

# **Performance of CMOS Circuits**

Instructed by Shmuel Wimer  
Eng. School, Bar-Ilan University

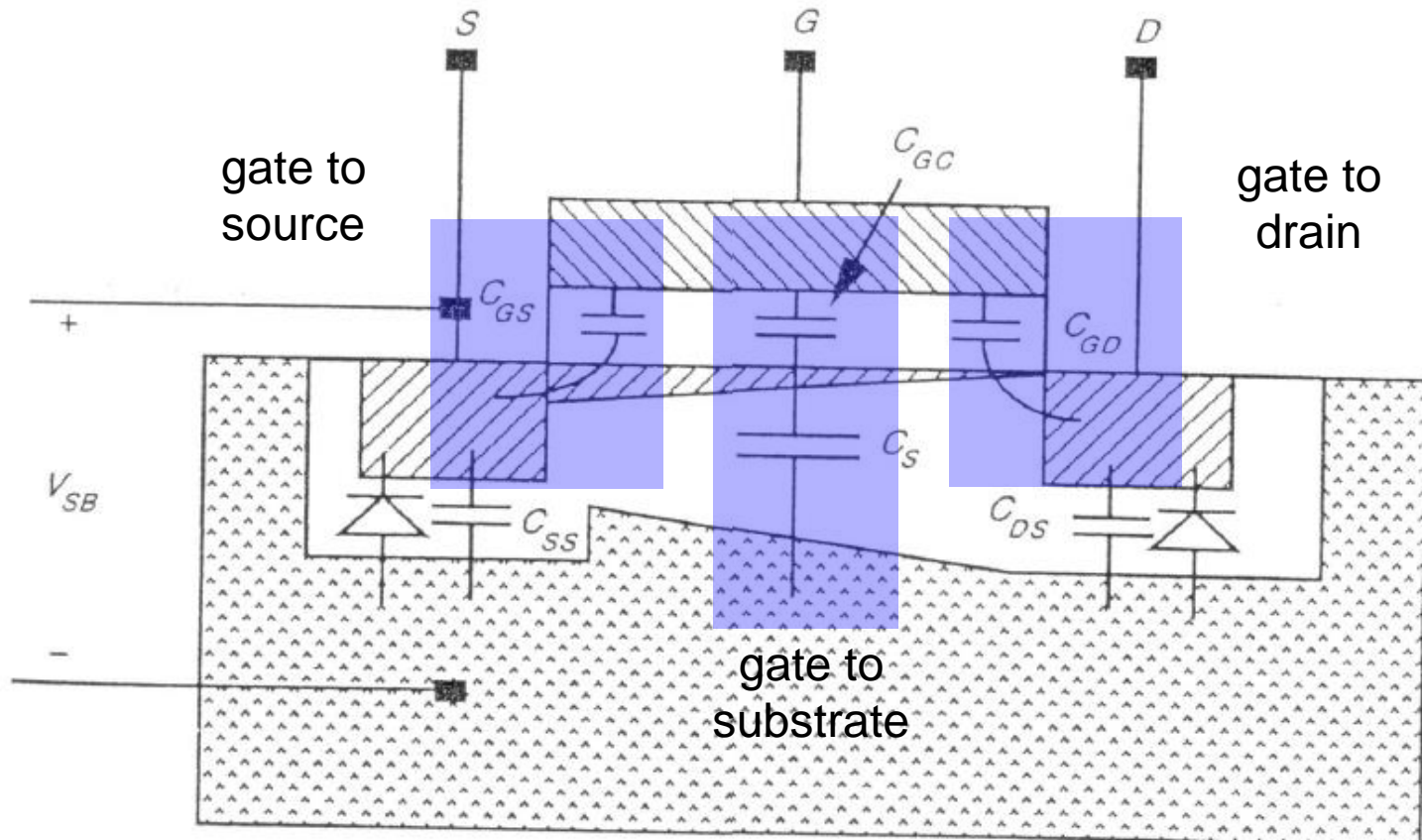
Credits: David Harris  
Harvey Mudd College

(Some material copied/taken/adapted from  
Harris' lecture notes)

# Outline

- ❑ Gate and Diffusion Capacitance
- ❑ RC Delay Models
- ❑ Power and Energy
- ❑ Dynamic Power
- ❑ Static Power
- ❑ Low Power Design

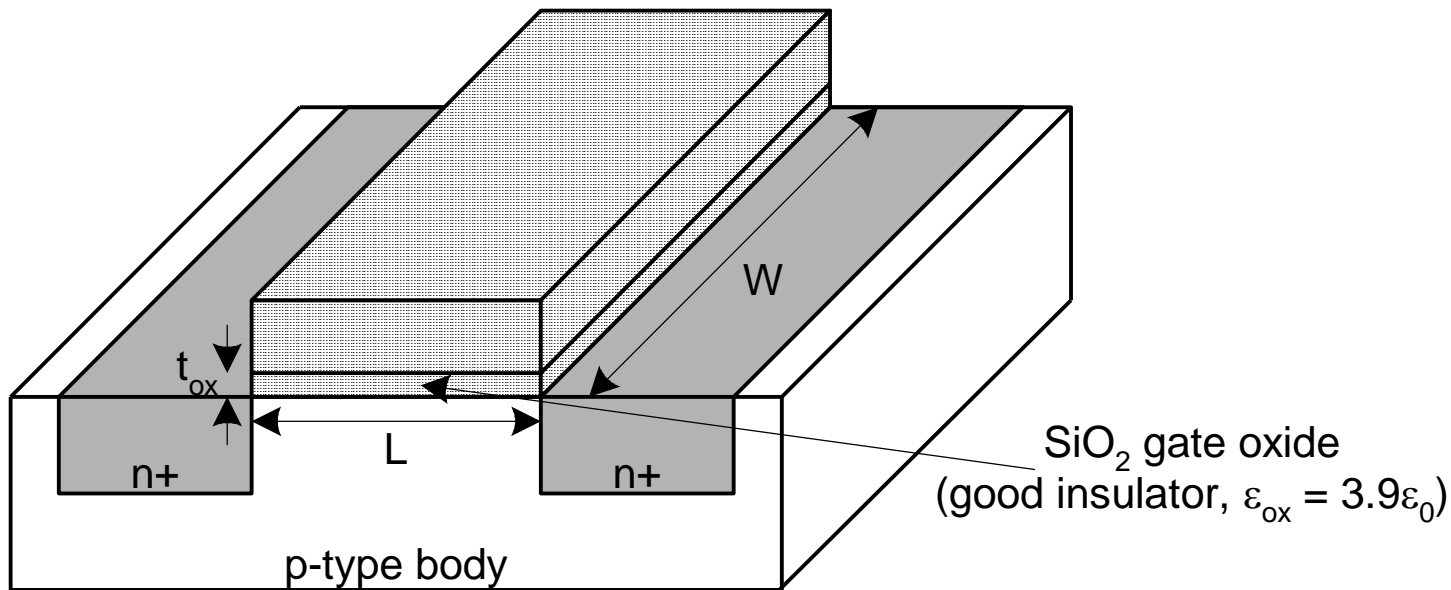
# MOSFET Capacitance

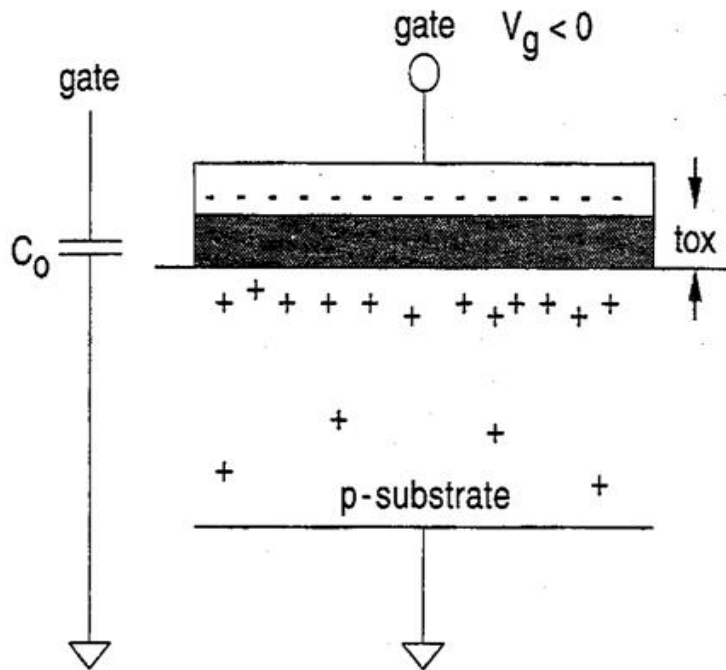


- ❑ Any two conductors separated by an insulator have capacitance
- ❑ Gate to channel capacitor is very important
  - Creates channel charge necessary for operation
- ❑ Source and drain have capacitance to body
  - Across reverse-biased diodes
  - Called diffusion capacitance because it is associated with source/drain diffusion

