

DIGITAL INTEGRATED CIRCUITS (ENCS333)

Dr. Khader Mohammad

TR : 8:30-9:45 : Mosri08

DIGITAL INTEGRATED CIRCUITS (ENCS333)

Dr. Khaider Mohammad

T.R : 8:30-9:45 : Masrif08

L1: Introduction

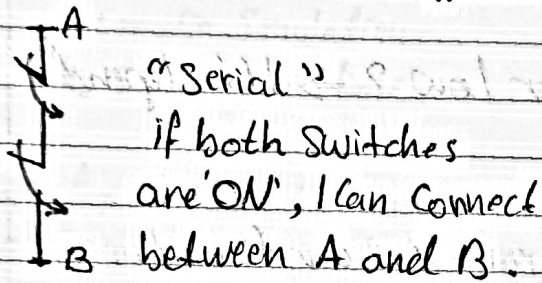
R/17-2-2019

- The smallest things that we have in ICs are transistors.
- In the IC design the main transistor we are using is CMOS technology.

* **CMOS** : Complementary metal Oxide Silicon Semiconductor

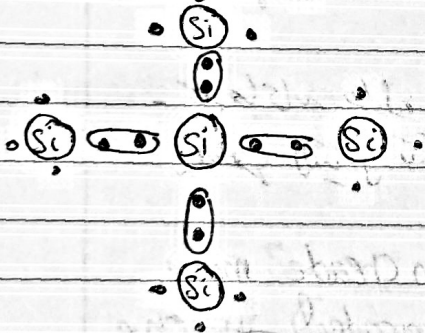
Transistor can use as switch \rightarrow ON / OFF

Combination with different switches \Rightarrow different functions.

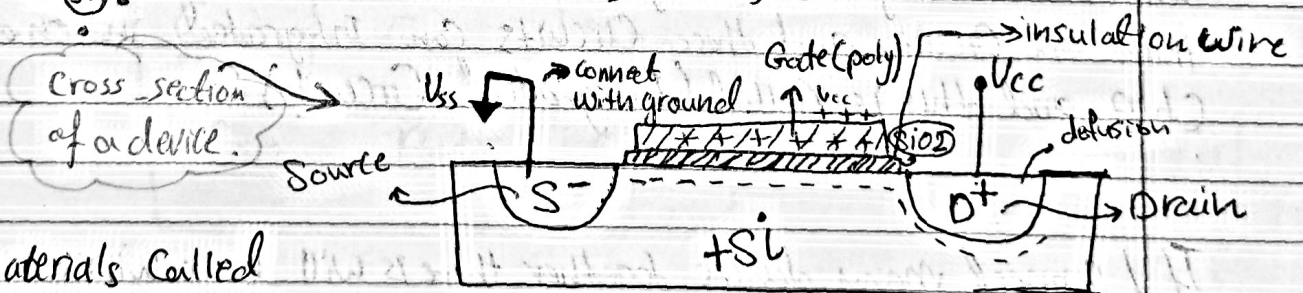


Sure, using Serial is different than using parallel!

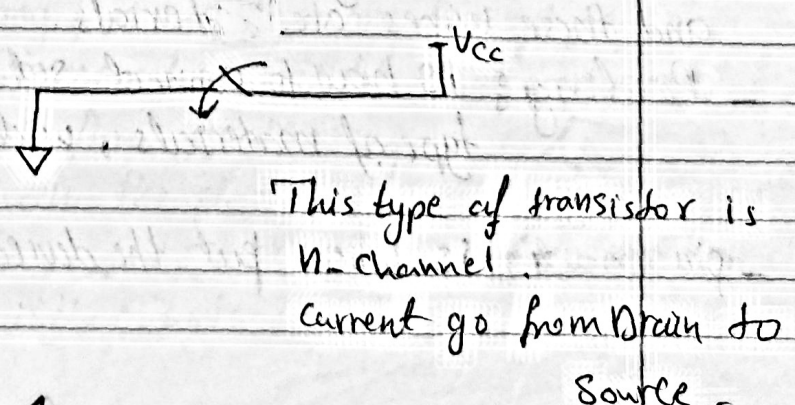
- Switch :

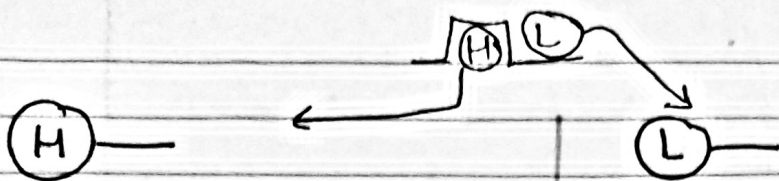


group of silicon atoms linked together.
the charge of silicon is positive (Si^+).
we build on it different materials.
materials \Rightarrow Drain, Source



These materials called diffusions with specific properties. each material have different characteristic.





The device on active high that's mean it works on high phase only (high signal).

If this device will be active-low that's mean it will work on low phase only.

- Some transistors work on high and another work on low. the transistors that work on high called \rightarrow (NMOS device), and transistors work on low called \rightarrow (PMOS device).

* How much the voltage is high or low? \Rightarrow It depends on process node

process node means the length of transistor on the device (Length of Gate).

Because this length controls the size of device and the speed of current.

كل ما كان طول الـ transistor في الـ device اقل يعني الـ transistor يكون اقل في الحجم والسرعة.

When I have huge number of transistors, I can create more than one circuit and these circuits can be integrated in one chip and this is called (Integrated Circuit).

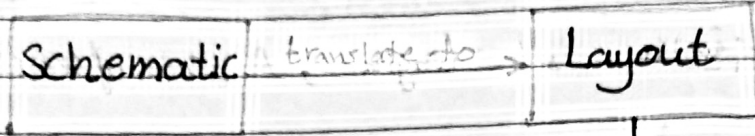
If I connect many devices together there will be a wires and these wires are materials.

- Routing: is how to connect set of wires with devices, and type of materials, width, spacing, ...

- placement: is where put the devices (cells).

Each wire or each device contains resistance and capacitance, Inductance. This is what we call "Parasitic Extraction".

(R, C, L) that's what we need to be able to run Timing.

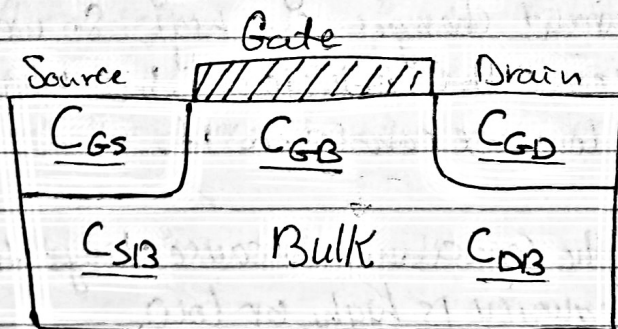


↓
 build the design then simulate it.
 we make simulation to check the time.
 ↓
 actual physical layer

- model for devices contains specific information about all these devices and wire that will be used.

If I work on 10 micro, the model of device should be 10 micro and so on.

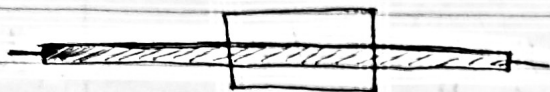
Timing: Delay between input and output.
 Slope → time to be high and time to be low.



* Capacitors: If I have any two wires or any two plates there's a capacitor between them. and it comes from the voltage difference between the plates.



* We connect chips to the outside world → By [packaging] → make an interface with outside world.



R??

Wire

* What we do exactly?

1. design combinational logic from transistors.
2. simulate it to insure that it's functionality correct - "malle"

3. layout

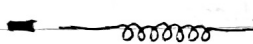
4. Timing Analysis \rightarrow Because now we have the actual design

* $V = IR$

* $R = \frac{\rho \cdot L}{A}$, ρ : resistivity , L : length , A : area of cross section.

* $C = \frac{\epsilon \cdot A}{d}$, ϵ : material between wires
 , A : Area between two wires ($L \times W$)
 , d : distance between wires.

* $\tau = R \cdot C$

 We ignore the inductors in chip level, but when we talk about top level design we need the inductance.

* $I_C = C \frac{\partial V}{\partial t}$

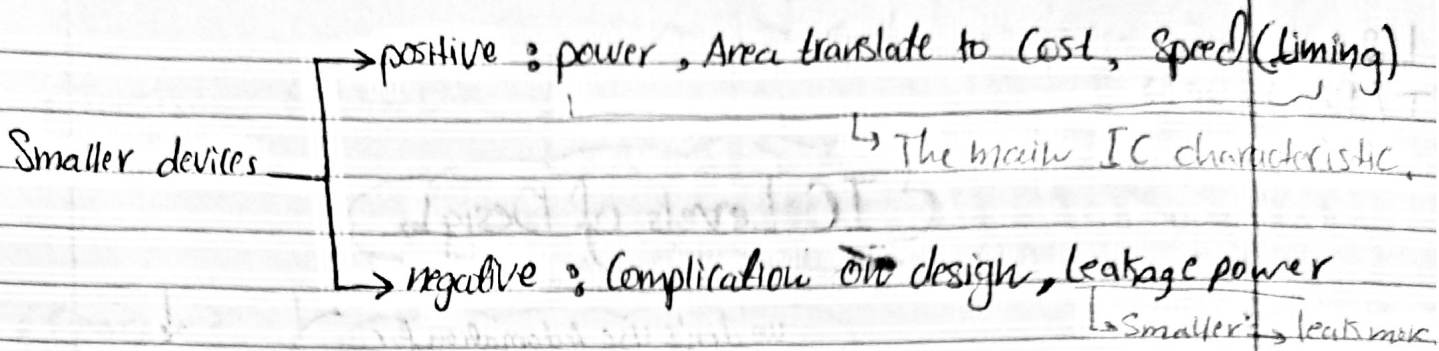
* $V_L = L \frac{\partial I}{\partial t}$

The relation between frequency and (Capacitance & Inductance) is:
 When the frequency is high then the current changes are high so we can not ignore the inductance, but when the freq. is low we can ignore it.

* But we cannot ignore the capacitance because always there is a voltage whether the frequency is high or low.

- DRC (Design Rules) : we take it from people who works on Fab \rightarrow (manufacture the machine).

∇ - Each device have numbers of transistors, and these transistors are doubles every 2 years.



- * **Wafer** : We can make thousand of chips on the same wafer.
- * **Die** : the smallest thing is on the wafer.

- The thickness of wafer depends on the process that it works on.

* number of layers that connect between devices are affect to thickness of wafer. How?

↳ answer

* When the thickness of wafer be small then it will be less expensive → less complexity because the connections between devices are few → less number of layers.

* When the thickness of wafer be large that is mean the complexity is increased because number of layers will be large to connect the devices together → more expensive.

* Design Abstraction levels

- System Architecture.
- Module
- Gate : build the gates we need.
- Transistor level
- device & layout design

- we cut the die from the wafer then we put it in package which has pins to connect with outside world, and those pins should be fast so we care about the material of wires.

[P should be very Low] like Gold.

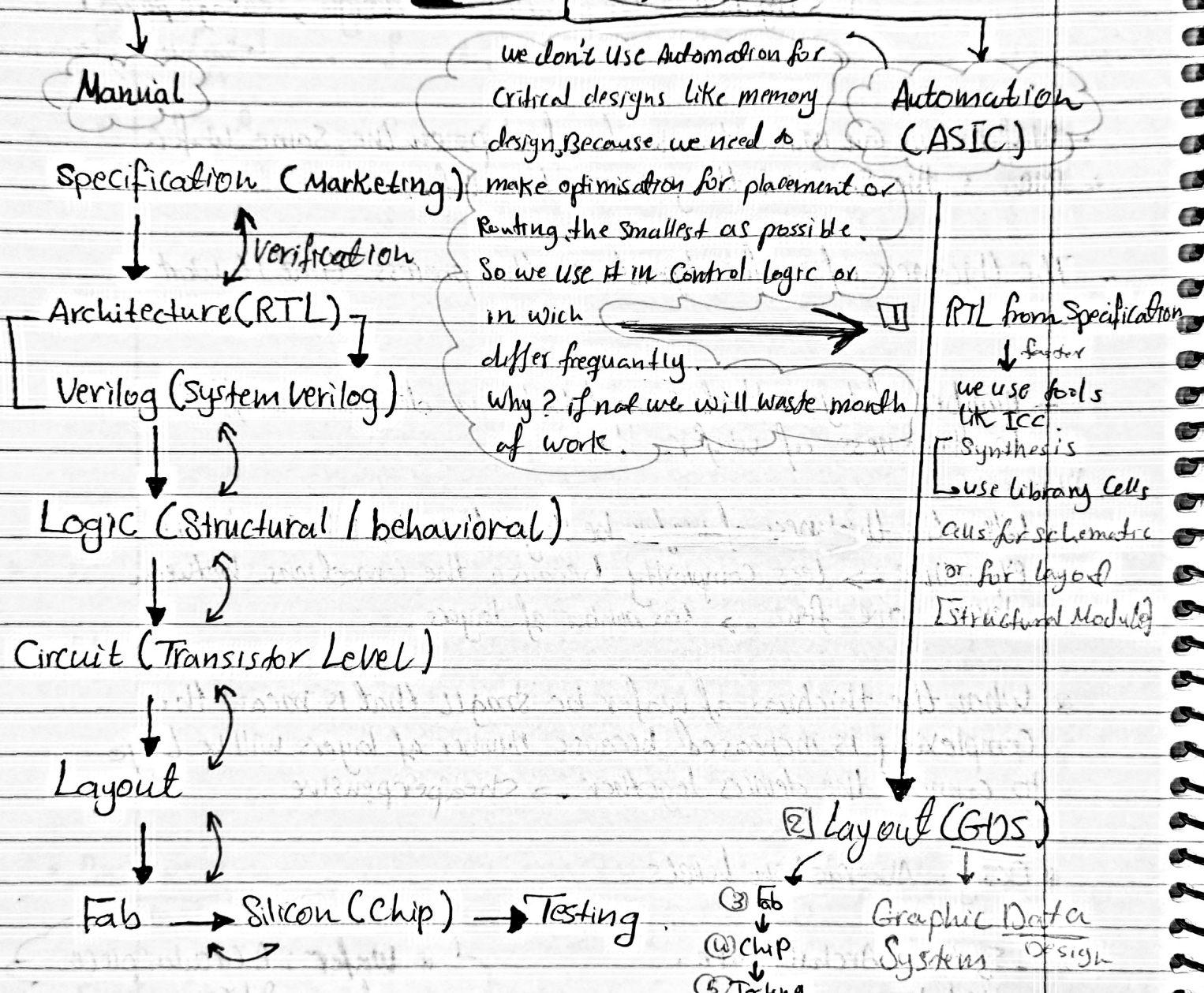
Less resistance

* **Wafer** : A circular piece of pure Silicon.

* **Die** : A rectangular piece of silicon that contains one IC design.

