

Layout power optimization

(129)

- Layout design is very important for performance and power optimization in deep submicron technology.

∴ it determines the lengths & widths of interconnect.

- In deep submicron design, 50% ~ 70% of clock cycle is consumed by interconnect delay.
- Significant amount of power is spent on charging & discharging interconnects.

∴ Interconnect capacitance ↑

short circuit power

static power

} secondary

factor when compared to total capacitive power

- Will cover:

— power & delay models ↔ interconnects & gates

— device layout optimization

— interconnect layout optimization

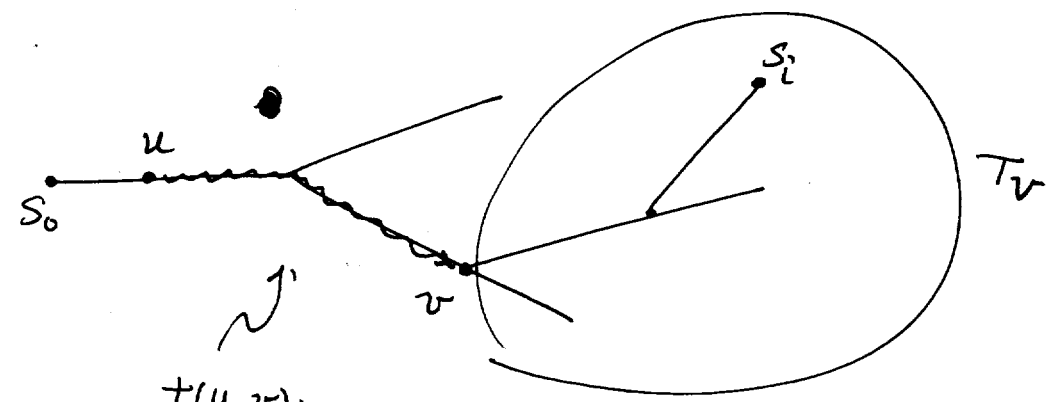
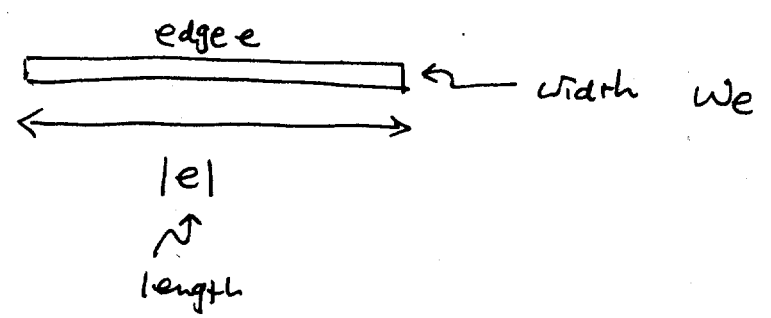
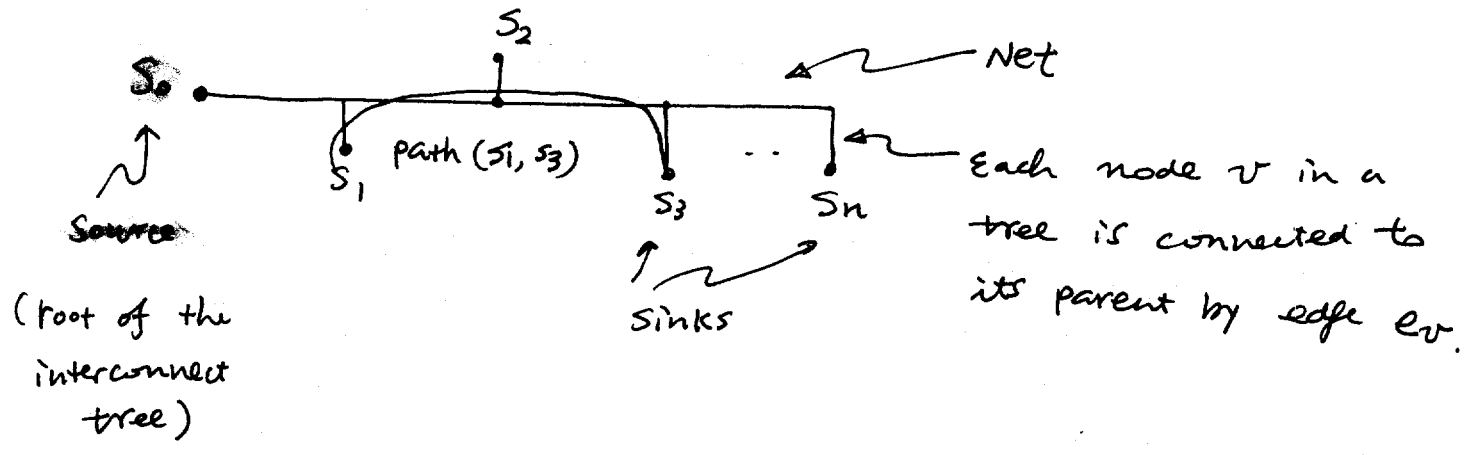
— device & interconnect simultaneous optimization

↙ must study this

∴ delay

must be maintained same level while optimizing for power

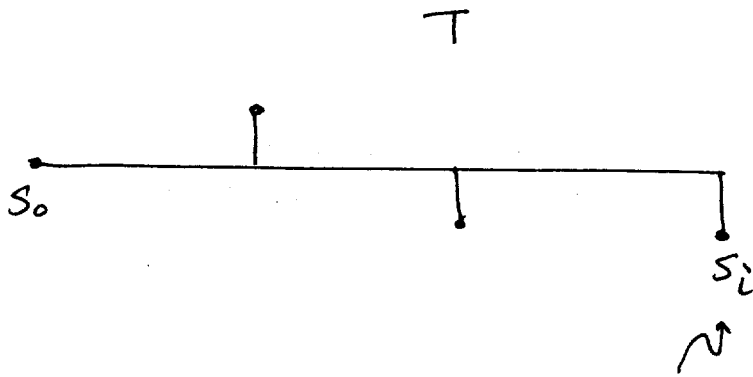
~~Definition~~



$t(u, v)$:
signal delay from u to v .

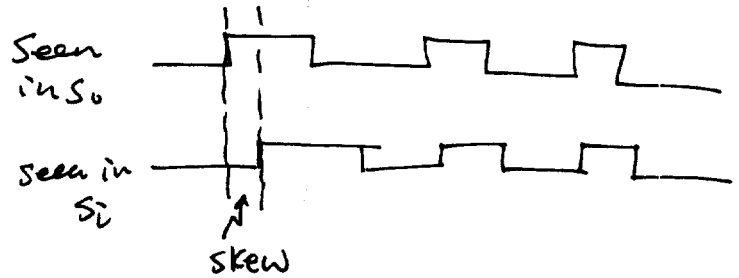
$t(s_0, s_i)$: delay from source to sink s_i

$t(s_0, s_i)$ is shortened by t_i .



clock skew of s_i :

The difference in the clock signal delay from source to s_i .



$$\text{skew}(T) = \max_{s_i, s_j \in S} |t_i - t_j|$$

$S = \{ \text{the set of all sinks} \}$.

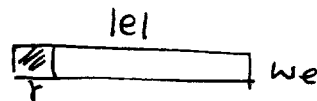


r : unit square wire resistance

C_a : unit area capacitance

C_f : unit length fringing capacitance. (both sides)

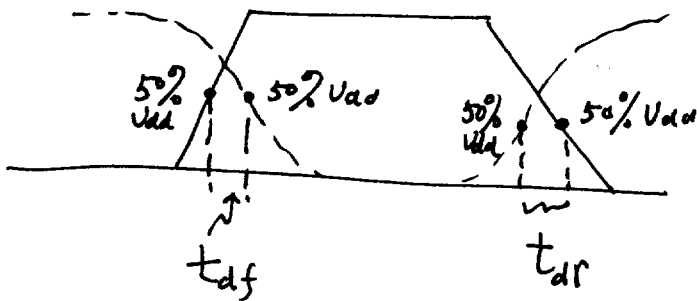
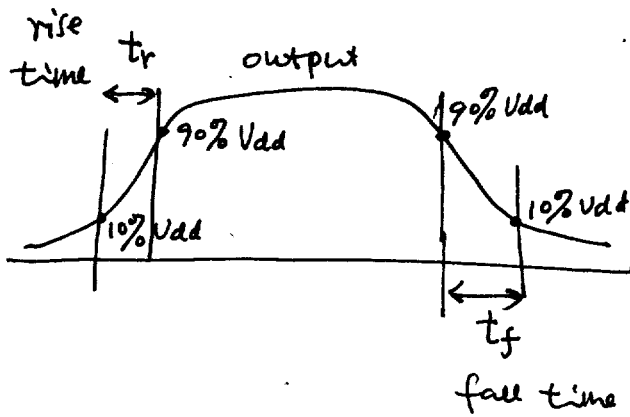
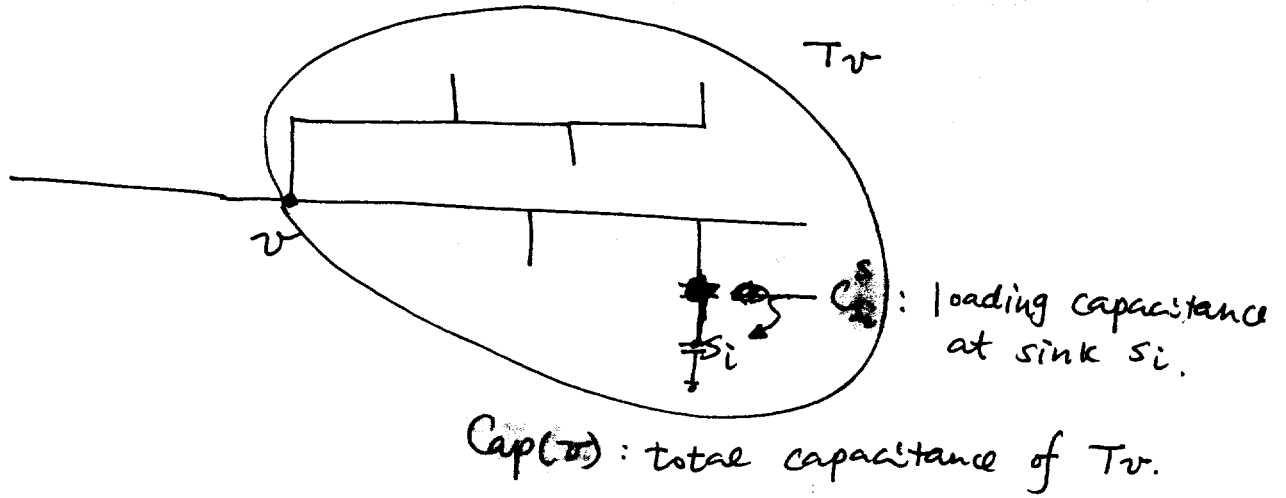
Wire resistance of edge e :