# Gate Capacitance

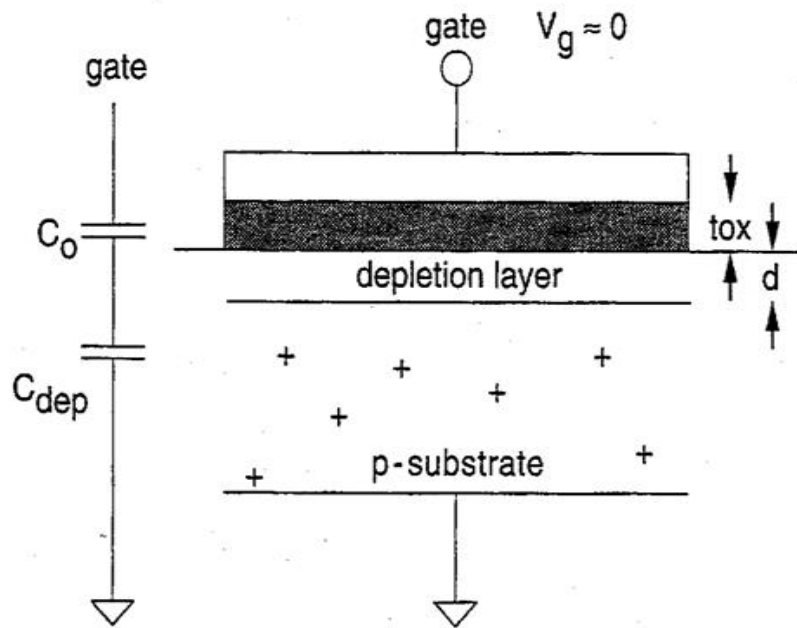
- Approximate channel as connected to source
- $C_{gs} = \epsilon_{ox} WL / t_{ox} = C_{ox} WL = C_{permicron} W$
- $C_{permicron}$  is typically about 2 fF/ $\mu\text{m}$



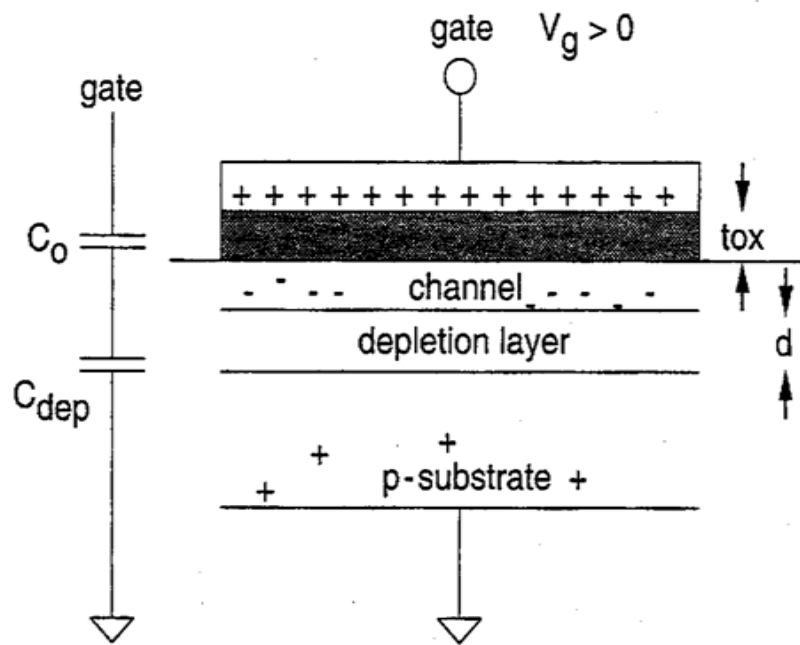


Accumulation occurs when  $V_g$  is negative (for P material). Holes are induced under the oxide.

$C_{gate} = C_{ox}A$  where  $C_{ox} = \frac{\epsilon_{SiO_2}\epsilon_0}{t_{ox}}$

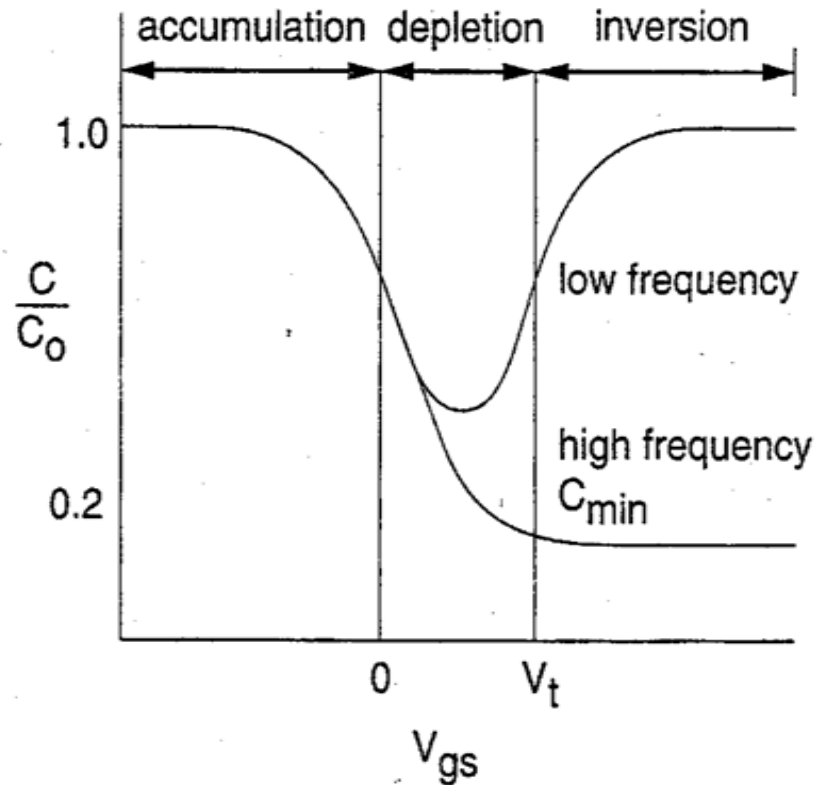


Depletion occurs when  $V_g$  is near zero but  $< V_{tn}$ . Here the  $C_{gate}$  is given by  $C_{ox}A$  in series with depletion layer capacitance  $C_{dep}$



Inversion occurs when  $V_g$  is positive and  $> V_{tn}$  (for P material). A model for inversion is comprised of  $C_{ox}$  A connecting from gate-to-channel and  $C_{dep}$  connecting from channel-to-substrate.



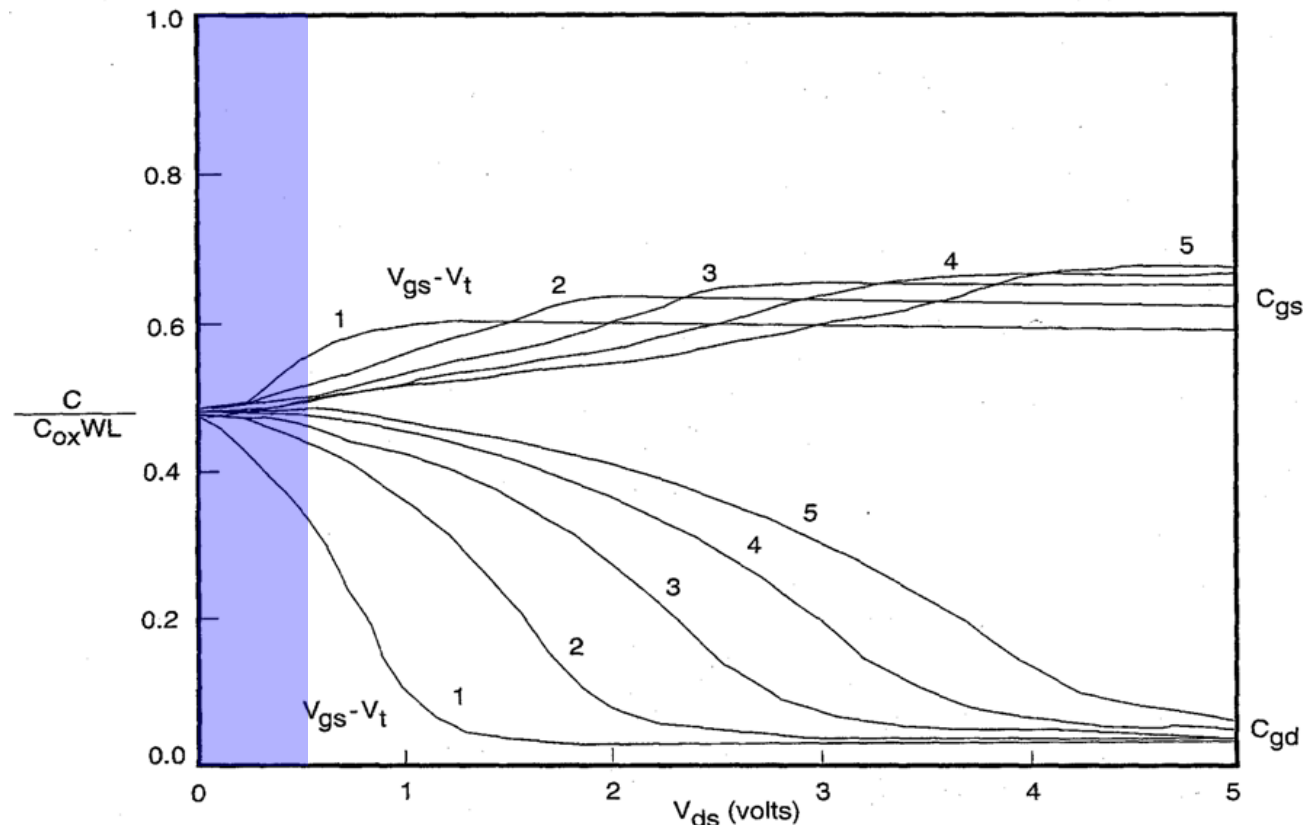


Normalized gate capacitance versus Gate voltage  $V_{gs}$ .