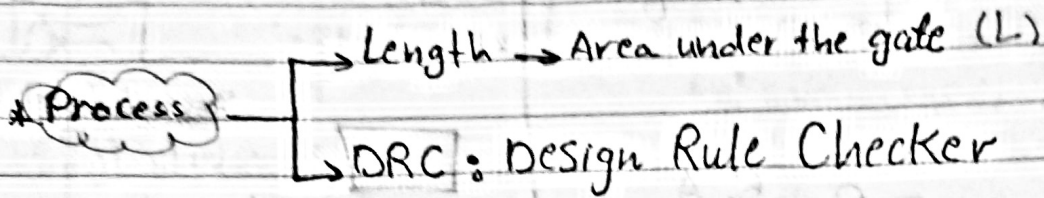
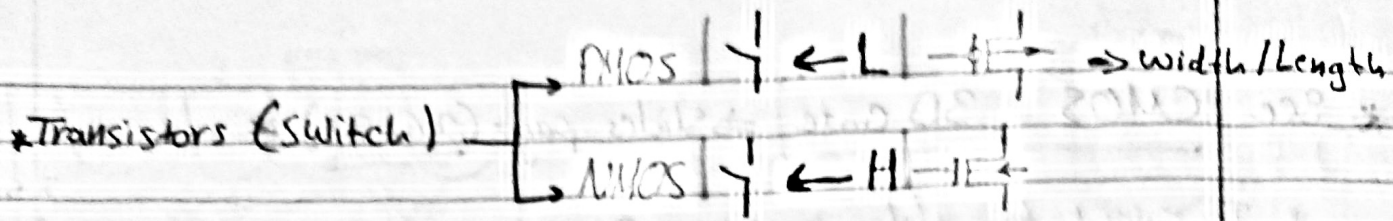
L2:
T/12-2-2019

IC Levels of Design



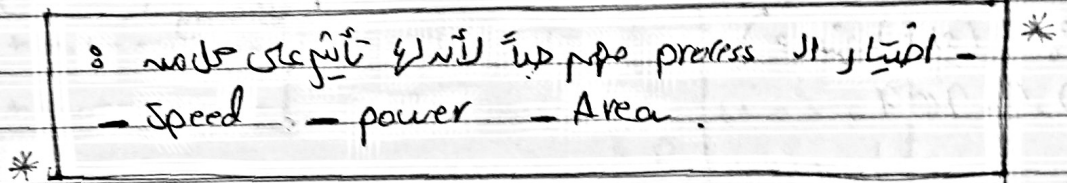
** The meaning of these steps are on slides

ICC ⇒ This is a tool takes RTL then do the Routing and placement.



* The three basic DRC checks: width, spacing, Enclosure.
 - Each process have different set of DRC.

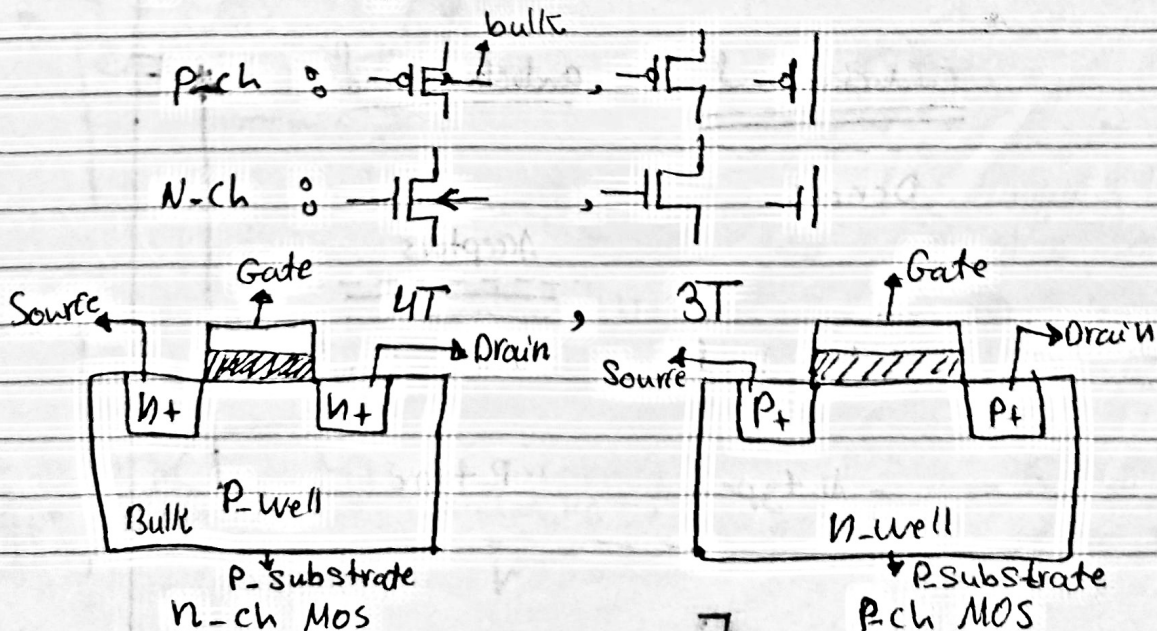
$L = 45 \text{ nm}$ \rightarrow easier to make a design than $L = 10 \text{ nm}$.



* When we choose $L = 10 \text{ nm} \Rightarrow$

- less area
- Speed: faster
- Cost: Less
- Complexity will increase

* we get the design rule from Fab.



* See CMOS, 3D Gate \Rightarrow Slides page (26+27) *

* To calculate delay in wire $\tau =$

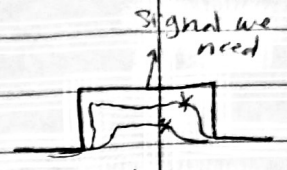
$\tau = R \cdot C$

$R = \frac{\rho \cdot L}{A}$, $C = \frac{\epsilon_0 \cdot A}{d}$

$\Rightarrow \tau = \frac{\rho \cdot L}{A} \cdot \frac{\epsilon_0 \cdot A}{d}$

$\Rightarrow \tau = \frac{\rho \cdot \epsilon_0 \cdot L}{d}$

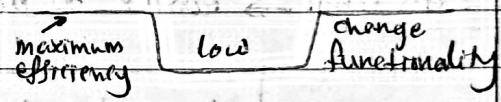
* Reliability:



So we do liability verification.

\rightarrow check if the device still on long time it will function the same [every transistor can switch

Complimentary] In Reliability we use "Bathtub" Curve.



describe how the device work

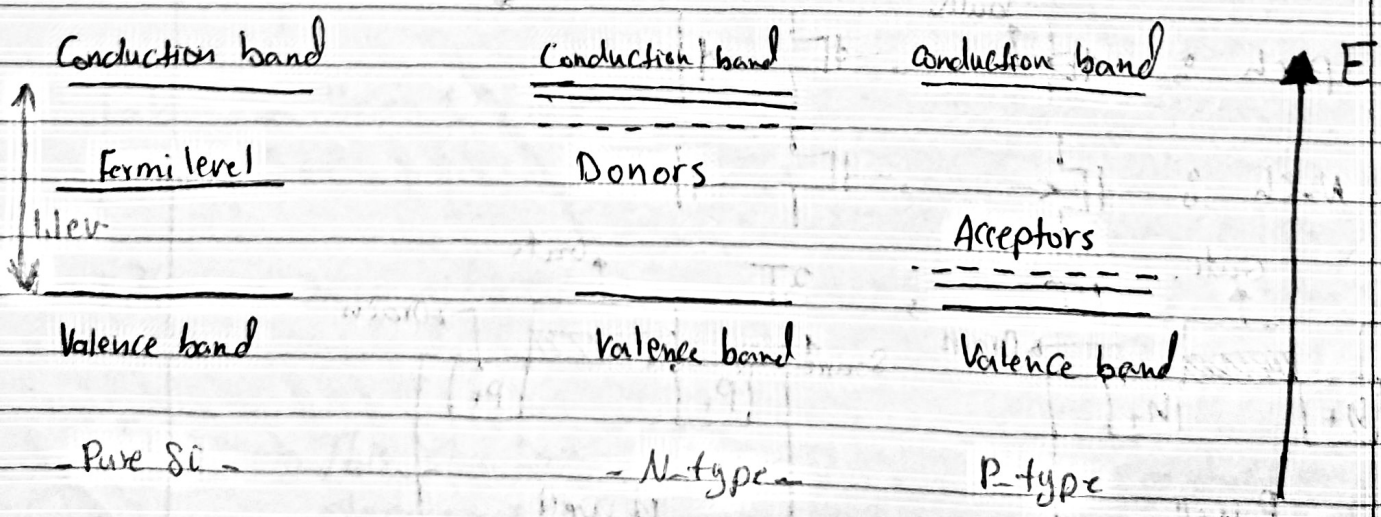
L3: (جهاز الذاكرة)
T/12-2-2019

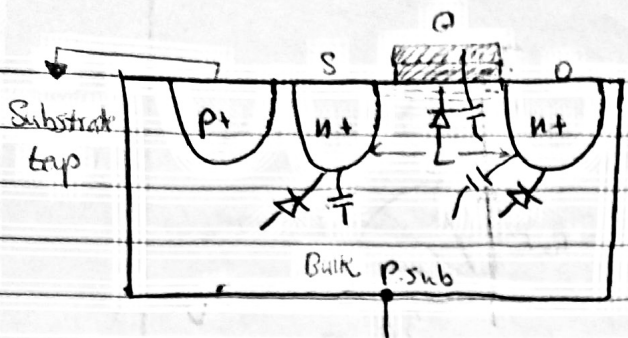
The Semiconductor

NMOS \rightarrow N-type \rightarrow free electrons (-ve charge)
PMOS \rightarrow P-type \rightarrow holes (+ve charge)

Slide 5
after Doping

* Energy gap: (slide 6 on lect. 2)



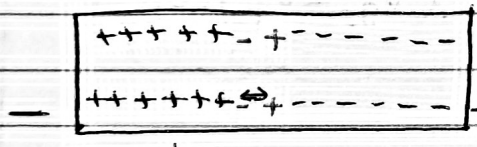
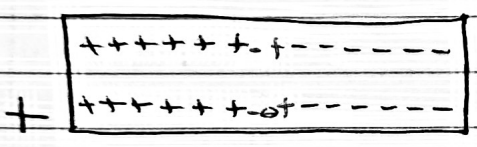
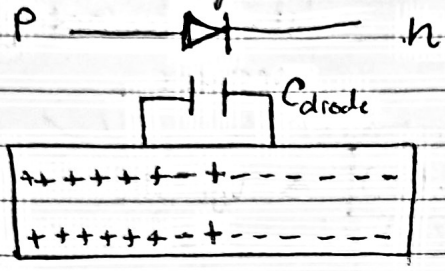


The Mos Transistor (Slide 11)

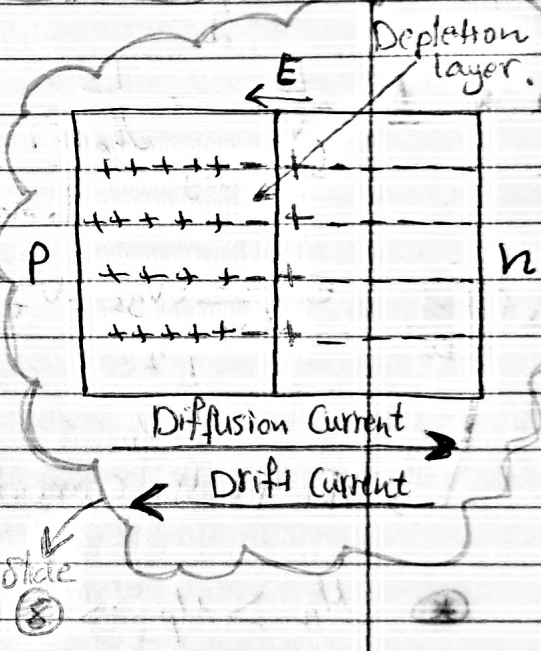
← n-channel device

slide 11

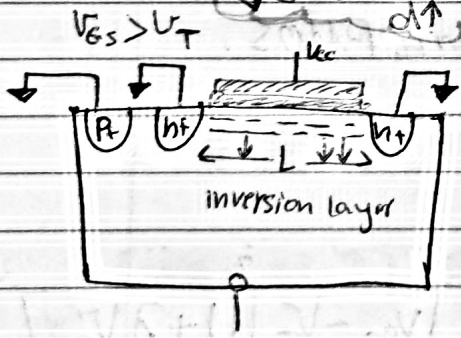
(Slide 10) Depletion Capacitance



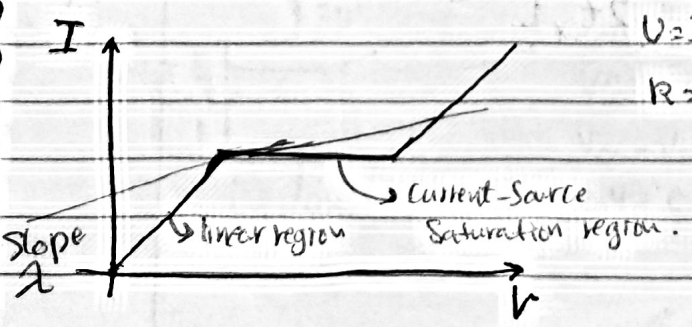
- The Semiconductor
- Silicon crystals is composed by atoms with 4 valance electrons.
 - Diamond lattice
 - No free charge
 - very high resistance.



$$C = \frac{\epsilon_0 A}{d}$$



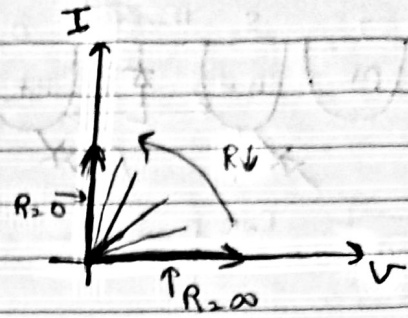
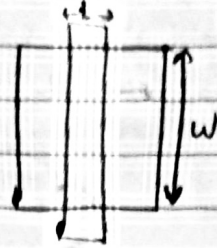
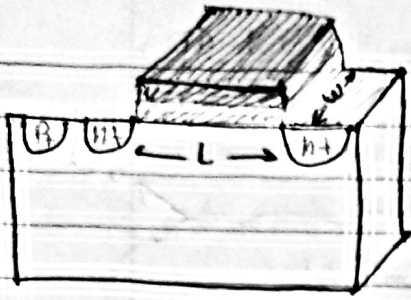
U_T is minimum Voltage we apply on gate to get inversion layer.



$$V = IR$$

$$R = \frac{L}{\mu}$$

→ $W, \mu, t_{ox}, L, \epsilon_0$
width mobility thickness of channel
these all affect the current in the channel.



→ L increases so R is decreased.

in linear region the device acting as resistor.

$$i = \frac{dQ}{dt} = C_{ox} \cdot \frac{dV}{dt}$$

$$C_{ox} = \frac{\epsilon (WL)^{area}}{t_{ox}}$$

$$V = (V_{gs} - V_t) \Rightarrow V_{gs} > V_t$$

$$t = \frac{L}{v} = \frac{L}{\mu_e \cdot E} = \frac{L^2}{\mu_e V_{Ds}}$$

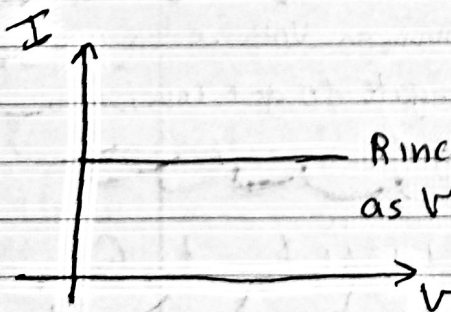
electric field = $\frac{V_{Ds}}{L}$

$$I_{DS} = \frac{\epsilon WL (V_{gs} - V_T)}{t_{ox} L^2} \cdot \frac{\mu_e V_{Ds}}{L} = \frac{\mu_e \epsilon W}{t_{ox} L} (V_{gs} - V_T) \cdot V_{Ds}$$

→ This is the current in linear region.

● $\mu_n > \mu_p \Rightarrow$ so always NMOS faster than PMOS.

