock signal;
[modulation + random noise]
- * clock jitter \Rightarrow within the same FF:
 - Due to variation in PLL or due to crosstalk.
 - In one cycle there was supporting crosstalk & in one cycle there was destructive crosstalk.
 - so b/w two cycles of the same FF itself,
 - one \Rightarrow ~~clock~~ edge \Rightarrow arriving little early due to constructive crosstalk.
 - one edge \Rightarrow arriving little late due to destructive crosstalk.
- \Rightarrow And this is called as jitter.

3. Variation of pulse width:

- important for level sensitive clocking.

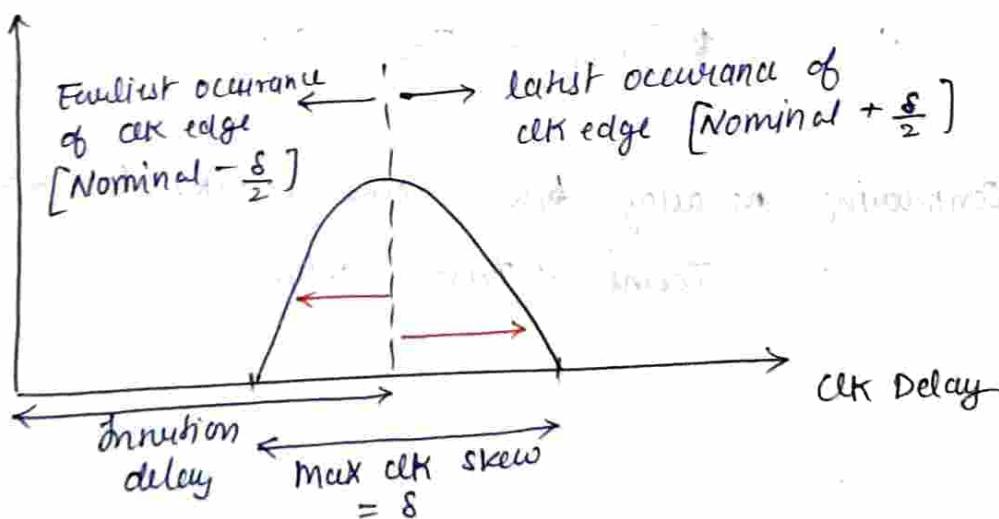


- skew is because of fixed + temporal variations
jitter is only because of temporal variation.

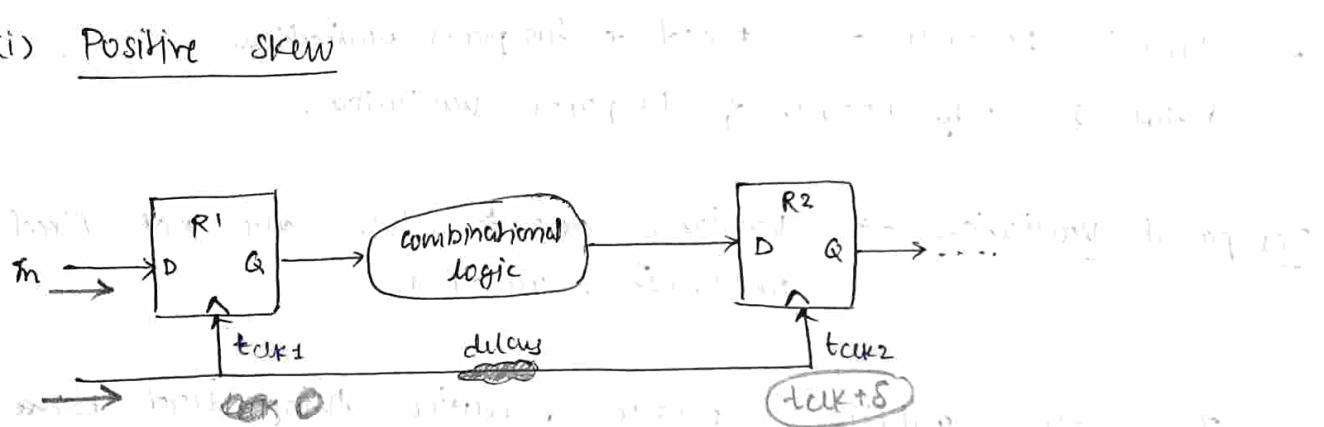
Temporal variations → Variation due to which are not fixed
unpredictable, random.

- Temperature variation, source & device during ~~lead to~~ ^(freq) both jitter & skew ~~versus~~ clock.
 - skew clock is also due to routing (ie. for different path)
 - skew clock is deterministic & jitter is ~~is~~ ^{not} deterministic.
 - therefore for designing the CRT we will ~~look to~~ leave the margin.
- (Q) A clk as input to the FF then there is some delay b/w the normal clk & inverted clk. Are these the clock skew?
- that delay will come in setup & hold time, because inverted clock is generated inside the FF itself.

Clock Skew:

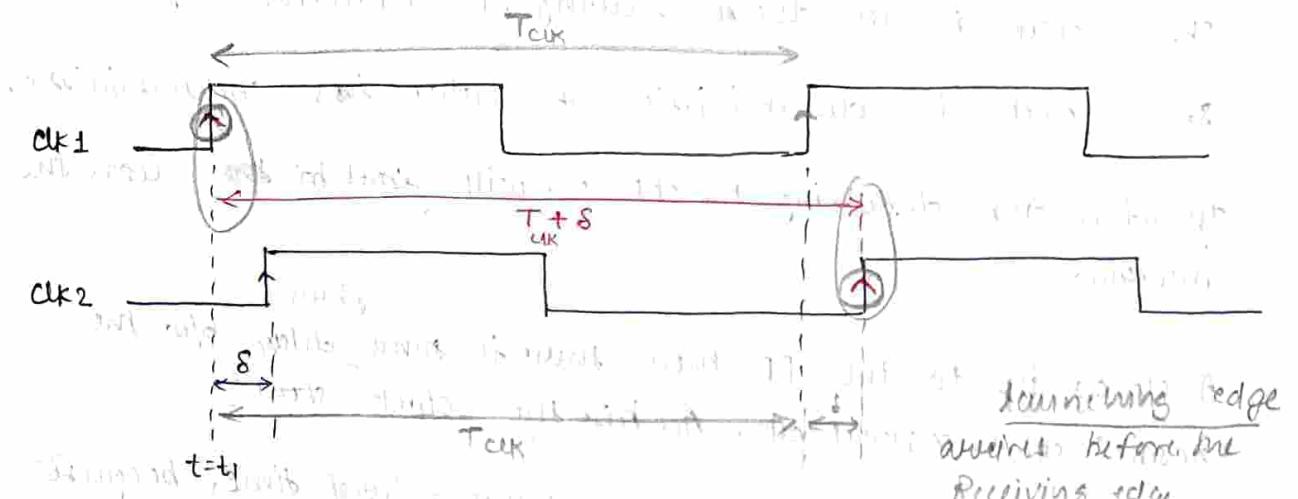


(i) Positive Skew



- If CLK-1 arrives at $t = t_1$ at R1 then CLK2 arrives at ~~$t = t_1 + \delta$~~

$$t = t_1 + \delta \text{ at } R_2.$$



- * CLK1 ka tre edge $t = t_1$ pe FFI me pahucha
CLK2 ka tre edge $t = t_1 + (T_{clk} + \delta)$ pe FF2 me pahuchega

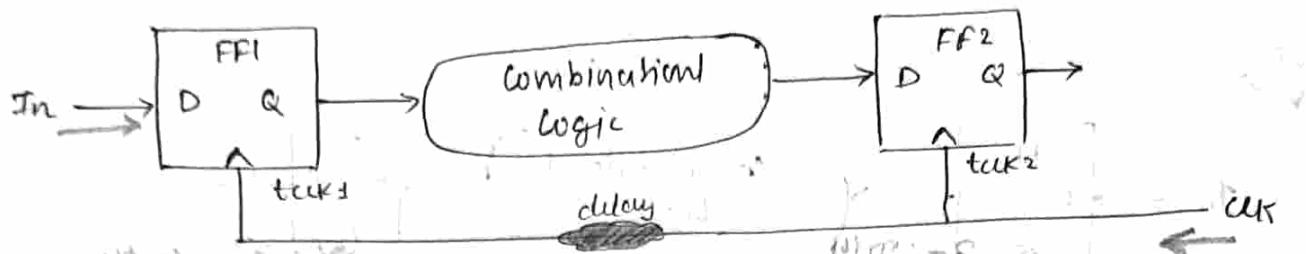
* therefore we can place bigger combinational path b/w FFI & FF2 when delay should be less than $T_{clk} + \delta$.

$$\therefore T_{comb} < T_{clk} + \delta$$

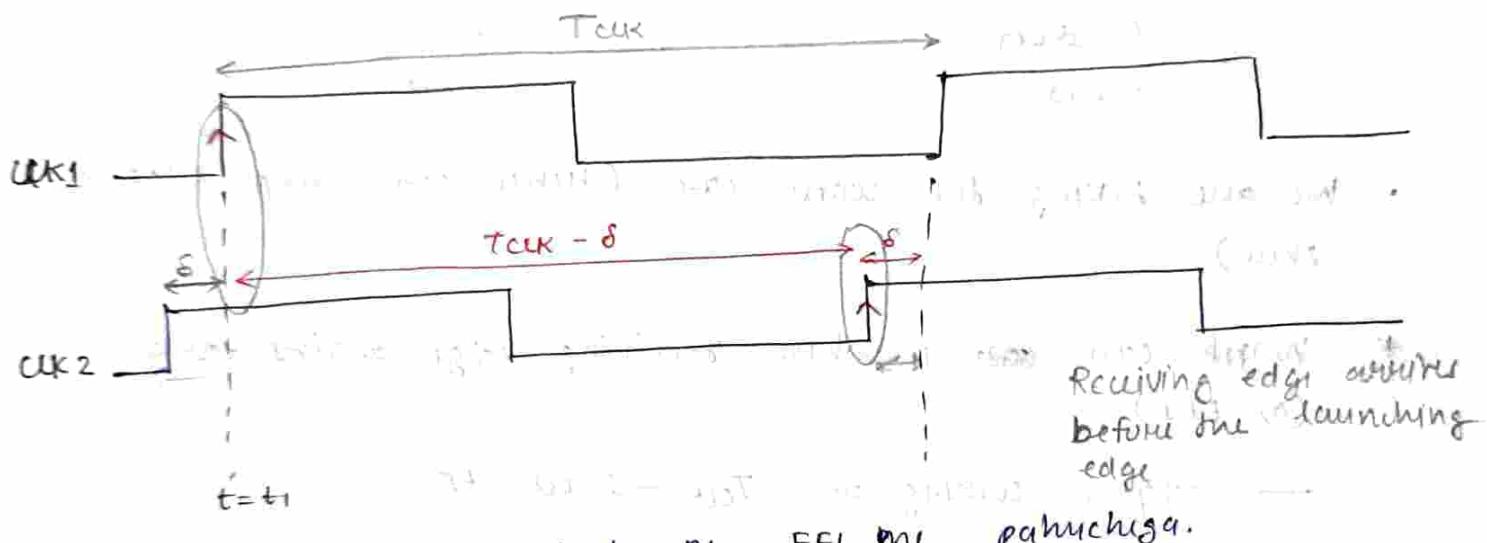
* Considering no delay b/w CLK1 & CLK2 (i.e. $\delta = 0$)

$$\therefore T_{comb} < T_{clk} \text{ only.}$$

(ii) Negative Skew



- If CLK1 arrives at $t=t_1$ at FF1
at $t=t_1 - \delta$ at FF2.



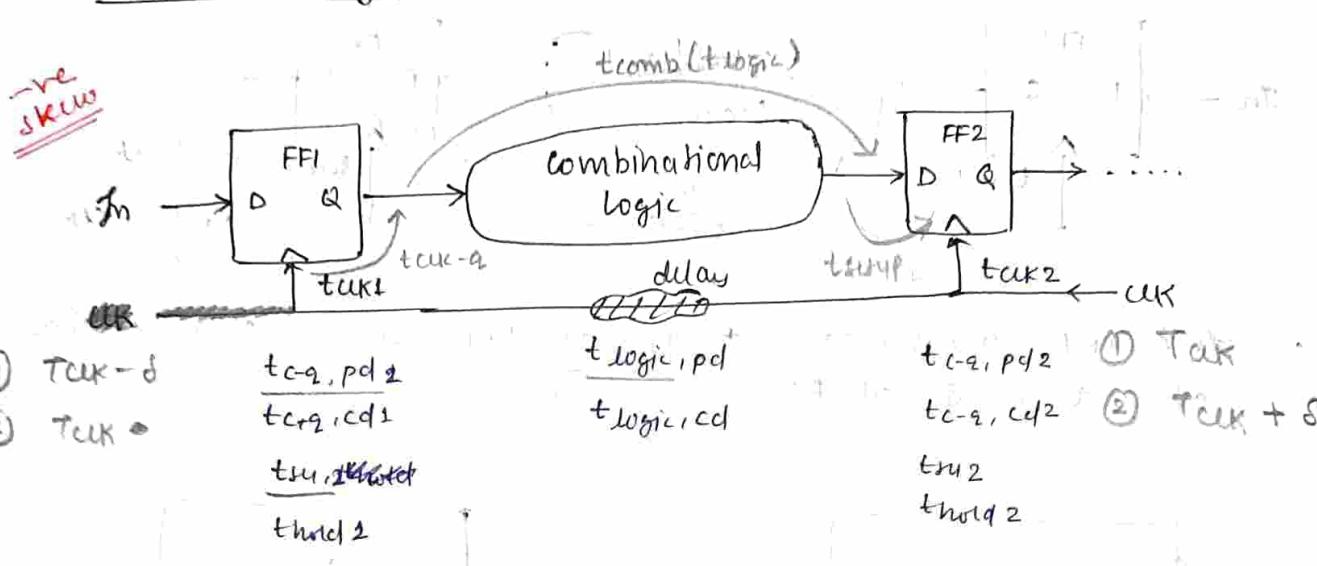
• $t=t_1$ पर लॉच का एजेंसी $t=t_1$ पर FF1 में पहुंच जाएगा।

• $t=t_1 + (Tck - \delta)$ पर FF2 में पहुंच जाएगा।

* therefore we have to plan combination logic of delay such that $T_{comb} < Tck - \delta$

Timing constraint: [Considering only the worst case]

i) minimum cycle time:



We are taking the worst case (Hence considering Negative skew).

* Worst case \rightarrow when receiving edge arrives early (at FF1)

\rightarrow edge \rightarrow coming at $T_{clk} - \delta$ at FF1.

Condition: t_{clk} should be greater than worst case path delay.

$$T_{clk} - \delta \geq t_{c-q, pd} + t_{su, setup_2} + t_{logic}$$

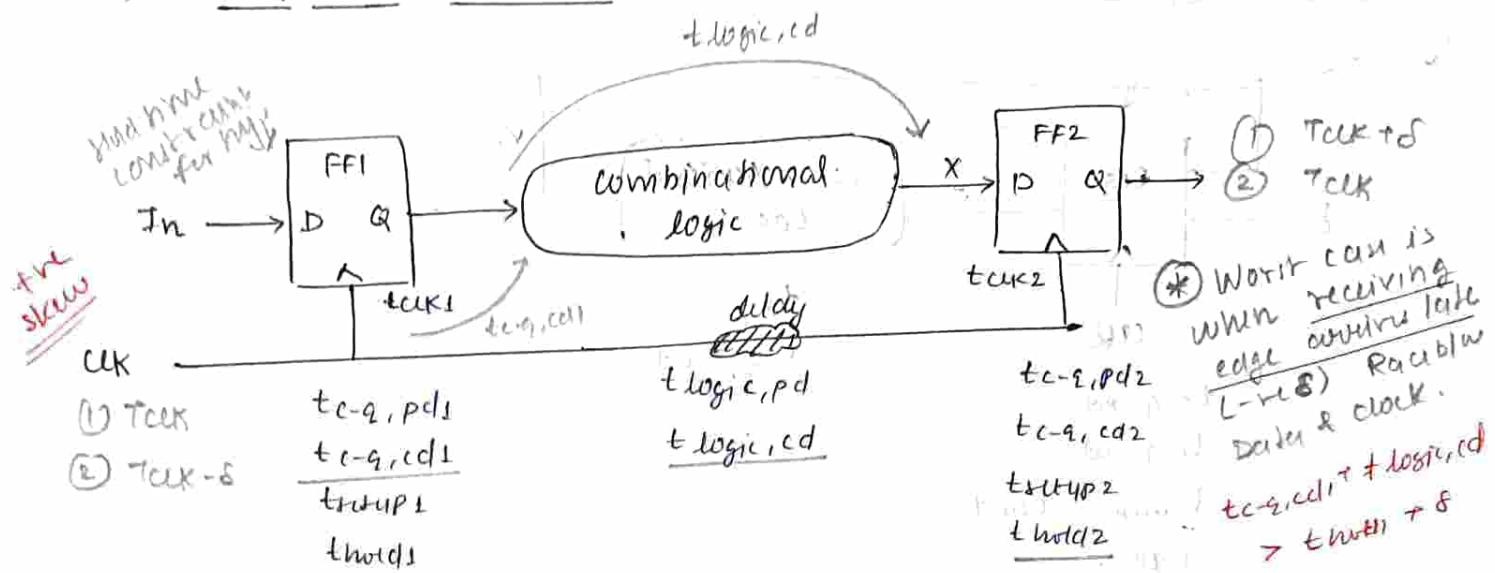
$$\underline{\text{skew}} \quad [T_{clk} \geq t_{su, setup_2} + t_{c-q, pd} + t_{logic} + \delta]$$

Max delay to reach FF2 & delay of clock is also considered

$$\underline{\text{skew}} \quad [T_{clk} + s \geq t_{c-q, pd} + t_{logic} + t_{su, setup_2}]$$

(ii) Hold time constraint

$$[t_{path} + \delta \geq t_{hold}]$$



To avoid hold time violation, we need to ensure that when the next clock edge comes (i.e. when t_{CK2} comes) after that my earlier output must have remained stable for the duration of t_{Hold} time.