$$r_e = \frac{r \cdot |e|}{w_e}$$

Wire capacitance of e

$$C_e = C_a |e| w_e + C_f |e|$$



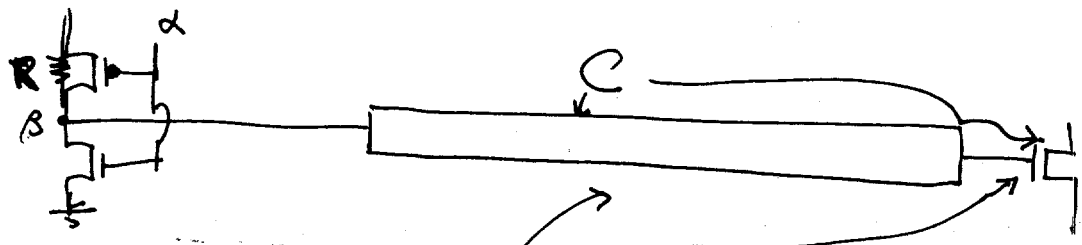
delay time for falling signal

delay time for rising signal

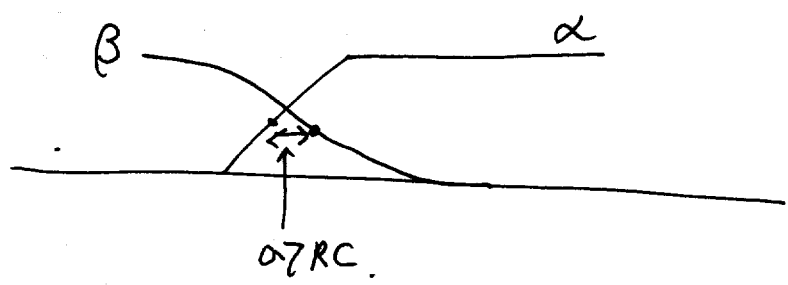
Interconnect Delay Model

Lumped delay model \rightarrow pathlength delay model \rightarrow Elmore delay model.

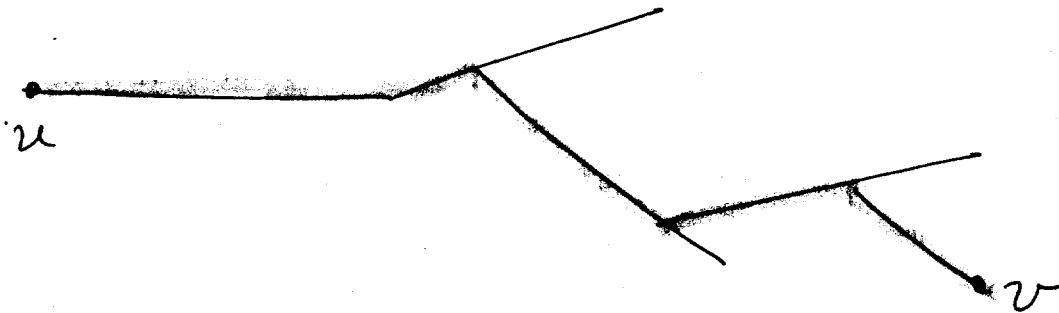
Lumped Capacitance Model



- The source driver sees C as a single load capacitance
- Ignore the resistance at the interconnect (!: driver $r \gg$ interconnect r)
- The delay for a ^(driver) gate to switch to 50% of its final value under step input $\approx 0.7RC$.
- True for feature size $\geq 1.2 \mu m$.



• pathlength (linear) model:



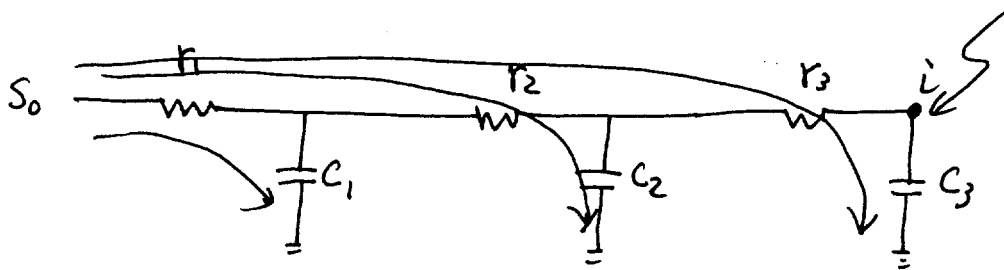
$$t(u, v) \propto \sum_{e \in \text{path}(u, v)} |e_w|$$

Limitations:

- ① It ignores the wire resistance, considers only wire capacitance along the path.
- ② It ignores the effect of edges not along the path.

• Elmore delay model.

- The most commonly used delay model.



$$t(s_0, i) = r_1 C_1 + (r_1 + r_2) C_2 + (r_1 + r_2 + r_3) C_3$$

either charging or discharging

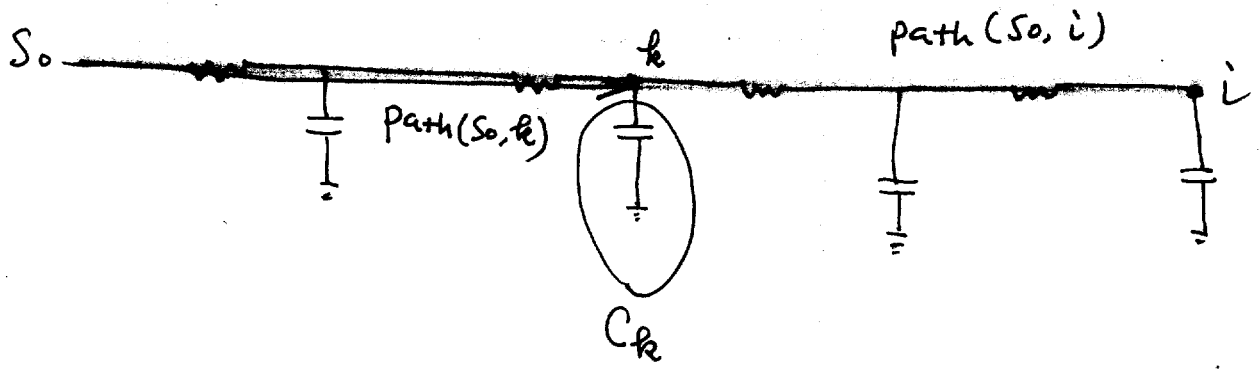
$$= r_1 (C_1 + C_2 + C_3) + r_2 (C_2 + C_3) + r_3 C_3$$

• General form

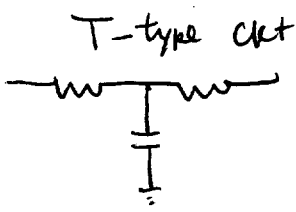
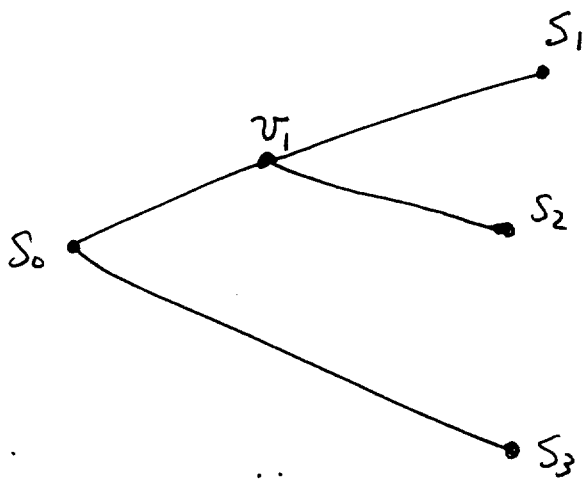
$$t(s_0, i) = \sum_{\text{all nodes } k} R_{ki} \cdot C_k$$

C_k : lumped capacitance at node k .