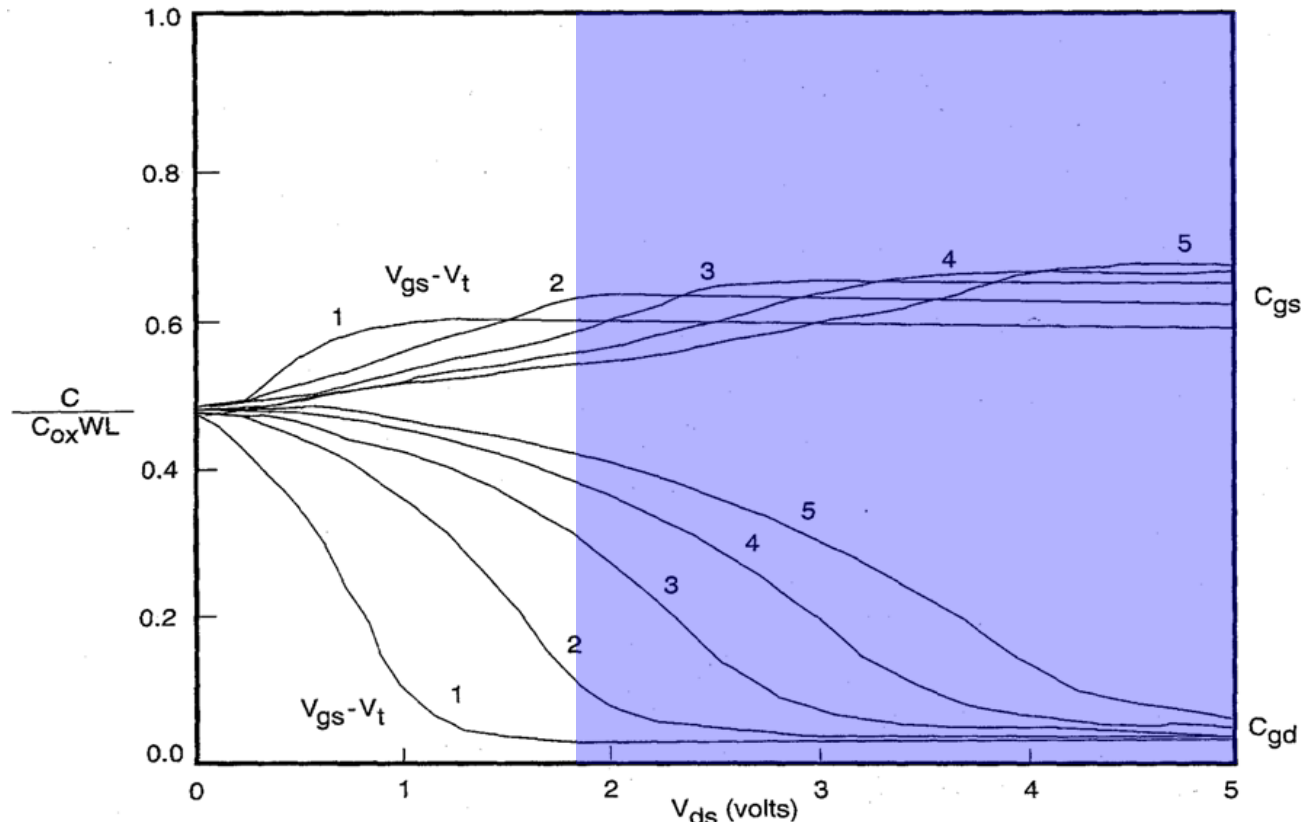
High freq behavior is due to the distributed resistance of channel

# Normalized Experimental MOS Gate Capacitance Measurements vs $V_{ds}$ , $V_{gs}$

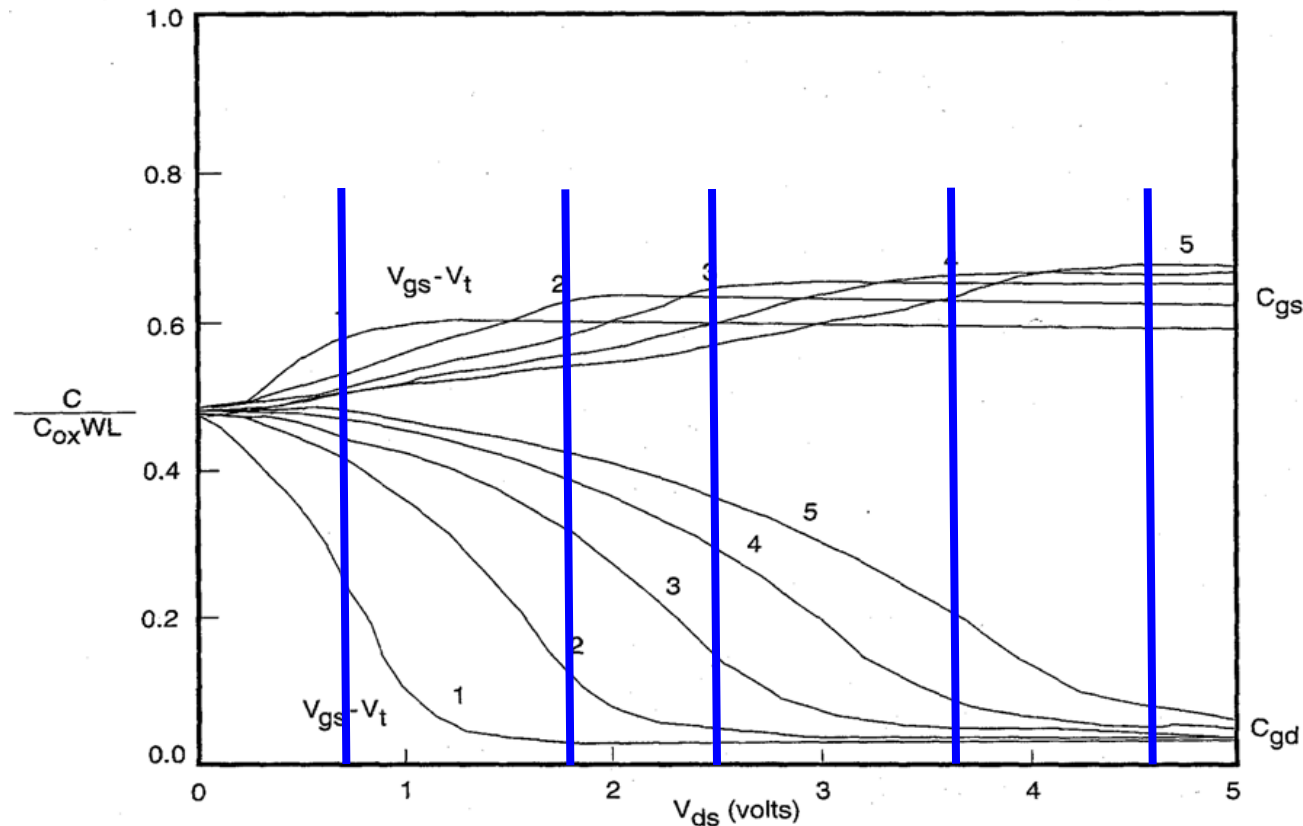
For  $V_{ds} = 0$ , the total gate capacitance  $C_{ox}$  splits equally to the drain and source of the transistor.



For  $V_{ds} > 0$ , the gate capacitance tilts more toward the source and becomes roughly  $2/3 C_{ox}A$  to the source and 0 to the drain for high  $V_{ds}$ .

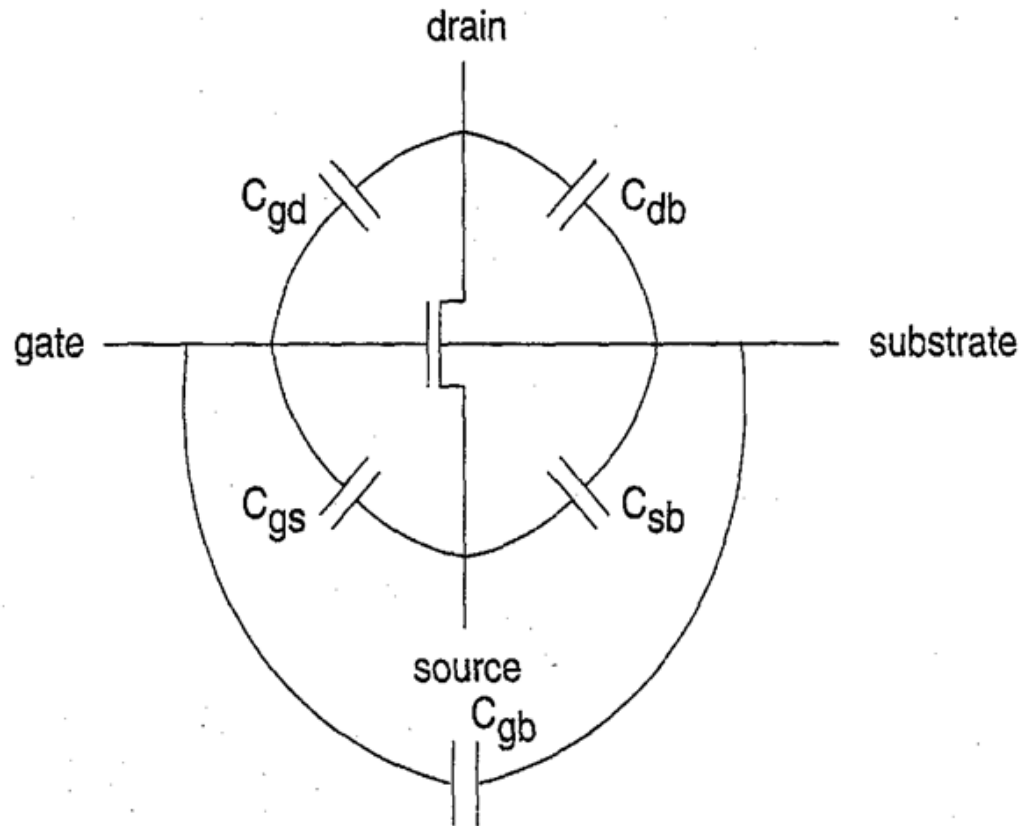


Higher  $V_{gs} - V_t$  forces this tilting to occur later, since the device is linear up to  $V_{gs} - V_t = V_{ds}$ .