↓
more current



$$I_{DS} = \frac{\epsilon W \mu}{2 t_{ox} L} [V_{gs} - V_T]^2 [1 + \lambda V_{Ds}]$$

in Saturation region.

* Leakage Current: $I_{DS} = I_0 e^{\frac{V_{GS}}{V_T}}$ @ cut-off also @ weak inversion.

[low process] \rightarrow Smaller device $\rightarrow I_0$ is higher, V_T is lower
 So I_{DS} is higher, which affect the power analysis.

* If we need to increase the speed of the device, we usually play with W and V_{GS} .

* $V_{GS} > V_{TH} \Rightarrow$ inversion layer

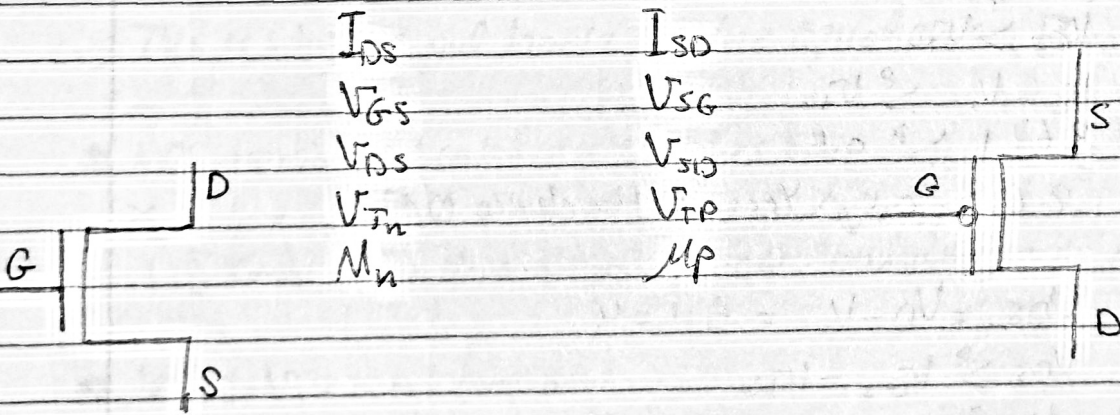
* $V_T(N) \Rightarrow '+'$, $V_T(P) \Rightarrow '-'$

L4:

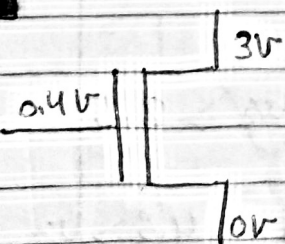
R/14-2-2019

NMOS

PMOS



* Examples: $[V_{TH} = 0.5]$

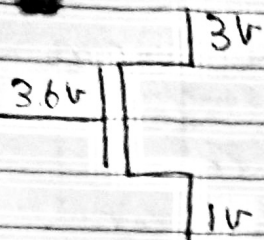


\Rightarrow ON/OFF?

$V_{GS} \stackrel{?}{>} V_{TH}$

$V_G - V_S \stackrel{?}{>} V_{TH}$

$0.4 - 0 \stackrel{?}{>} 0.5 \Rightarrow$ No, so the device off.



→ ON/off?

$$V_{GS} \stackrel{?}{>} V_{th}$$

$$V_G - V_S \stackrel{?}{>} V_{th}$$

$$3.6 - 1 \stackrel{?}{>} 0.5$$

$2.6 > 0.5 \Rightarrow \text{yes} \Rightarrow \text{So the device ON.}$

$$V_{GS} = 2.6$$

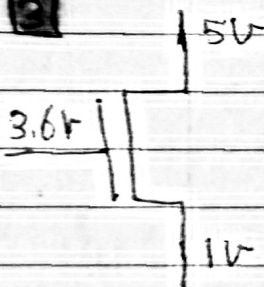
$$V_{th} = 0.5$$

$$V_{DS} = V_D - V_S = 3 - 1 = 2.$$

$$V_{DS} \stackrel{?}{>} V_{GS} - V_{th}$$

$$2 \stackrel{?}{>} 2.6 - 0.5$$

$2 \not> 2.1 \Rightarrow \text{No} \Rightarrow \text{device in linear mode.}$



→ ON/off?

$$V_{GS} \stackrel{?}{>} V_{th}$$

$$V_G - V_S \stackrel{?}{>} V_{th}$$

$$3.6 - 1 \stackrel{?}{>} 0.5$$

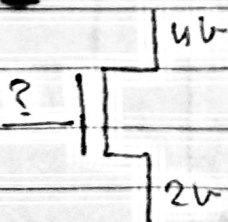
$2.6 > 0.5$. Yes. So the device ON.

$$V_{DS} = V_D - V_S = 5 - 1 = 4$$

$$V_{DS} \stackrel{?}{>} V_{GS} - V_{th}$$

$$4 \stackrel{?}{>} 2.6 - 0.5$$

$4 > 2.1 \rightarrow \text{Yes, device in Saturation.}$



what is V_G make
the device in
Saturation?

$$V_G - V_S > V_{th}$$

$$V_G - 2 > 0.5 \quad \text{--- (1)} \Rightarrow V_G > 2.5$$

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

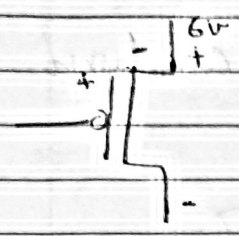
$$2 > V_G - 2 - 0.5 \quad \text{--- (2)} \Rightarrow 4.5 > V_G$$

$$\Rightarrow 4.5 > V_G > 2.5$$

-12-

* What is V_G make the device on edge?

$V_G = 4.5$

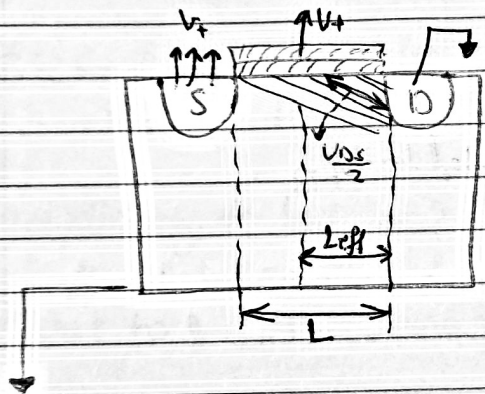


$V_{SG} > |V_{th}|$ check out/ off

If on's

check $V_{SD} > V_{SG} - |V_{th}|$

Yes \rightarrow Saturation, No \rightarrow linear



"Pinch-off"

channel length modulation

\Rightarrow This is why we put λ in the equations [slope]

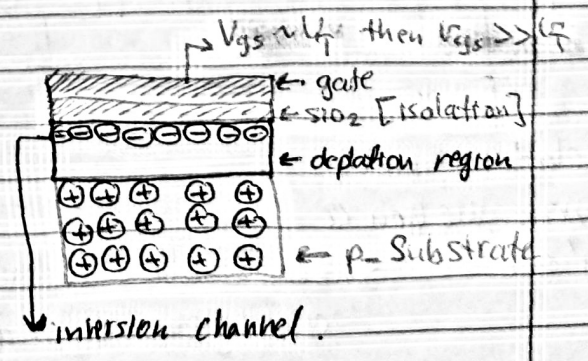
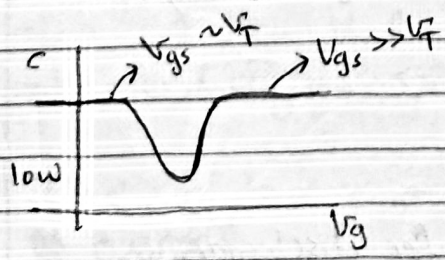
* MOS transistors treated as switches

\swarrow or switch with resistance \searrow

[Come from linear region since transistor acts as resistor]

• $V_{GS} > V_{th} \Rightarrow$ inversion.

• Capacitance in depletion region varies due to V_{GS} and its area [channel]



* Why it is limited?
 this limit based on V_{is}

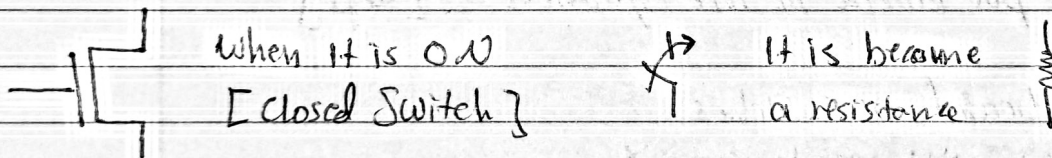
* V_T :- the minimum voltage you need to apply so you can turn the device ON

* Depends on :- - Bulk charge [voltage on substrate]
 - Bulk doping - energy [P or N device] - surface
 - Charge [SiO₂] - flux

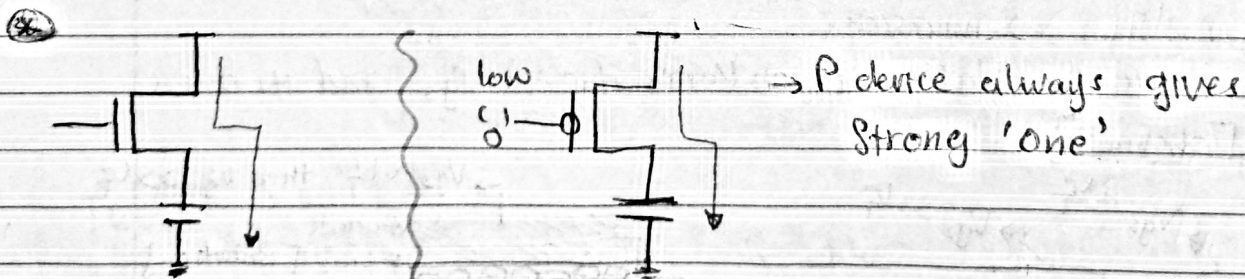
* How to measure?

we need :- V_{FB} "forward bias"
 Q_B "charge" → related to V_{FB}
 C_{ox} "capacitance"
 ϕ "of doping"

* What is R_{on} [R_{off}]?



* in Saturation → $R = \infty$ "the device become current source"

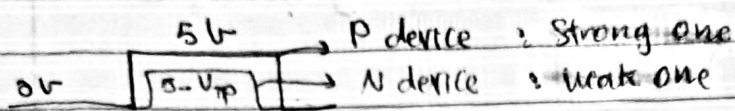


the electrons which move never give you a strong 'one'

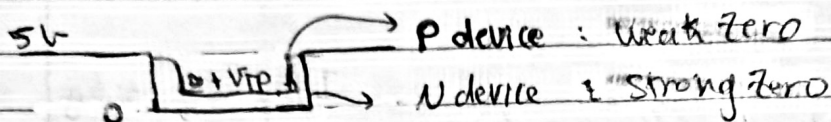
→ so N device always used for discharging and gives Strong 'Zero'

* What this mean?

T and we need to charge it.



V_{TP} : $V_{\text{threshold}}$ for p device



* AS the device grew in size [higher speed] But [rising time & discharge time differ]

* What is gate leakage?

Leff will be smaller from the one we draw in layout due to (1) the channel (2) the overlap with D&S.

* As V_D is higher the drain depletion region increases causing a decrease in $V_T \Rightarrow V_T(V_D)$: V_T function of V_D .

* low $V_T \Rightarrow$ good for speed.
 \Rightarrow bad for leakage.

* How does the temperature affects the performance of devices?

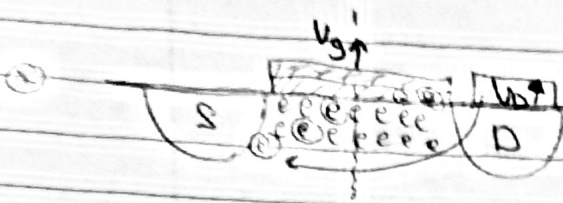
μ_n, μ_p are functions of temperature $T \uparrow$ then $\mu \downarrow$ [slower] also $T \uparrow$ then $V_T \uparrow$ [more slower] & we need more voltage to make the device ON, so we need more power & more area.

\rightarrow So we need to consider & -

- ON/off
- α
 - \rightarrow linear
 - \rightarrow Saturation

which all affect the design matrices.

● T, V_T, μ .

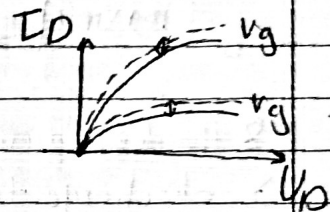


"Reliability"

⊗ ⊙ "Hot electrons". It has energy so instead of going to Drain it goes to the oxide. [Effect V_T] [slower] affect the Mobility (μ).

- ⊗ Bias Temperature "Pch BT"
- ⊗ V_T stability.
- ⊗ gate stress.

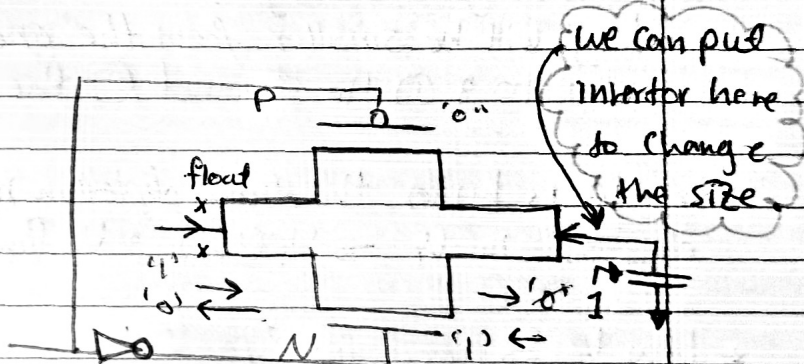
→ all these variables make the current less.
 → and as V_g higher → differences is greater.



L5:

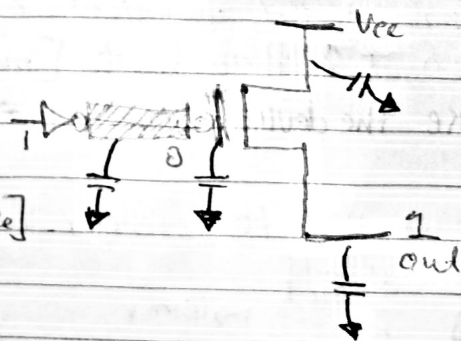
T/19-2-2019.

* Pass-gate :-
 2-switches
 PMOS & NMOS to
 Pass Strong one &
 Strong zero to the
 device.

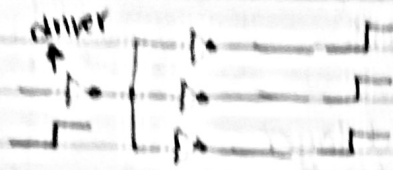


Problems :- ① if I have a float node, it will charge & affect the input. to solve this we can put inverter to prevent the returning data from influencing the input.

* Cloud :- These capacitors consider My load for the inverter load → [wire, gate, drain, source]



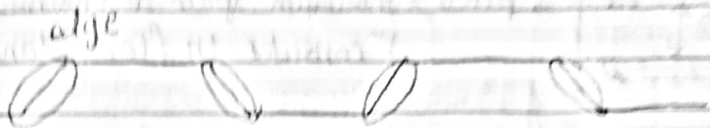
* fanout 2.



It shows how many outputs we can connect to the driver, and the outputs with the same slope as the driver.

