

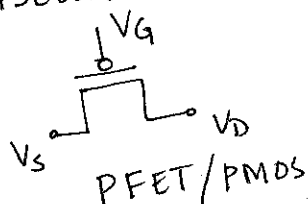
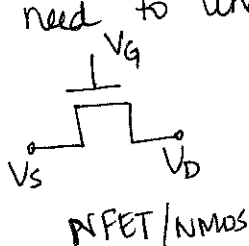
- 1) Technology scaling :- Every generation, shrinking transistor sizes leads to an improvement in their performance and reduces energy dissipation. At least, this was true for classical transistor scaling.

Scaling a technology reduces the gate delay by 30% and the lateral and vertical dimensions are also scaled down by 30%.

The die area decreases by 50%.

The areal & fringing capacitances all decrease by 30%.

- 2) To understand the design of integrated circuits, you need to understand transistors.



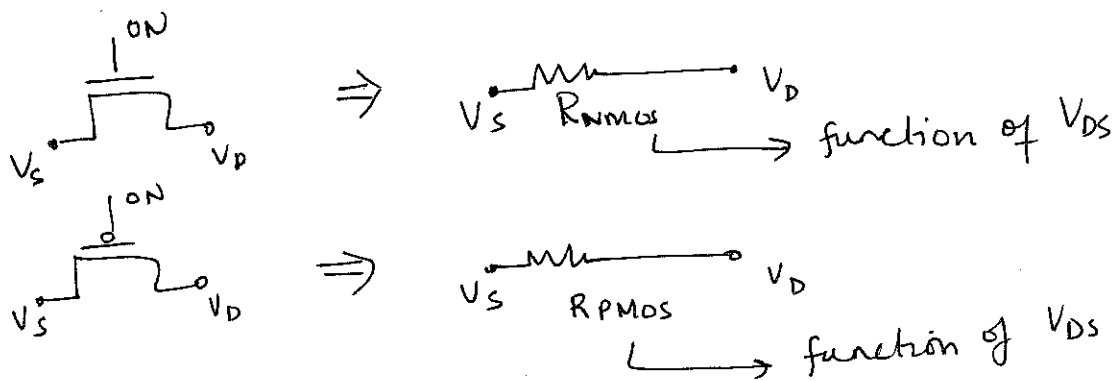
⇒ Essentially a switch that is controlled by the gate voltage V_G

$\left\{ \begin{array}{l} \text{NFET is good For transmitting "0" when it is ON.} \\ \text{PFET is good For transmitting "VDD" when it is ON.} \end{array} \right.$

NFET is ON when $V_{GS} > V_{Tn}$

PFET is ON when $V_{SG} < |V_{Tp}|$. The actual value of V_{Tp} is negative.

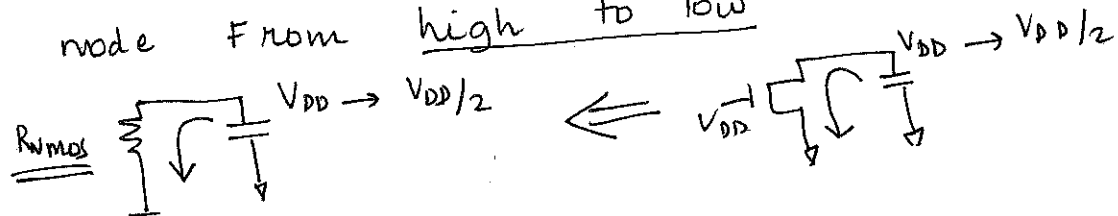
However when the transistor is ON it essentially is not a perfect on switch. That is, it has some finite resistance. Unfortunately, the resistance of the switch is dependent on the drain to source voltage. That is, the transistor is a non-linear resistor.



This makes life a bit harder to work with transistors. But we can often make simplifications to evaluate R_{NMOS} and R_{PMOS} . When we make some simplifications, we call these resistances as "equivalent" resistances. In that case, it really depends on the voltage swings across the transistor terminals.

We will look at a few specific cases to evaluate equivalent R_{NMOS} and R_{PMOS} .