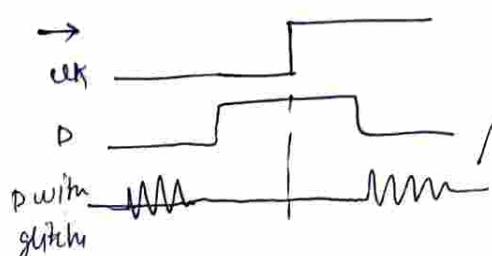
∴ the max. delay from 'In' to 'X' has to be greater than t_{Hold} i.e. $t_{C-Q,CD1} + t_{Logic,CD} > t_{Hold1} + \delta$

Note: (i) For min. cycle time, we are considering clock edges which are t_{CK} time apart.
(ii) For hold time constraint, we are considering same clock edges.

* FFI me data fab tak hold kewna hai jab tak hold time khatam nahi ho jata.

$$t_{C-Q,CD1} + t_{Logic,CD} > t_{Hold1}$$

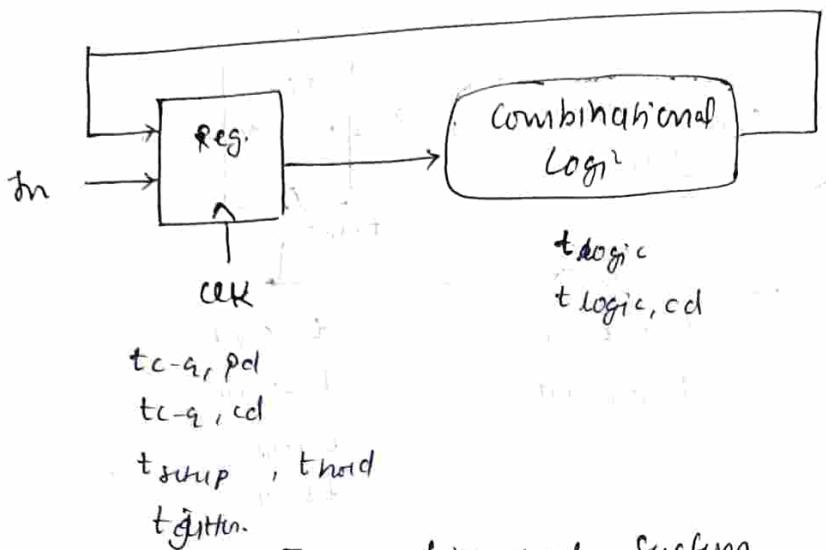
Q. Why are we taking contamination delay?



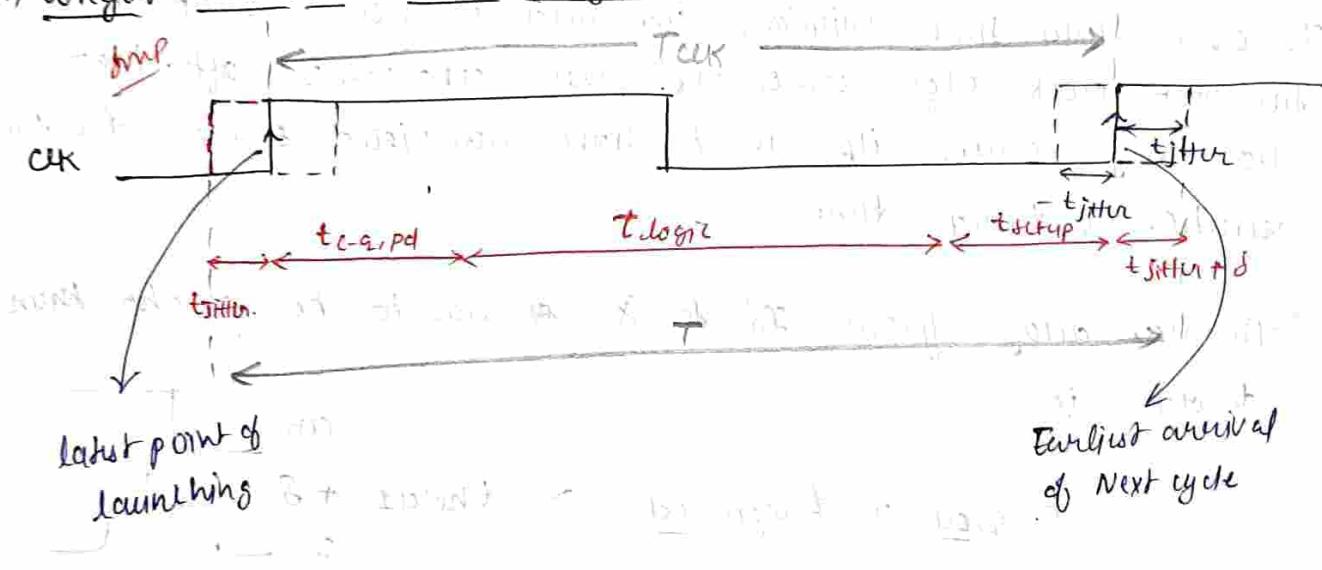
→ the first glitch should happen after the hold time

→ the first glitch should not arrive

Impact of Jitter



(i) longest Path in Edge triggered System



- (ii) If launching edge and date of receiving edge is early, the data will not be too late if:

~~total jitter < trap~~

$$[t - t_{j,i,1} - t_{j,i,2} - \delta > t_{c-q,pd} + t_{logic} + t_{sjup}]$$

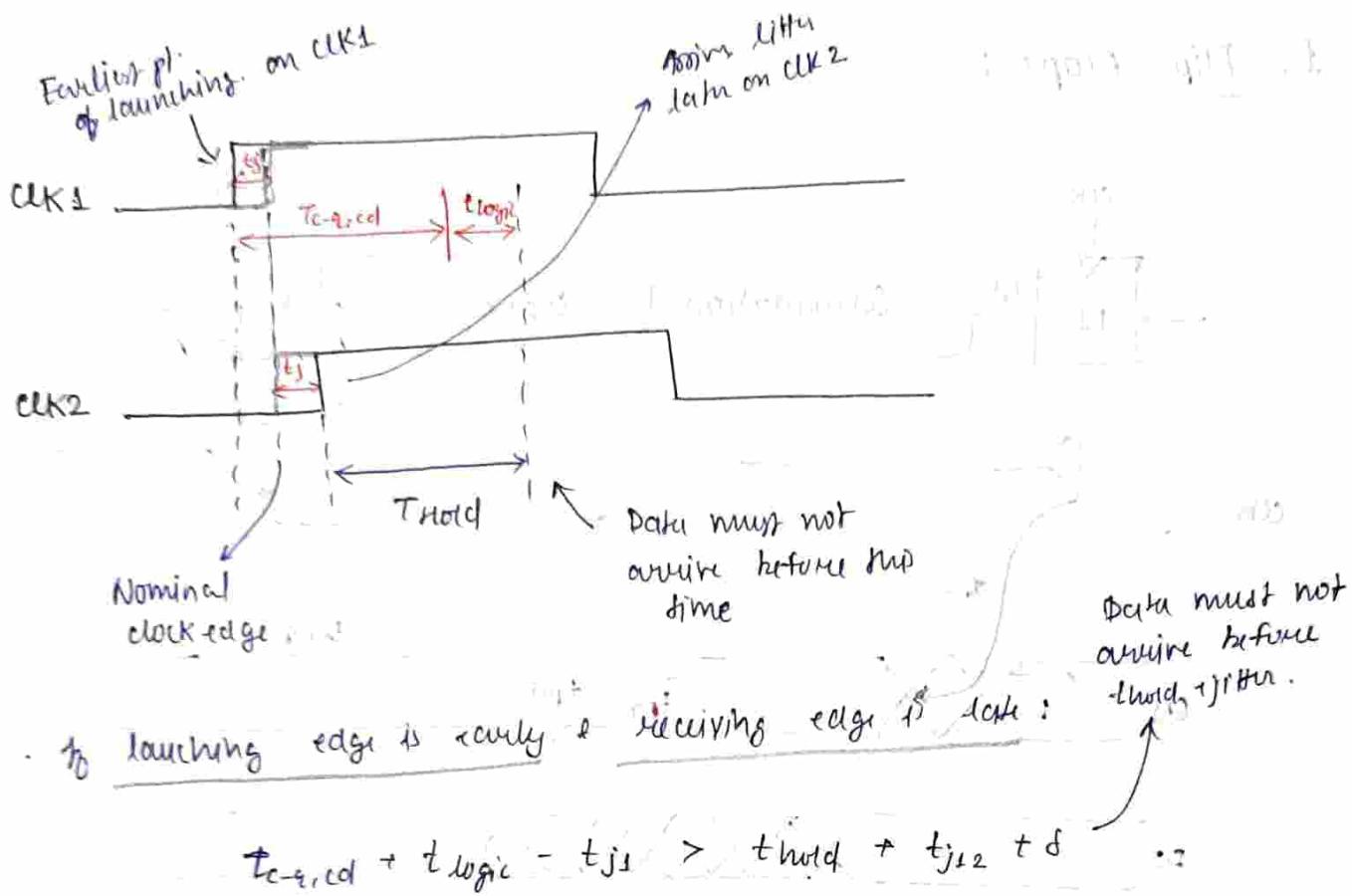
* minimum cycle is defined by the maximum delay through logic,

$$[t > t_{c-q,pd} + t_{logic} + t_{sjup} + \delta + 2t_{jitter}]$$

* skew can be either +ve or -ve.

(iii) Shortest Path

Yates CSE5770M
Date: 10/10/2017

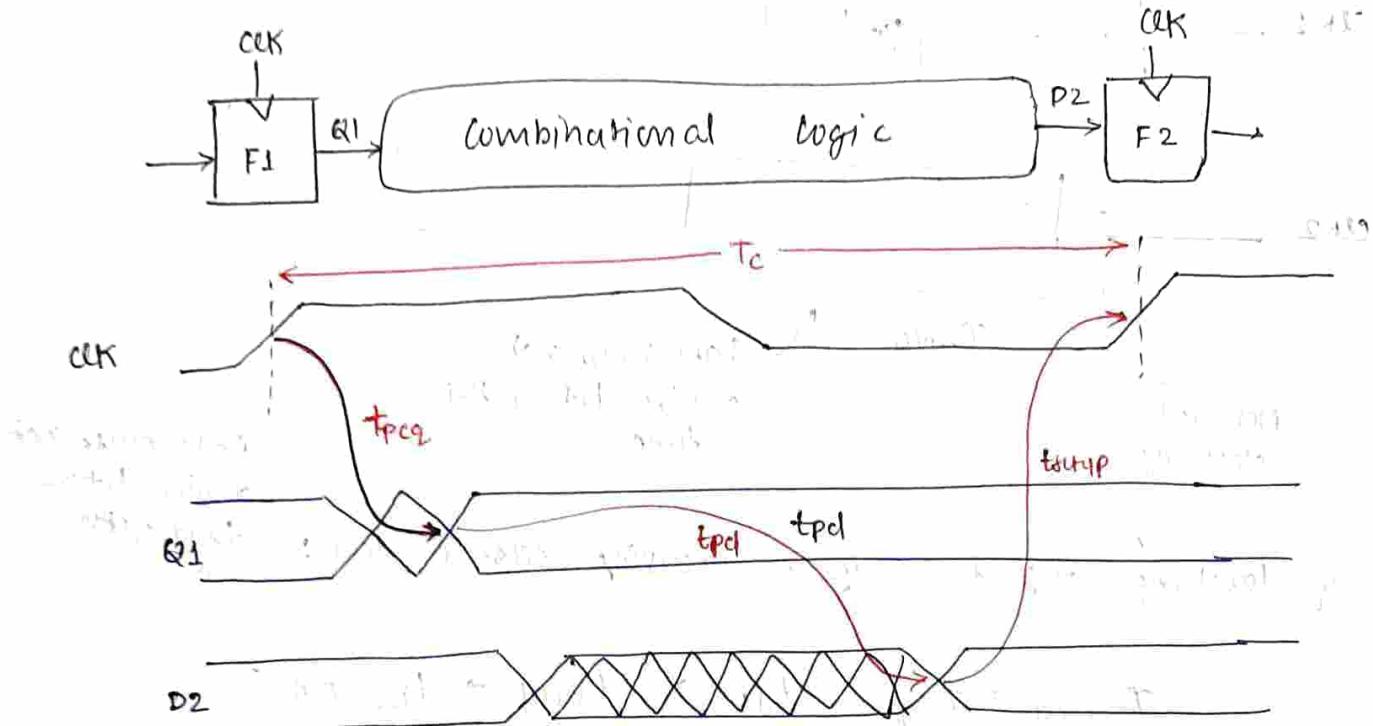


∴ Minimum logic delay, $t_{logic} = t_{j1} - t_{hold} - \delta$

$$t_{c-q,cd} + t_{logic} = t_{j1} - 2t_{j2} - \delta > t_{hold}$$

Maximum Delay : [Neglecting skew & Jitter] (in picoseconds)

1. Flip-Flops :



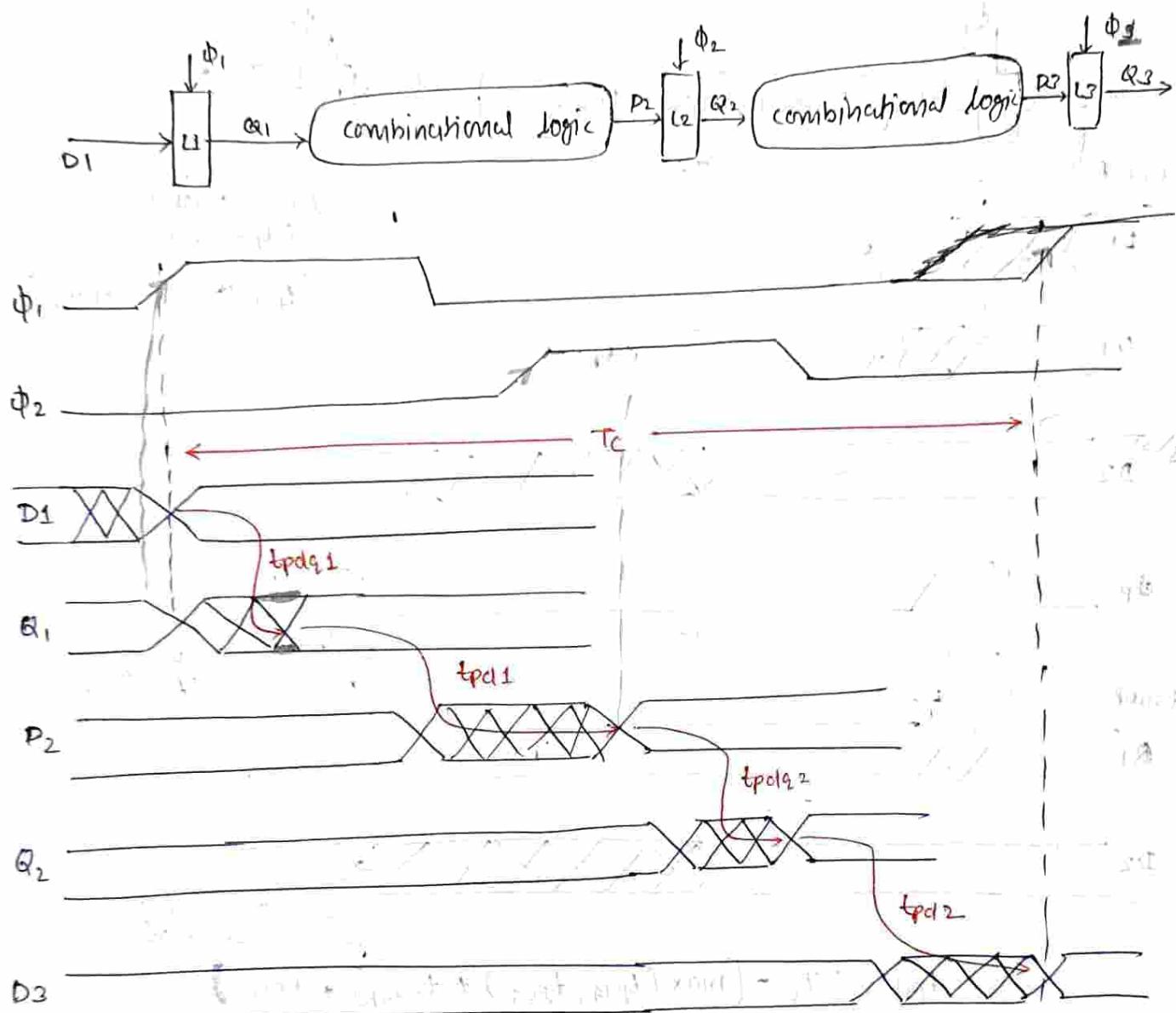
$$t_{PD} \leq T_c - [t_{Setup} + t_{PCQ}]$$

assuming no skew

- If FFs are not inserted in the combinational path then only t_{PD} will be the total delay.
- But now due to FF we have t_{Setup} & t_{C-Q} delay introduced in the path.