The driver depends on the load which we need to decrease it

* edge rate 2

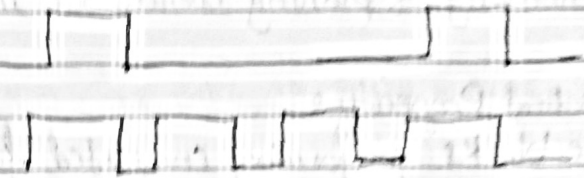


The time it needs to pull-up or pull-down. also this time increases when the load is higher

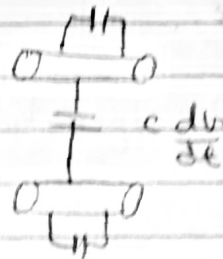
* AF "activity factor":

How the signal change with respect to the clock

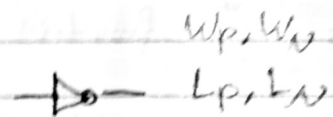
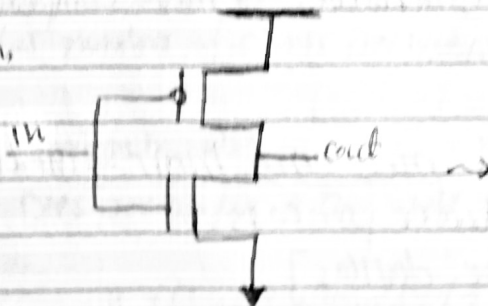
→ I need to reduce that.



* Coupling Capacitor: magnetic field about wires.

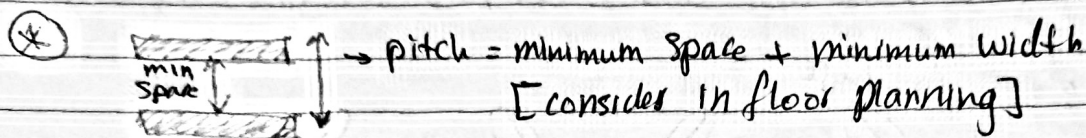


* CMOS: "inverter"



we need them to find $[R, L, C]$ → find load.

- ⊗ Size of a Component depends on W, L .
- ⊗ as the process is smaller \rightarrow more layers.
- ⊗ Via : to connect metal layers.
- ⊗ Contact : to connect diffusion with a metal layer.
- ⊗ as we go up in layer structure, the width of a layer is higher. [Less resistance].



⊗ Why we do Standardization?

- Utilization for the area.
- Routing & planing would be much easier.

⊗ Intellectual Property :

large blocks performing completed functions.

⊗ PLL \rightarrow generate CLKs.

L6 :

R121-2-2019.

"Productivity"

* Technology shrinks by 0.7 / generation. \rightarrow more complexity
 \rightarrow cost of a function decrease.

* How to design chips with more and more functions?

\rightarrow we go to lower and lower process
[flexibility to add more devices].

⊛ What is foundry?

Company that do the manufacturing not do the design but some do the both.

- * mixed signal \Rightarrow digital + analog
- * photo-lithograph: printed light \Rightarrow they use light to print materials.
- * masks of metal layers.



Scaled factor (2): if I want to change the scale of the design instead of changing all the design. I change lambda only.

- * wiring \Rightarrow micro
- * process \Rightarrow nano
- * active material \Rightarrow diffusions.

* non-recurrent cost: cost you once.

* The cost related to wafer and single die [Area]

- Yield & Defects -

$$Y = \frac{\text{number of good chips per wafer}}{\text{Total number of chips per wafer}} \times 100\%$$

$$\text{Die cost} = \frac{\text{wafer cost}}{\text{Dies per wafer} \times \text{Die yield}}$$

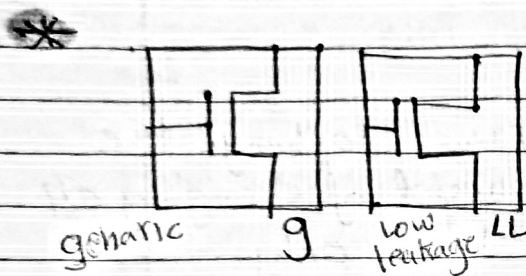
$$\text{Dies per wafer} = \frac{\pi \times (\text{wafer diameter} / 2)^2}{\text{die area}} = \frac{\pi (\text{wafer diameter})}{\sqrt{2 \times \text{die area}}}$$

$$\text{die yield} = \left[\frac{1 + \text{defects per unit area} \times \text{die area}}{a} \right]^{-a}$$

⇒ a is approximately 3

$$\text{die cost} = f(\text{die area})^4$$

⇒ Circularity wafer ⇒ Cost less [Lower waste]



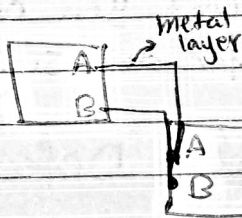
* for same process Ex 90nm

- lower leakage power
- ↳ higher V_T
- ↳ slower [long critical path]

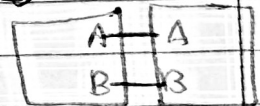
- * always in speed path devices, we don't use LL
- * low voltage devices used in memory

Power ⇒ (1) leakage · (2) dynamic · (3) short circuit [static]

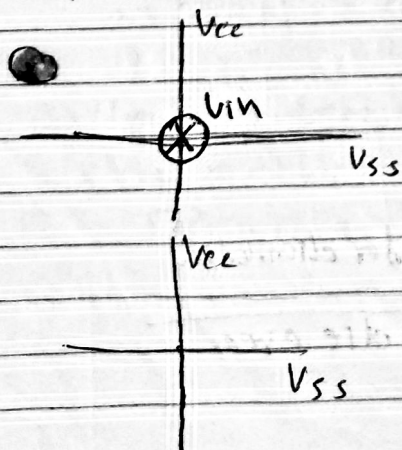
alignment & instead of → So we can save metal layers.



we do ⇒



poly head : It is the metal to poly contact.

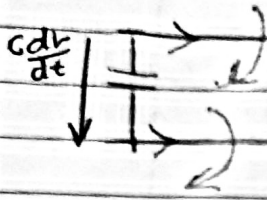


you short Vee! you need a hard time to find it.

this called "open", [not connected].

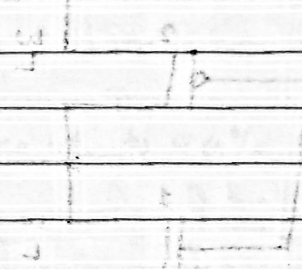
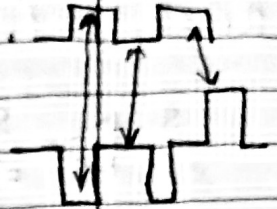
electric field

$dv \rightarrow$ Voltage difference between wires

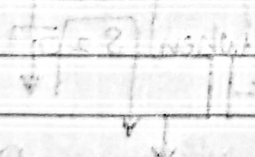


← mutual coupling
 \equiv noise

\Rightarrow we put shield [usually for dlc signals]



$2\pi \times 10^{-7} \times I$
 $2\pi \times 10^{-7} \times I$



L7: Inverter Design

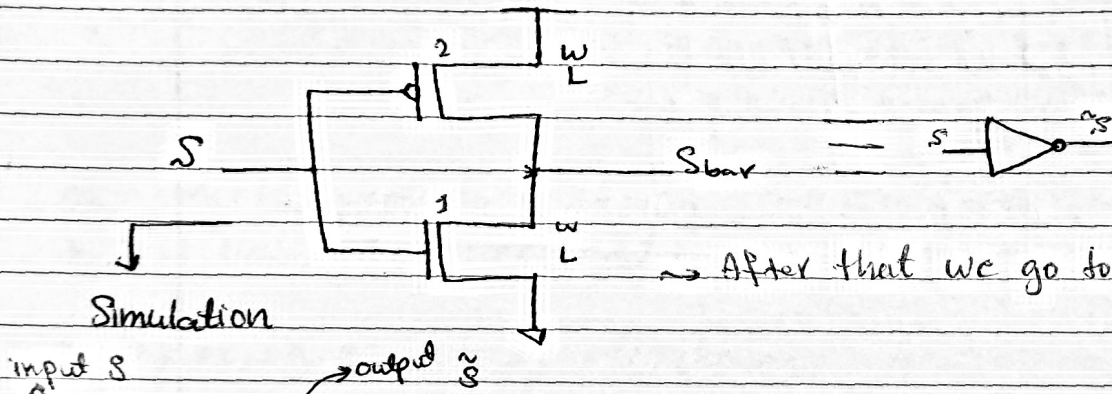
26-2-2019 / T

$S \rightarrow \text{Inverter} \rightarrow \begin{matrix} S_{bar} \\ \approx \sim S \end{matrix} \Rightarrow$ we need to model it in Behavioural mode.

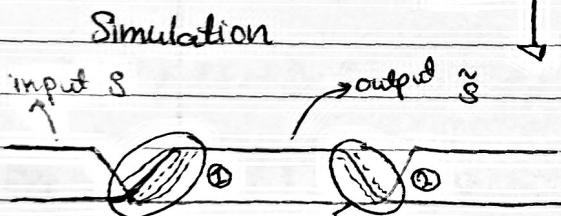
[How does the logic behave]

\rightarrow RTL done by Architect or Logic Designer then it converted into Schematic done by Circuit designer.

\rightarrow we have transistors inside Schematic level [NMOS + PMOS].



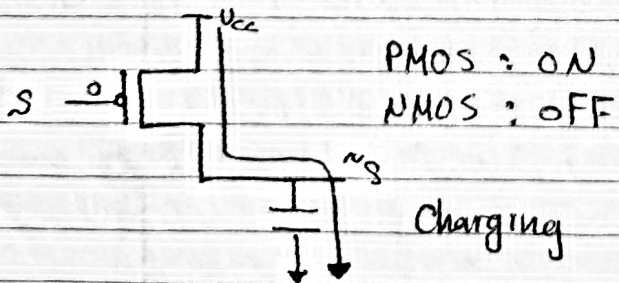
\rightarrow After that we go to the layout



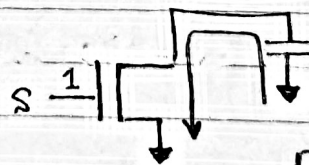
\rightarrow what determine this slope?

N \Rightarrow take the signal to ground

P \Rightarrow take the signal to High [it's connected with V_{cc}]



PMOS : ON
NMOS : OFF



PMOS : OFF
NMOS : ON

Discharging

$I_0 \propto \frac{W}{L}$ \rightarrow only change this
 \rightarrow its fixed.

\rightarrow When $S=0$

The Capacitor starts charging and its charge depends on W & R of the device.

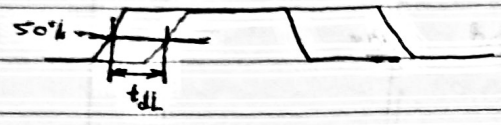
\rightarrow When $S=1$

The fast of discharging depends on the width of the device.

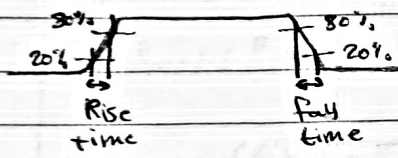
- How fast I can charge the capacitor?
Depends on W for PMOS.
- How fast I can discharge the capacitor?
Depends on W for NMOS.

* W_{10} for N device take the signal to low faster than W_5
 W_{10} for P device take the signal to high faster than W_5 .

• How do we measure delay between input and output?



• How do we measure slopes?



* We can take 10% to 90%.

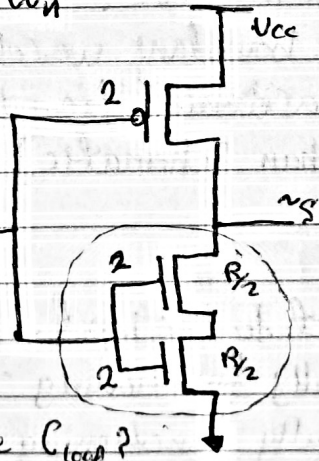
* If I want rising slope to be equal to falling slope??

→ make the width of PMOS larger [NMOS is faster]. So, if the

ratio $\frac{M_n}{M_p} = 2 \rightarrow W_p = 2W_n$

* L is fixed.

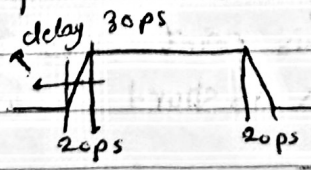
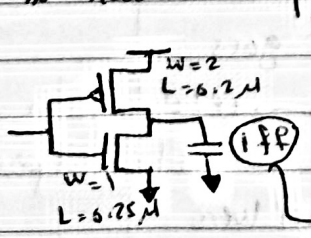
* we still need to have the same rise and fall slope.



This is what we call sizing the device.

$I \propto W, V = RI$

* What will happen if we increase the C_{load} ?

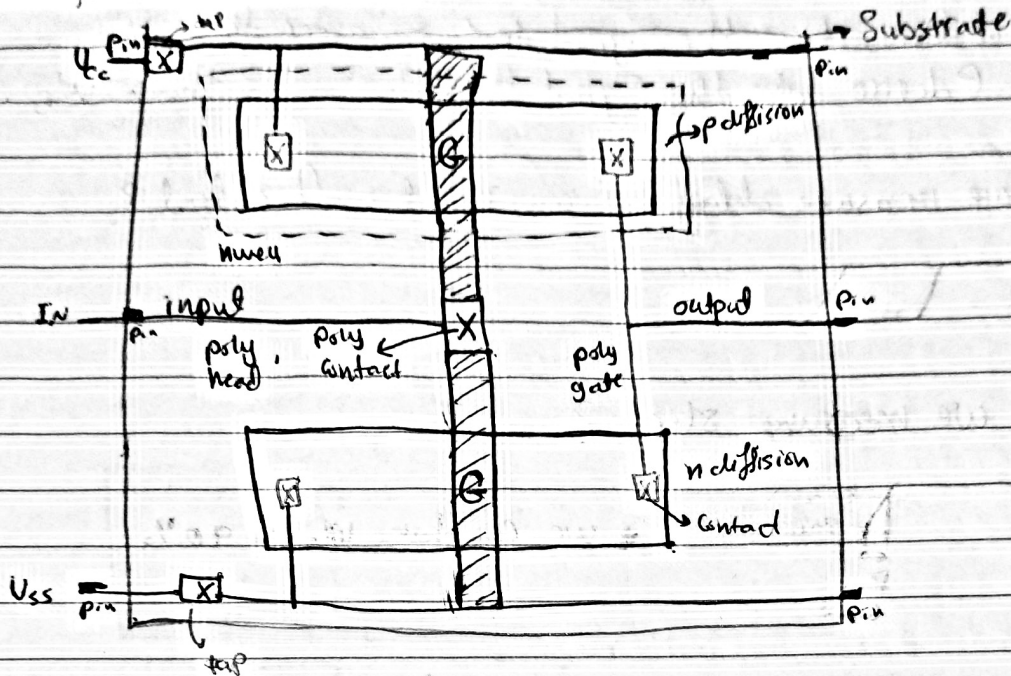


change to 10ff??

⇒ The delay will increase to 50ps like example, and slopes could be 30ps.

- ⇒ But if we don't want to affect the delay?
 - increase the width of devices [Strong driver]
 - The area will increase ↑. delay → constant.

* layout for the inverter :-



- * input came from metal layer 1
- * contacts usually high resistive.
- every thing touch the boundary we called it a pin.
- * we extract the capacitance from layout. [more accurate than schematic].

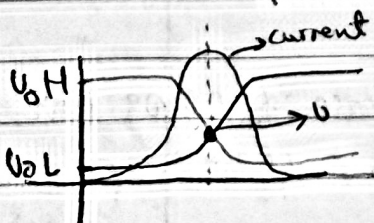
* VTC Curve :-

The output with respect to the input.

V_{OH} → when PMOS fully functioning.

V_{OL} → when NMOS fully functioning.