If an NMOS is being used to pull down a node from high to low

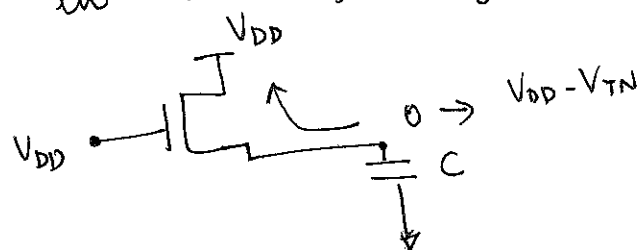


We evaluate R_{Nmos} at the beginning of the transition & then we evaluate R_{Nmos} half-way through the transition and take the average of the two resistances.

$$R_{eq} = \frac{1}{2} (R_{Nmos}(V_{DD}) + R_{Nmos}(V_{DD}/2))$$

Also note in this case, NMOS is not under any body effect, since $V_{SB} = 0$.

However, if NMOS were used to pull-up a node in the following configuration.



The problem is slightly harder because now $V_{SB} \neq 0$. At the beginning of the transition

$V_{SB} = 0$ and $V_{TN} = V_{TO}$.

At Half way through the transition,

$V_{SB} = \left(\frac{V_{DD} - V_{TN}}{2} \right)$, where V_{TN} must be

obtained self consistently by solving the equation

$$V_{TN} = V_{TO} + \gamma \left[\sqrt{2\phi_f + \frac{V_{DD} - V_{TN}}{2}} - \sqrt{2\phi_f} \right]$$

This is only a quadratic equation in V_{TN} that has two roots. Select the root that is positive and greater than V_{TO} .

NOTE:- whether the NMOS is pulling a node from V_{DD} to 0 or its pulling up a node from 0 to $V_{DD} - V_{TN}$.

$V_{DS} \geq [V_{DSat} = V_{GS} - V_{TN}]$ so just use the equation for resistance in saturation.

In general (a) $R_{eq} \propto 1/(W/L)$

If the transistor size becomes bigger then R_{eq} becomes smaller.

(b) $V_{DD} \uparrow \quad R_{eq} \downarrow$ (c) $V_T \uparrow \quad R_{eq} \uparrow$ (d) $(\mu_{ox}) \uparrow \quad R_{eq} \downarrow$

These scaling trends are important to consider.

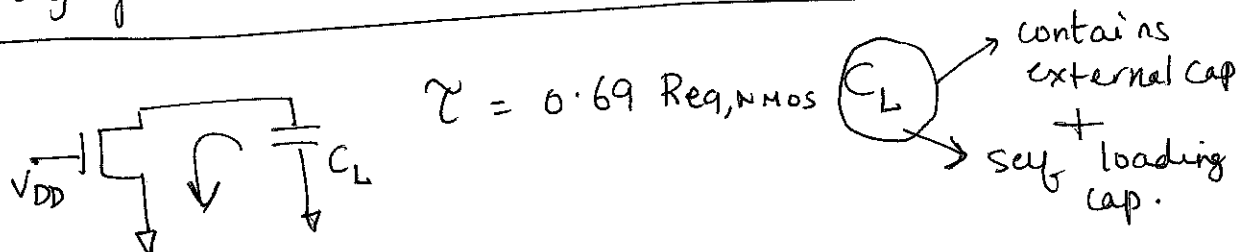
If all process parameters are same for PMOS and NMOS except the mobility such that $\mu_n > \mu_p$ then $R_{eqp} > R_{eqn}$ For $(\frac{W}{L})_n = (\frac{W}{L})_p$.

∴ If you want the same equivalent resistance

$$R_{eqn} = R_{eqp} \Rightarrow \left(\frac{W}{L}\right)_p > \left(\frac{W}{L}\right)_n$$

$$\frac{\left(\frac{W}{L}\right)_p}{\left(\frac{W}{L}\right)_n} = \frac{\mu_n}{\mu_p} \quad \text{if all other process parameters are identical. That is, } \frac{C_{ox,p} = C_{ox,n}}{V_{TN} = |V_{TP}|}$$

Charging and discharging through transistors



Calculate $R_{eq, NMOS}$ as explained earlier.

$$\left\{ \begin{aligned} R_{eq, NMOS} &\cong \frac{3}{4} \frac{V_{DD}}{I_{Dsat}} \left(1 - \frac{5}{6} \lambda V_{DD}\right) \\ I_{Dsat} &= \frac{1}{2} \mu_n C_{ox} \left(\frac{W}{L}\right)_n [V_{DD} - V_{TN}]^2 \end{aligned} \right\}$$

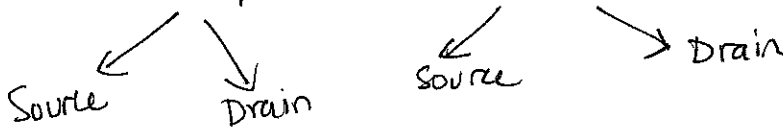
Note, the above is true only when we are discharging a node from V_{DD} to 0 using an NMOS.