$$T_c \geq t_{PCQ} + t_{PD} + t_{Setup} + \delta$$

2. ~~V~~ Two-phased latches



$$T_c \geq t_{pdq1} + t_{pdq2} + t_{pdq2} + t_{pd2} (+\delta)$$

$$t_{pd} = t_{pd1} + t_{pd2} \leq T_c - (t_{pdq1} + t_{pdq2})$$

$$t_{pd1} + t_{pd2} \leq T_c - (2t_{pdq})$$

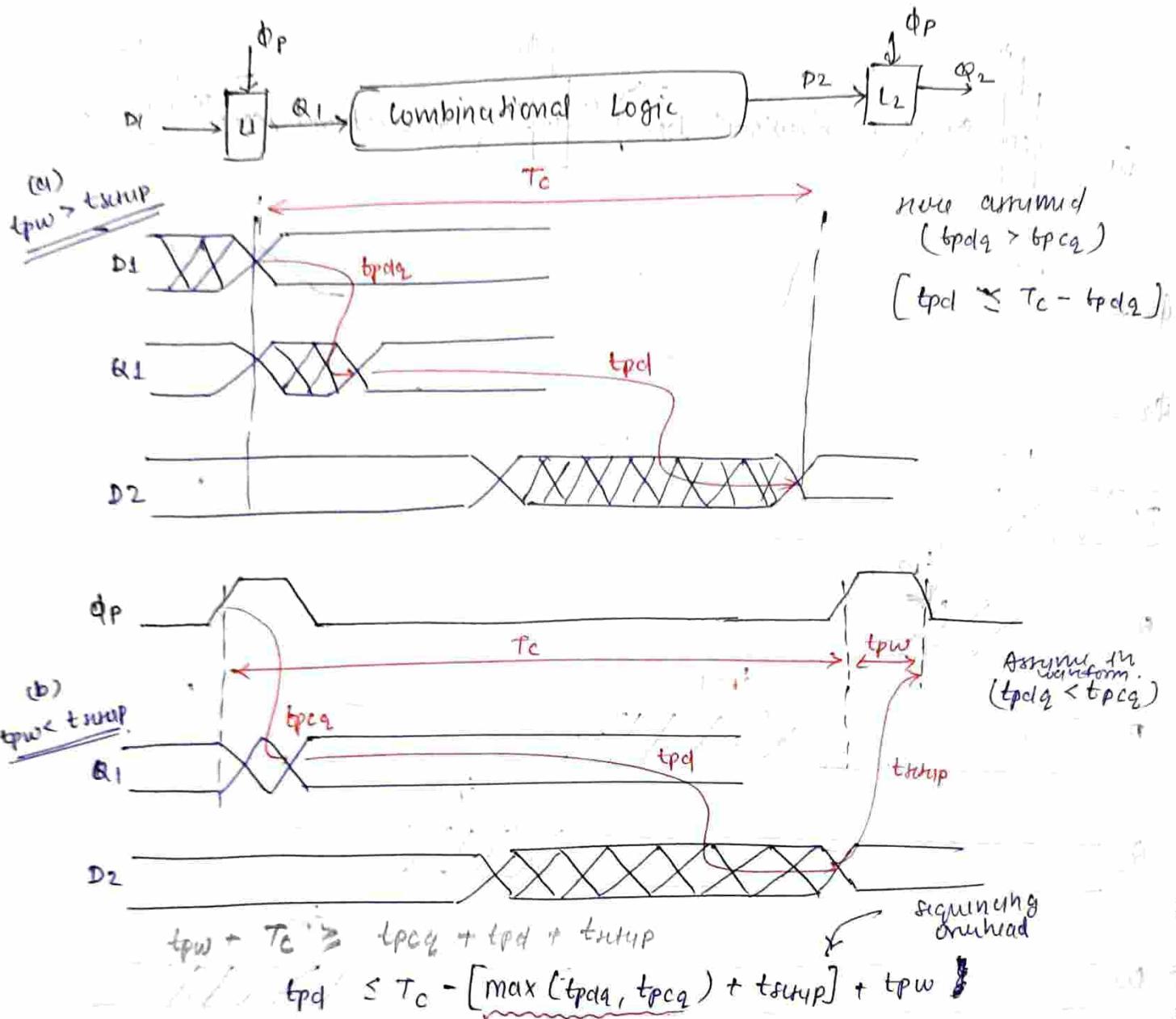
\nwarrow sequencing or overhead

~~We are not considering the setup time here because setup time is taken at sampling edge and for the latches it is the falling edge.~~

Q. Is total delay be more longer if we consider the setup time also?

→ It will be discussed in the Machine Borrowing.

3. Pulsed Latches



$$t_{pd} \leq T_c + t_{pw} - [max(t_{pdq}, t_{pcq}) + t_{qup}]$$

- When $t_{pw} \gg t_{qup}$ then we ignore the setup time.

- When $t_{pw} \ll t_{qup}$ then we can't ignore the setup time and we have to consider setup time in max delay.

- For a latch maximum of D-Q delay or CLK-Q delay is considered.

Q. When $t_{pw} > t_{qup}$ then why are we considering the setup time?

(d) Why we are subtracting the tpw from the sequencing overhead?

- In a latch you launch your data from rising edge & sample it at falling edge because setup time is from falling edge.
- But time period is from rising-to-rising edge of CLK. So the tpw has to be reduced from the sequencing overhead.

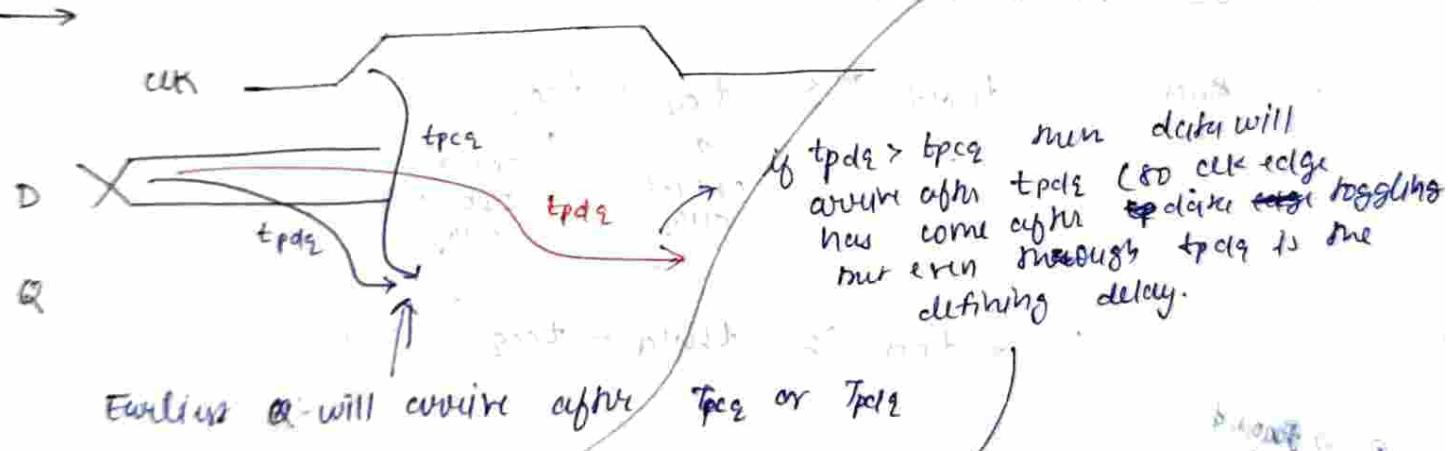
$$* T_c + t_{pw} > t_{pd} + \max[t_{pdq}, t_{pcq}] + t_{setup}$$

* Timing constraint window is $T_c + t_{pw}$
But my time period is only T_c .

* $t_{pdq} + t_{pcq}$ sum is e2.

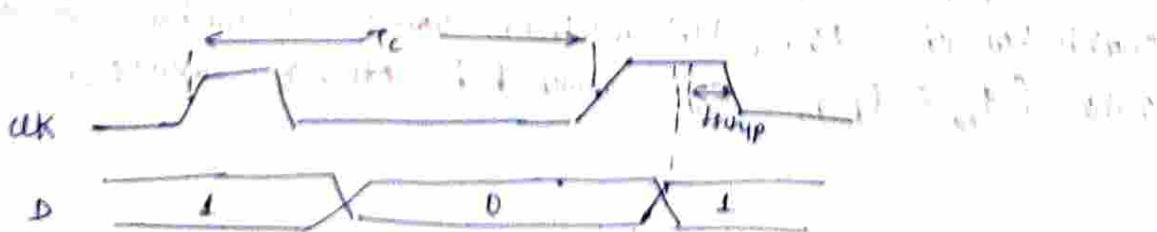
* $\max(t_{pdq} + t_{pcq})$ will be taken.

Ques * Why max of t_{pdq} & t_{pcq} ?



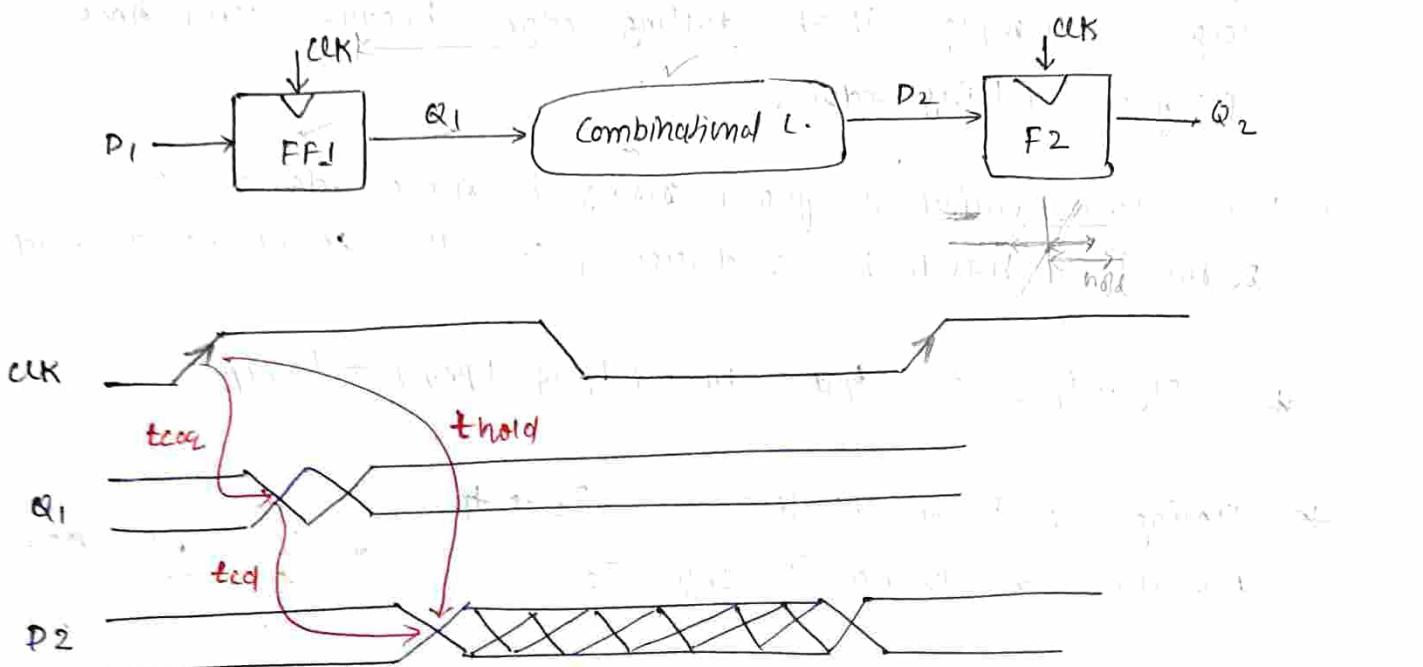
* For this can the toggling be glued to CLK-edge but the delay would be due to t_{pdq} because $(t_{pdq} > t_{pcq})$. Example

Q If pulse width (tpw) is large enough & clock is toggling at the middle of the PW + tsetup time is not violating so the difference of w/ the clock edge how come & data is toggled, then wouldn't that difference be taken in our timing calculation?



Minimum Delay :-

1. Flip-Flop



- D_2 should hold its value for ' t_{hold} ' time.

And $t_{hold} \leq t_{cd} + t_{ccq}$

(why MS. more sign?)

critical delay

clock to Q delay

considering skew

$$\therefore \{ t_{cd} \geq t_{hold} - t_{ccq} \} + \delta$$

for 2 phase
latch

skew

- 2). The Q_1 triggered op should reach Q_2 after hold time requirement has met. (Q_2 or P_2)

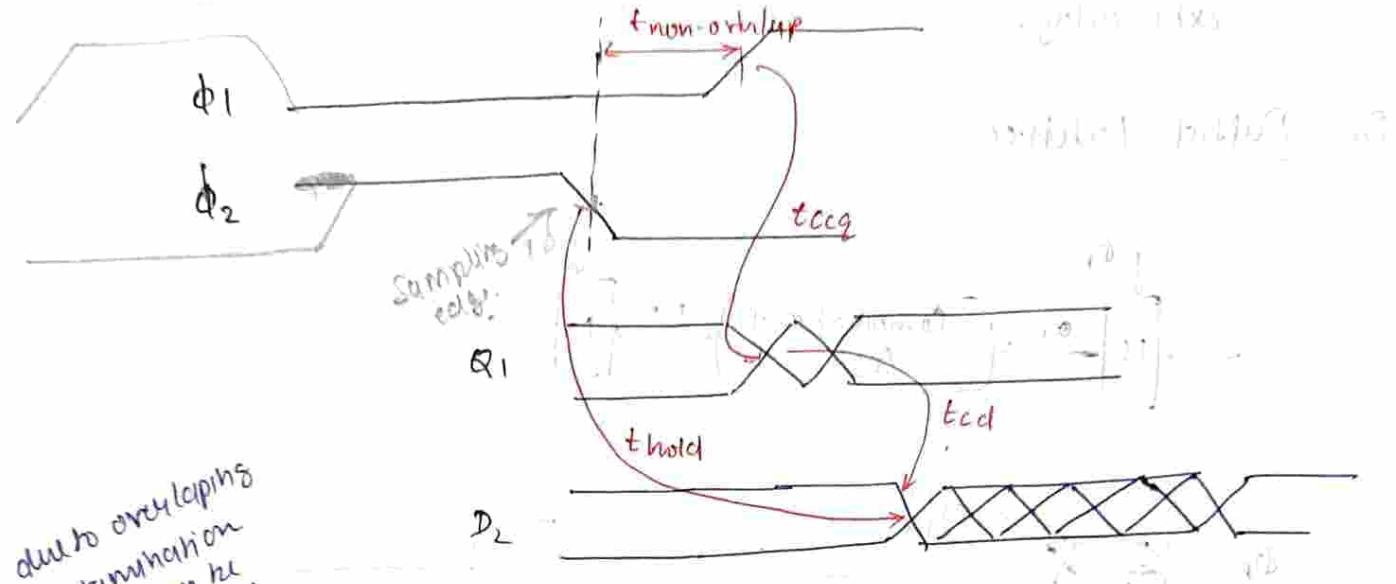
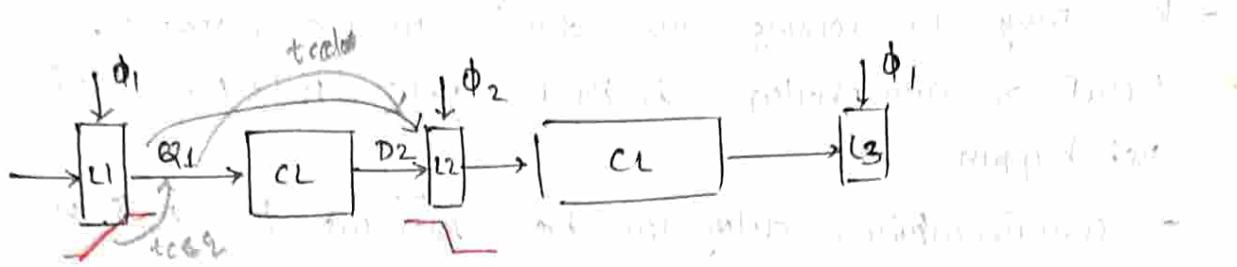
— And there is some non-overlapping window b/w Q_1 & Q_2 .

$$t_{cd} \geq t_{hold} - t_{ccq} - \text{non-overlap} + \delta$$

Note: for taking skew, add δ in t_{hold} kuki hum chahte hai ki t_{hold} kisi minimum value se ~~by~~ by minimum rakhne [$t_{hold} \leq t_{pwm} - \delta$] would be in effect of skewing].