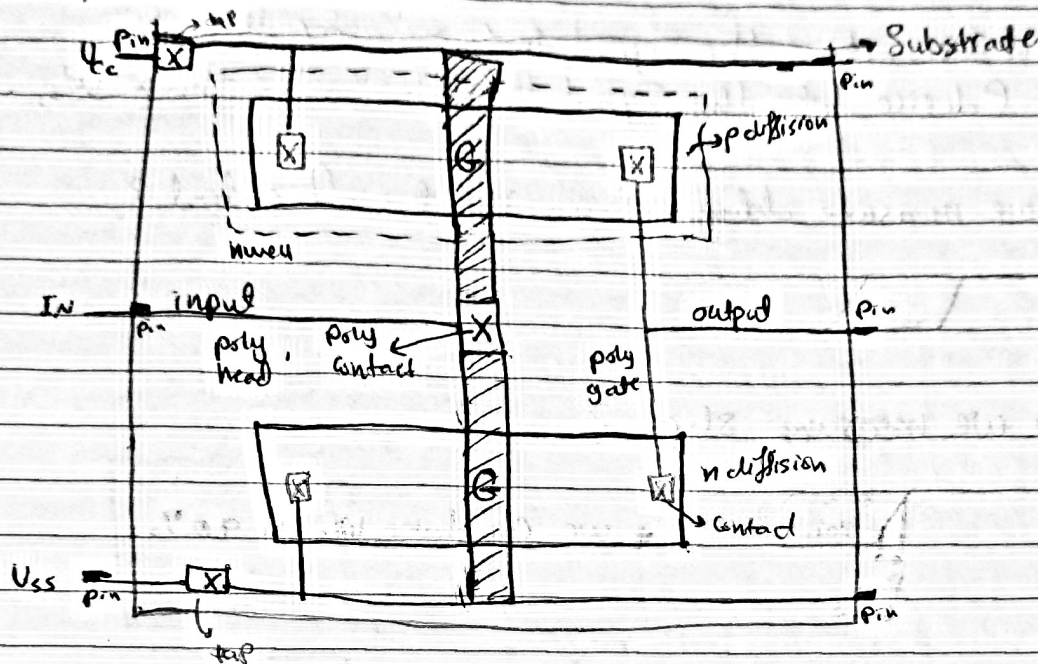
V_M → Both of them are ON. [in this point the current goes from V_{cc} to V_{ss} and this is short circuit].



is better → the two devices together were ON for a long time → consumed more power

- ⇒ But if we don't want to affect the delay?
 - increase the width of devices [Strong driver]
 - The area will increase ↑, delay → constant.

* layout for the inverter :-



- * input came from metal layer 1
- * Contacts usually high resistive.
- every thing touch the boundary we called it a pin.
- * we extract the capacitance from layout. [more accurate than schematic].

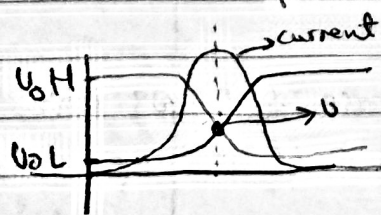
* VTC Curve :-

The output with respect to the input.

V_{OH} → when PMOS fully functioning.

V_{OL} → when NMOS fully functioning.

V_M → Both of them are ON. [in this point the current goes from V_{cc} to V_{ss} and this is short circuit]. → related with Short circuit power.

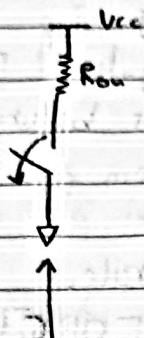
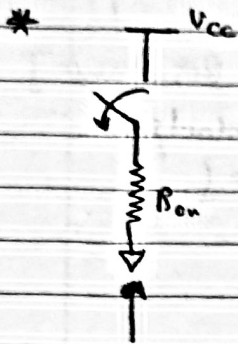


is better → the two devices together were ON for a long time → consumed more power

* If I have the current at mid point and I have W_n and need to find W_p ?

→ $I_{D_n Sat} = I_{D_p Sat}$

$W_n \Rightarrow W_p$?



$V_{OH} = V_{cc}$

$V_{OL} = 0$

$V_M = f(R_{on_n}, R_{on_p})$

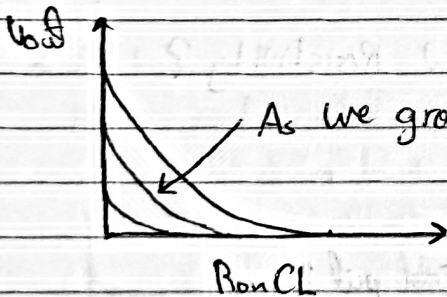
affect fall time

affect Rise time

$R_{on} \uparrow$ faster fall time.

more power ($I^2 R$).

* $t_{PHL} = \ln \frac{1}{2} R_{on} C_L$

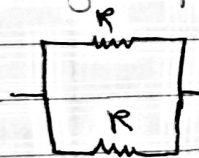


LS: Layout

T/5-3-2019

- * Paint from set of different colors. [based on the circuit].
 - The designer try to build a library cell of the inverter.
 - make sure that everything does not break [cannot correctly]
- * metal 1 \rightarrow Connection with diffusion or poly
 - contact \uparrow poly-head \rightarrow [high Resistance]
- \rightarrow from DR we get the dimensions and distances between contacts, in addition of the amount of overlap. [or from process file]
- \rightarrow Everything I put in the layout it has a rule.
 - \Rightarrow these rules programmed in the tools [DRC].

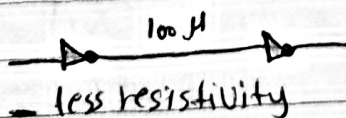
- any movement could make a short.
 \Rightarrow why to put two contacts instead of one?
reduce the resistance to $R/2$, and this is limited to the DR \rightarrow same to poly-head.



- Why we don't use another material with low resistivity?
 \Rightarrow more expensive.

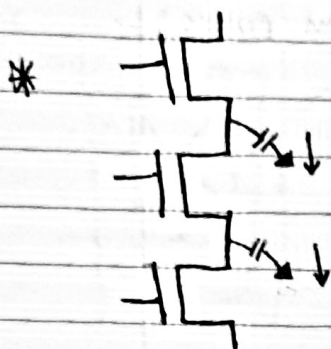
- * In the process file :-
 - number of layers: How many metal layer I can put?
 - poly Spacing: the distance between two polys.

\Rightarrow The higher I get in metal layer, the resistance decreased
 \rightarrow So, that why we use higher metal layer to do routing between far blocks.



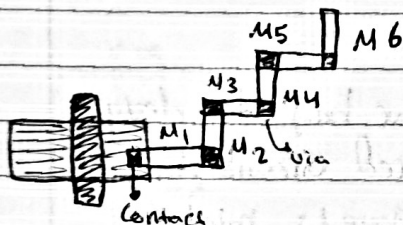
$$* \frac{R \ell^2}{2}$$

* Metal layers divided into binary groups to enable vertical and horizontal connection.



I need to reduce this capacitance.
 → So, we share the diffusions.

* How to connect diffusion with M6?



* geometry class to generate the tools based on DR.

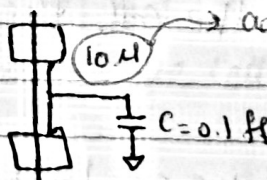
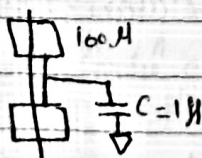
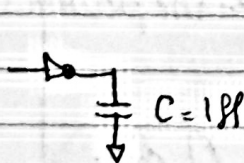
* Fram View → to deal with pins not all things inside.

* process scaling → using lambda [scaled factor]

* Antenna:

Wires not needed, sometimes we need it [to add capacitance]

→ match the delay:



→ add wires to make it local to make C = 1pF

* Density Rule:

amount of metals, poly there is ratio between all the materials
 → this is important for the machines that would cut.

* Why we have two input for inverter?

→ to connect horizontally or vertically. [to get the flexibility how to connect]

L9: Design Rules

R1 7-3-2019

→ we do Design Rules to be able to tolerate fabrication errors; -

- mask misalignment
- dust
- process parameters
- Rough Surfaces

تفاوت

* masks → metal layers $[M1, M2, \dots]$.

* Dust can make short circuit so we make cleaning rooms.

* half design Rules:

anything on edges [excludes pins] must be far half the design rules, because it might have another cell placed around it.

→ outside the cell, because there is might be another cell share it.

Shared metal
or shared diffusion



↓
make width smaller.

* you should make the size of cell smaller as possible, because the bigger you make the cell, the more areas, the more the power and the more delay.

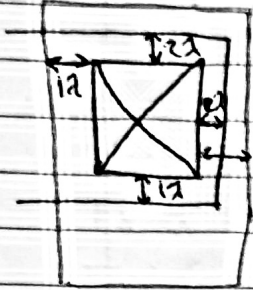
* to minimize the effort we can build a library of Standard cells.

- ▶ to optimize cells very well.
- ▶ to make the tool make the layout for me.
- that's what we call the synthesis flow.
- ⇒ we make all the cells with the same height.

* We should always put PMOS at the top to connect the V_{cc} in one line.

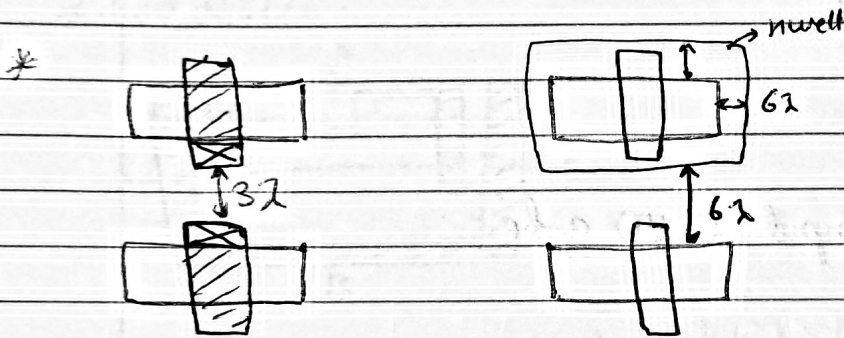
= one substrate for whole the cell.

*



→ the more the area you put, you increase the Capacitance.

* poly uses a width of 2λ → if I have 10nm process, then $\lambda = 5nm$.



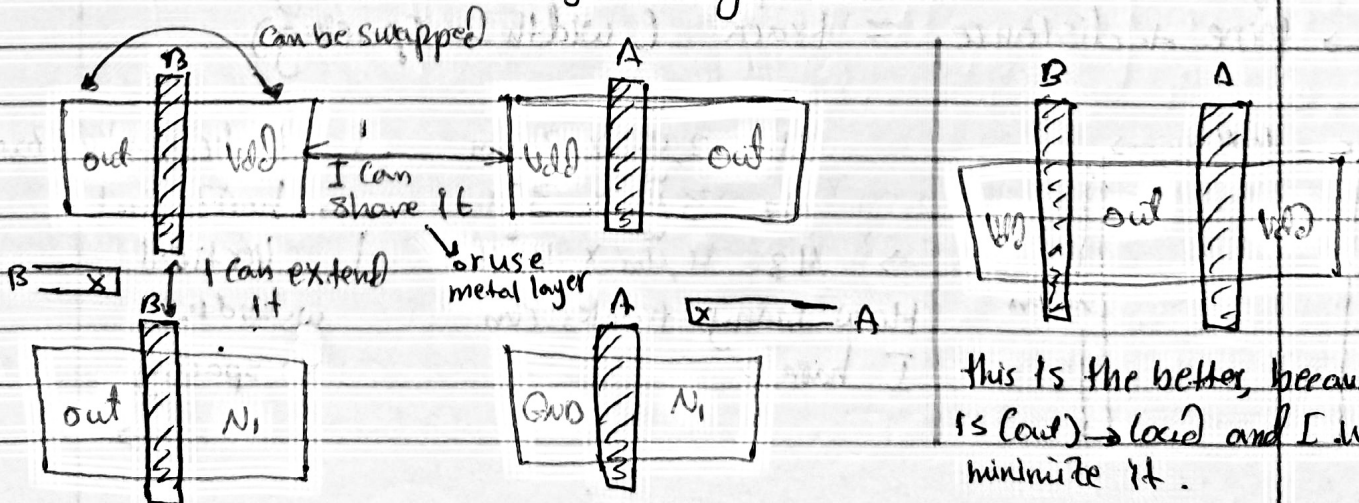
*



Why so long? you can change the size as the load of the driver can be changed.

→ another reason, to do standardization.

* Limitation: ⇒ ① Device limited. ② Metal limited → metal layers that could go through the cell.



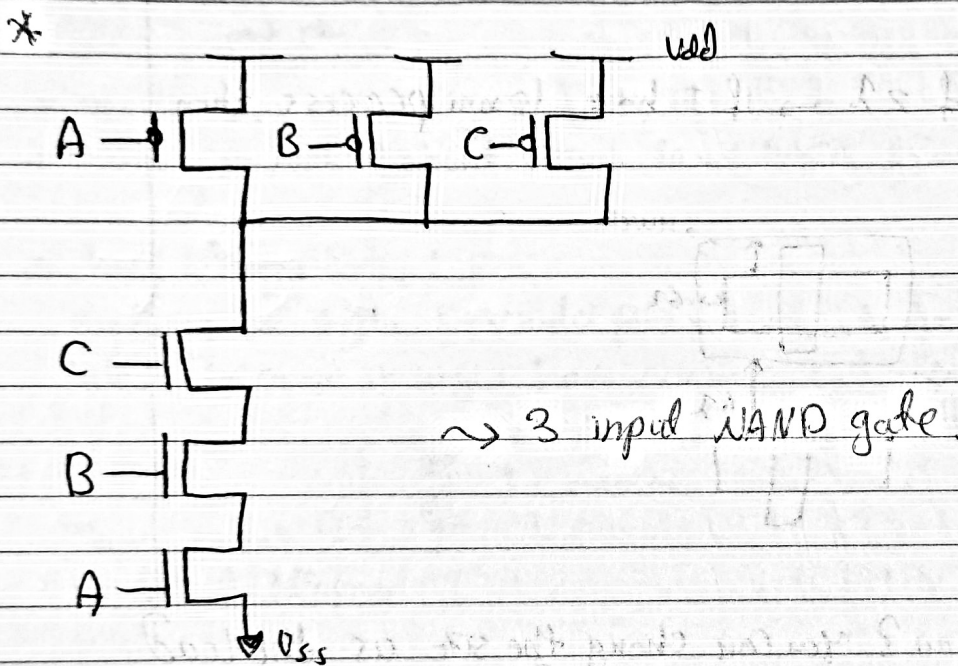
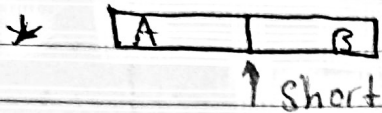
this is the better, because this is (out) → load and I want to minimize it.

L10:

T/12-3-2019

* What does it mean the width of the devices?

→ As I did the simulation, this is what will work for me.



* Wiring Tracks:

Space required for a wire, and it used to calculate height of the cell [metal limits]

→ depends on the number of tracks and metal layers.

→ wire + distance = track (width + spacing)

* 

	w	s	
M1	$2 \mu m$	$2 \mu m$	Vertical + horizontal
M2	$3 \mu m$	$3 \mu m$	horizontal
⇒ M3	$4 \mu m$	$4 \mu m$	Vertical

How many tracks can I have?
w: width
s: spacing

→ each metal layer has a track for its own.

m1: In horizontal $\Rightarrow \lfloor \frac{20}{4} \rfloor = 5$ tracks and I need to make Spacing from the edges (2 μ m from above and 2 μ m below)

m3: In vertical $\Rightarrow \lfloor \frac{16}{8} \rfloor = 1$

Vertical we do 2.