Go through Hw 2 where we were charging a node using PMOS to calculate $R_{eq, PMOS}$ the proper way.

Concept of self-loading

Self loading capacitance is the capacitance at the output node due to the transistor structure itself:

Components include :- overlap + Junction capacitance



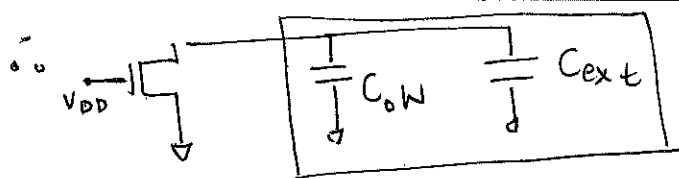
Typically all these capacitances are collectively termed as 'Parasitic capacitances'. All of these scale according to $C_{par} \propto W$

$$\left\{ \begin{array}{l} C_{ov,s} = C_{ov,s} W \\ C_{ov,d} = C_{ov,d} W \end{array} \right. \quad \text{where} \quad \left. \begin{array}{l} C_{ov,s} = \frac{\epsilon}{t_{ox}} X_{ov,s} \\ C_{ov,d} = \frac{\epsilon}{t_{ox}} X_{ov,d} \end{array} \right\} \text{ typically same}$$

$$\underline{C_{jn} = C_j W L_s + C_{jsw} (W + 2L_s)}$$

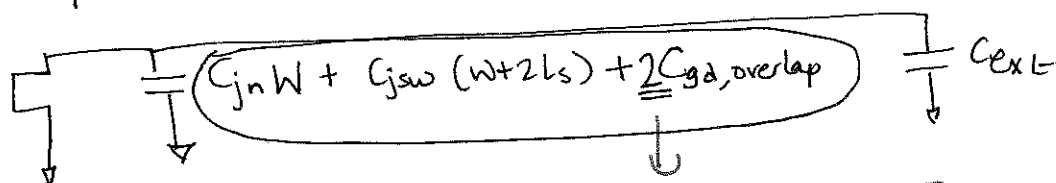
To remember:- parasitic capacitances arise due to the transistor structure & scale proportionally w/ device width.

If we lump all of the parasitics into one component $C_{par} = C_0 W$
↓
per-unit-width capacitance.



total capacitance
at the output node
of the transistor

C_0 comes from overlap and junction of the NMOS. However, the overlap capacitance must be doubled because it is a miller capacitance. And remember, when capacitance connects input node with the output node then it is essentially a miller capacitance and its effect is doubled.



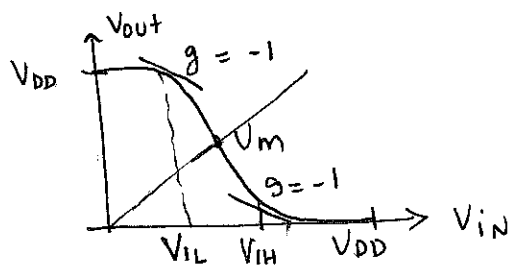
VERY IMPORTANT.

DO NOT FORGET !!

INVERTER ANALYSIS

- a) Inverter consists of PMOS in pull-up and NMOS in pull-down network.
- b) Figures of merit
 - (a) Switching threshold
 - (b) Low to high delay
 - (c) High to low delay
 - (d) Driving another inverter
 - (e) inverter gain \rightarrow useful concept for noise margins.