# MOS Transistor Gate Capacitance Model

Gate capacitance has different components in different modes, but total remains constant.



Gate capacitance has different components in different modes, but total remains constant.

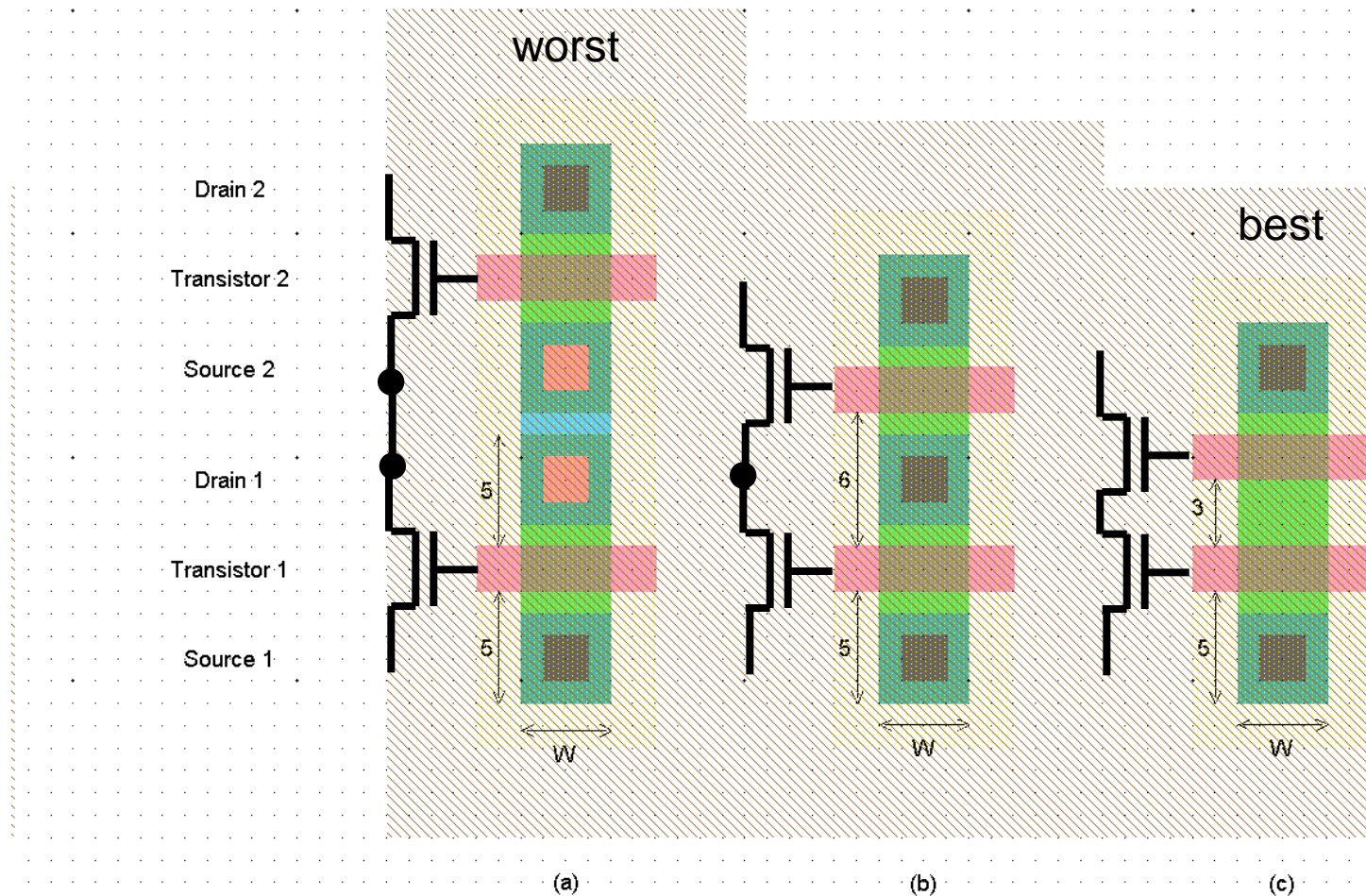
**TABLE 4.3 Approximation of intrinsic MOS gate capacitance**

Parameter	CAPACITANCE		
	Off	Non-saturated	Saturated
$C_{gb}$	$\frac{\epsilon A}{t_{ox}}$	0	0
$C_{gs}$	0	$\frac{\epsilon A}{2t_{ox}}$	$\frac{2\epsilon A}{3t_{ox}}$
$C_{gd}$	0	$\frac{\epsilon A}{2t_{ox}}$	0 (finite for short channel devices)
$C_g = C_{gb} + C_{gs} + C_{gd}$	$\frac{\epsilon A}{t_{ox}}$	$\frac{\epsilon A}{t_{ox}}$	$\frac{2\epsilon A}{3t_{ox}} \rightarrow \frac{.9\epsilon A}{t_{ox}}$ (short channel)

# Diffusion Capacitance

- ❑  $C_{sb}$ ,  $C_{db}$
- ❑ Undesirable, called *parasitic* capacitance
- ❑ Capacitance depends on area and perimeter
  - Use small diffusion nodes
  - Comparable to  $C_g$   
for contacted diff
  - $\frac{1}{2} C_g$  for uncontacted
  - Varies with process

# Diffusion Capacitance (Cont'd)



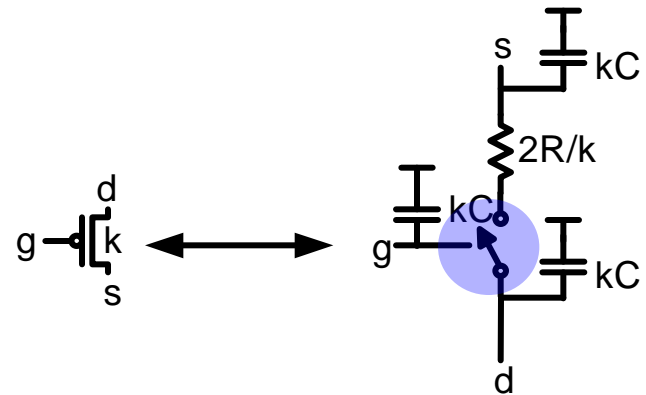
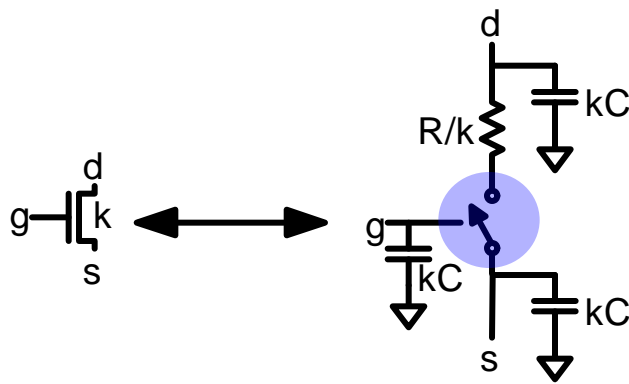


# Effective Resistance

- ❑ Shockley models have limited value
  - Not accurate enough for modern transistors
  - Too complicated for much hand analysis
- ❑ Simplification: treat transistor as resistor
  - Replace  $I_{ds}(V_{ds}, V_{gs})$  with effective resistance  $R$ 
    - $I_{ds} = V_{ds}/R$
  - $R$  averaged across switching of digital gate
- ❑ Too inaccurate to predict current at any given time
  - But good enough to predict RC delay

# RC Delay Model

- Use equivalent circuits for MOS transistors
  - Ideal switch + capacitance and ON resistance
  - Unit nMOS has resistance  $R$ , capacitance  $C$
  - Unit pMOS has resistance  $2R$ , capacitance  $C$
- Capacitance proportional to width
- Resistance inversely proportional to width