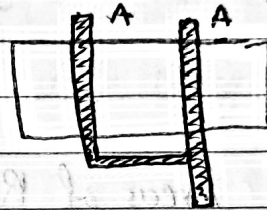
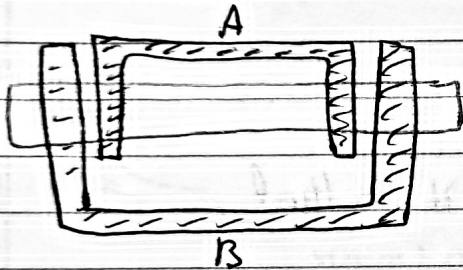
* Area estimation based on number of tracks.

* Stick diagrams:

Is a way of sketching the layout [with out doing it].

* Legging \rightarrow high width into two legs but save the performance.

* two devices:



L11:

R/14-3-2019

* How to evaluate our design?

- cost
- Reliability
- Speed / performance [delay / freq]
- power.

* Reliability :-

- real world is analog and we try to process things into digital. [we care about the signal is high or low].



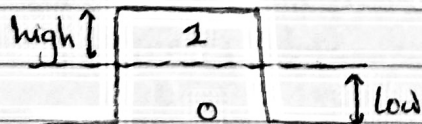
Areas of Reliability

- * How Signal Switch and how frequently it do that.
- * Also Routing

* Things that affects the Reliability :-

- inductive coupling: we have 2 signals, next each other and switch very often $[\frac{di}{dt}]$.
- capacitive coupling: depends on $\frac{dv}{dt}$
- power and ground noise.

* Noise and Digital Systems :-

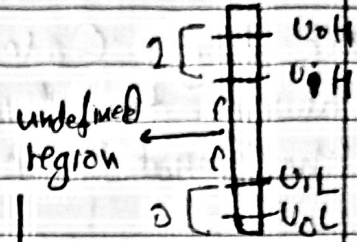
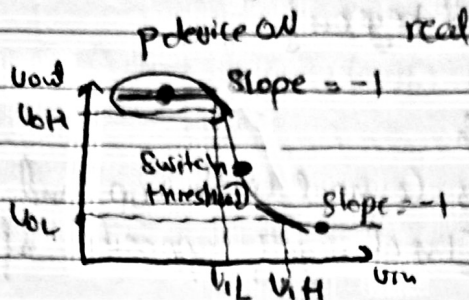
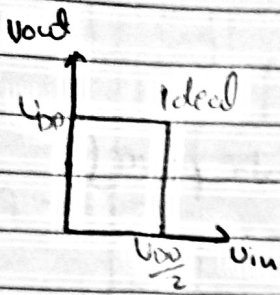


* you need to set up your design Consistance. [where high and where low]
→ where the noise.

* digital signal is discrete value of analog signal, but I set a threshold [where consider high and where low].

→ analog is accurate but we do not take.

* for small amount of noise, out put noise is less than input noise.



* noise margin

high $\rightarrow NM_H = V_{OH} - V_{IH}$

low $\rightarrow NM_L = V_{IL} - V_{OL}$

both transistors are ON in this region.

→ does the size of the device affect the noise margin?

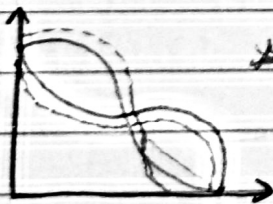


∴ size affected.
also the slope of the input.

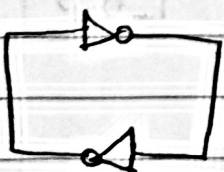
* nmos: act normally but pmos not charge to highest value

and so on... the noise amplified.

→ the size of them should be proportional to each other.



* I want it to be bigger → the smaller → more noise can't be rounded.



→ keep change the size until I see the eyes.

* NM: the bigger the noise margin the better.
 [less Undefined region]

* Static Logic Gates:

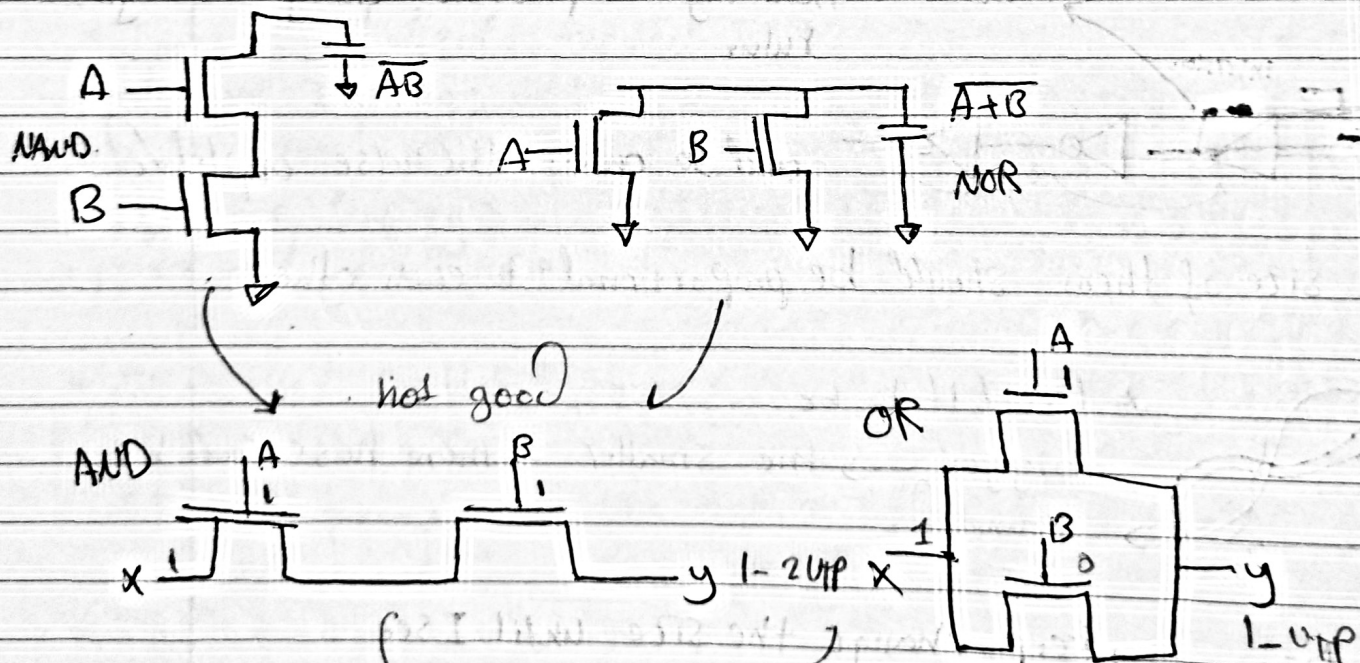
- always the devices connect with V_{DD} and V_{SS} . [static power]
- the signal static [unchange] opposite of dynamic.

→ CMOS: we have two networks, one is connected with PMOS and the other with NMOS.
 high ← Complement to each → low
 other

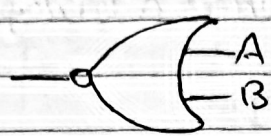
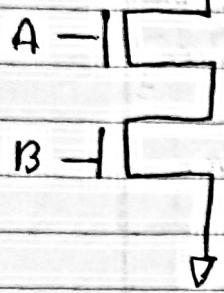
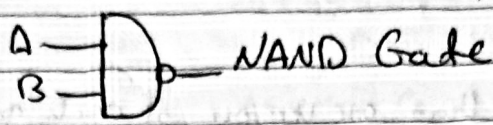
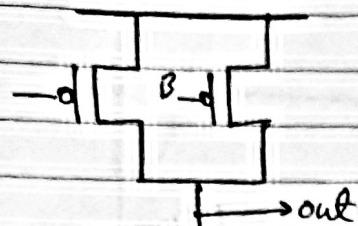
(PON) pull up pull down (PDN)

→ So I build the Pnetwork then mirror it to Nnetwork or the opposite.

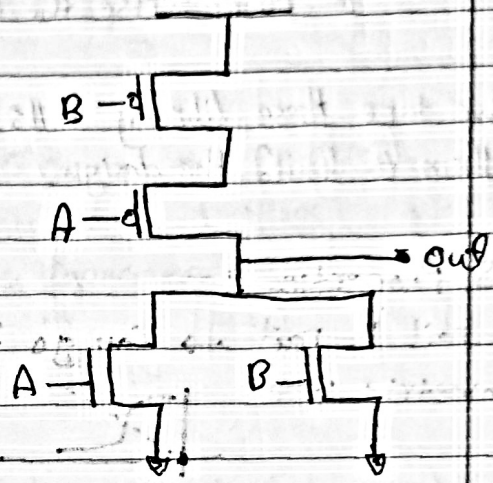
→ If I use NMOS to pull up → V_{OH} will decrease → less NM



not good → does not give strong one.



NOR Gate



* If you have to choose between NAND, NOR what to choose? Why?

NAND \rightarrow nmos is faster, so I prefer the nmos to be in series [long path].

* to design CMOS complex gates \rightarrow make pull up and pull down networks

L12:

T/19-3-2019

* Library cells: more than one version of each gate.

→ First we do schematic then layout then do simulation to characterise the cell. [know slopes, delays, Capacitance]

→ we take these library cells, write a verilog code which use them to build the design.

* CMOS Properties :-

full rail to rail → Can go fully to up or low [1 or 0]

Symmetric VTC

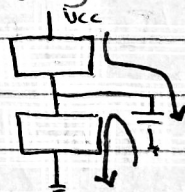


Propagation delay is function of load capacitance & resistance.

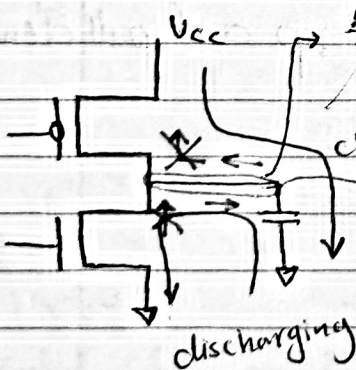
↳ between change in output due to change in input.

No static power dissipation just charging and discharging [Power Stored]

Direct path current during switching.



just these ways
Self heat



this is a problem
In design.

charging [self heat] temperature ↑ R ↑ delay ↑

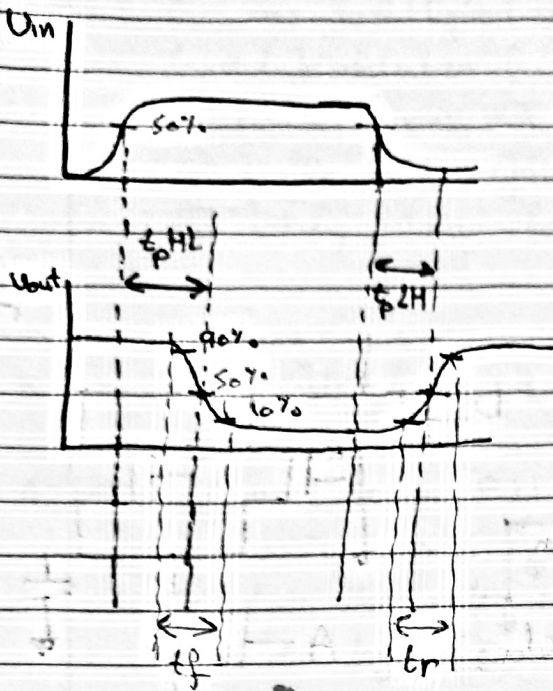
the width of this wire must be larger.

* Example: let $f = (A+B) \cdot C$

$$\text{pull up} = \overline{((A+B) \cdot C)} = \overline{(A+B)} + \overline{C} = (\overline{A} \cdot \overline{B}) + \overline{C}$$

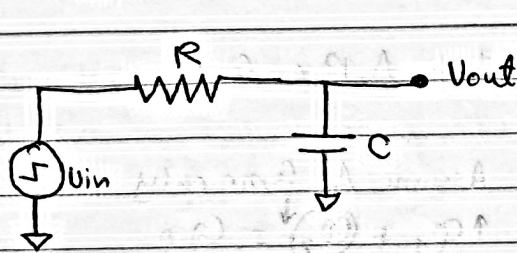
$$\text{pull down} = \overline{f} = (A+B) \cdot C$$





* Delay Definitions:
 We use CMOS, we don't use dynamic but in certain packages

* $t_r = 10\text{ps}$
 $t_r = 5\text{ps}$ \rightarrow this is better
 • [larger width of PMOS] or
 • [less capacitor load]



$$V_{out} = (1 - e^{-t/\tau}) V$$

$$t_p = \ln(2) \tau = 0.69 RC$$

* Power dissipation :-

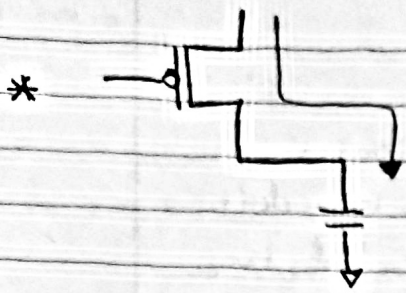
- instantaneous power : $p(t) = v(t) i(t) = V_{supply} \cdot i(t)$
- Peak Power : $P_{peak} = V_{supply} \cdot i_{peak}$
- Average power : $P_{av} = \frac{1}{T} \int_t^{t+T} p(t) dt = \frac{V_{supply}}{T} \int_t^{t+T} i_{supply}(t) dt$

\rightarrow When we change the width we must look at both power + delay
 \rightarrow Power-Delay : we do optimization.

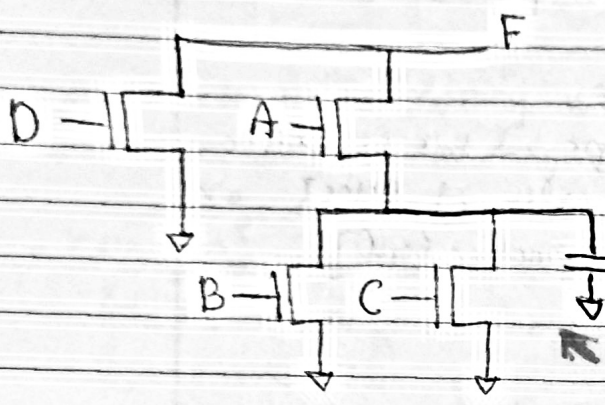
* Energy = Power \times time

\rightarrow $Q_{charged} = C_L V_{DD}$ so the energy = $Q V_{DD} = C_L V_{DD}^2$

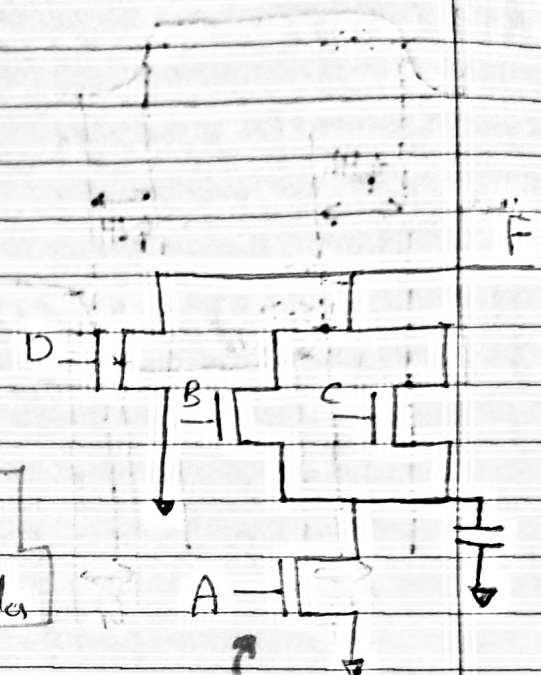
$E_C = \frac{1}{2} C_L V_{DD}^2$: half of the time I do charging and another half do discharging.