Inverter can work even when supply voltage $< V_{TN}$ because it still conducts in sub-threshold. As long as $|gain| > 1$, inverter will continue to work.



$$\left. \begin{aligned} NM_L &= V_{IL} \\ NM_H &= V_{DD} - V_{IH} \end{aligned} \right\} \text{noise margins}$$

$$g = \text{gain} = \left(\frac{dV_{out}}{dV_{in}} \right)$$

\rightarrow We want $V_{IL} \uparrow$ and $V_{IH} \downarrow$
 Max. value of $V_{IL} = V_m$ and min. value of $V_{IH} = V_m$.

$$\underline{V_m = V_{IL} = V_{out}} \leftarrow \text{switching threshold}$$

V_m is calculated when both NMOS & PMOS are in saturation.

$$V_m = \frac{V_{DD} - |V_{TP}| + \sqrt{\frac{\beta_n}{\beta_p}} V_{TN}}{1 + \sqrt{\frac{\beta_n}{\beta_p}}}$$

A general rule is that when PMOS becomes stronger then V_m goes toward V_{DD} , while when NMOS becomes stronger then V_m goes toward GND.

What does a transistor being "strong" mean?

A transistor is said to be strong when its ON current is high. By ON-current, we mean the current in saturation.

A transistor can be made strong by increasing its size or by increasing $(V_{GS} - V_t)$ drop.

$$\text{For } V_m = \frac{V_{DD}}{2} \Rightarrow \sqrt{\frac{\beta_n}{\beta_p}} = \frac{0.5 V_{DD} - |V_{TP}|}{0.5 V_{DD} - V_{TN}}$$

General equation for $V_m = V_{DD}/2$

Switching Delay of the inverter

$$t_{PLH} = 0.69 R_{eq,p} C_L$$

$$t_{PHL} = 0.69 R_{eq,n} C_L$$

$$t_p = \frac{1}{2} (t_{PLH} + t_{PHL})$$

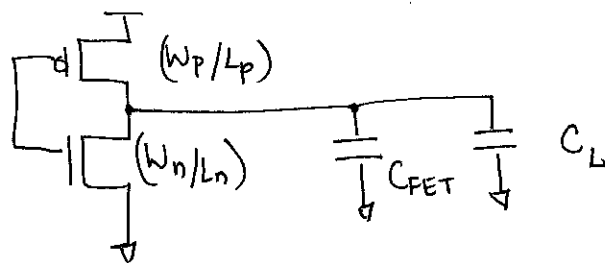
For making $t_{PLH} = t_{PHL}$

$$\frac{R_{eq,p}}{R_{eq,n}} = \frac{\mu_n}{\mu_p}$$

$$R_{eq,p} = \frac{3}{4} \frac{V_{DD}}{I_{Dsatp}}$$

$$R_{eq,n} = \frac{3}{4} \frac{V_{DD}}{I_{Dsatn}}$$

What is C_L ? (IMPORTANT)



Assume $L_p = L_n = L_{min}$

C_{FET} = PARASITIC

C_L = Fan-out or a constant load + wire.

Lets compute C_{FET}

$$C_{FET} = C_{FET-nmos} + C_{FET-pmos}$$

$$C_{FET-nmos} = C_{jn} W_n + 2 C_{ov} W_n$$

$$C_{FET-pmos} = C_{jn} W_p + 2 C_{ov} W_p$$

$$C_{FET} = C_{FET-nmos} \left(1 + \frac{W_p}{W_n} \right) \quad \text{V.V. Important.}$$

$$t_{PLH} = \frac{0.693 V_{DD}}{4 I_{Dsatp}} \left(C_{FET-nmos} \left(1 + \frac{W_p}{W_n} \right) + C_L \right)$$

Assume that there is no extra load at the output from C_L but that the inverter is self loaded.

$$t_{PLH} = \frac{0.693 V_{DD}}{4 I_{Dsatp}} C_{FET-nmos} \left(1 + \frac{W_p}{W_n} \right)$$

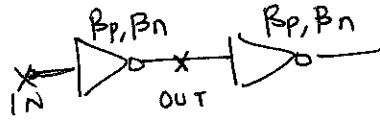
$$I_{Dsatp} = \frac{1}{2} \beta_p (V_{DD} - |V_{TP}|)^2$$

$$t_{PLH} = \frac{V_{DD}}{\beta_p (V_{DD} - |V_{TP}|)^2} C_{FET-nmos} \left(1 + \frac{W_p}{W_n} \right)$$

As $\left(\frac{W_p}{L_p} \right) \uparrow$ $\beta_p \uparrow$ $R_{eq,p} \downarrow$ but parasitic

capacitance increases. Hence, self-loading ~~does~~ means that increasing the device size does not affect the delay.