2. X two-phased ladder

W. Smp



$$t_{cd1}, t_{cd2} \geq t_{hold} - t_{ccq} - t_{non-overlap}$$

- We want ϕ_1 , triggered $\phi_0/p.$, should reach L_2 after hold time requirement has met.

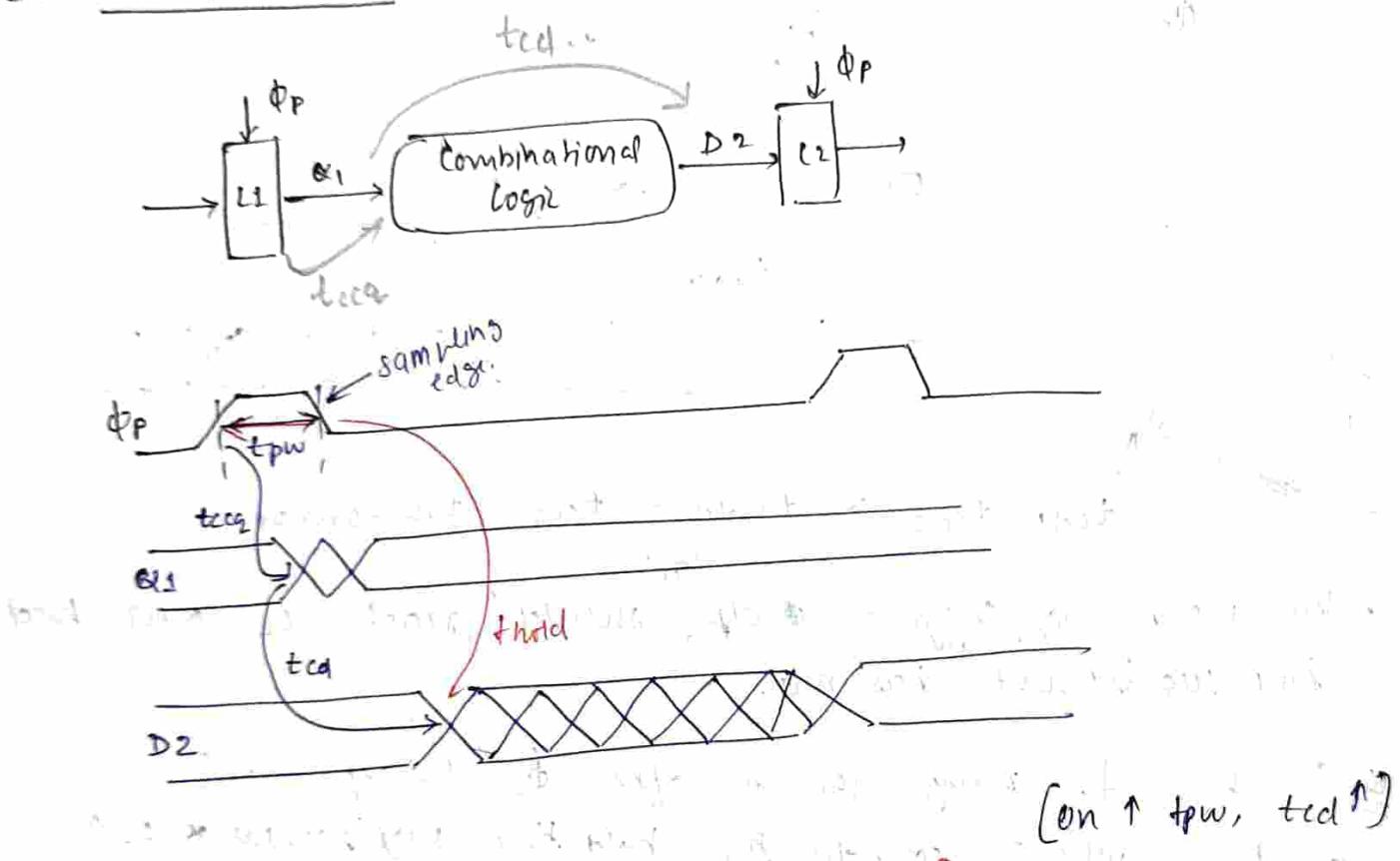
- * why ϕ_1 - rising edge is after ϕ_2 falling edge?
→ now, will I consider the hold time requirement of L_2 ?
→ when L_2 sampling edge is closed, now we want ip of L_2 to be stable and for that there should not be any rising edge of ϕ_1 (which would change the value of ip of L_2).

Therefore now let's assume rising edge of ϕ_1 is coming. We are interested on that rising.

After sampling edge of ϕ_2 it will stop sampling but it's enough to meet the hold time requirement.

- * You can't change non-overlap time period of the clk.
 - We may be losing data due to non-overlapping.
- * Benefit of non-overlap is that your hold time failure will not happen.
 - contamination delay can be reduced to ms. non-overlap extra delay.

3. Pulsed Latches



$$t_{cd} \geq t_{hold} - t_{cs} + t_{pw} + \delta$$

- For entire duration of pulse, no change in Q_2 over here.
 - Data is triggered from the edge of ϕ_P .
 - And hold time is defined from this rising edge of ϕ_P .
- * We characterize setup & hold time at the original clock of the system.
 → Here we have pulsed clock in our system

- Q. What is the requirement for L2?
- Data should be stable for ~~threshold~~ ~~before the~~ falling edge of the CLK.
- * And O/P \Rightarrow biggest clocking edge of the CLK at L1
 $\therefore t_{cqq} + t_{col} \leq t_{hold} + t_{pw}$ will work if t_{pw} is ~~between~~ \Rightarrow with ~~clock~~ ~~edges~~ ~~are~~ ~~in~~ ~~it~~ because ~~t_{hold}~~ \Rightarrow coming from the falling edge of O/P.
~~(as falling edge of O/P is sampling edge of CLK)~~
- * Note that system setup-time is more relaxed + hold time very tight so and very sticky because you can't recover hold time constraint — therefore may not very commonly used.
- Q. Since we dedicate t_{pw} to get better hold-time requirement. Then how much t_{pw} can be?
- Impulse & system behaves as FF.
- Q. Does $t_{pw} > t_{cqq}$? Is it necessary?
 → There is no connection b/w them
- would not suggest

Time Borrowing

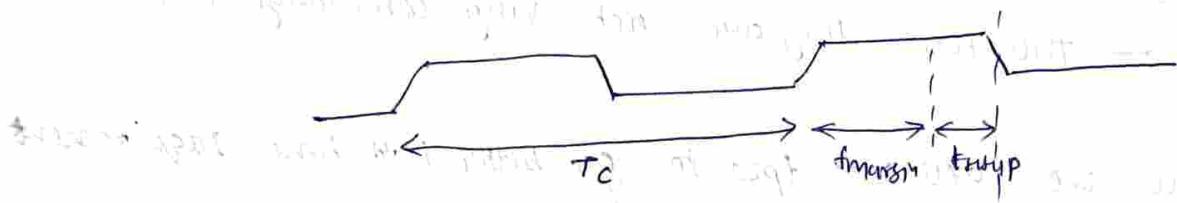
(Week Pg: 387) timing and initial

in FF based system:

- Data launches on the rising edge
- must setup before the next rising edge
- If it arrives late, system fails
- If it arrives early, time is wasted.

in latch based system:

- Data can pass through while it is transparent if it does not limit itself for any rising edge
- long cycle of logic can borrow time into next.

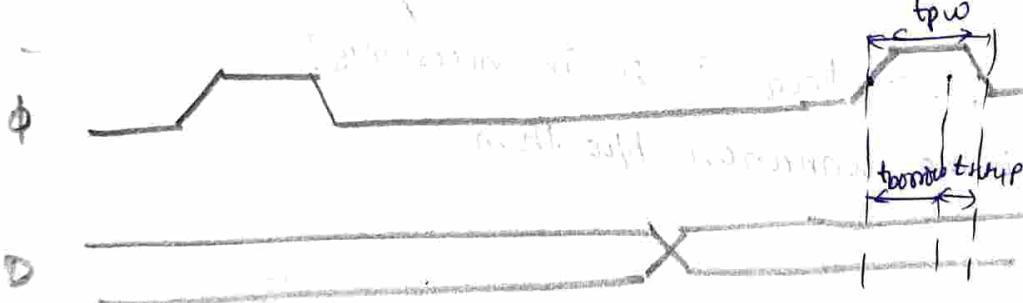


(i) two-phase latch

$$t_{borrow} \leq \frac{T_c}{2} - (t_{startup} + t_{holdup})$$

(ii) pulsed latch

$$t_{borrow} \leq t_{pw} - t_{startup}$$

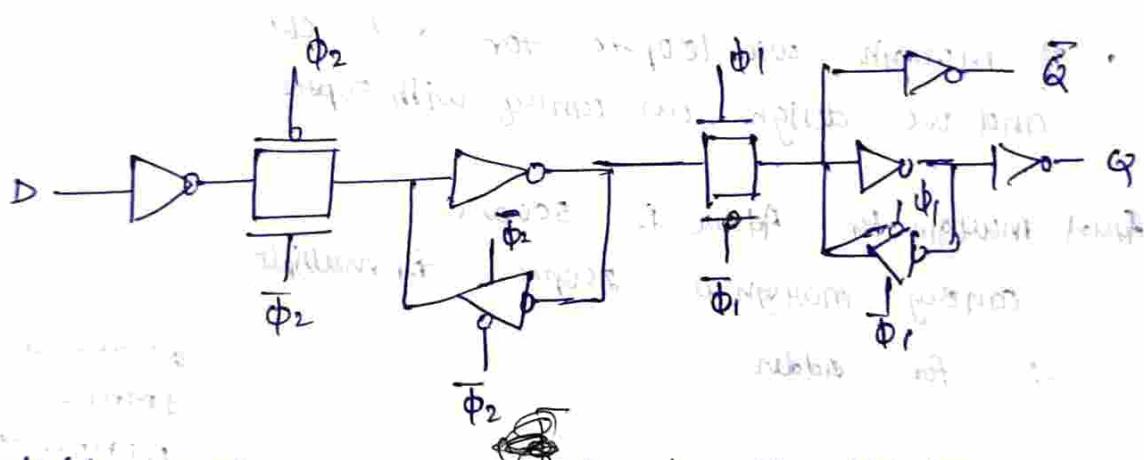


$$t_{borrow} \leq t_{pw} - t_{startup}.$$

- * If setup time is getting violated you reduce the clock speed. (inuman Take)
- * If hold times are violated, chips for any speed (ie, you can't improve this by changing Take you have to change your CKT) → But with two phased latches you can avoid such type violations.
 - To get working chip → the most imp requirement

**** A way to guarantee hold time is to use 2-phase latches with big non-overlap time**

Safe Flip-Flop:



- To avoid overlapping of ϕ in the ff we can use two different Φ_1 & Φ_2 within the same FF.
- When we design the FF with non-overlapping clocks
 - it is robust in the presence of hold time failure because there is no overlap
 - It will slow down the system, because it adds to the setup requirement.

Adaptive sequencing

below printing in dark ink
(top row) page 2000

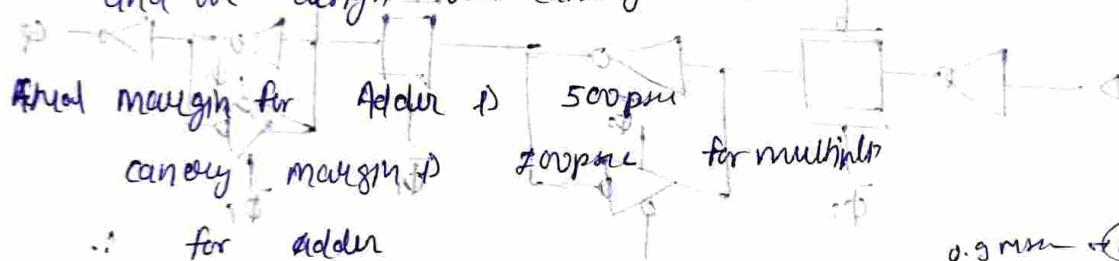
~~Designs include~~

- if a silicon is manufactured at typical lot size
 - then we have margin for SS lot for the setup time.
 - then we also have margin for FF lot for the hold time.

* your designed system can operate at 450 MHz as well as you can operate at 950 MHz while taking the margin into account.

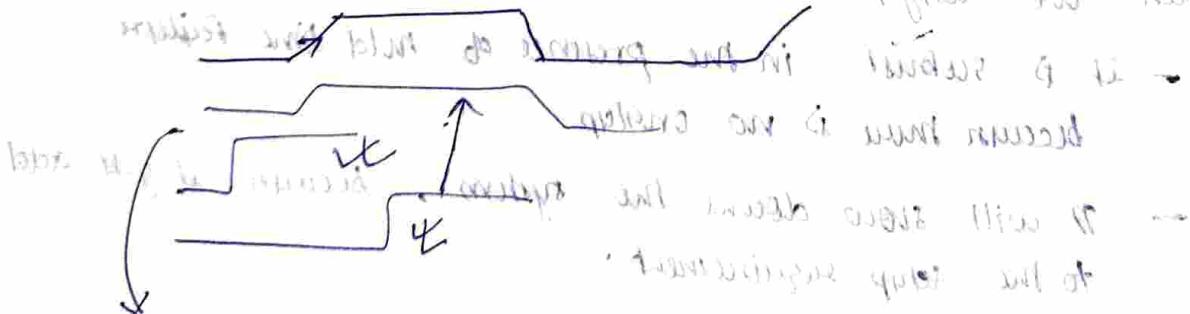
Eg. of Adder and Multiplier. (11:50 min)

- If malogn was 100sec for out $\frac{d}{dt}$
and we design our convey with 50sec



#Cancery and the ^{FG} principles of employment \rightarrow
#Razor ^{FG} and the ^{FG} principles of employment \rightarrow

1950s 1960s 1970s 1980s 1990s 2000s 2010s



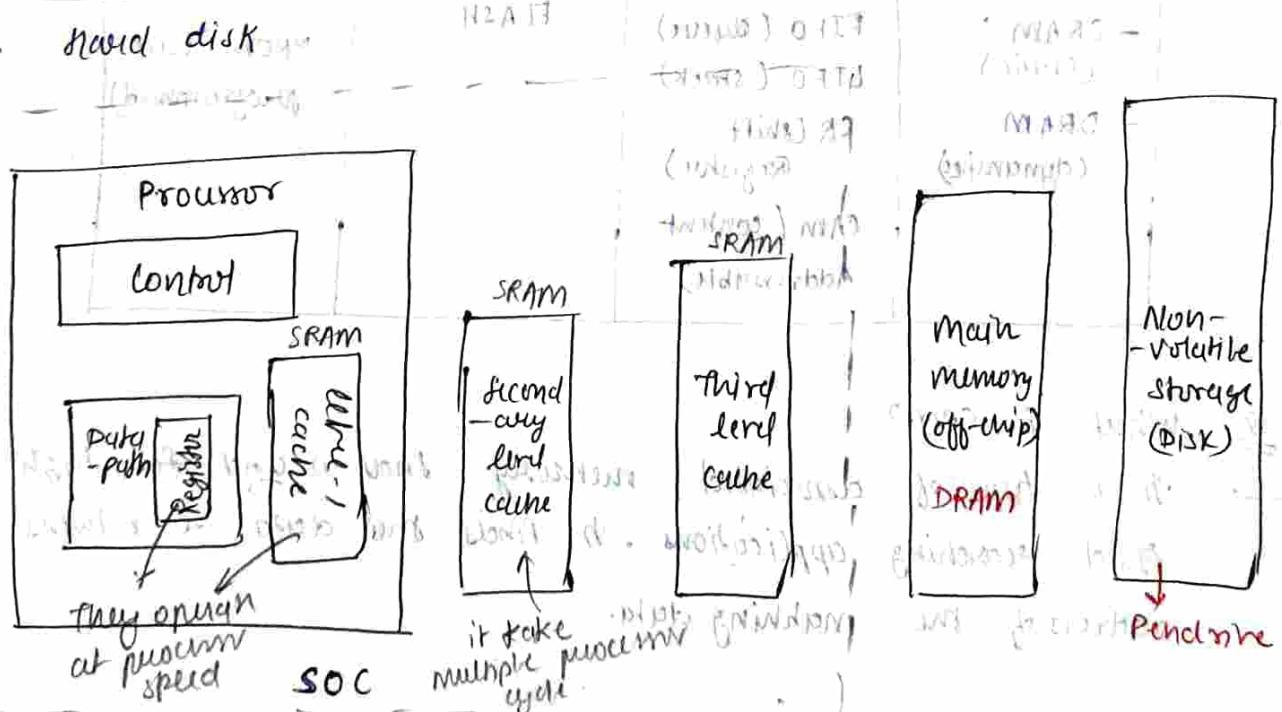
Memory Design

- Memory are high density storage elements.
- Memories are minimalistic storage elements which provides high density. As we take FF, register or latch which involves large no. of transistors even though they provide storage functionality but they don't provide high package density. hence we use memory models. (FF takes more area as compared to standard - 60 memory cell)
- Memory cover 70-80% of chip area in advanced SoCs.