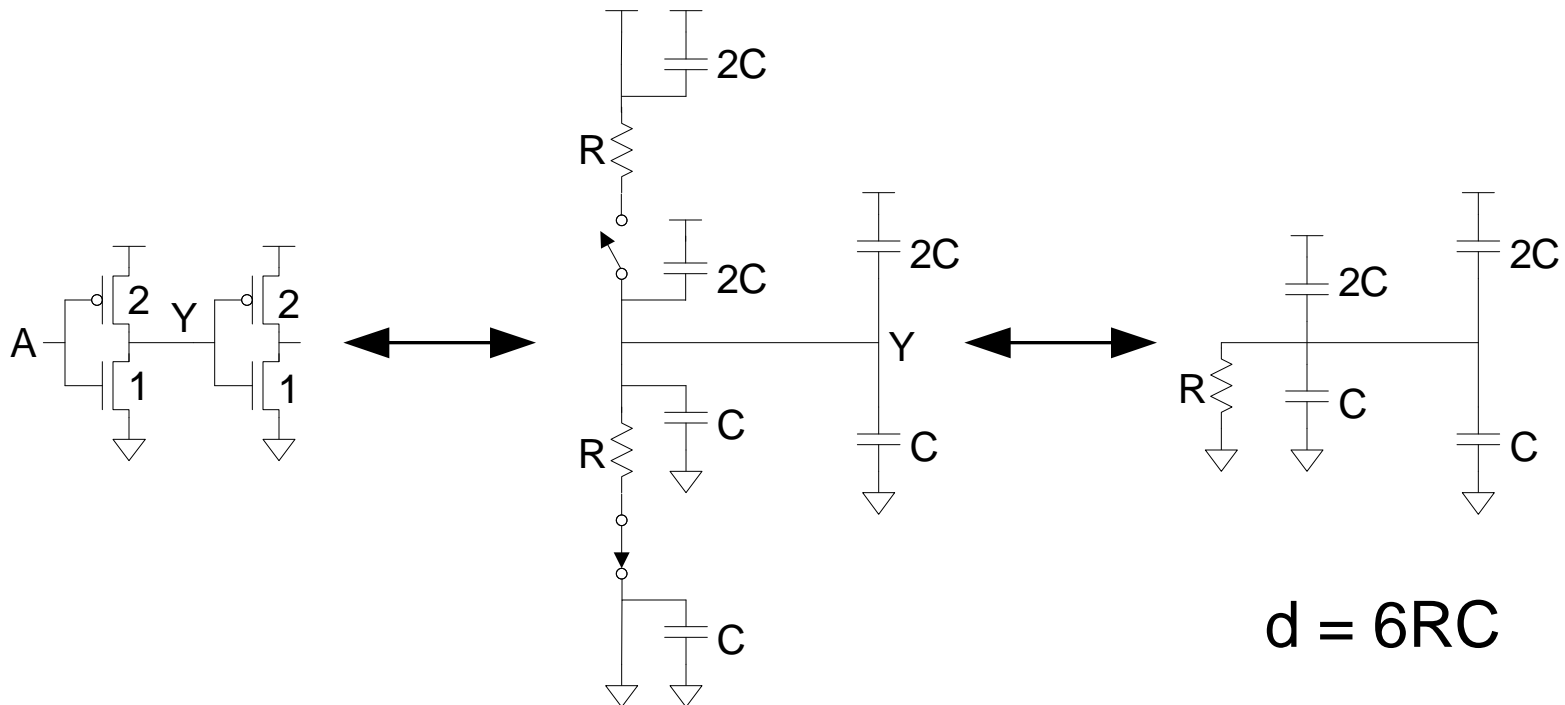


# RC Values

- ❑ Capacitance
  - $C = C_g = C_s = C_d = 2 \text{ fF}/\mu\text{m}$  of gate width
  - Values similar across many processes
- ❑ Resistance
  - $R \approx 6 \text{ K}\Omega \cdot \mu\text{m}$  in  $0.6\mu\text{m}$  process
  - Improves with shorter channel lengths
- ❑ Unit transistors
  - May refer to minimum contacted device ( $4/2 \lambda$ )
  - Or maybe  $1 \mu\text{m}$  wide device
  - Doesn't matter as long as you are consistent

# Inverter Delay Estimate

- Estimate the delay of a fanout-of-1 inverter



# Transient Response

- ❑ *DC analysis* tells us  $V_{\text{out}}$  if  $V_{\text{in}}$  is constant
- ❑ *Transient analysis* tells us  $V_{\text{out}}(t)$  if  $V_{\text{in}}(t)$  changes
  - Requires solving differential equations
- ❑ Input is usually considered to be a step or ramp
  - From 0 to  $V_{\text{DD}}$  or vice versa

# Inverter Step Response

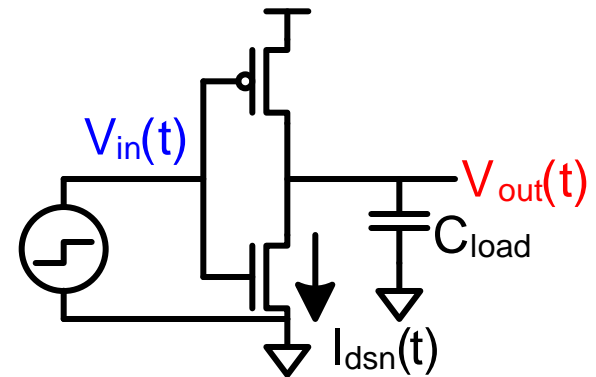
□ Ex: find step response of inverter driving load cap

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

$$V_{out}(t < t_0) = V_{DD}$$

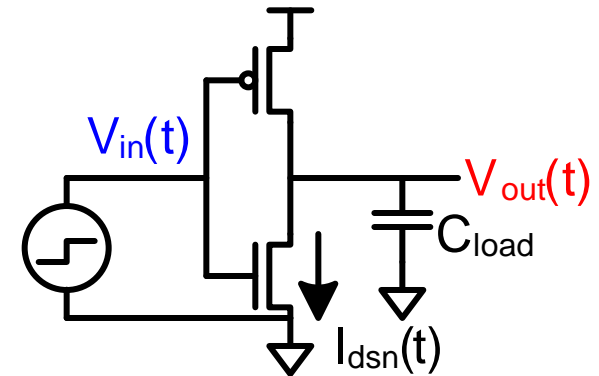
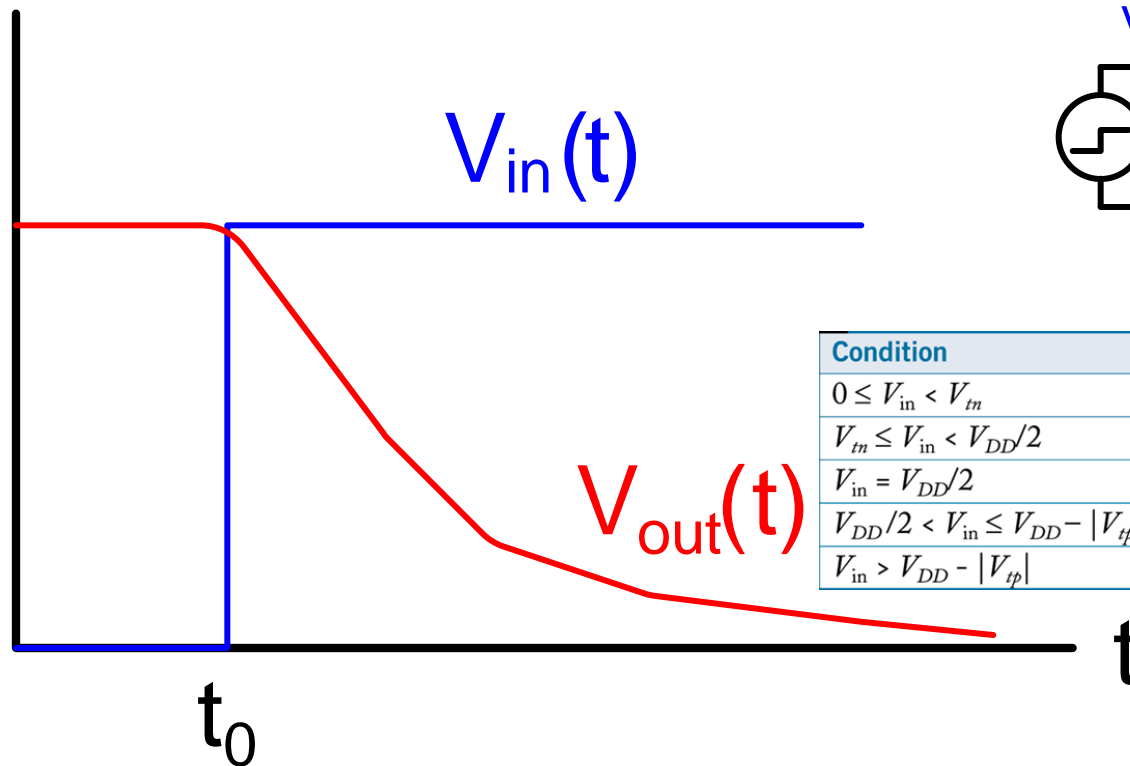
$$\frac{dV_{out}(t)}{dt} = -\frac{I_{dsn}(t)}{C_{load}}$$

$$I_{dsn}(t) = \begin{cases} 0 & t \leq t_0 \\ \frac{\beta}{2}(V_{DD} - V_t)^2 & V_{out} > V_{DD} - V_t \\ \beta\left(V_{DD} - V_t - \frac{V_{out}(t)}{2}\right)V_{out}(t) & V_{out} < V_{DD} - V_t \end{cases}$$