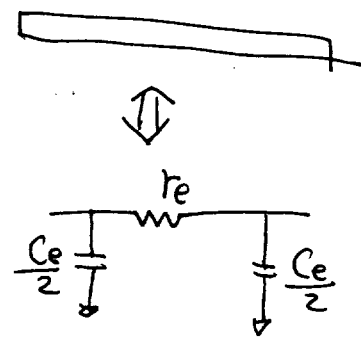
R_{ki} : The resistance of the portion of path(s_0, i)
Common with path(s_0, k)

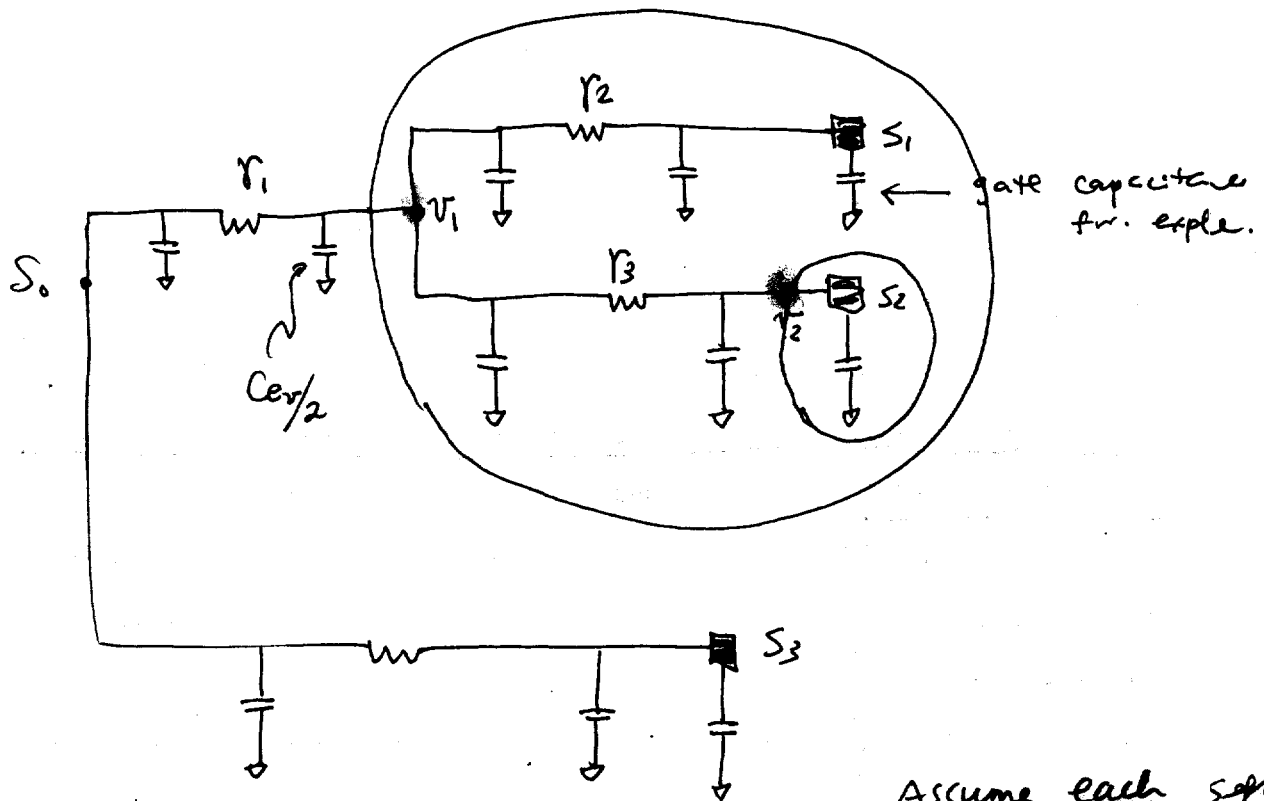


Example:



Each edge can be modeled by Π -type circuit





Assume each segment has cap = C_{ev} .

$$t(S_0, S_i) = \sum_{e_v \in \text{path}(S_0, S_i)} r_{e_v} \cdot \left(\frac{C_{ev}}{2} + \text{cap}(v) \right)$$

$$t(S_0, S_2) = r_1 \left(\frac{C_{ev}}{2} + \text{cap}(v_1) \right) + r_3 \left(\frac{C_{ev}}{2} + \text{cap}(v_2) \right)$$

- Good things for Elmore model:
 - * An optimal or near optimal solution using this model is also nearly optimal using spice for routing, wiresizing, clock skew.

Disadvantages of Elmore model:

① Absolute value of Elmore delay may not be accurate.

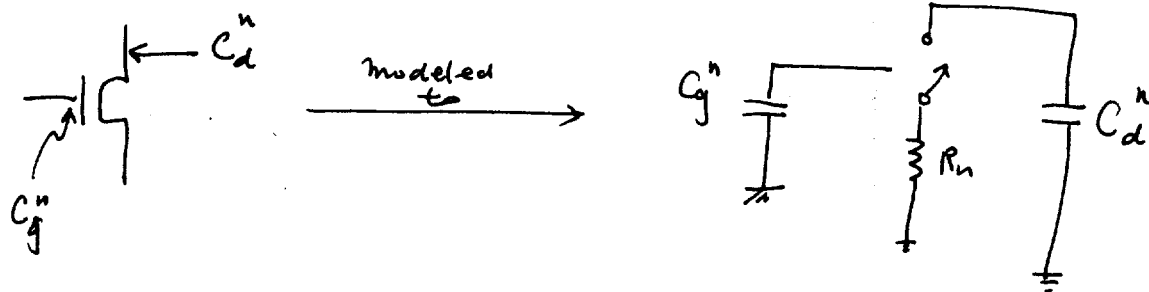
② It can't handle inductive effects

∴ More accurate delay can be modeled by RLC tree for interconnects.

$$\text{Delay} \cong k l \sqrt{L(\omega) C}$$

\uparrow \uparrow
 constant length

Different delay model



R_n : resistance of minimum transistor n .

C_g^n : gate cap. of min n transistor

C_d^n : output diffusion cap. of min n transistor.

• For an arbitrary transistor with size $d \geq 1$.

$$R_d = R_n / d$$

\uparrow
Width

$$\text{gate cap} = C_g^n \cdot d$$

$$\text{diff cap} = C_d^n \cdot d \leftarrow \text{part of } C_L$$

fixed channel length

$$t_f = K \cdot \frac{C_L}{\beta_{min}^n \cdot d \cdot V_{DD}}$$

$$\beta_{min}^n = \frac{W}{L} \cdot \frac{\epsilon_{ox} \cdot \mu}{t_{ox}}$$

for step input.

C_L includes diff cap & load capacitance.

$$\beta_{min}^n = \frac{W}{L} \cdot \frac{\epsilon_{ox} \cdot \mu}{t_{ox}}$$

$$t_{df} = t_f / 2$$

↑ delay time for the falling signal.

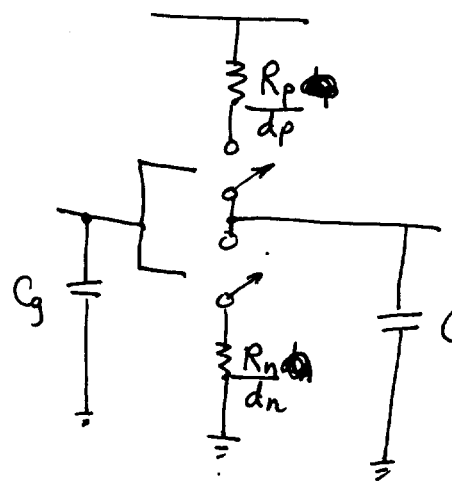
∴ $R_d \propto \frac{1}{\beta_{min} \cdot d}$

$$t_{df} = t_f / 2 = \frac{\text{constant}}{2 V_{DD}} \left(\frac{1}{\beta_{min} \cdot d} \right) C_L$$

∴ $t_{df} \propto R_d \cdot C_L$

t_{df} can be approximated by $R_d \cdot C_L$

Inverter



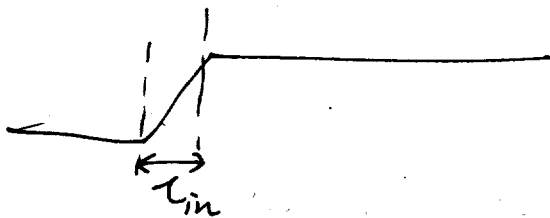
← sum of diff cap. due to the p-, n-transistors.

$$t_f = K \cdot \frac{C_L}{\beta_{min}^n \cdot d_n \cdot V_{DD}}$$

$$t_r = K' \cdot \frac{C_L}{\beta_{min}^p \cdot d_p \cdot V_{DD}}$$

• Shortcoming of the above simple RC model:

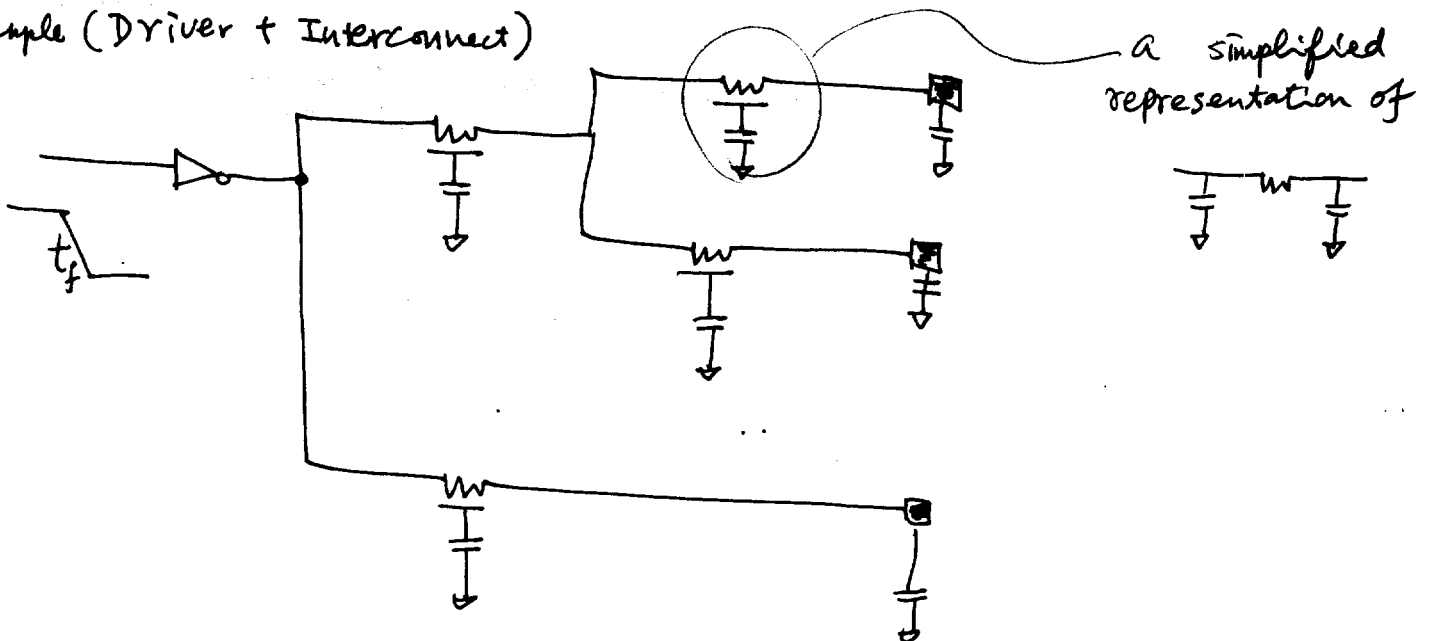
Can't deal with the shape of input waveform.



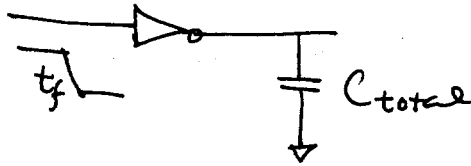
$$t_{df} = t_f/2 + \frac{t_{in}}{6} \cdot \left(1 + 2 \frac{V_{Th}}{V_{DD}}\right)$$

threshold voltage of n-transistor.

• Example (Driver + Interconnect)



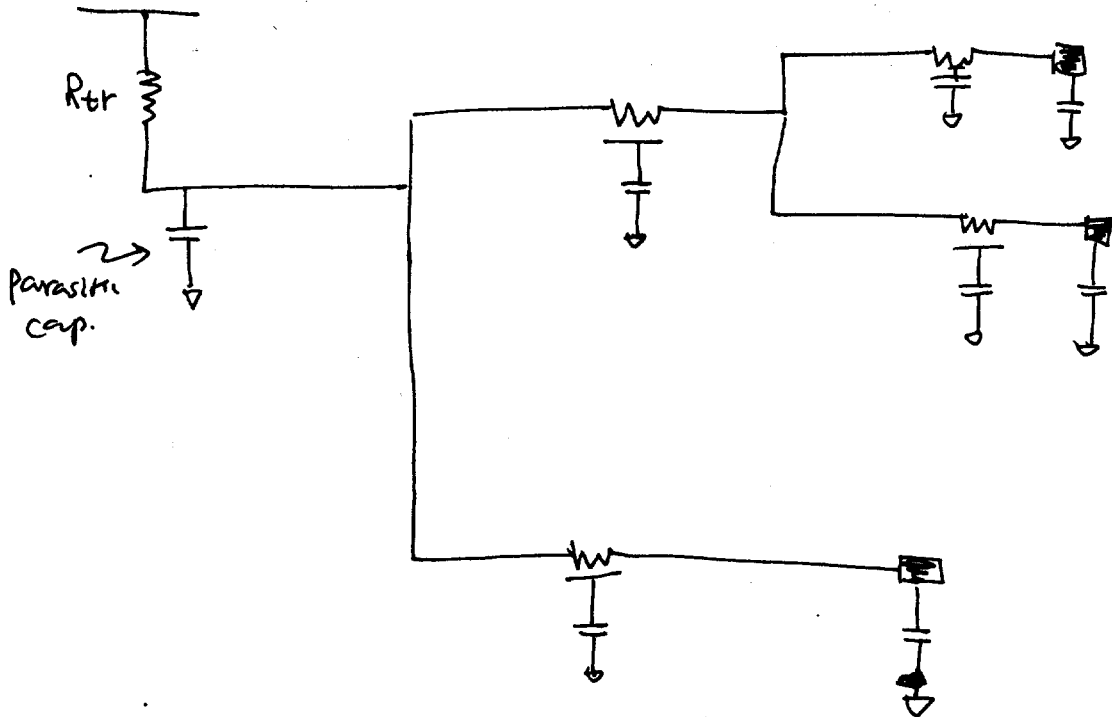
Can be simplified to



$$t_r = k' \frac{C_{total}}{\beta_{min} \cdot d_p \cdot U_{DD}}$$

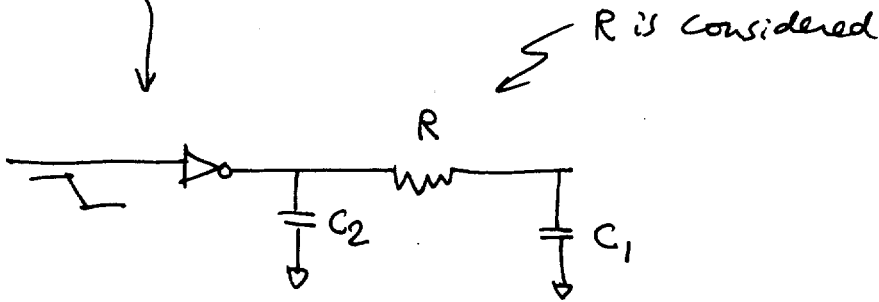
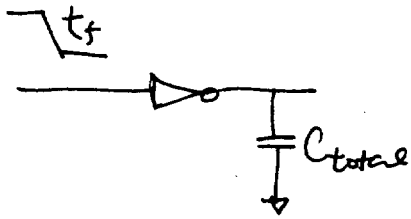
← include $C_{parasitic}$
 $(C_{diff(P)} + C_{diff(N)})$
 $+ C_w + C_{driven-gate}$

• A better estimation

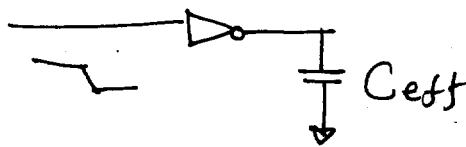


Then, use Elmore model to estimate the delay.

Another Estimation



Compute iteratively the effective capacitance seen by the driver.



power model

Short Circuit

input transition time

$$\bar{I}_{sc} = \frac{\beta \cdot t_{in}}{12 \cdot V_{DD}} \cdot (V_{DD} - 2V_T)^3 \cdot f$$

$$P_{sc} = \bar{I}_{sc} \cdot V_{DD} = \frac{\beta \cdot t_{in}}{12} (V_{DD} - 2V_T)^3 \cdot f$$

$$\beta = \frac{W_{eff}}{L_{eff}} \cdot \frac{\epsilon_{ox} \cdot \mu}{t_{ox}}$$

~~Assume $\beta_n = \beta_p = \beta$~~

The gain factor

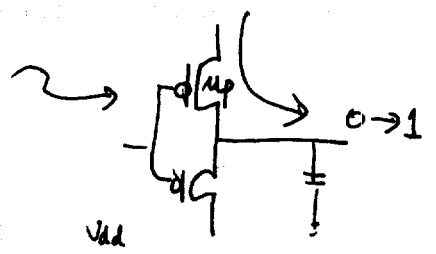
$$\beta = \frac{W}{L} \cdot \frac{\epsilon_{ox} \mu}{t_{ox}}$$

is determined by width of transistor W &

mobility of carrier responsible for the transition

(μ_p for low-to-high transition)

(μ_n for high-to-low transition)



is P_{sc} for output

$\hookrightarrow P_{sc}$ can be reduce to

dominated by pull-up transistor.
 ∴ pull-down was on

$$P_{sc} = K \cdot \mu W L_{in} f$$

↑
comes from
 β .

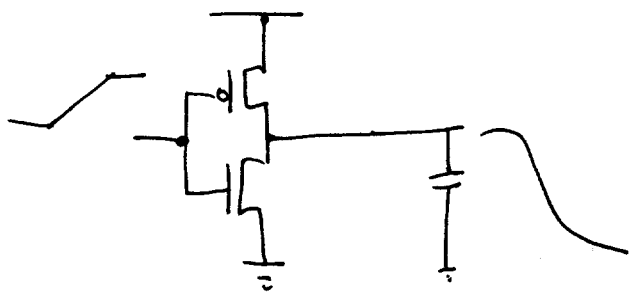
P_{sc} for output
 dominated by pull-down transistor
 ∴ pull-up was on.

process & voltage dependent.

$$P_{sc} = \frac{W_{eff}}{L_{eff}} \frac{\epsilon_{ox} \mu}{t_{ox}} \frac{\tau_{in}}{12} (V_{DD} - 2V_T)^3 \cdot f$$

$$= \left[\frac{\epsilon_{ox} \cdot (V_{DD} - 2V_T)^3}{L_{eff} \cdot t_{ox} \cdot 12} \right] \mu \cdot W_{eff} \cdot \tau_{in} \cdot f$$

↑
K



When input transition time \gg output transition time

\Rightarrow short ckt path achieved for a longer time

\rightarrow incur more power dissipation

\therefore For power minimization, try to achieve equal input and output transition times. &

Keep both as small as possible.