Q: Why so much of area is given to memories?
 → Because in typical computer architecture, memories ~~consume less~~ are more at various levels.

Memory Hierarchy: * Memory management

- * Typical system has multiple level of storage systems.
1. Very fast registers in the Data path.
 2. Processor speed driven L1-cache
 3. Slower but dense on-chip secondary (& tertiary) memories.
 4. Off-chip, but on-board storage with latency of 100s of cycle.
 5. Hard disk



- * As we move away from processor they become slower.

Classification of Memories

- Read-Write Memory → memory in which we can read the data from memory & write data to memory.
(RW memory)
- Read only memory → memory in which data can't be written
(ROM) & we can only read data from it.
- Non-volatile RW memory → memory that can retain stored information even after power is removed.
(NVRWM)
- Volatile memory → memory that don't retain stored information when power is removed.
- Random Access memory → RAM is a computer's short term memory, where the data that the processor is currently using is stored.

RW memory		NVRWM	ROM
Random Access	Non-Random Access	EPROM EEPROM FLASH	MASK Programmed PROM (fuse programmed) Erasable
- SRAM (static) - DRAM (dynamic)	FIFO (queue) LIFO (stack) SR (shift Register) CAM (content Addressable)		

volatile

Q. What is CAM?

→ It is a type of cache memory that is used for high speed searching applications. It finds the data & returns the address of the matching data.

- Typically in memory, you give address location & you would like to find out what is written on particular address.
- * In CAM, depending on the data, you get a different address to send your data.
 - CAM are tightly coupled to processors, so that when particular instruction is given, the follow-up instructions will be brought up quickly based on the content by doing the quick content search, as to which location is this particular instruction stored & then I pull out remaining instructions from that region quickly.

- * EEPROMs are used in Boot sequence.
 - Also used in setup Box of Airtel, Jata sky.

- * Flash are used in memory card, cell-phones.
 - All solid state devices (SSD) are flash.

- * ROM are used in Boot sequence.
 - Once it is programmed you can't change it.

Feature comparison in memory Type

- * ~~speed~~ ON-chip memories typically SRAMs, whether we talk about L1-cache, L2-cache, L3-cache. They are very fast but they are not very dense.

- It's 6-transistor per SRAM.

- Q: What is 6T-SRAM cells?

- It has two-transistors to control bit-lines and two back-to-back coupled inverters.

- Density
- In SRAM we need 6-transistors to save 1-bit of data.
 - DRAMs use one-transistor + 1-capacitor.
 - Flash is effectively less than one transistor per bit of data.
 - On one transistor you would be storing two bits or 3-bits of data. Therefore flash are much more dens.

Memory Type	SRAM	DRAM	Flash
✓ Speed	Very fast	Fast	Very slow
✓ Density	low	high	Very high
✓ Endurance	Butter	Butter	Poor
Power	low	High	Very low
✓ Refresh	No	Yes	No
Retention	Volatile	Volatile	Non-volatile
Scalable	Good	Bad	Good
✓ Mechanism	Bi-stable Latch	Capacitor	FN Tunnelling @ Hot carrier my electrons

Endurance

- Endurance → No. of times you can read & write on a particular memory.

- In SRAM, it is 6-transistors fully static CMOS logic + regular voltage operation & endurance is maximum.
- In DRAM, it has capacitor which can leak the data but even though endurance is very high
- In Flash we use very high voltage & therefore endurance goes bad. (No. of reads & writes are very less)

Refresh

- Refresh → Restoring the data that was stored on the cell.
 - * if you don't restore the data on regular basis then you end up losing data. Hence you have to refresh.
- DRAMs store data on capacitor hence data can leak over time. Therefore we need to refresh the DRAM.

Q. how we can store two or three bits on 1-transistor
in a Flash?

→ We have floating gate in flash.

In Flash, we use special kind of circuit, where you have
(gate-1 \leftrightarrow silicon-dioxide \leftrightarrow gate-2 \leftrightarrow SiO_2 \leftrightarrow substrate)

Just like FGMOs \rightarrow floating gate mos.

in which you can trap charges & hence V_t would change
and then you would be able to change say what is the
kind of data stored here.

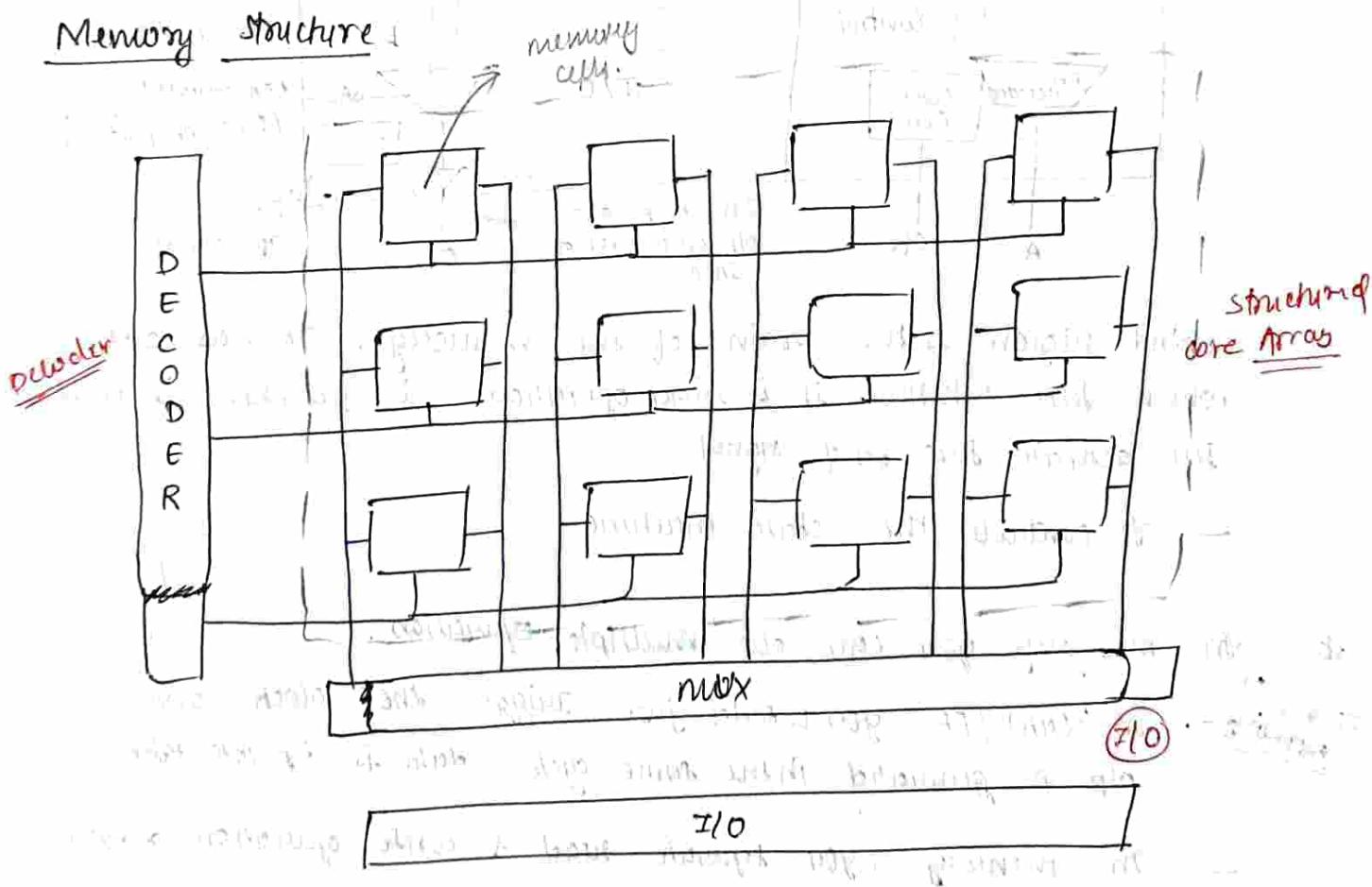
Q. Why we need to refresh in DRAM?

→ If we don't refresh then we have kip 2 transistors + 1 cap
(I keep it) this will increase memory.

420
1730133

refresh + 2 transistors + 1 cap

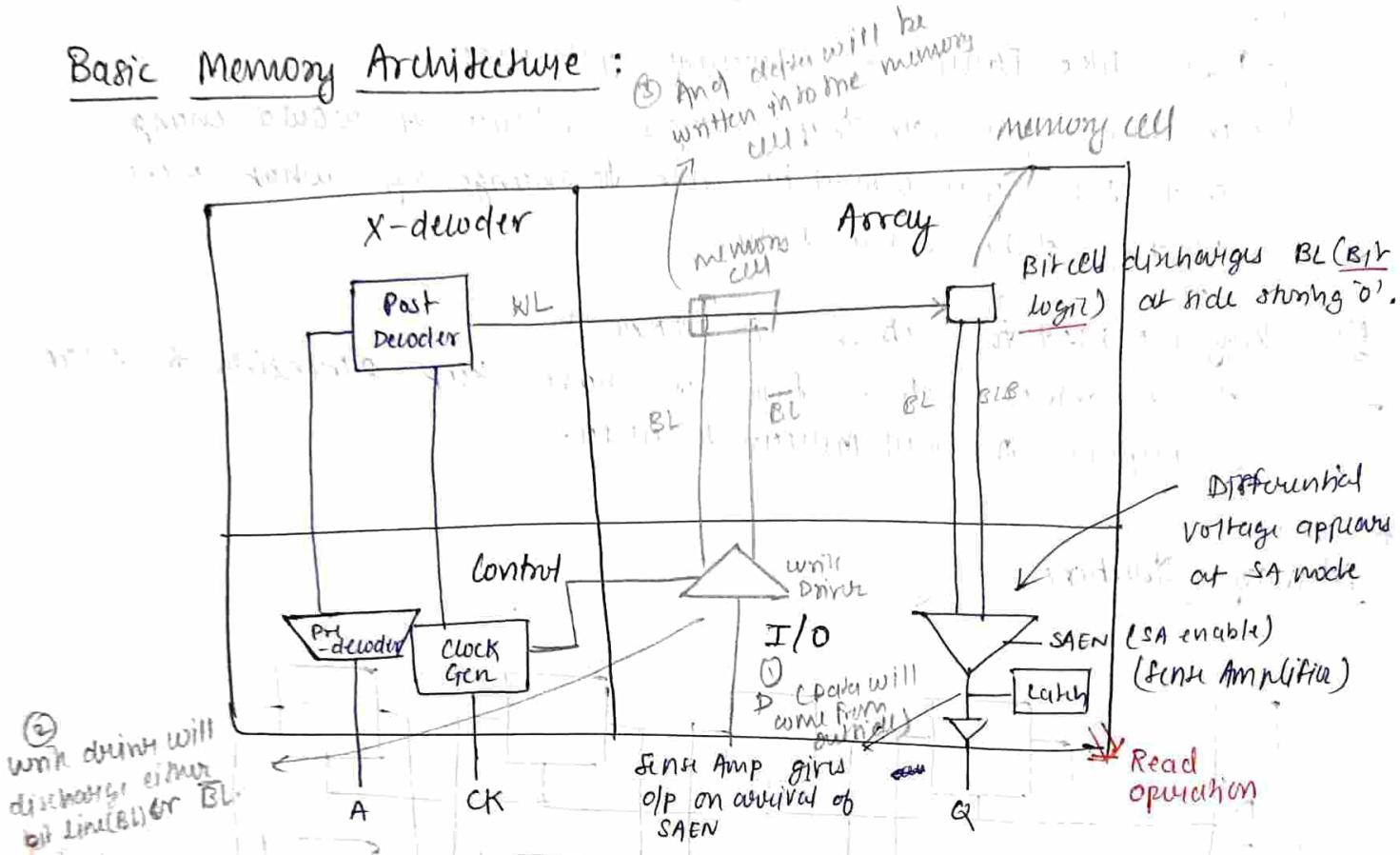
Memory Structure



- * There is a cluster of storage elements which are closely coupled not only through word line but also through bit line, that unit remains constant.

- there is some queuing architecture that remains constant
 - there is I/O which have got to buffer the old ~~old~~ demands same
 - some control region in memory remains constant.

Basic Memory Architecture:



- Control region is the brain of the memory. It has cell which tells whether it is exact operation, if yes then it will tell activate this set of signal.

- It controls the state machine.

it in memory you can do multiple operation.

~~FF vs memory~~ In Latch/FF, you would just trigger the clock and o/p is generated in the same cycle. Data is latched.

- In memory, you switch read + write operation. occurs cycle.

* for ~~ip~~ we use diff set of ips for with operation are not diff. set of ips

* There are many states in which the memory can be controlled by control unit.

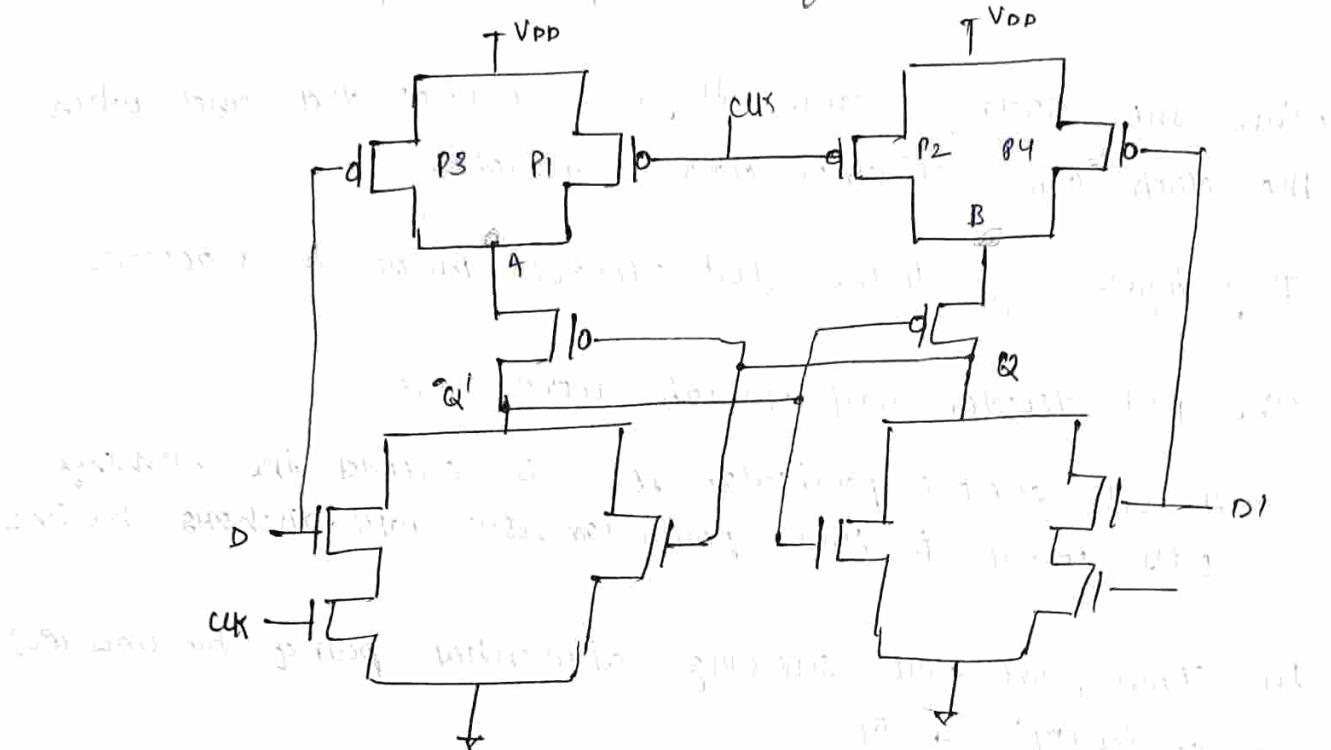
Control Region: It has clock generator and address pre-decoding.