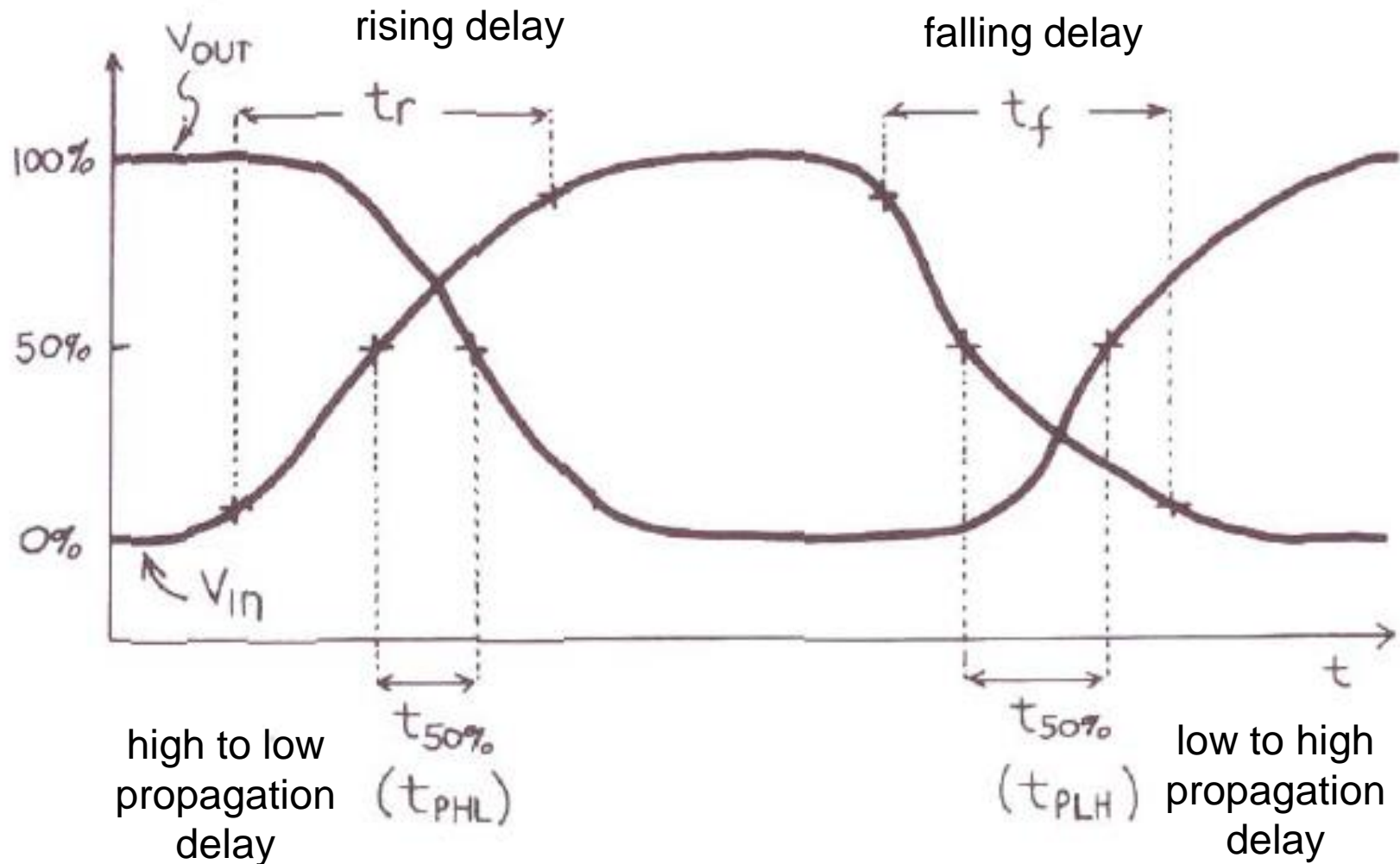
# Inverter Step Response

- Ex: find step response of inverter driving load cap

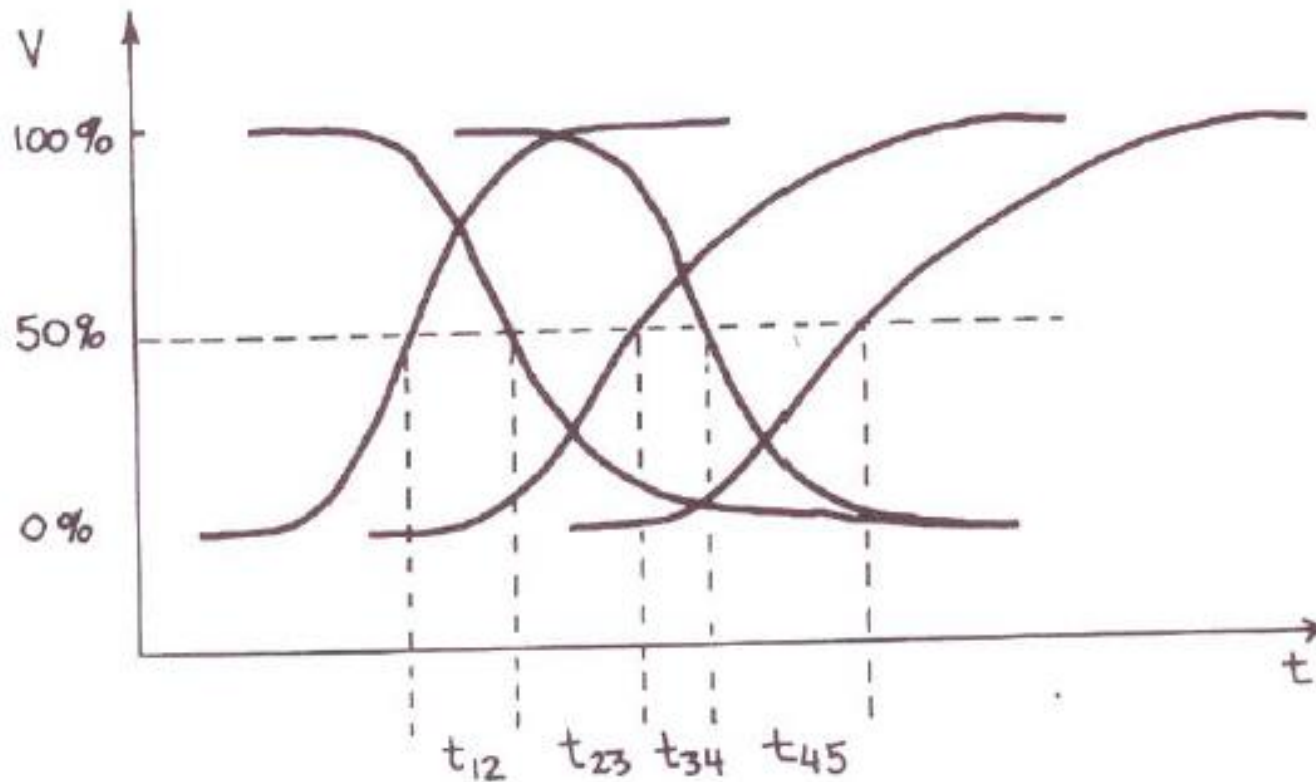
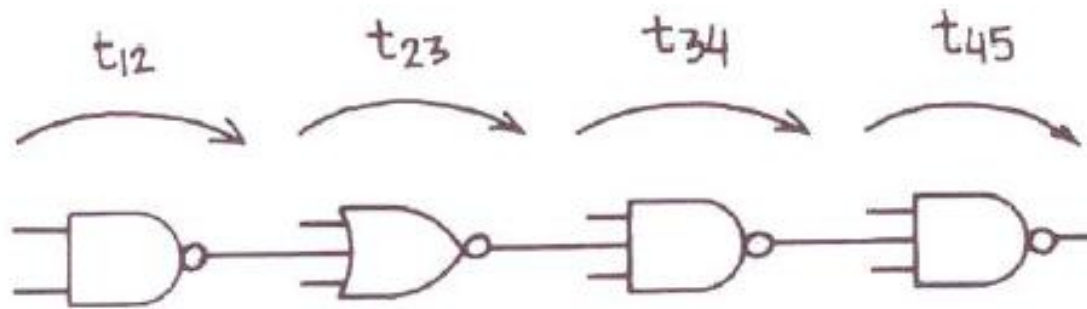


Condition	p-device	n-device	Output
$0 \leq V_{in} < V_{tn}$	linear	cutoff	$V_{out} = V_{DD}$
$V_{tn} \leq V_{in} < V_{DD}/2$	linear	saturated	$V_{out} > V_{DD}/2$
$V_{in} = V_{DD}/2$	saturated	saturated	$V_{out}$ drops sharply
$V_{DD}/2 < V_{in} \leq V_{DD} -  V_{tp} $	saturated	linear	$V_{out} < V_{DD}/2$
$V_{in} > V_{DD} -  V_{tp} $	cutoff	linear	$V_{out} = 0$