When a CMOS inverter is loaded by an identical inverter, choose $\frac{W_p}{W_n} = \alpha$



How do you minimize the delay from IN to OUT?

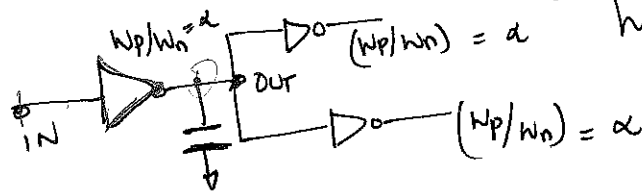
Calculate the net capacitance at the output.

$$C_L = \underbrace{C_{FET,n}}_{\text{parasitic}} (1 + \underbrace{\alpha}_{(W_p/W_n)}) + C_{\text{wire}} + (1 + \alpha) \underbrace{C_{g1}}_{\text{input gate cap}}$$

$$t_d = 0.5 \frac{C_L V_{DD}}{(V_{DD} - V_T)^2} \left(\frac{1}{\beta_n} + \frac{1}{\beta_p} \right)$$

Substitute C_L in the above equation & minimize the delay $\frac{\partial t_d}{\partial \alpha} = 0 \Rightarrow \alpha_{\text{opt}} = \sqrt{\frac{\mu_n}{\mu_p}}$ For neg. C_w .

Now if the inverter is driving two identical inverters.



how should you size the first inverter

$$C_L = C_{FETn} (1 + \alpha) + C_{\text{wire}} + \underbrace{2(1 + \alpha)}_{\text{extra due to two inverters as the Fan-out}} C_{g1}$$

$$t_d = \frac{0.5 C_L V_{DD}}{(V_{DD} - V_T)^2} \left(\frac{1}{\beta_n} + \frac{1}{\beta_p} \right) \Rightarrow \frac{\partial t_d}{\partial \alpha} = 0$$

So the general principle of obtaining delay is the same. All you need to know is two things: (a) Equivalent resistance (b) Load Capacitance.

For parasitic capacitance of a transistor, I always use the symbol C_{FET} .

For input capacitance of a transistor, I always use the symbol C_g .

$$d = p + g + b \quad \begin{matrix} g=1 \\ b=1 \end{matrix}$$

$$d = p + h \quad h = \frac{C_{out}}{C_{in}} = 2$$

$$d = 2 \text{ for } p = 0$$

$d \Rightarrow$ normalized delay.

$$\tau_{ref} = 0.69 R_{ref} C_{ref}$$

$$C_{ref} = (1 + \alpha) C_{in}$$

$$D = 0.69 \underline{R_{ref}} (1 + \alpha) C_{in} \times 2$$

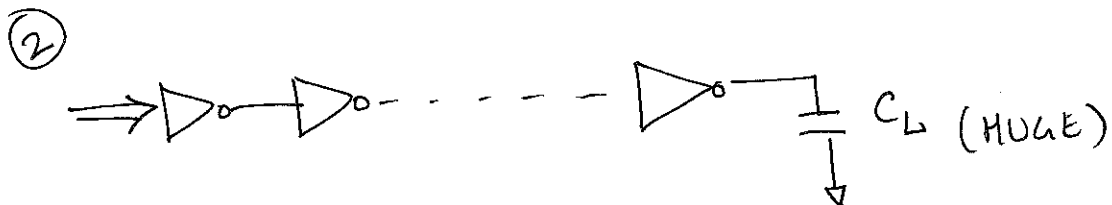
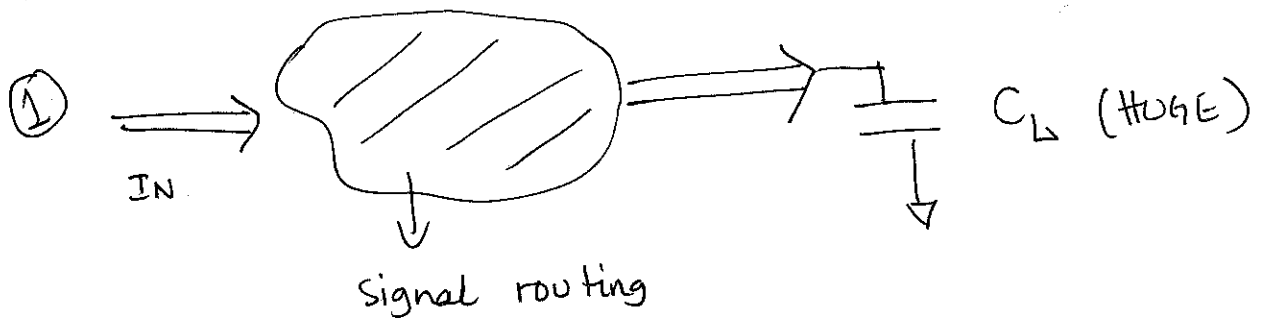
$$R_{ref} = \frac{3}{4} \frac{V_{DD}}{I_{sat}}$$