- When the address comes they are pre-decoded and when the clock comes internal clock is generated.
- They together go to the post decoder inside the X-decoder.
- The post decoder will generate word line.
 - It will select a particular row & selected the memory cells stored in that particular row will discharge bit lines.
- In SRAM, we are showing differential pair of bit lines (BL).
 - Bit line (BL) is \bar{BL} .
- In SRAM also we can have single ended i.e. 1-bit per memory cell.
- In PROM, DRAM, flash memories, you would typically one-bit line per memory cell.

Q: Why we don't have 2-bit lines in DRAM?
→ If we have to route 1-metal vs 2-metal which is denser?
→ Using 1-metal will be denser. As we want to keep density very high hence we don't use 2-bit lines in DRAM.

* We you select some bit line would sense a memory cell & if it discharge you would measure using the sense Amplifier.

D-FF Implementation This D-FF will be used to select only one word in memory block.



- Total = 14 transistors $\equiv 12 + 2$, for generating D'.

Q: Out of these 14 transistors, how many transistors which are not really needed or more are free transistors which could actually be shared across memory cells.

→ For example, consider the case where word line is not selected, where the clock is not coming (i.e. $CK=0$) means P_1 & P_2 are ON. Therefore VDD comes directly on node A & B.

→ means P_{10}, P_3, P_4 are not required 99% of the time, because for memory clock would come rarely.

• Suppose you have 1024 words in memory, even if you access memory in every cycle, then also every word is accessed in every 1 out of 1024 times.

→ so more than 99% of times P_1, P_2, P_3, P_4 are not getting selected (then P_{10} are simply waste)

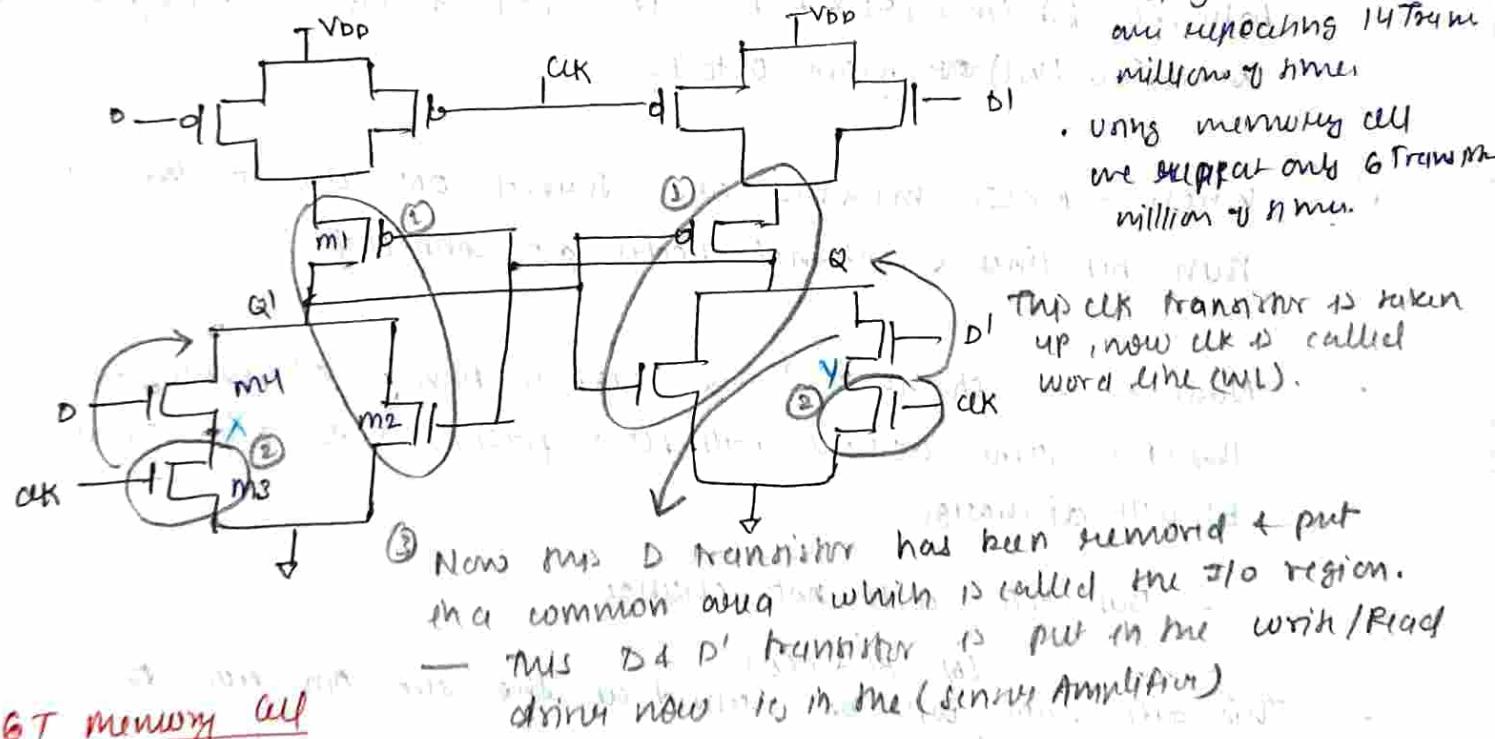
Q. What will be the alternative way to implement more 1% cells where P_1, P_2, P_3, P_4 are needed?

→ Instead of using full NMOS static ckt we can use pseudo NMOS.

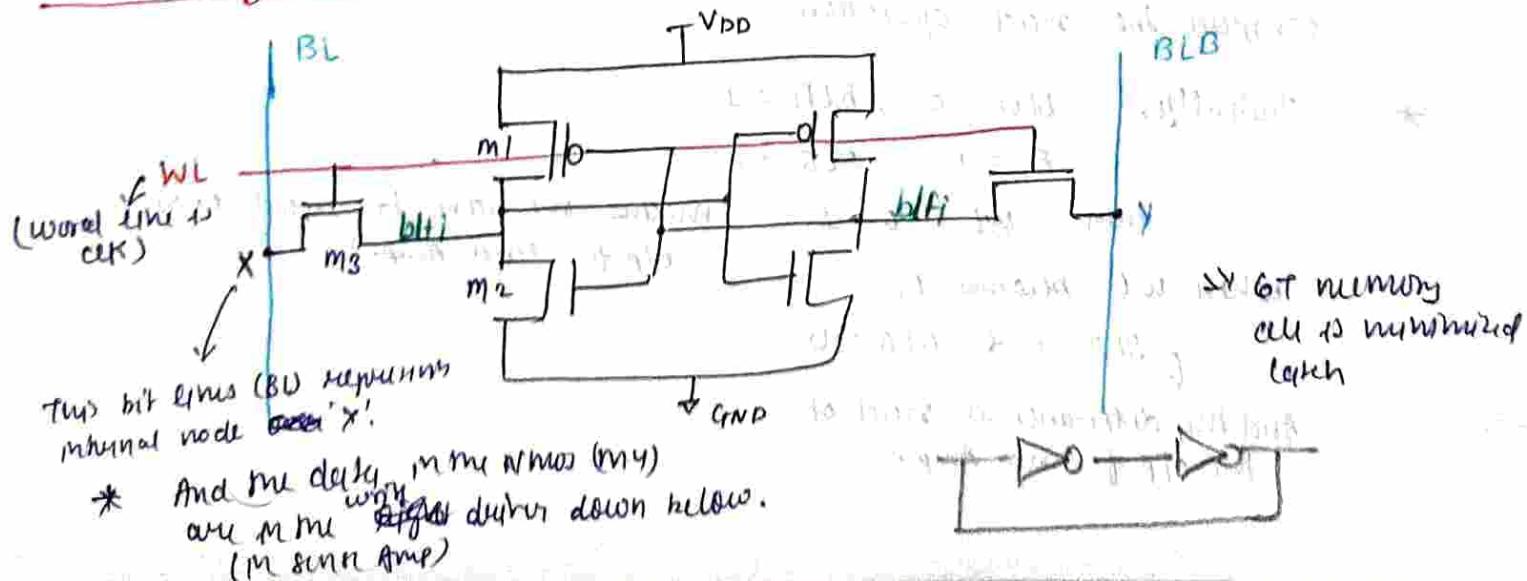
* But the challenge in designing pseudo NMOS is "sizing".
∴ we have size our MOSFET very carefully.

* When data is read from one memory cell then cell is ON for that cell & off for all other memory cell, therefore pmoses of other memory cell are not used for that time.

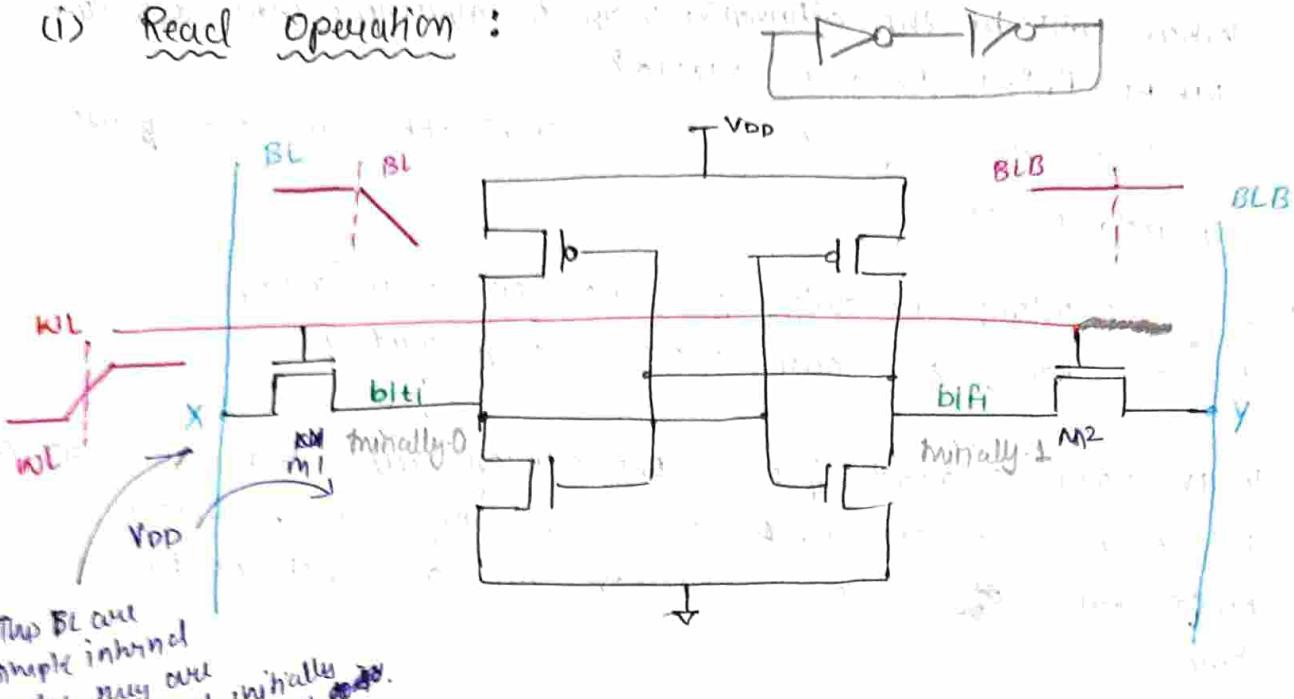
Cell operation: (Memory cell operation)



GT memory cell



(i) Read Operation:

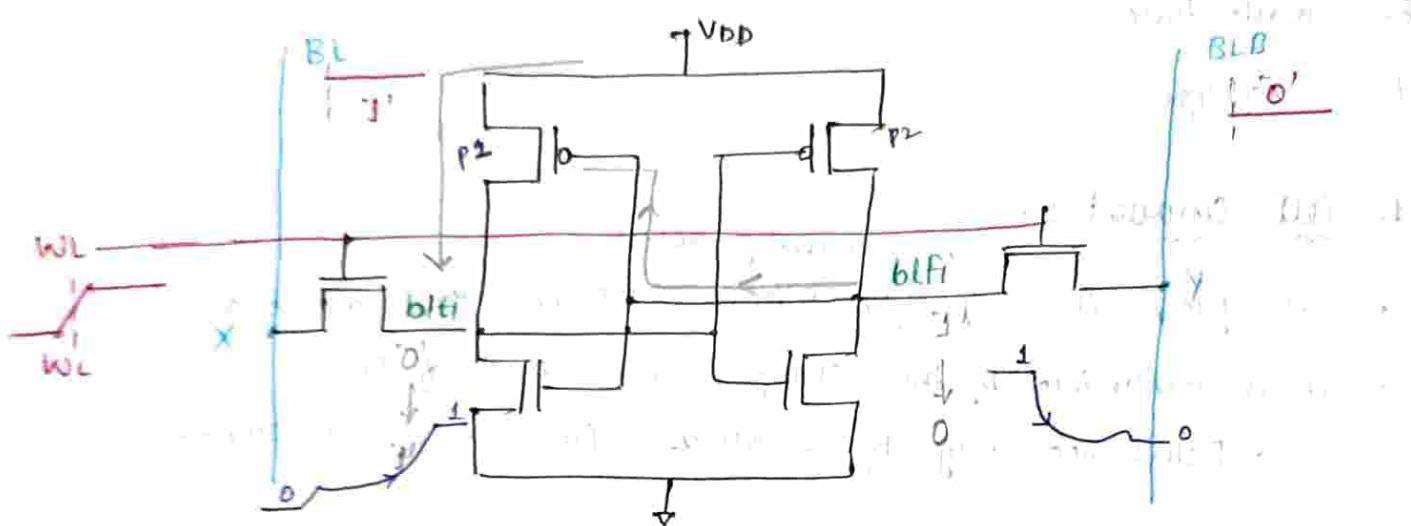


- To read into MS 6T memory cell, we first precharge both the bit lines (BL & BLB) i.e. $BL=1$ & $BLB=1$ and take one word line (WL) from 0 to 1.
- when NMOS M_1 & M_2 are turned ON due to $WL=1$. Both bit lines & internal nodes get connected.
- Node X is showing '1' & bli is having '0' initially. Therefore some current will flow from X to bli. Then BL will discharge
 - But BLB will not change
- (b/w BL & BLB)
 - This differential will be sensed at the sense amplifier to complete the read operation
- * Initially, $bli = 0$, $blfi = 1$
 $BL = 1$, $BLB = 1$
 Then $WL = 0 \rightarrow 1$ (means we have to read & at o/p to sense Amp.)
 when WL becomes 1,
 $BLB = 1$ & $BLB = 0$
 And this difference is read at the o/p of sense Amp.

Note: The side on which $(\text{BL}) = 0$, '0' was stored. But particular side, the bit like would discharge.