$T_p = 0.69 R_p C$
 $p = C V_{DD}^2$
 change of C affect both

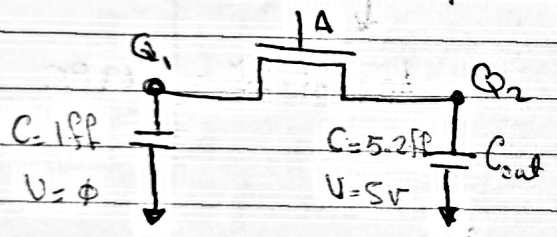


or
 Charge when A=1
 $C_{dc} + C_{db} + C_{da}$



Charge when B or C = 1
 $C_{dc} + C_{db} + C_{da}$
 Charge Sharing Problem

→ So it depends when A, B, C come



Assume A come only
 $\uparrow Q_1 + \downarrow Q_2 = Q_p$
 $C_1 V_1 + C_2 V_2 = C_p V_p$
 $0 + 5.2 pF * 5 = 6.2 pF * V_f$ a big deal.
 $V_f = \frac{5.2 pF * 5}{6.2 pF}$ ↓ not good

→ make the signal that will take the out to zero the late.

L13:

R/21-3-2019.

* sharing metal tracks to minimize its number and so minimize the width of cell.

* Sometimes we go to another level of metal layer instead of add more tracks to avoid the increase in the width.

[see stick diagrams on slides]

* LVS: Be sure that layout match the circuit built in schematic

logical verification for schematic.

- 1 check connectivity
- 2 width and length.
- 3 verify timing delay → not very accurate in schematic.

Compare the netlist of each of them

* In Spice netlist → mt1 or mt2 or ... ; depends of the process.

* EDA: Convert from RTL / Verilog to layout [Synthesise flow]

what is the advantage and disadvantage?

in control logic [change very frequently] → don't do reduction or optimization.

→ each process has its own libraries, and in each library we have different families

for: Speed power Area } depends on what we do optimization for

* Test Creation: test the logic not Schematic

→ do different test patterns and insure that the output match our expectation. Hard work testing.

L14:

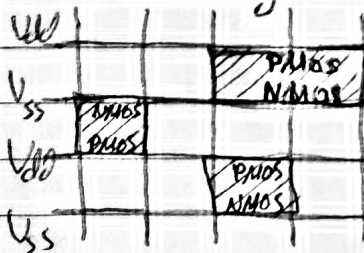
T/26-3-2019

* Physical Synthesis Steps:

- floorplanning: interface and shape of design + supply network [Vdd & Vss]

- placement: minimize the total area and interconnect length.

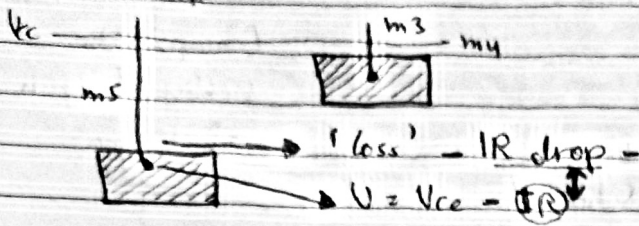
- Routing → we do it based on how we do the floorplanning.
→ we use grid, 'unit tile' of floorplanning build this grid



Cells in library designed to be multiple to unit tile.

→ make the design compact more

→ Power grid: make power network [power delivery]



- * Formal Verification: just test the functionality of the design.
- * timing after Pn layout: STA Static timing Analysis.

→ physical Synthesis: we optimize the cell as required by the designer.

* clock Delay Problems.

- all clock pins are driven by a single clock.

→ we need it to reach all components without delay.

→ Sometimes we allow some clock skew [no problem]

* clock distribution [clock tree].

↳ buffer to compensate the delay } we have special algorithms.
↳ how to connect this tree.

* Digital Standard Cell Library (DSSL)

AND, OR, NAND, NOR, ...

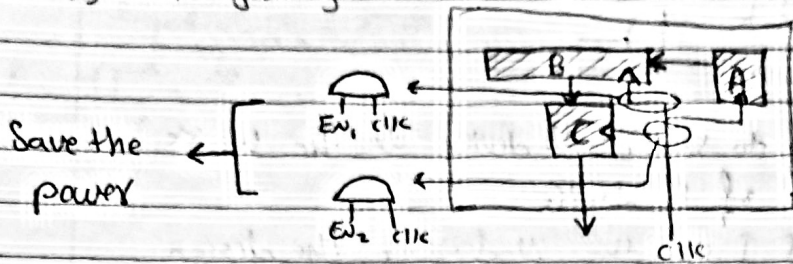
→ multivoltage and multithreshold libraries ⇒ to do optimization

high/low

Ex. low leakage: pick high V_T but slower.

* low power design techniques

↳ clock gating



at first clk, B & C no need to work, so I can turn them off by Enablers.

↳ Static multi-voltage

Same as clock gating but enable on voltage [switching of voltages]

L: 15

R128-3-2019

→ put power switch but it will cost you.

* Design corners :

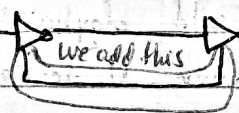
at the time we run the characterization for any cell, we want to check the speed, slope, capacitance ...

Design Setup

[min, Typical, Max] ← Voltage transfer curves reflected to what you choose.

* PD: Power times delay

- Specifications
- data sheet
- circuit design w/ do simulation
- layout check DRC
- extract parasitics
- we generate netlist
- library characterization. find slopes, delay, power delay, noise margin ...

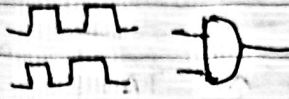


$$\downarrow R = \frac{PL}{AF} \rightarrow C \uparrow, P \uparrow$$

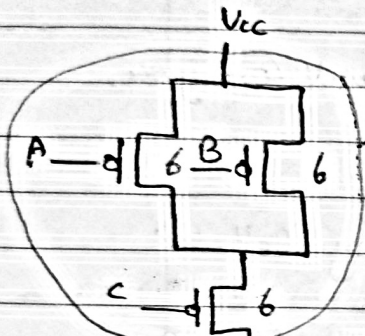
slides-18) * Dynamic Gates.

*

Static design



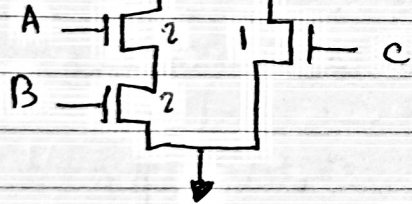
Dynamic design → we try to make N-device do the work instead of P-device
 So, if I use the N-devices, the design should be faster and less area.



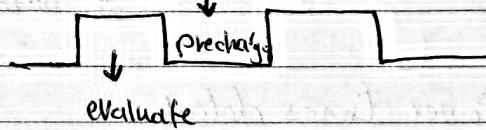
we try to eliminate the P-devices

Ex: Sizing → $\frac{N}{P} = \frac{1}{3}$

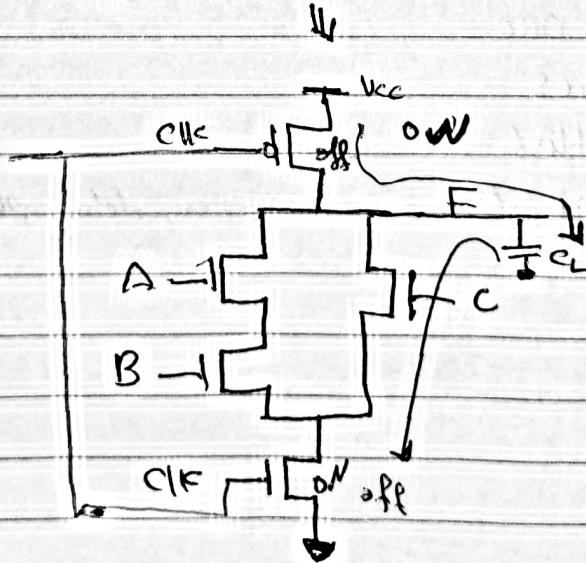
$F = \overline{AB + C}$



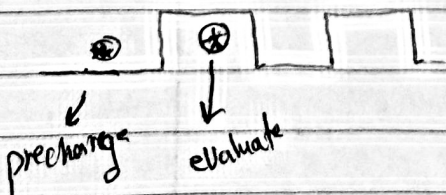
Precharge, evaluate

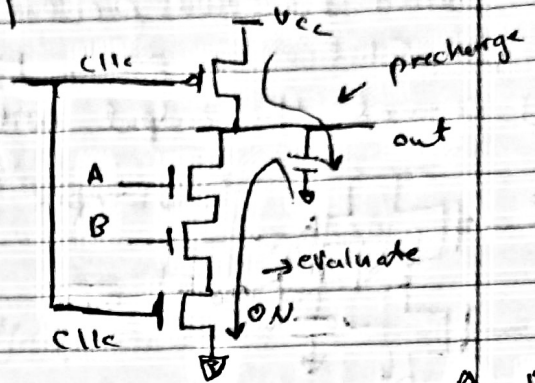
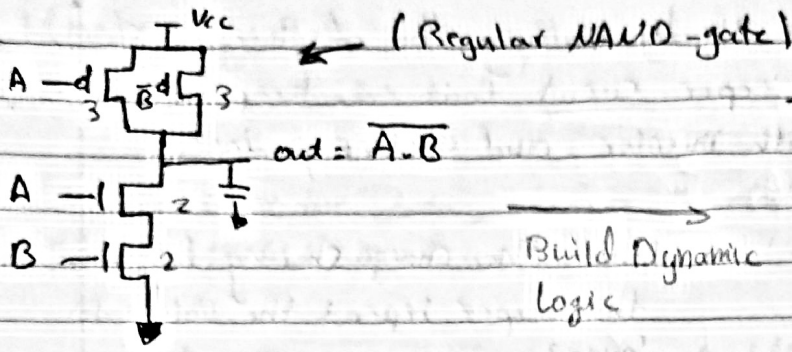


* to build dynamic logic, we use (pull down N-devices) and remove (pullup-P) and put P-device & N-device with clk.



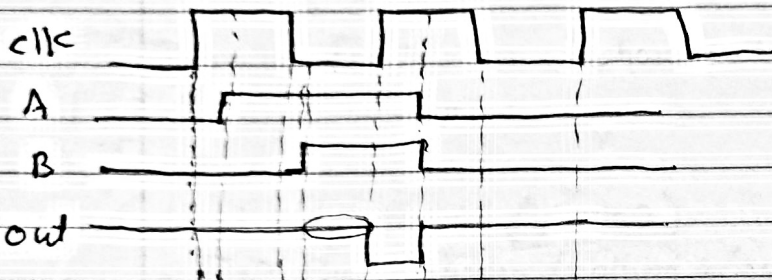
precharge: charging the capacitor before starting the evaluation.





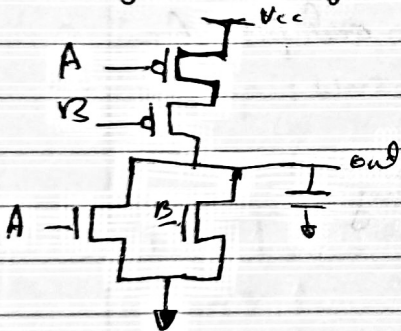
Build Dynamic Logic

(*) Timing diagram

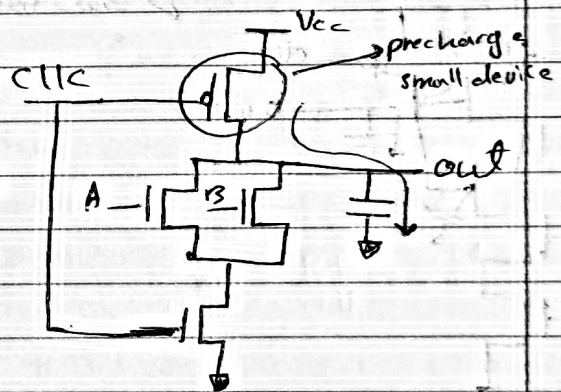


A	B	out
0	0	1
0	1	1
1	0	1
1	1	0

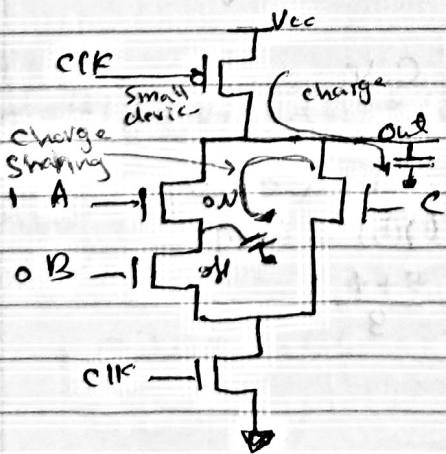
* Regular NOR-gate



Dynamic NOR-gate



(*)

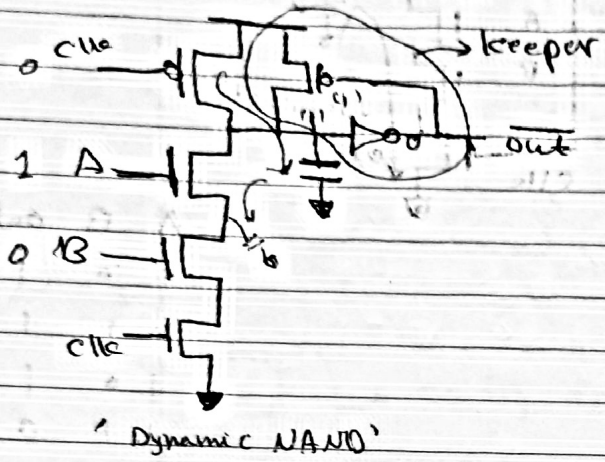


This is happen for each precharge and this is the problem in dynamic logic.

(*) The problem is: after each precharge, the charge sharing happens and we loss the charging

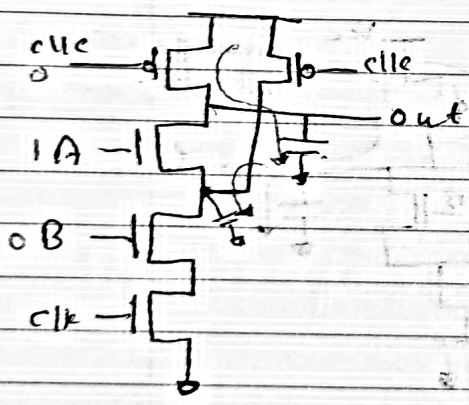
or prevent it

- ① The solution of previous problem is by making the discharging slow.
 So, the solution is put a 'Keeper circuit' that can keep the value.
 ② by connect the output with inverter and another p-device.

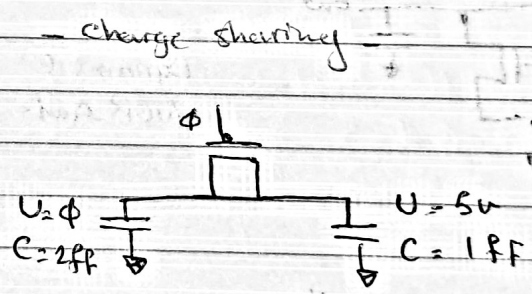


* when the charge sharing happens, the keeper repeat the value to 'one'

② → another way :



← but this increase the area and power



$$Q = C \cdot V$$

$$Q_F = Q_1 + Q_2$$

$$(C_1 + C_2) \cdot V_F = C_1 \cdot U_1 + C_2 \cdot U_2$$

$$(3) V_F = (1)(5) + (2)(\phi)$$

$$V_F = \frac{5 + \phi}{3}$$

L16: slides - L9

* Sequential Gates

① - Use to store the data, because the data is changing. If we don't store it, the data will loss.