Dynamic power

$$P_{cap} = \frac{1}{2} E(sw) C_L V_{DD}^2 f$$

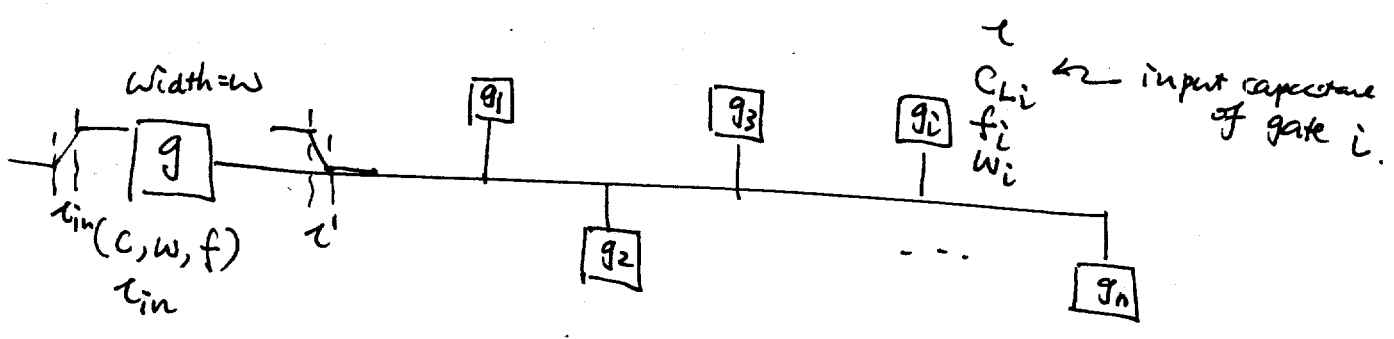
\uparrow
includes device capacitance &
interconnect capacitance

\uparrow
critical for deep submicron technology.

~~Transistor and Gate Sizing~~

- Transistor sizing: Determine the optimal width for each transistor
 \Rightarrow minimize delay or power under the delay constraint.

• Gate sizing: All transistors in a gate are scaled by a common factor to perform power/delay optimization.



C : input capacitance for a minimum-width transistor

W : size of switching transistor (width)

f : switching frequency of g .

t_{in} : input transition time of g .

t : output transition time of g , i.e.,
input transition time of g_1, g_2, \dots, g_n .

pull-up w
pull-down
depend on
output w

- If output w is for pull-down

- If output w is for pull-up.

⑥

~~ACET~~ Project III assigned

Due 3/17/03

Interconnect Layout optimization.

- ckt partitioning

* Gate or Transistor
size

- placement

- Routing

* Transistor reordering

- Wire sizing

definitions

properties: Monotone.

Separability

Dominance

Sizing method

$$P = V_{dd}^2 (C_{gate} \cdot w \cdot f + \sum_{i=1}^n C_{Li} \cdot f_i) + \dots$$

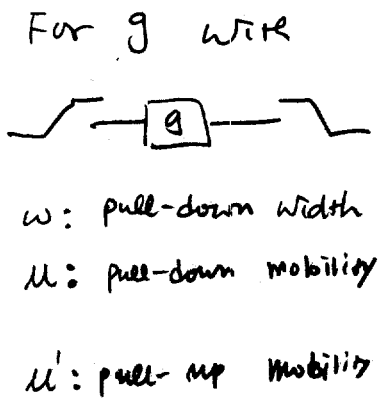
Capacitive power for g (gate) → $C_{gate} \cdot w \cdot f$
 Power of gate cap of driver
 Capacitive power for g_1, g_2, \dots, g_n
 → $\sum_{i=1}^n C_{Li} \cdot f_i$
 → $C_{L1}, C_{L2}, \dots, C_{Ln}$
 → can be ignored.

Let $P_1 = V_{dd}^2 \sum_{i=1}^n C_{Li} \cdot f_i$

$$k (\mu \cdot \tau_{in} \cdot w \cdot f + \sum_{i=1}^n \mu' \cdot \tau \cdot w_i \cdot f_i)$$

mobility of carrier for g → μ
 mobility of carrier for g_i → μ'
 Comes from β of P_{sc}

Constant based on technology & voltage
 (comes from β & V_{DD}) of P_{sc} .



The delay output transition time of g:

$$\tau = \phi \frac{C_L}{\mu W_0} + \frac{\tau_{in}}{6} \left(1 + \frac{2V_T}{V} \right)$$

process dependent constant → ϕ
 load cap. of g → C_L

$$P = P_1 + (k_1 \mu \tau_{in} + k_2) w \cdot f + k \mu' \sum_{i=1}^n w_i f_i \left(\phi \frac{C_L}{\mu} + k_2 \tau_{in} \right)$$

$k_1 = V_{DD}^2 \cdot C$
 g short-ckt → k_1
 g source → k_2
 $g_1 \dots g_n$ → $\sum_{i=1}^n w_i f_i$
 Switching power

↑
 $g_1 \dots g_n$
 short-ckt power

→ $p = \alpha * W + \frac{\beta}{W} + \sigma$
 ← Increasing w → increase power of g
 reduce power of $g_1 \dots g_n$
 ∴ trade-off

$\frac{\partial P}{\partial W} = 0$ gives the condition where P is minimum on W .

~~for~~

$$W = \sqrt{\phi \frac{\mu'}{\mu} C_L \left(\sum_{i=1}^n W_i f_i \right)}$$

$$\sqrt{\left(\frac{k_p}{k_n} + \mu \tau_{in} \right) f}$$

the power optimal p-transistor size:

$\mu \leftarrow \mu_p$

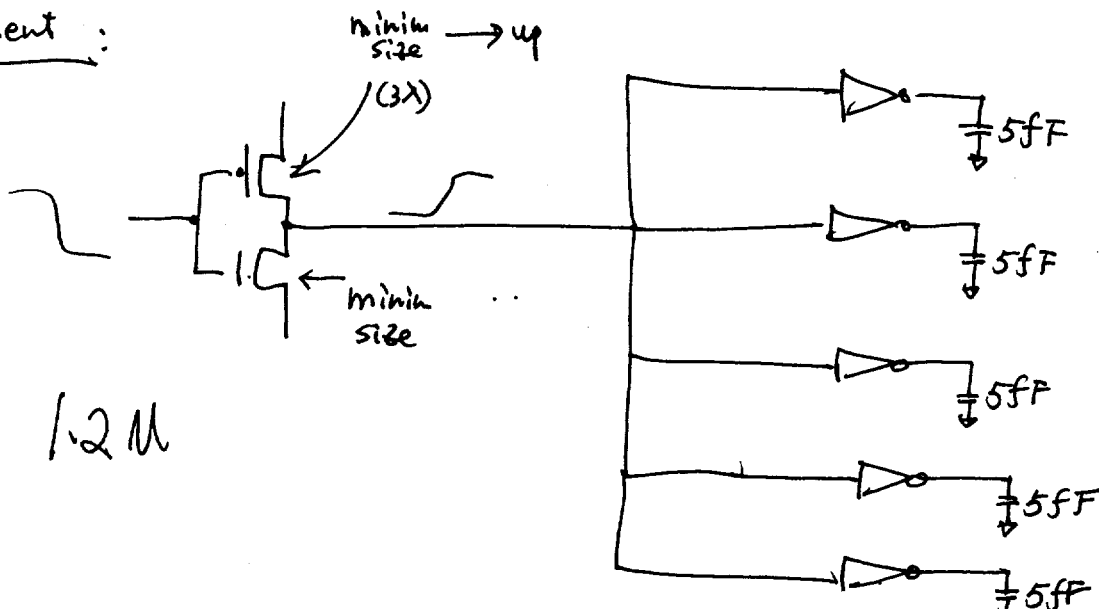
$\mu' \leftarrow \mu_n$

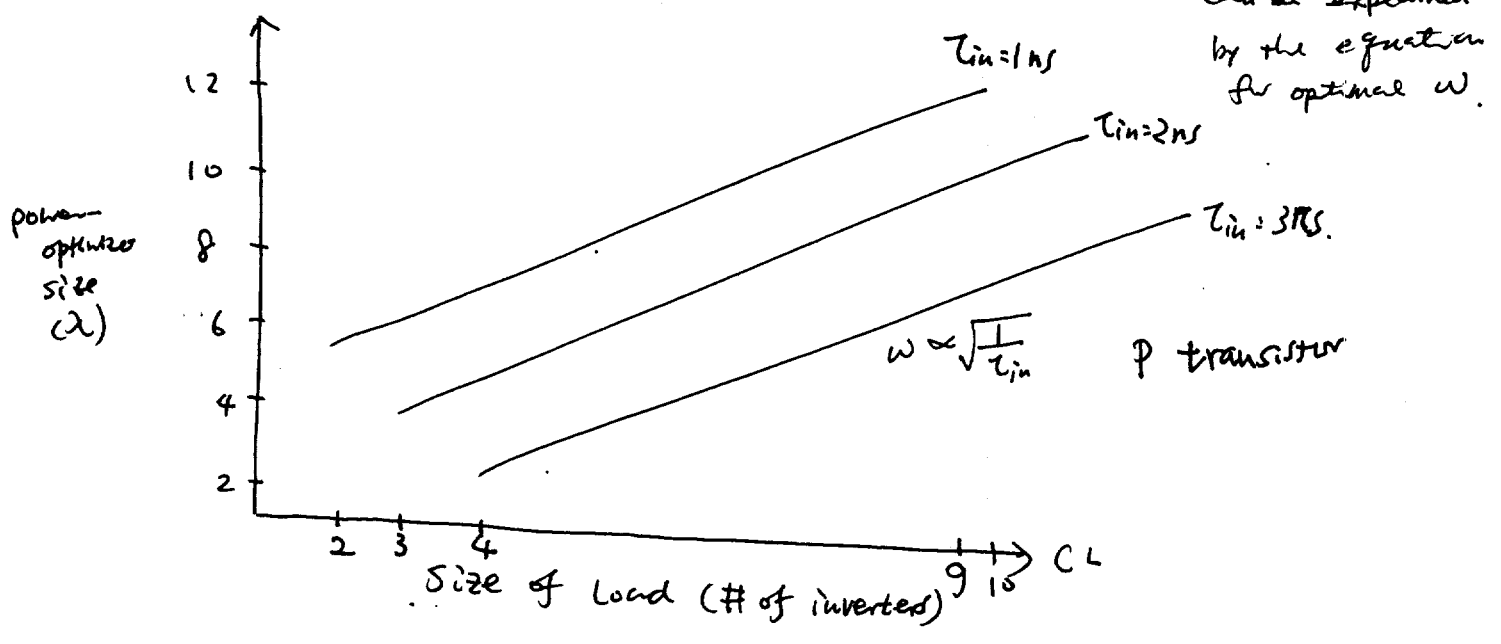
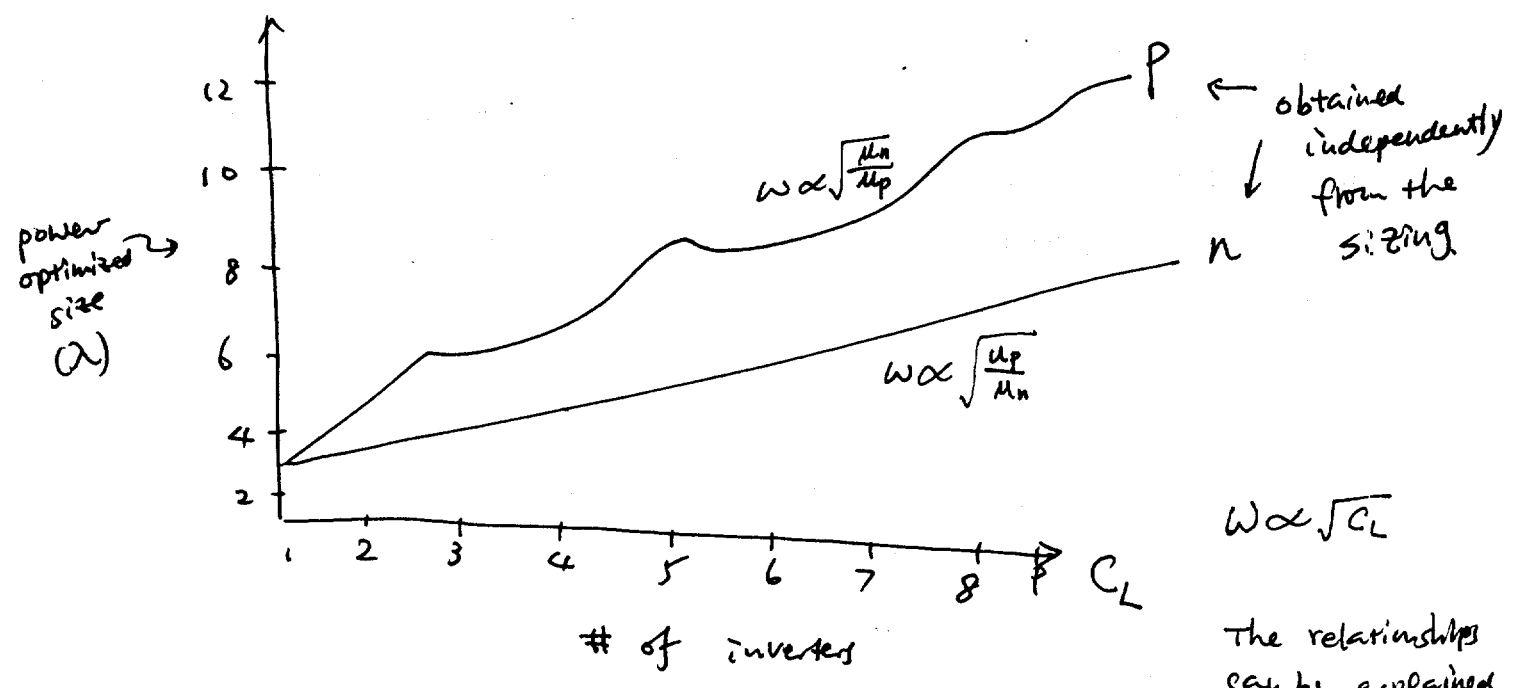
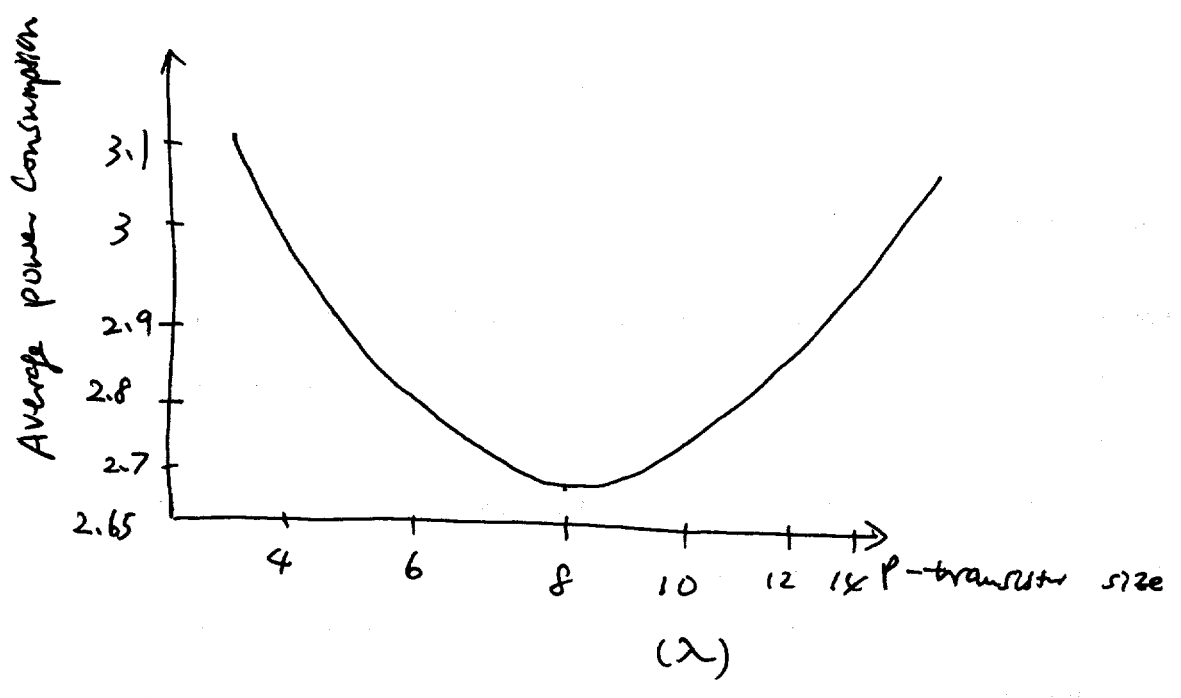
the power optimal n-transistor size:

$\mu \leftarrow \mu_n$

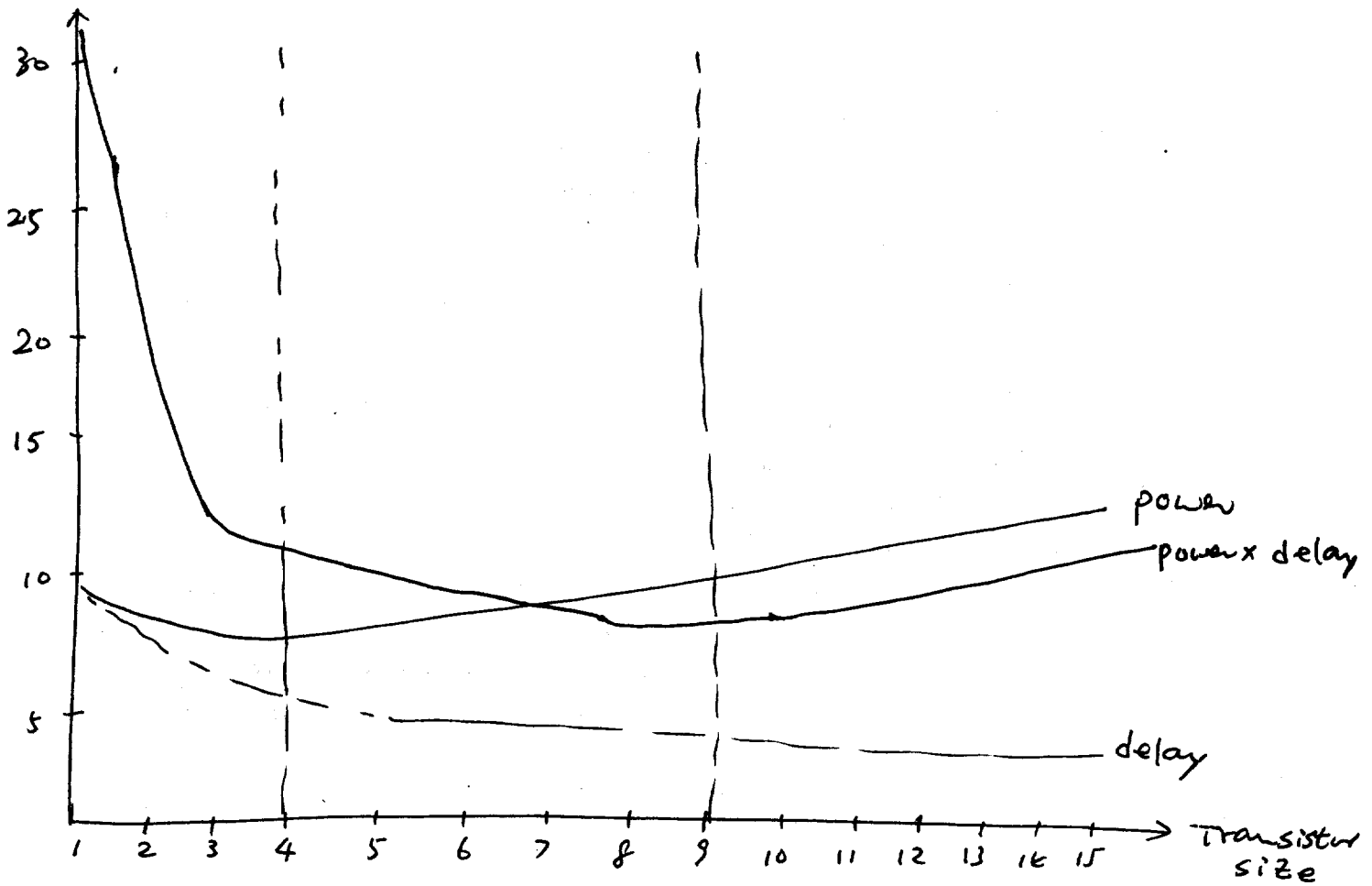
$\mu' \leftarrow \mu_p$

Experiment :





The relationships can be explained by the equation for optimal w .



- ↑
- ① Decrease in power & delay
- ② Rapid decrease in power-delay curve

- ↑
- ① Decrease in delay
- ② increase in power
- ③ Slow decrease in power-delay curve

- ↑
- ① increase in power-delay curve
- ② Δ (increase in power)
- rate $>$ Δ (decrease in delay)

P transistor of driver

Transistor sizing method for power optimization

① All transistors are assigned power-optimal sizes.
(region 1)

② Timing analysis is performed.

③ While (\exists path with delay $>$ constraint)
power-delay ^{optimal} ~~optimal~~ sizes are used for all transistors
in the critical path.

④ If (\exists path with delay $>$ constraint)
the transistor with least slope on the power-delay curve
is sized up

⑤ If (\exists path with delay \ll constraint)
the transistor with steepest power-delay curve
is sized down.

⑥ Repeat ④ and ⑤ until delay constraint is satisfied.

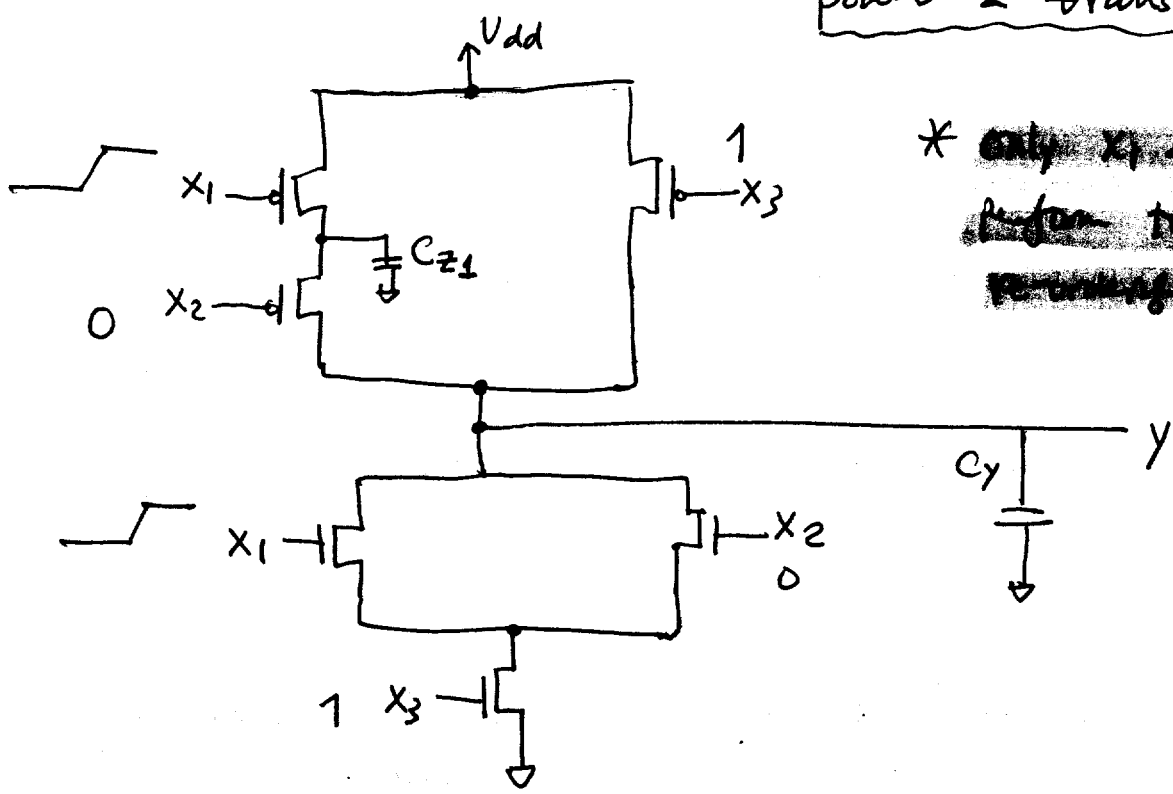
Circuit	minim area power	minim-power power	power saving
C432	306.4 MW	298.3 MW	2.64%
C490	227.12 MW	220.10 MW	3.09%
			0.74%
			1.57%
			6.4%
			4.35%
			5.05%
			9.53%
			4.71%
			10.91%

↳ It is hard to optimize power dissipation in layout level.

Transistor Re-ordering

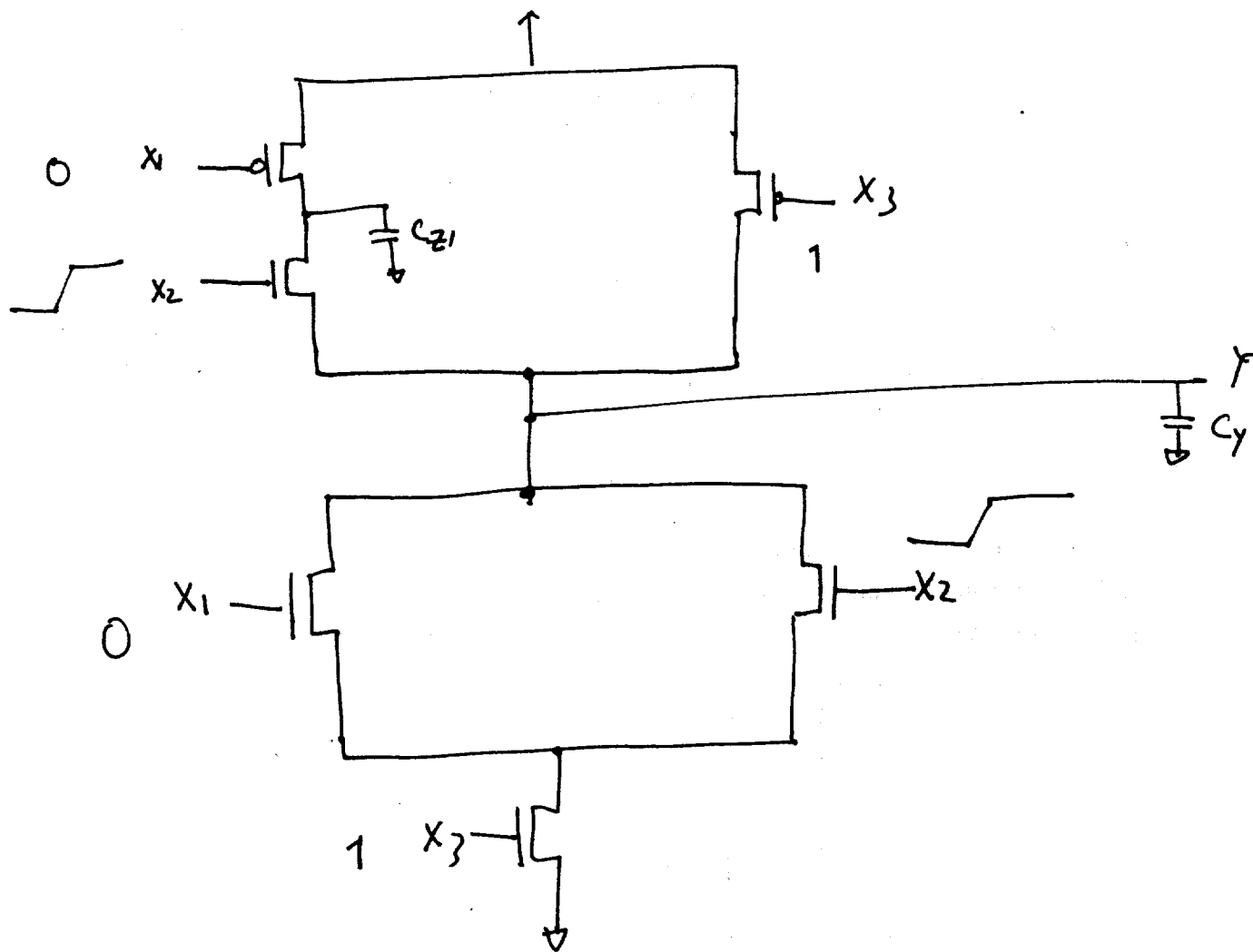
- Find the best ordering of series connected transistors in each gate to minimize the delay or power
- No area penalty when compared with device sizing techniques.