- * If BL discharges, '0' is read at the OLP.] or will
- If BLB discharges, '1' is read at the OLP.] never

(ii) Write Operation :



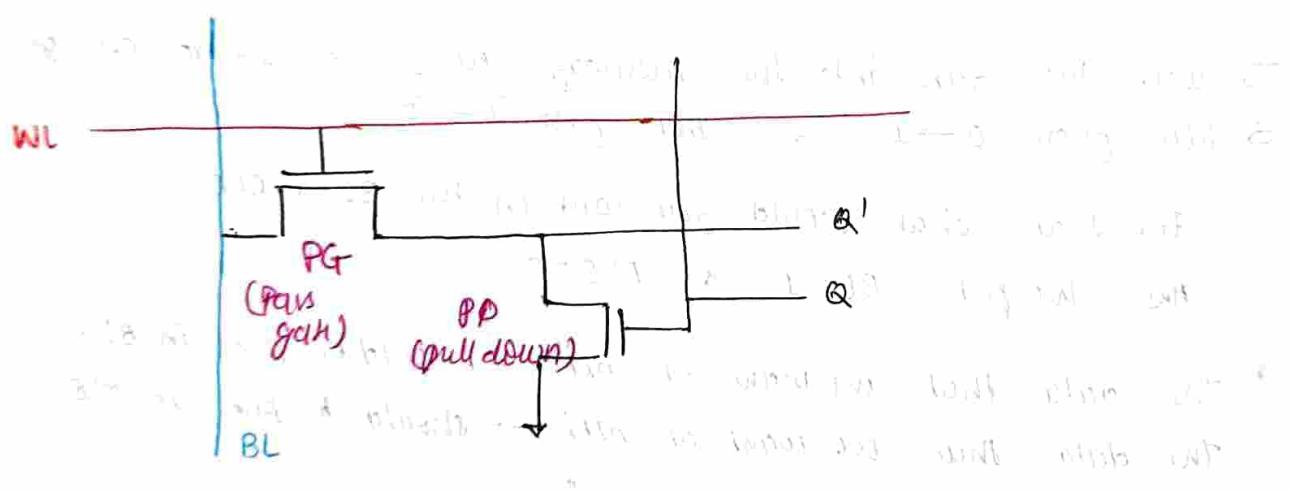
- To write the data into the memory cell, we have to change
→ bLTI from '0' → '1' + bLFI from '1' → '0'
- For that what would you work on the BL + BLB.
Ans - We put $\text{BL} = 1$ + $\text{BLB} = 0$
- * the data that we want at bLTI, should be put in BL.
The data that we want at bLFI, should be put in BLB.
- The node bLFI discharges fast and reaches '0' and on other side node bLTI will only charge a little and it will completely gets charged when bLFI becomes '0' & P2 turns on then only bLTI becomes completely '1'.
- * Writing '0' on '1' would take more time than writing '0' (on bLTI)
- Because to write a '1' you have to write a '0' (on bLFI) & then only '1' will be written (on bLTI)

Bitcell - Figure of Merits

1. Cell current \rightarrow It tells how quickly my bit line will change.
(It controls Read speed)
 2. Bit line leakage \rightarrow Power & Read speed.
 3. Cell stability
 4. Retention voltage
 5. Write margin
 6. Write time
 7. Leakage



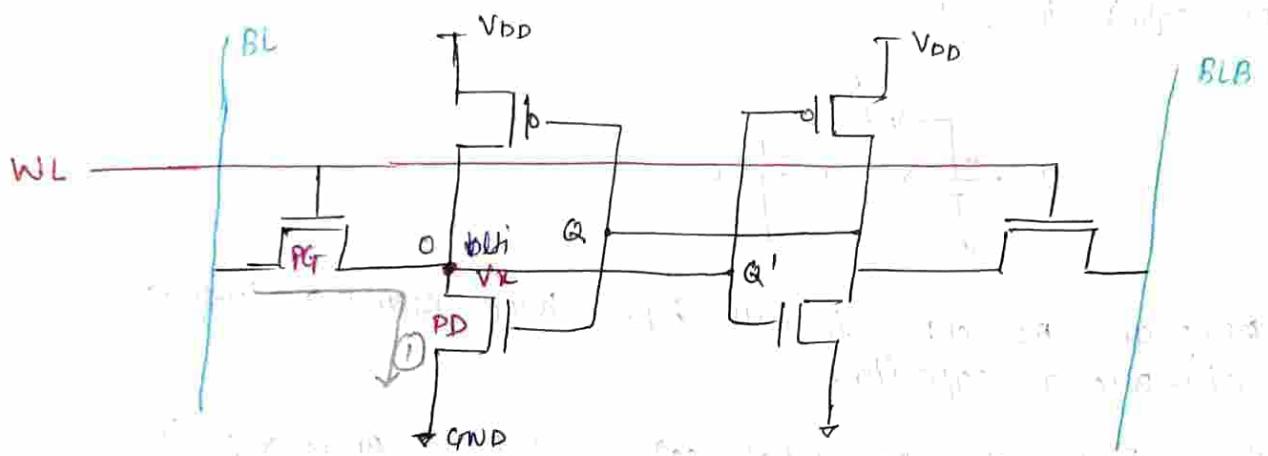
1. Cell Current :- of memory access
 - It defines the speed at which Bit line can discharge.
 - It is determined by the sizing of the stack of PMOS.
 — Determined only by 2 devices PG & PD of side storing.



- Bigger the PG & PD, the bigger the current would be.
 — hence leakage also increases, more the capacitance of BL also.
 — so we can't keep on increasing the size.

2. Cell Stability

- It is the most important specification of a storage element.
- Challenge is higher during read operation. Because during read operation you expose your latch (your storage element) to external noise.



- When we start the read operation, current starts to flow like ①.
- Initially bitⁱ was storing '0', but to ensure the current flows from 'bitⁱ' to 'GND', the value of 'bitⁱ' goes higher than '0' (let V_x).
- This V_x , turns OFF + turns ON on Rns and which makes the system unstable and lisk starts to appear.
- hence cell stability starts to degrade
- therefore, we try to ensure that V_x is as low as possible. For that we get another sizing constraint.
- PG must be much more resistive than the PD (pull down)

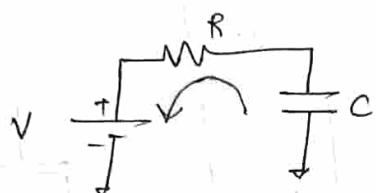
$$PG \text{ (size)} < PD \text{ (size)}$$

- Q. What will happen if PD is more resistive than PG?
- Then node 'bitⁱ' @ a will rise above $\frac{VDD}{2}$ and it will kind of flip the content of memory cell. And we don't want that.
- * When we remove the pre-charge then to do the memory cell operation we change the size of other transistors.

Positive F/B → Data flip → ERROR

3. Write Margin → metric of speed,写入速度.

- * when BL goes to '0', min node will discharge.
- Bit line is an RC circuit and you ~~can't~~ discharge a capacitor to '0' in finite time.
- How much time is required to discharge a capacitor (through RC path) to '0'?

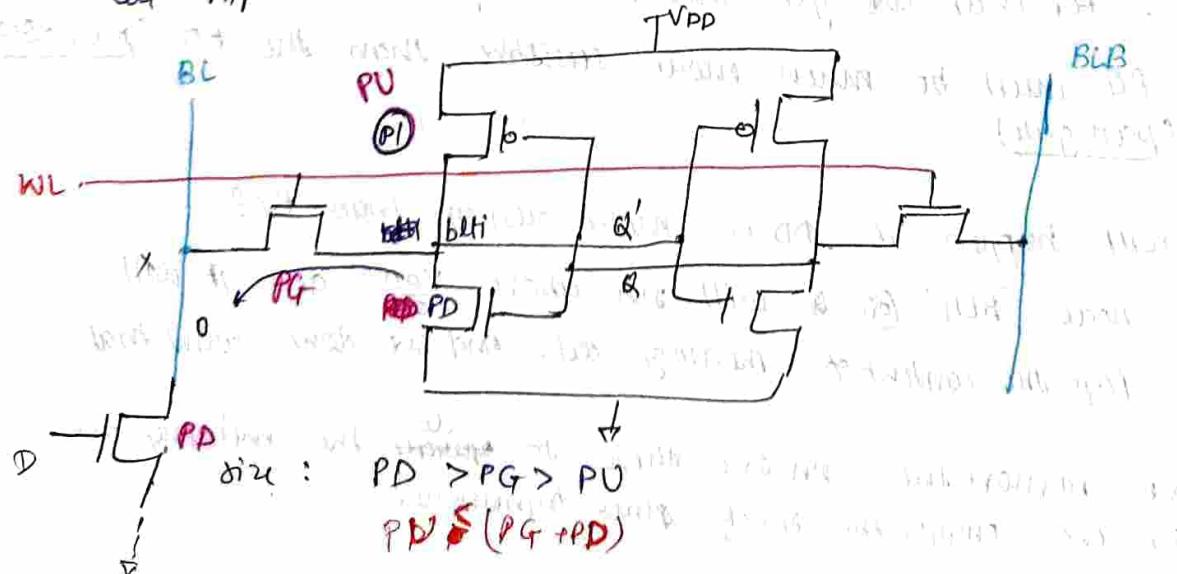


- Being an RC circuit it will take huge amount of time to discharge a capacitor.

- It is first order differential eqⁿ, therefore practically it takes infinite amount of time to discharge. But we don't have that amount of time to read & work a memory.

- * Since we can't discharge the capacitor to '0', then what is the higher level of BL at which I will still be able to write into the memory cell? → called write margin

- Now with this write margin we don't have to discharge the cap. to '0'. We have to discharge it to the point where it ensures memory cell flip with WL at V_{PP} .



- * Highest level of BL to ensure cell flip with WL at V_{PP}.
- * OR Lowest level of WL to ensure cell flip with BL at gnd.

Q: What is meant by -ve write margin?

- We take value of Bit-line voltage lower than V_{DD} so that it can't write into the memory cell.
- it Negative Bit-line write assist scheme

(*) When you want to write into the memory, that the pull up PMOS (PU) should not be able to drive & should draw lesser current than what this stack of PG + PD can sink.

- series of PG & PD should be stronger than the pull up (PU)

$$(PG + PD) > PU$$

series of PG (pull up) & write driver (PD) should be stronger than

PU.

PD should be resistive to give higher bump on other side.

PD should be resistive to give higher bump on other side.

so, write ability of the cell is inverse to its stability.

So, write ability of the cell is inverse to its stability.

4. Leakage

All off cells leak.

Even ON devices have junction leakage.

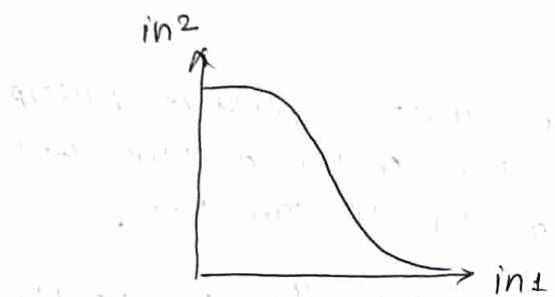
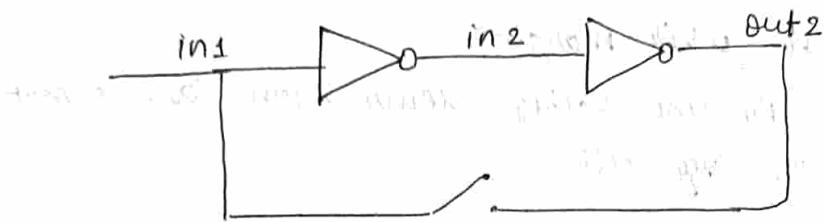
we should have small PG, PP + PU for low leakage

long channel

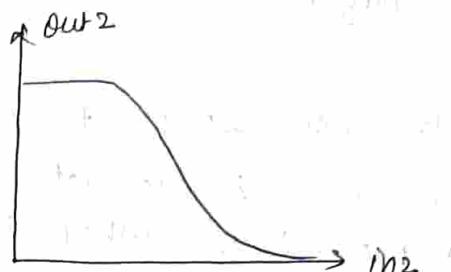
- min. length of transmission is \propto min. in memory cell length and will be 10μ @ 20μ @ 1.2μ @ 1.5μ .

thick gate oxide, less doping on source/drain.

Butterfly curve with cell stability

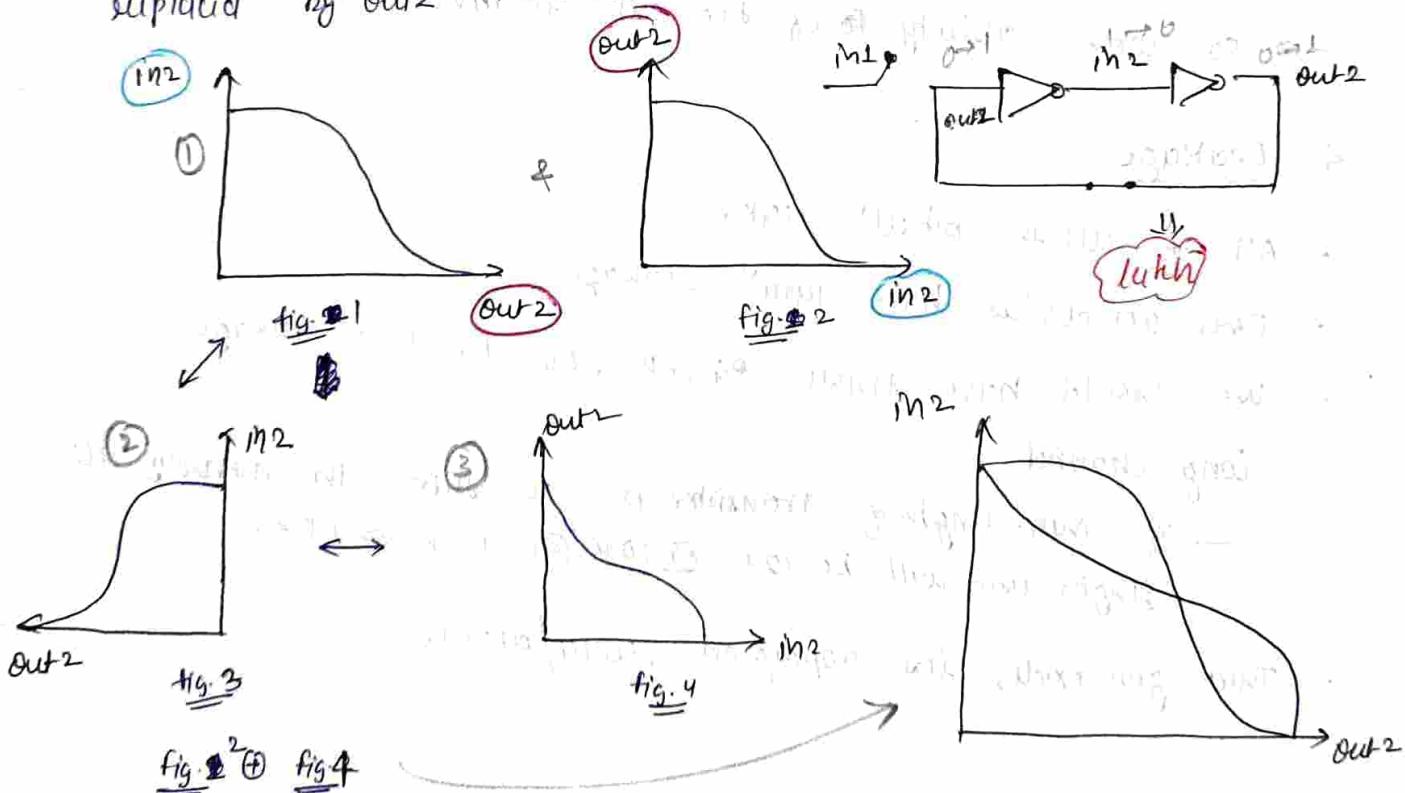


As $m_1 \oplus in_2 \oplus$



As $m_2 \oplus out_2 \oplus$

- In memory cell these two switches are connected back to back.
- Therefore we close the switch & don't give any ip to in_2 replaced by out_2 in a f/b.

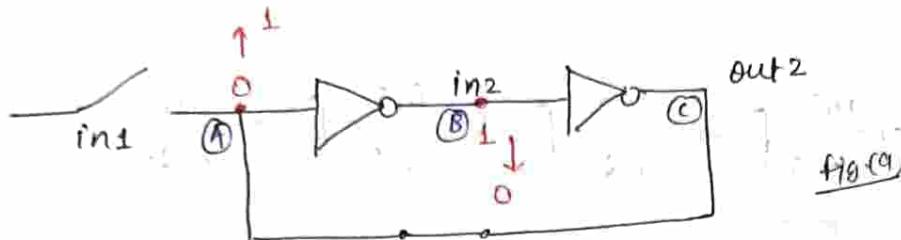


* Transfer funcⁿ of the latch is
Butterfly curve.

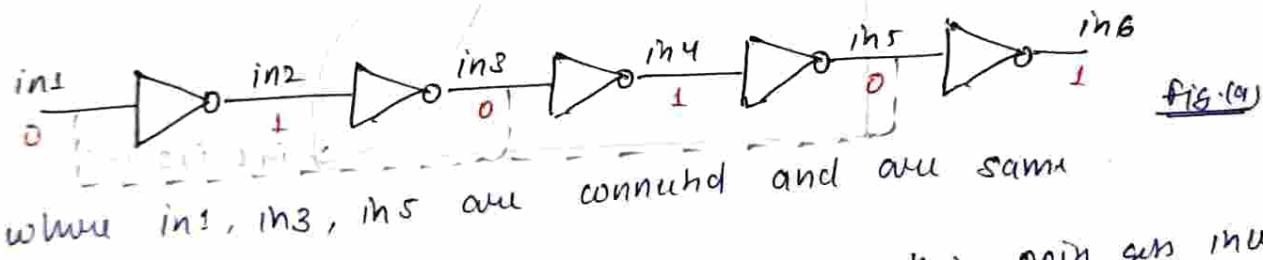
Butterfly curve
"Transfer funcⁿ of the
latch"

- Now the important thing is how much noise we can put in there ie. the cell stability.
 - the noise for which latch is immune to.