# Delay Definitions







# Delay Definitions (Cont'd)

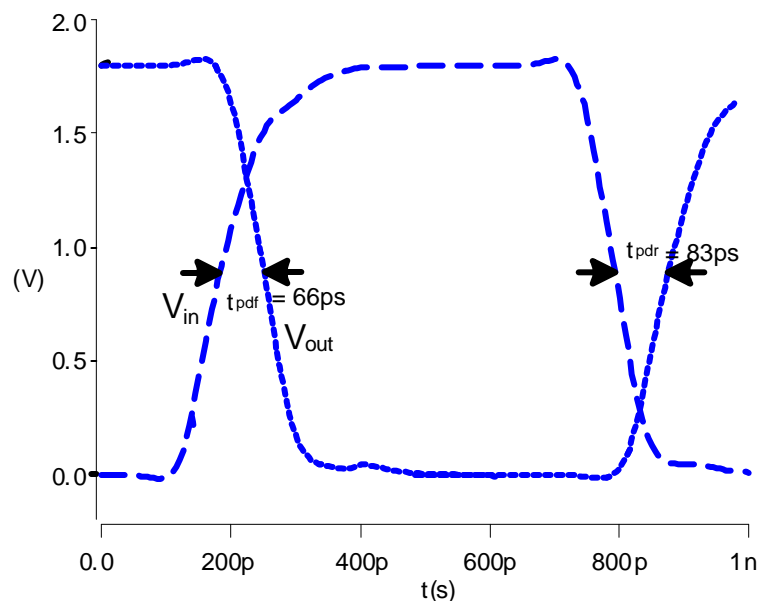
- ❑  $t_{pdr}$ : *rising propagation delay*
  - **Maximum** time from input crossing 50% to rising output crossing 50%
- ❑  $t_{pdf}$ : *falling propagation delay*
  - **Maximum** time from input crossing 50% to falling output crossing 50%
- ❑  $t_{pd}$ : *average propagation delay*
  - $t_{pd} = (t_{pdr} + t_{pdf})/2$
- ❑  $t_r$ : *rise time*
  - From output crossing  $0.2 V_{DD}$  to  $0.8 V_{DD}$

# Delay Definitions (Cont'd)

- ❑  $t_f$ : *fall time*
  - From output crossing  $0.8 V_{DD}$  to  $0.2 V_{DD}$
- ❑  $t_{cdr}$ : *rising contamination delay*
  - **Minimum** time from input crossing 50% to rising output crossing 50%
- ❑  $t_{cdf}$ : *falling contamination delay*
  - **Minimum** time from input crossing 50% to falling output crossing 50%
- ❑  $t_{cd}$ : *average contamination delay*
  - $t_{cd} = (t_{cdr} + t_{cdf})/2$

# Simulated Inverter Delay

- ❑ Solving differential equations by hand is too hard
- ❑ SPICE simulator solves the equations numerically
  - Uses more accurate I-V models too!
- ❑ But simulations take time to write

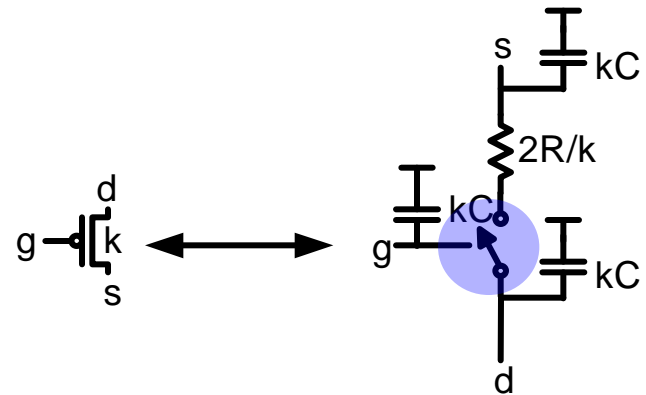
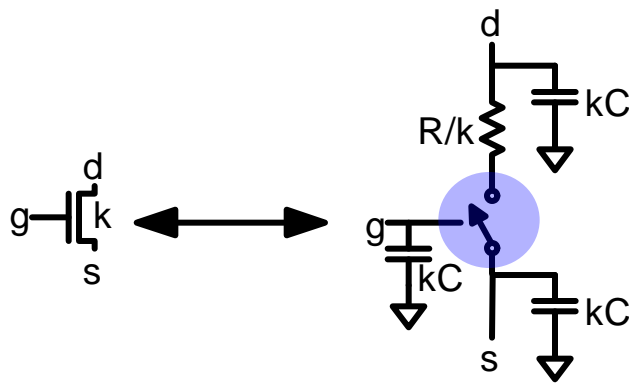


# Delay Estimation

- ❑ We would like to be able to easily estimate delay
  - Not as accurate as simulation
  - But easier to ask “What if?”
- ❑ The step response usually looks like a 1<sup>st</sup> order RC response with a decaying exponential.
- ❑ Use RC delay models to estimate delay
  - $C$  = total capacitance on output node
  - Use *effective resistance*  $R$
  - So that  $t_{pd} = RC$
- ❑ Characterize transistors by finding their effective  $R$ 
  - Depends on average current as gate switches

# RC Delay Model

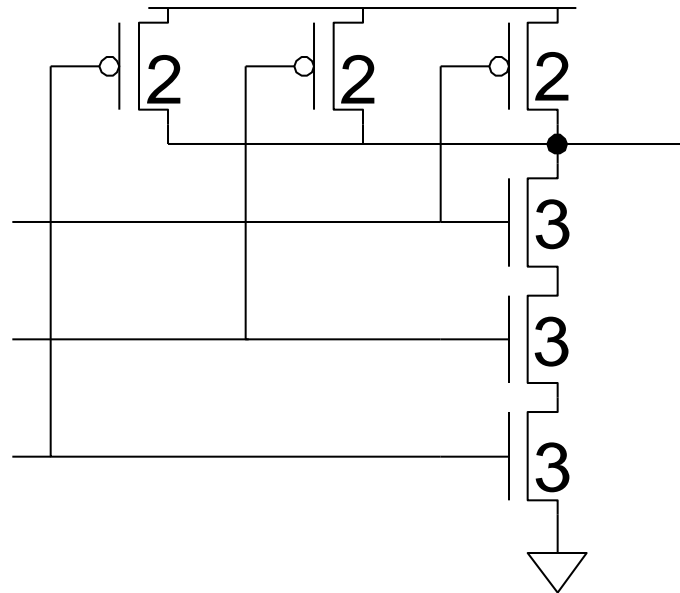
- Use equivalent circuits for MOS transistors
  - Ideal switch + capacitance and ON resistance
  - Unit nMOS has resistance  $R$ , capacitance  $C$
  - Unit pMOS has resistance  $2R$ , capacitance  $C$
- Capacitance proportional to width
- Resistance inversely proportional to width



# Example: 3-input NAND

- Sketch a 3-input NAND with transistor widths chosen to achieve effective rise and fall resistances equal to a unit inverter ( $R$ ).

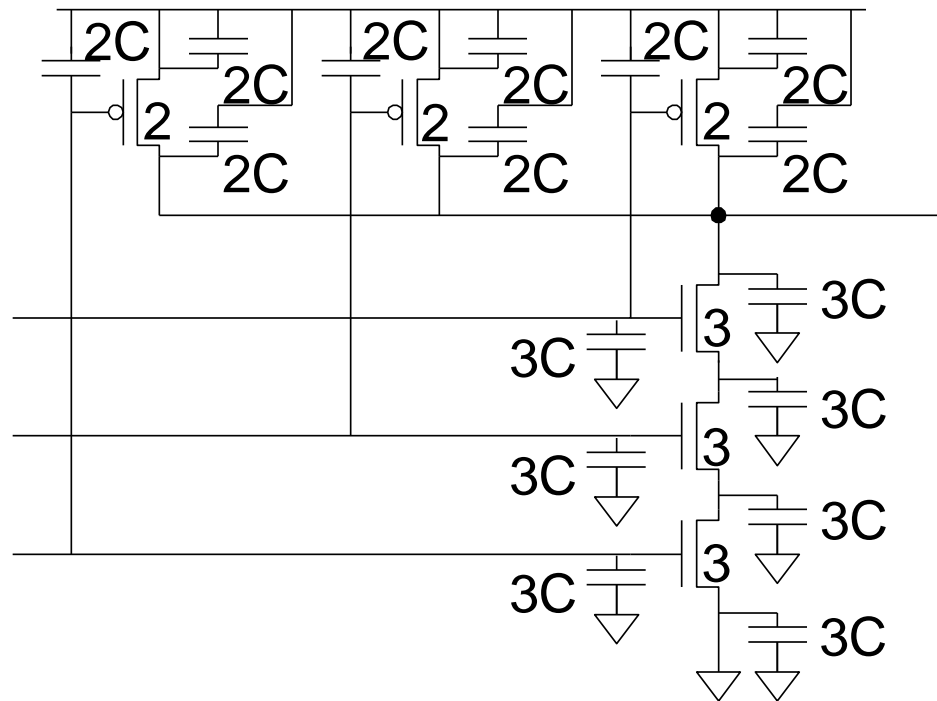
In worst case of P  
only one device is  
opened.



all N devices must  
be opened.