ERROR: undefined
 OFFENDING COMMAND: F3
 STACK:
 /F3_84



* ~~only x_1 & x_2 are~~
~~input transistors~~
~~ordering!~~

- When x_1 , C_y and C_{z1} are both discharged.
- With $x_1=0$ and x_2 , only C_y is discharged.
- If x_2 has a higher switching prob. than x_1 , the transistor driven by x_2 in the series-connected p-transistor chain should be placed closer to the output, to reduce the power dissipation.
- Based on switching activity, we can find the best ordering for power dissipation.

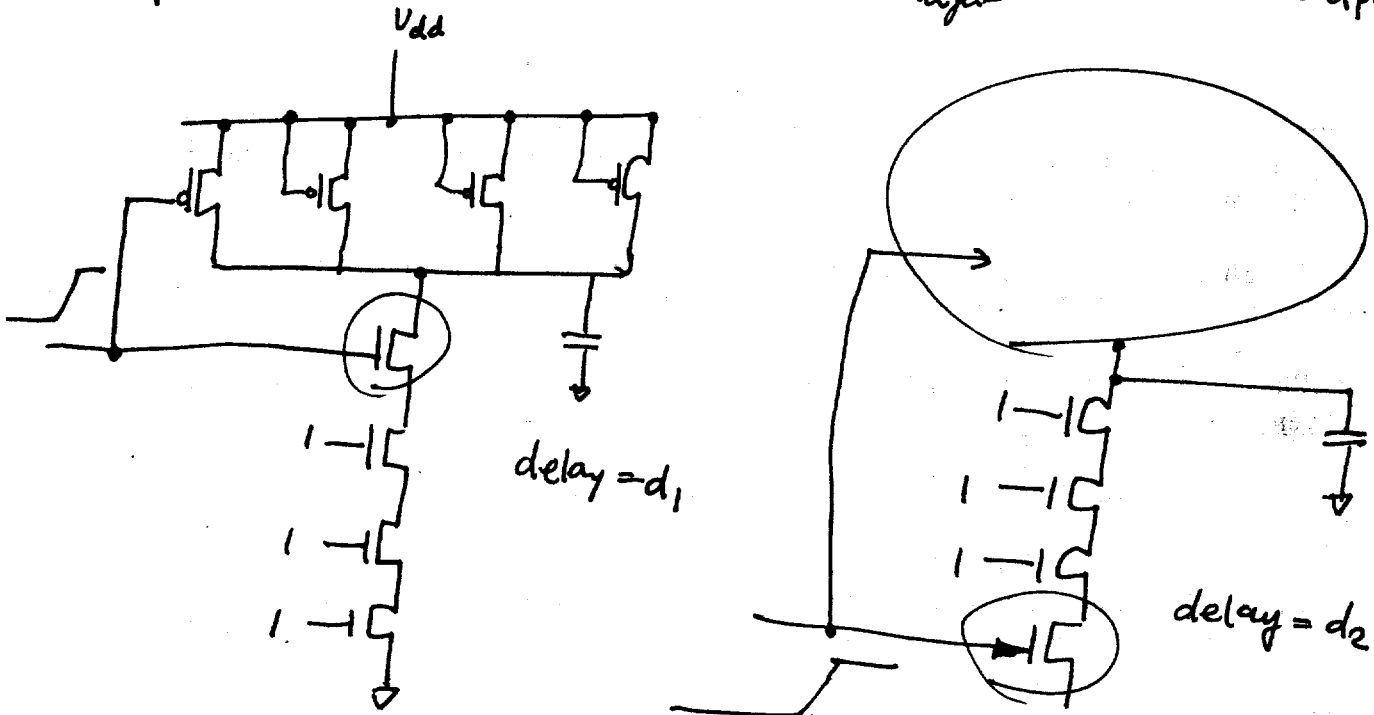


- C_{z1} will not be discharged

Delay and Transistor ordering

- Difficult to determine the order for delay minimization.
- Depends on input arrivals time and transition times.
- The example at page 150 suggests to move ^{high} transition signal close to output.

Exeple:



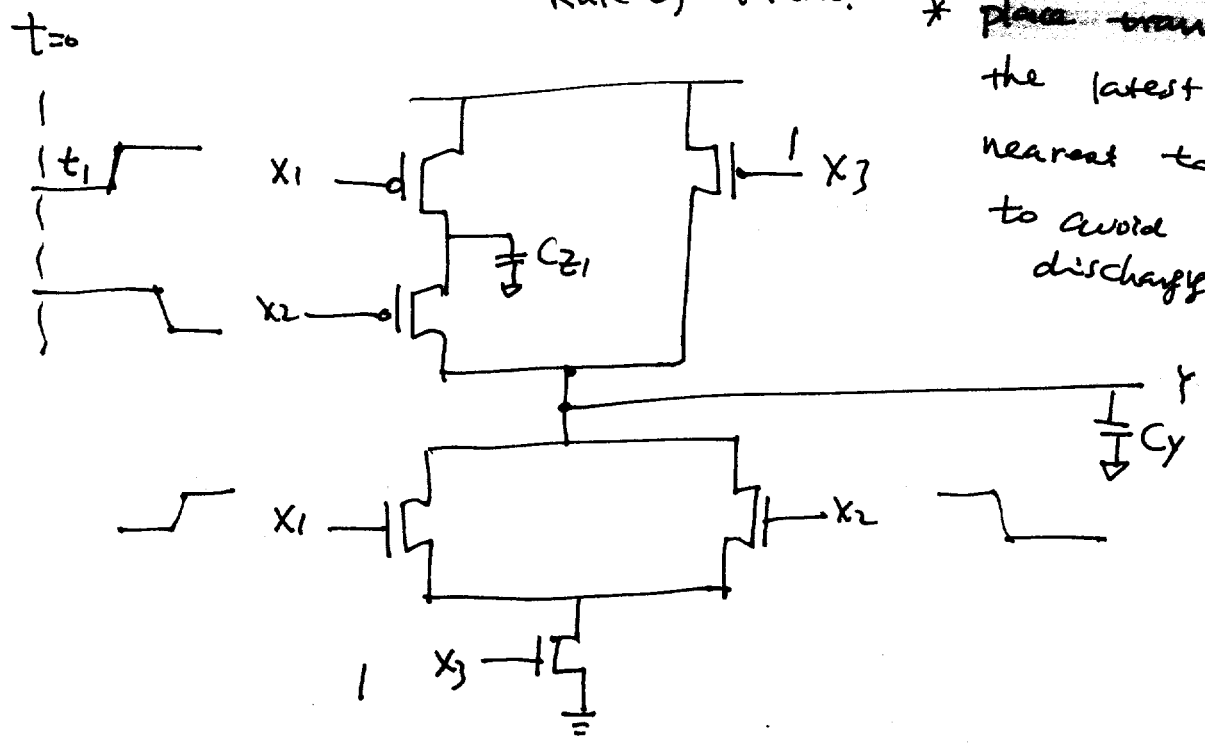
If rise time of input = 5 ns, $d_1/d_2 = 1.23$

If rise time of input = 1 ns, $d_1/d_2 = 0.92$

- spice simulation must be used to determine the optimal ordering for delay minimization.

Rule of thumb.

* place transitions with the latest arrival signal nearest to the output to avoid charging or discharging the internal capacitances to reduce gate delay.



True or False.

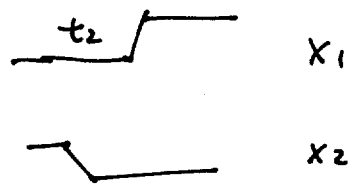
Case I.

X1 charge Cz1
 X2 discharge Cz1 + Cy

$$t(X2 \text{ transition}) > t(X1 \text{ transition})$$

Case II.

$$t(X2 \text{ transition}) < t(X1 \text{ transition})$$



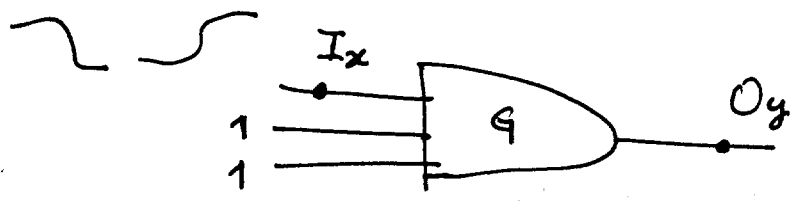
$$\therefore t_2 > t_1$$

\therefore X1 charge Cz1 more than above
 then X2 discharge Cz1 + Cy.

A possible solution

WEN B JONE
ASSOC PROFESSOR
M I : 0030

• Delay characterization for all gates in the library.



It appears that t_f and t_r are not considered.

$T_{xy}^i(G)$: intrinsic delay (without output load)

$R_{xy}(G)$: additional delay per unit fanout load.

→ Total propagation delay through a gate G from I_x to O_y is

$T_{xy}^i(G) + R_{xy}(G) \cdot C_y(G)$

↙ total output load.

- Create characterization for each input
- Sort the delay according to their ordering