Noise injection



- If I am storing '0' at node (A) & '1' at node (B) then worth cell noise will be ~~value~~ injected when '0' goes to '1' on node (A) & '1' goes to '0' on node (B). (Worst case noise pair).
- This latch can be considered as the long chain of inverters storing 010101. [Just for sake of understanding]

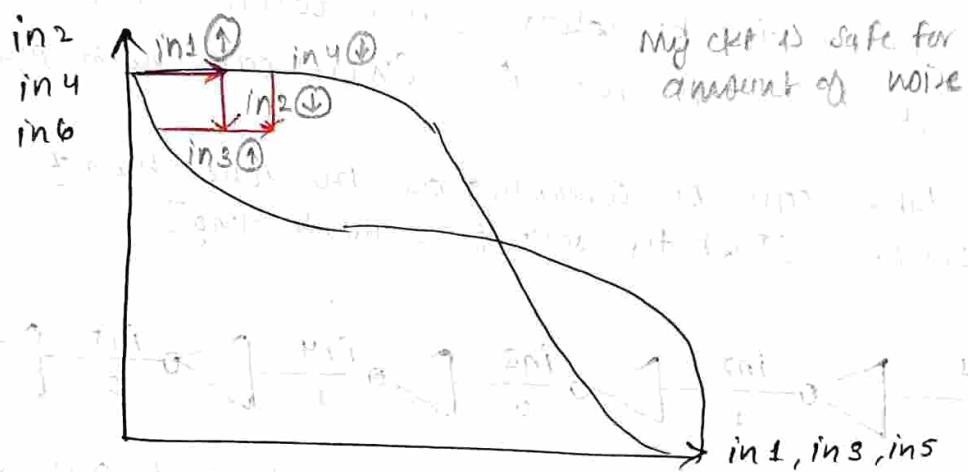
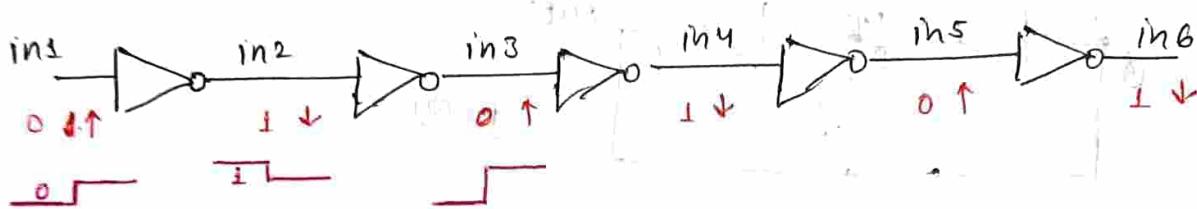
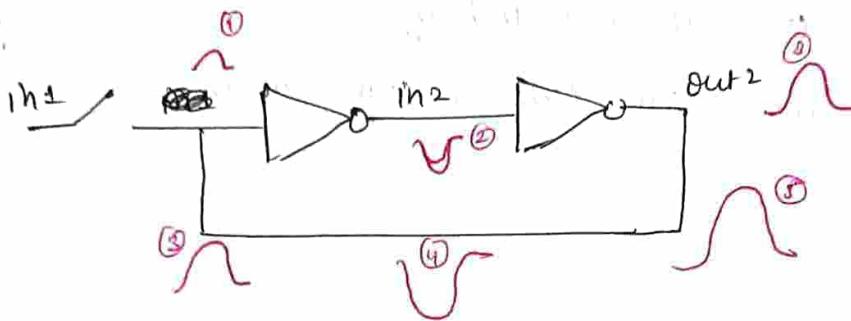


- If small noise is given at node (A) then this noise gets injected when it goes to node (B) then node (C) will have higher noise and again this noise is fed to node (A). ~~Because~~
- Because latch with amplifying action.

- So to overlap the big noise at node (C) & small noise at node (A) we use the fig.(b). Just for sake of understanding.

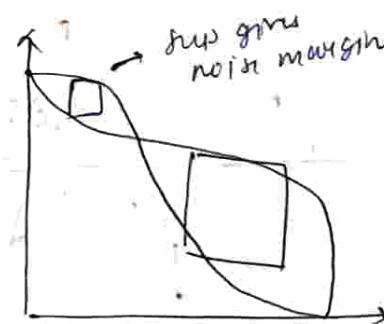
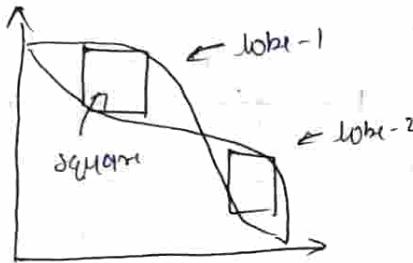
Butterfly curve with noise injection [ε Noise margin]

for noise margin



For max
Noise margin

* The largest square that we can fit into both lobes gives us the noise margin.



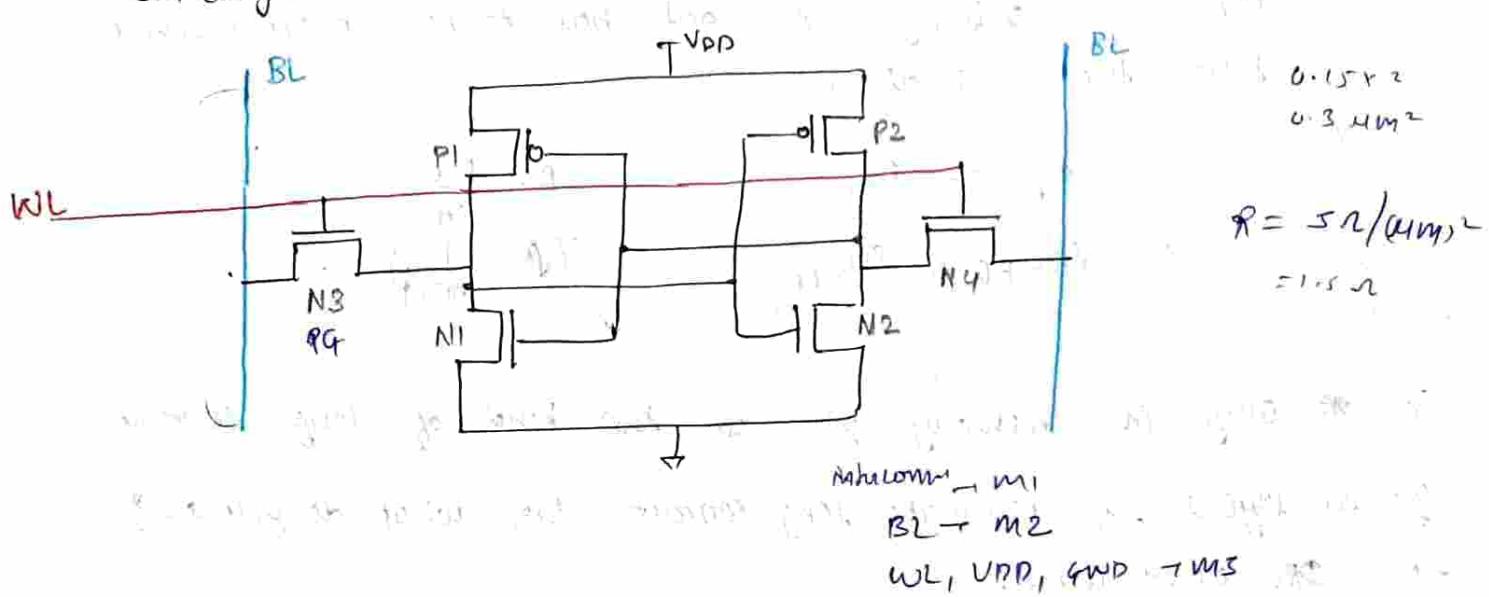
* Stability of memory cell can be measured by fitting the largest square in the smaller lobe of the hysteresis curve.

Bit-cell Layout

In standard cell layout, we have to keep a distance of $\frac{1}{2}$ DRC from the PR boundary.

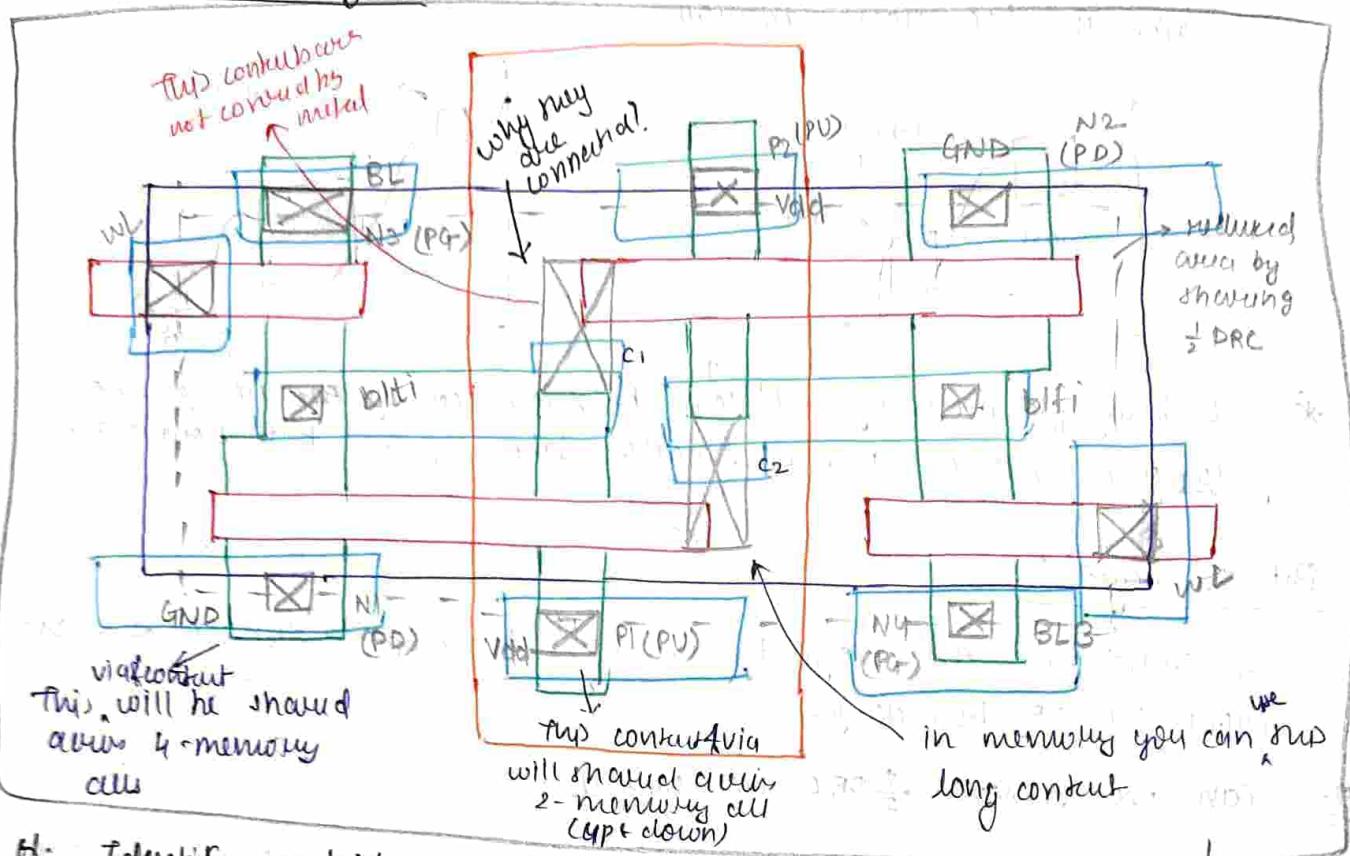
Q. Can we remove $\frac{1}{2}$ DRC in contacts?

* In memories you want to share as much around the cell boundary as possible, because you have a structured array. You know that your memory cell is connected with another memory cell only.



Advantages of this layout:
1. Reduces the number of contacts.
2. Reduces the area of the cell.

6T-SRAM Layout



Q: Identify which one is PU, PD & PG?

- the one in the Nwell is PU (P1 & P2)
- ~~pull up~~ pull down

For cell stability pu's gate has to be more resistive then the pull down

$$R_{PG} > R_{PD}$$

$$\therefore \text{size } PG < \text{size } PD$$

$$R = \rho \frac{l}{A}$$

$$(R \propto \frac{1}{\text{size}})$$

* Only in memory you can have kind of long contact.

Q: In layout, if you make long contact then what do you see?

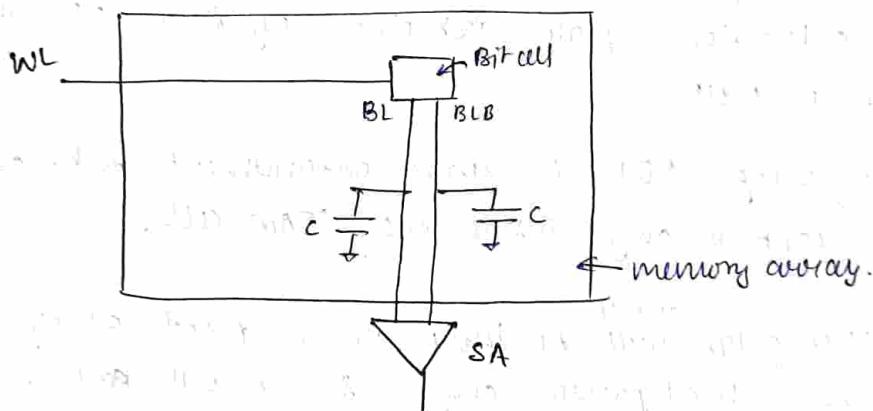
→ DRC error comes.

→ Because there is minimum contact to contact spacing for DRC. Hence it is increasing the area of layout.

- * In memory, instead of putting two contacts (to decrease the area of layout) [ie, putting one contact on poly & one contact on ~~the~~ active and maintaining the DRC b/w them].
 - We actually kind of merge two contacts and form a long contact. And this long contact is allowed to be used inside the memory.
 - Q: Why has rectangular contact only allowed in memory?
→ In Memory this is the only configuration where you allow the long contacts.
 - * The DRC is because ophs don't allow to fabricate two things.
 - In ophs you can predefine patterns then you can make the diffraction plate in such a way that you make the structure well.
 - So by using RET (Resolution enhancement technique), we allow long contacts only inside the SRAM cells.
 - * This rectangular contact will be there in the fixed array only and also in the fixed pattern only. So we will have a diffraction pattern for which I have to optimize my mask making process.
 - Then you put some metal layers above the contact to connect it to BL, BLB, VPP, GND, WL
 - (Vias & contacts) coming in the corner are shared across 4-memory cells & (vias & contacts) coming in the up-side & down-side are shared across 2-memory cell.
 - * Eg. of 2x2 array ~~and~~ (ie 4 memory cells shared across PR boundary).
 - WLI
 - WL2

Sense Amplifier (why?)

- Bit line capacitance becomes significant for large array.
 - If each cell contributes $2fF$.
 - For 256 cells: $\Rightarrow 512fF + \text{wire capacitance}$
 - pull-down resistance is about $15K$
- $t = \frac{RC\Delta V}{V_{dd}}$ [assuming $\Delta V = V_{dd}$] (ideally $t < 1\text{nsec}$)



- WL will come & select a Bit cell and that Bit cell will discharge the BL when we get the O/P.
- But this BL has very huge capacitance.
- If we want to discharge this BL which is already charged to V_{dd} . Therefore to read the data we have to discharge the BL upto $\frac{V_{dd}}{2}$ so that it can read the O/P.
- means we have to discharge $\frac{CV_{dd}}{2}$ charge from the BL to read '0' or the O/P.
- And depending on the cell current, there is a time it will take to discharge particular bit line.

* Solution: We can't easily change R, C or V_{DD}, but we can change ΔV.

$$T = \frac{RCAV}{V_{DD}}$$

• ΔV → smaller time
voltage

- It can reliably sense ΔV as small as < 50mV
the lower I discharge, the faster I can be.
- If we don't discharge it to $\frac{V_{DD}}{2}$ saves time by only discharging the BL upto $\frac{V_{DD}}{8}$ @ $\frac{V_{DD}}{10}$ (to get $T \approx 1\text{sec}$).
— So we put an amplifier down in the way, which would amplify $\frac{V_{DD}}{12}$ discharge into $\frac{V_{DD}}{2}$ full V_{DD} discharge at the o/p.
* Therefore we use sum Amplifier.