where $I_{satp} = I_{satn}$
so that pull-up & pull-down are the same.

The whole deal about inverter sizing :-

often times we want to route a signal from one point on the chip to another and the load capacitance is large. It always makes sense to route the signal via a chain of inverters to minimize the propagation delay.

Consider the scenario :-

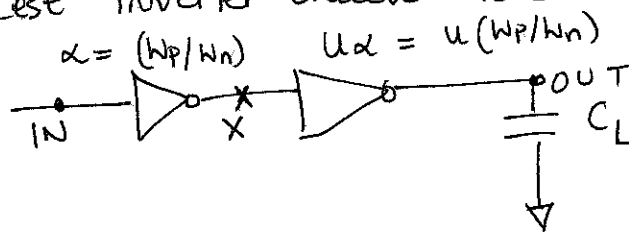


It was shown a long time ago that if you were to use ② to drive the large load and sized the inverters appropriately, the delay of ② $<$ delay of ①.

So the question becomes a) how do you size this chain of inverters b) how large should this load capacitance be for us to know that scenario ② makes sense.

First, let's address point (b) above. That is, how large should C_L be for us to consider the inverter chain.

The smallest inverter chain will be



$R_{eq} \rightarrow$ 1st inverter equivalent resistance.
assume $R_{eqn} = R_{eqp}$
For a given α .

Let's say the first inverter is sized ' α ' and second one is sized ' $U\alpha$ '.

$$t_{\text{delay}} = t_{\text{INV1}} + t_{\text{INV2}}$$

$$t_{\text{INV1}} = 0.69 R_{eq} C_x = 0.69 R_{eq} (C_{\text{INV1}} U)$$

$$t_{\text{INV2}} = 0.69 \frac{R_{eq}}{U} C_L$$

$$C_x = \text{input capacitance of inverter 2} \\ = C_{\text{INV1}} * U$$

Now we have the delay:-

$$t_d = 0.69 R_{eq} C_{INV,1} u + 0.69 \frac{R_{eq}}{u} C_L$$

$$\frac{\partial t_d}{\partial u} = 0 \Rightarrow u_{opt} = \sqrt{\frac{C_L}{C_{INV,1}}}$$

$$\begin{aligned} t_{d,opt} &= 0.69 R_{eq} \sqrt{C_L C_{INV,1}} + 0.69 R_{eq} \sqrt{C_L C_{INV,1}} \\ &= 2 \times 0.69 R_{eq} \sqrt{C_L C_{INV,1}} \end{aligned}$$

Without the second inverter, the delay of driving the load will be :-

$$t_d = 0.69 R_{eq} C_L$$

∴ Adding the second inverter makes sense

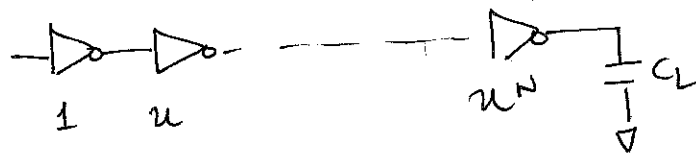
only when $2 \times 0.69 R_{eq} \sqrt{C_L C_{INV}} < 0.69 R_{eq} C_L$

$$\left[\frac{C_L}{C_{INV}} > 4 \right]$$

Makes sense to add only when the load is atleast "4" times bigger than the input capacitance of the 1st inverter.

Now we know when we must add a chain of inverters to drive a large load.

Key idea is :-



Make the 1st inverter in the chain as the "unit" inverter or the reference inverter. Upsize each following inverter by a factor of u .

$$N_{opt} = \ln \left(\frac{C_L}{C_1} \right) \Rightarrow u^{N_{opt}} = (C_L / C_1)$$

$$u_{opt} = e$$

In a general scenario, first calculate N_{opt} and round it off to the nearest integer.