② - use in finite state machine to represent different states of the machine.

③ - used in pipelined machine to designate pipe stages.

→ Instruction Fetch → Instruction Decode → Execution → STOP

* Types

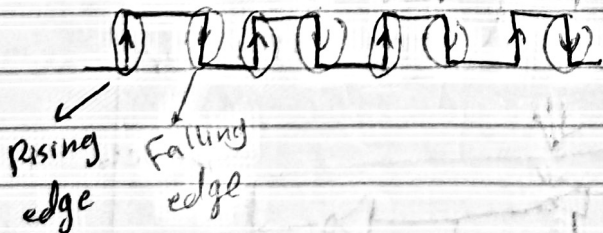
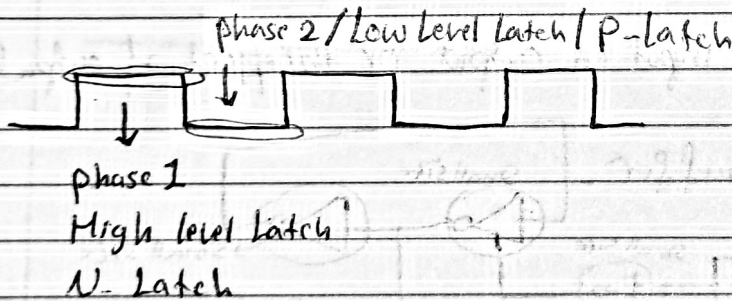
- Latch (phase 1 / phase 2)

- Flip Flop (Rising / Falling edge)

① Latch : works on Level

② Flip Flop : works on edges

③ In design we use Flip Flop, if we want a sequence on the edge not in all period.



- Works on rise edge.

- works on low edge.

- works on two edges (high/low at the same time).

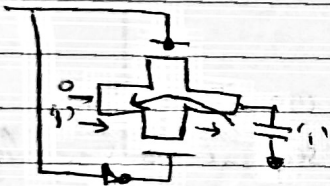
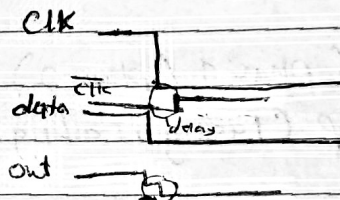
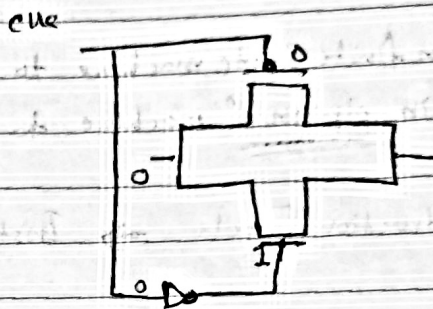
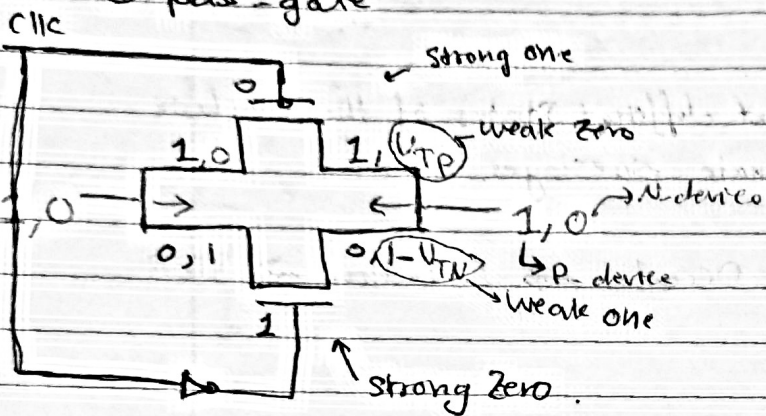
* Enabled Latch

→ (1) → work
0' → not work

① Synchronous Set / Reset → based on clk

② Asynchronous Set / Reset → not based on clk.

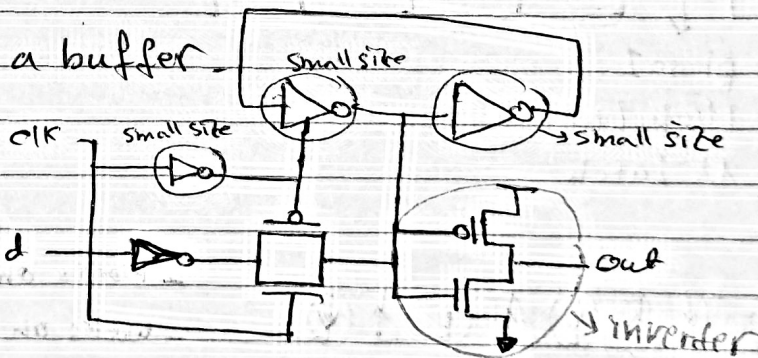
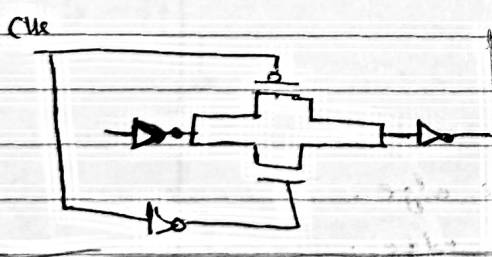
③ pass-gate



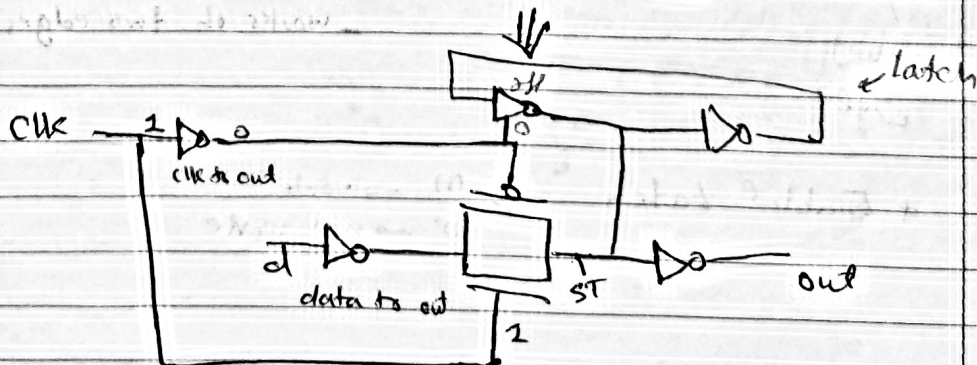
(glitch) delay from inverter

④ How I prevent the effect of input on output or the effect of output on input ??

The solution: put a buffer

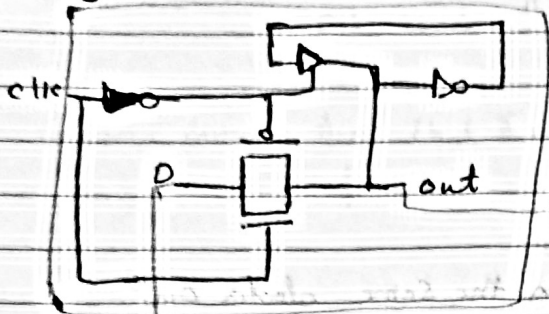


three state buffer



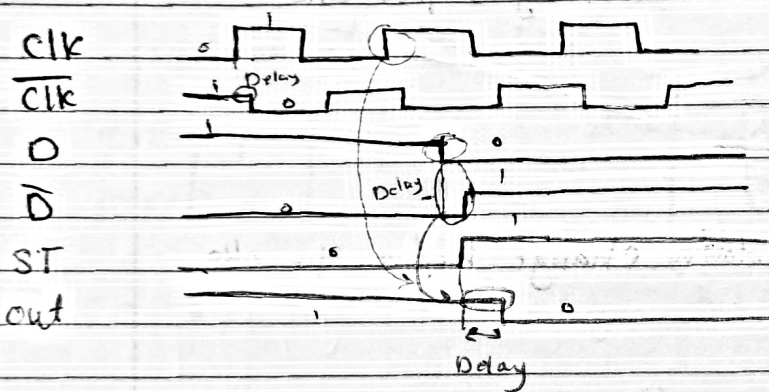
(N-Latch) → clk high - phase 1

Sometimes, we build the latch without the inverters on input & output. Why?



We use this design, when the data comes from another logic, and the output goes to another logic. In this case the inverters are already exist.

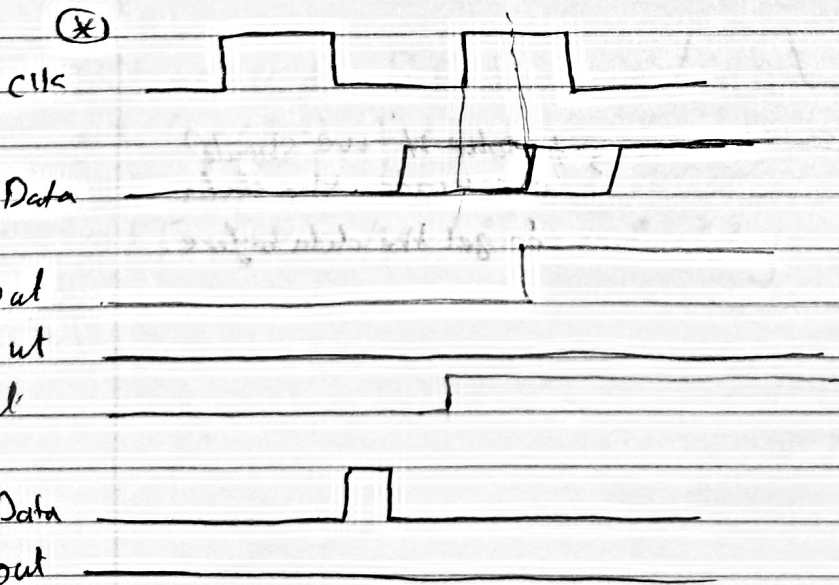
Timing Diagram of D latch



hold-time

after the latch closed and get the data, I don't want the data change after it. Because the variation (delay) on clk.

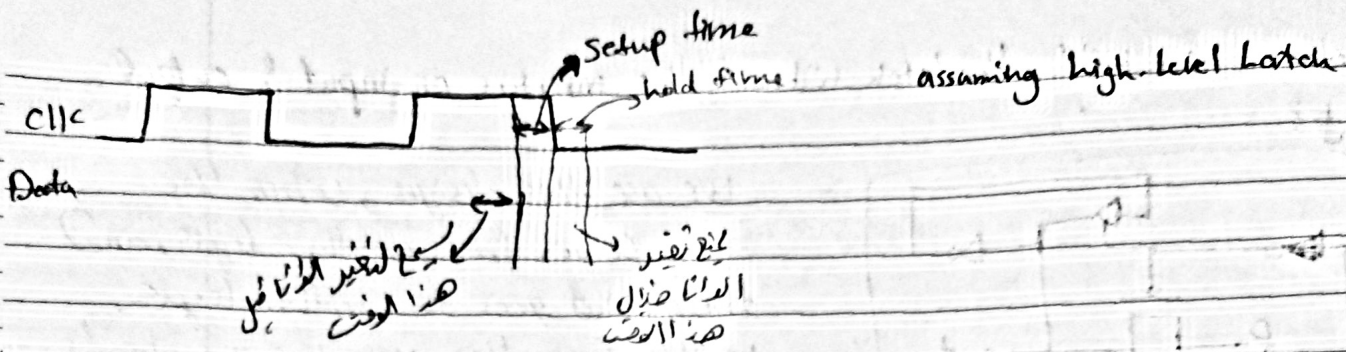
How to calculate the delay?
 → CLK to out delay
 → Data to out delay.



Setup time: the data be available before the clk closed.

phase 1: I take the variation of clock (h-L) or (L-h) at

phase 2: I take the variation of clock (L-h) or (h-L) at



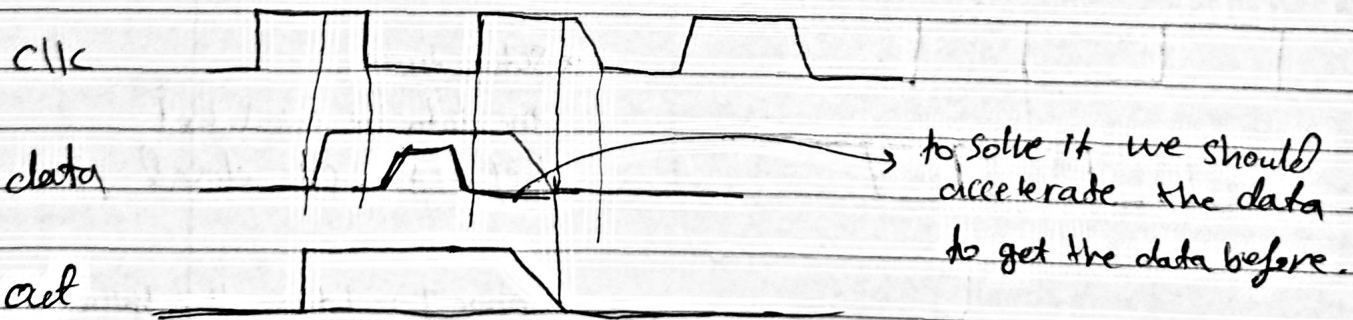
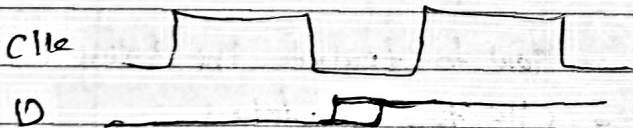
⊛ Non-Inverted Latch → The data in the same data out

⊛ Inverted Latch → The data in is not the same data out
(There is an inverter on input but there is no inverter on output)

⊛ Data to out delay → when the data change after clk



⊛ clk to out delay → when the data change before the clk




L17:

① Setup time : time that I need data to be ready before the clk is closed in latch.

② hold time : time that I need the data stable in a specific period of time.

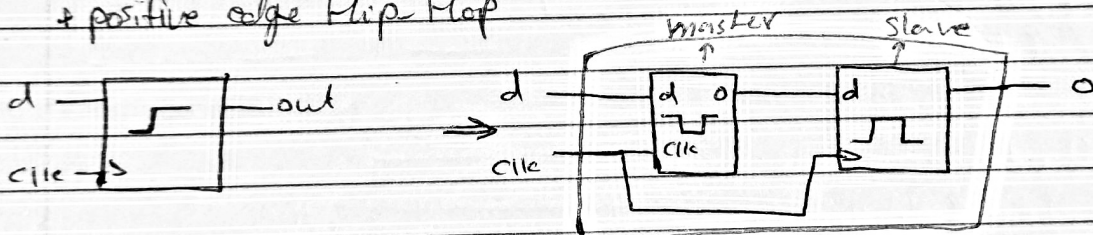
* P-Latch

works on clk low 

Flip Flop :-

two-latches (master & slave), one high and another low.

+ positive edge Flip Flop



* If last latch on high then the F.F in rising edge.

+ The disadvantage of Flip Flop

That is used two latches \rightarrow more area & power

But we use the F.F because we know when the data will get (on edge). but in latch there is a period to get the data.

in F.F there is no data to out delay