Then calculate $u_{opt} = (C_L / C_1)^{1/[N_{opt}]}$ ~~where $[N_{opt}]$~~

~~is~~
greatest int.
less than or
equal to.

→ This is equivalent to saying let's make the electrical effort of each stage the same.

Electrical effort is (output cap)/(input cap)

ENERGY DISSIPATION

$$E_{VDD} = \int_0^{\infty} i_{VDD}(t) V_{DD} dt$$

↓
current drawn from supply

$$E_c = \int_0^{\infty} i_c(t) V_{out} dt$$

$$E_R = \int_0^{\infty} i_R^2(t) R dt$$

BASIC
DEFINITIONS.

AVERAGE POWER DISSIPATION

$$P_{avg} = \alpha C_L V_{DD}^2 f_{clk} + \alpha E_{sc} + V_{DD} I_{static}$$

↓ ↓ ↓
Dynamic Short static power
power circuit Leakage

α = Activity Factor

Fraction of the times, the output node transitions
From "0" to "1"

STATIC CMOS GATES

Consist of pull-up network w/ PMOS

Pull-down network w/ NMOS.

- While sizing the gates, figure out the worst case path for low to high & high to low transition & make it equal to the reference inverter.
- When considering dependence of input pattern on the delay, pay special attention to the internal node capacitance.

Large Fan-in and large FO effect on delay can be represented as :-

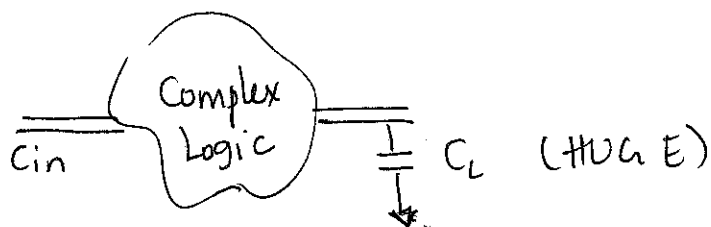
$$t_p = a_1 \underset{\substack{\downarrow \\ \text{Fan in}}}{FI} + a_2 FI^2 + a_3 \underset{\substack{\rightarrow \\ \text{Fan-out}}}{FO}$$

Remember techniques to size gates to minimize the delay.

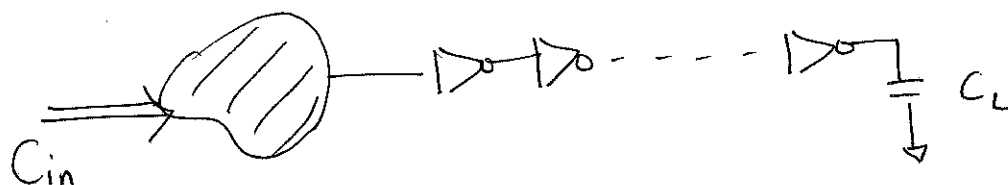
- progressive sizing.
- Input re-ordering.
- alternative logic structures

d) Isolating Fan-in from Fan out.

Let's look at the last point in more detail



If we have a large capacitance to drive then we can do something like



Remember this technique of buffering makes sense only when $F = \frac{C_L}{C_{in}} > 4$.

Also correlate this with logical effort discussion.

→ Size each one of these

$$\left. \begin{array}{ll} p = e & \text{For } \frac{p=0}{p=1} \\ p = 3.6 & \text{For } \underline{p=1} \end{array} \right\} \text{parasitic}$$



How should we size each gate in this logic path.

The idea is to make the stage delays equal
That is, $f_1 = f_2 = \dots = f_N = F^{1/N} = (BGH)^{1/N}$

$$\left[\begin{array}{l} f_i = b_i g_i h_i \\ b_i = \text{branching at } i^{\text{th}} \text{ stage} \\ g_i = \text{logical effort of } i^{\text{th}} \text{ stage} \\ h_i = \text{electrical effort of } i^{\text{th}} \text{ stage} \end{array} \right]$$

$$B = \prod_i b_i$$

$$G = \prod_i g_i$$

$$H = \prod_i h_i = (C_L/C_1)$$

$$D_{\text{opt}} = N F^{1/N} + \underbrace{P}_{\text{total parasitic delay}}$$