# 3-input NAND Caps

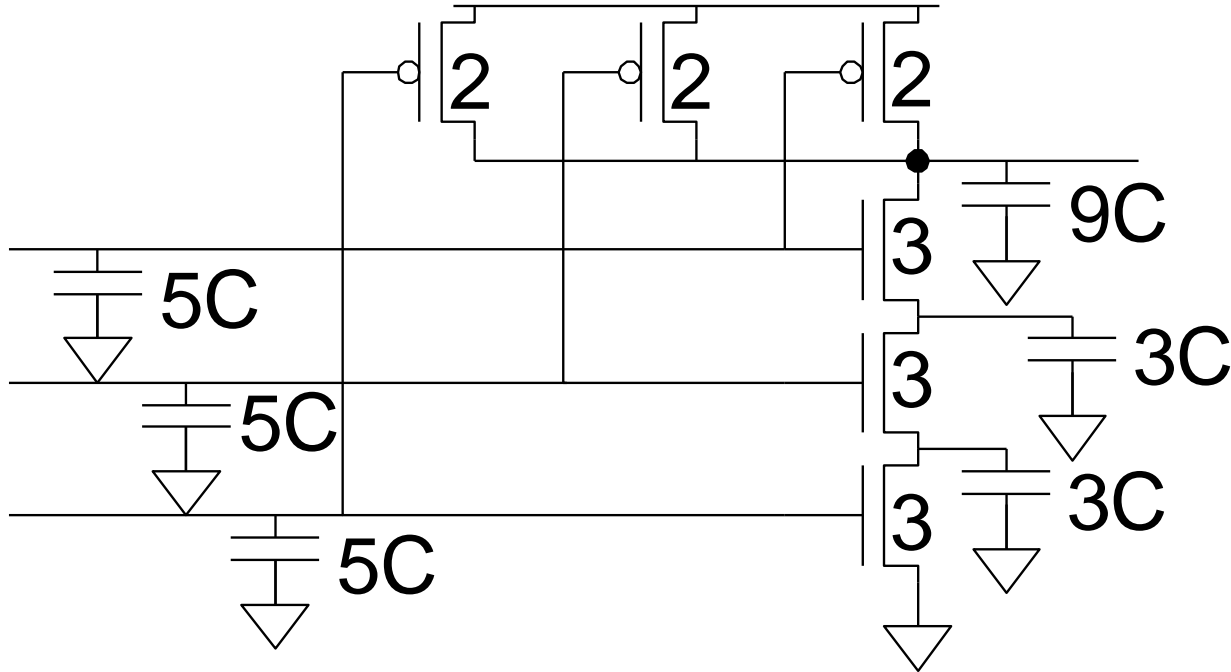
- Annotate the 3-input NAND gate with gate and diffusion capacitance.





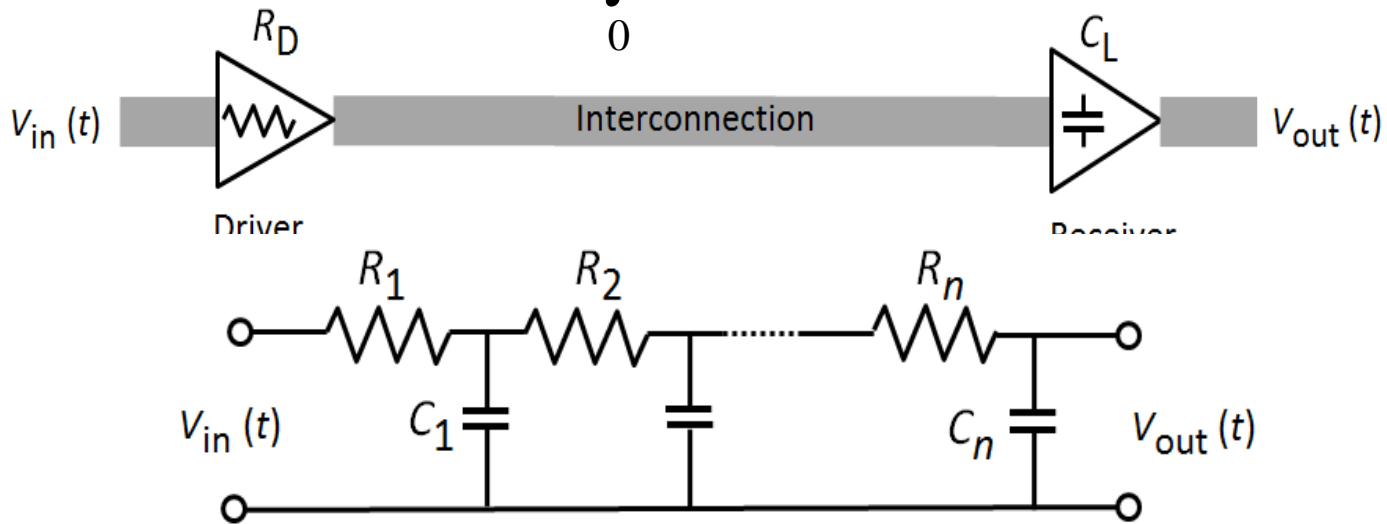
# 3-input NAND Caps (Cont'd)

- Annotate the 3-input NAND gate with gate and diffusion capacitance.



# Elmore Delay

$$\delta \triangleq \int_0^{\infty} t V_{\text{out}}(t) dt$$



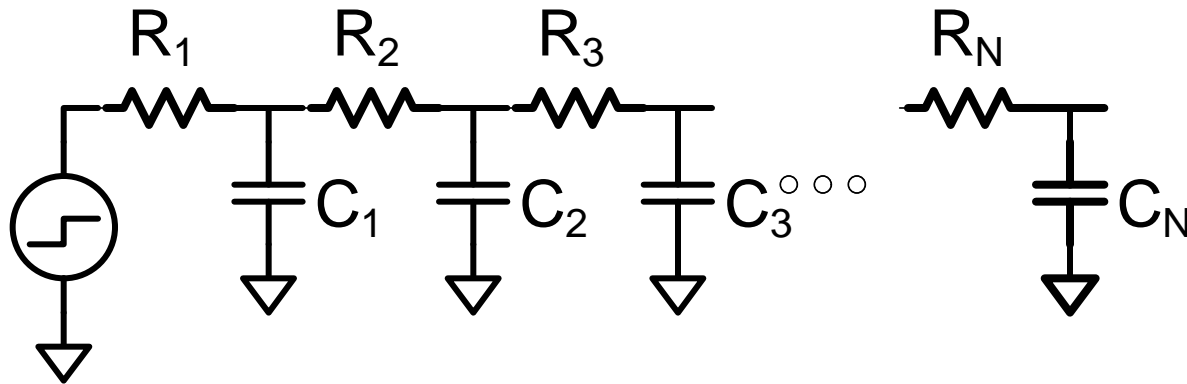
$$\delta \triangleq \int_0^{\infty} t V_{\text{out}}(t) dt \cong \sum_{i=1}^n R_i \sum_{j=i}^n C_j = \sum_{j=1}^n C_j \sum_{i=1}^j R_i$$

# Elmore Delay

- ❑ ON transistors look like resistors
- ❑ Pullup or pulldown network modeled as *RC ladder*
- ❑ Elmore delay of RC ladder

$$t_{pd} \approx \sum_{\text{nodes } i} R_{i\text{-to-source}} C_i$$

$$= R_1 C_1 + (R_1 + R_2) C_2 + \dots + (R_1 + R_2 + \dots + R_N) C_N$$



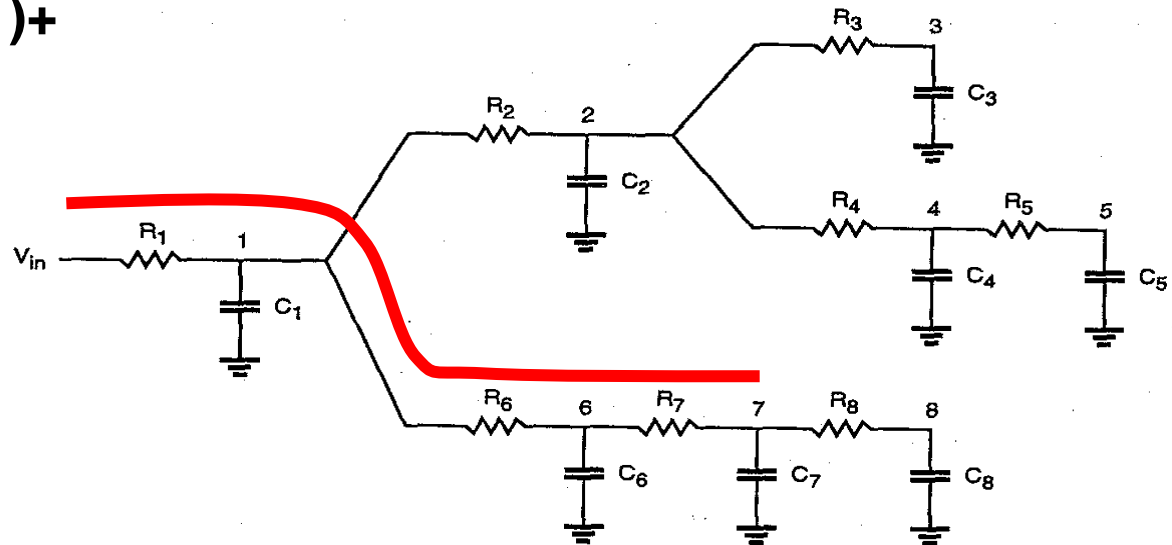
For a step input  $V_{in}$ , the delay at any node can be estimated with the Elmore delay equation  $t_{Di} = \sum C_j \sum R_k$

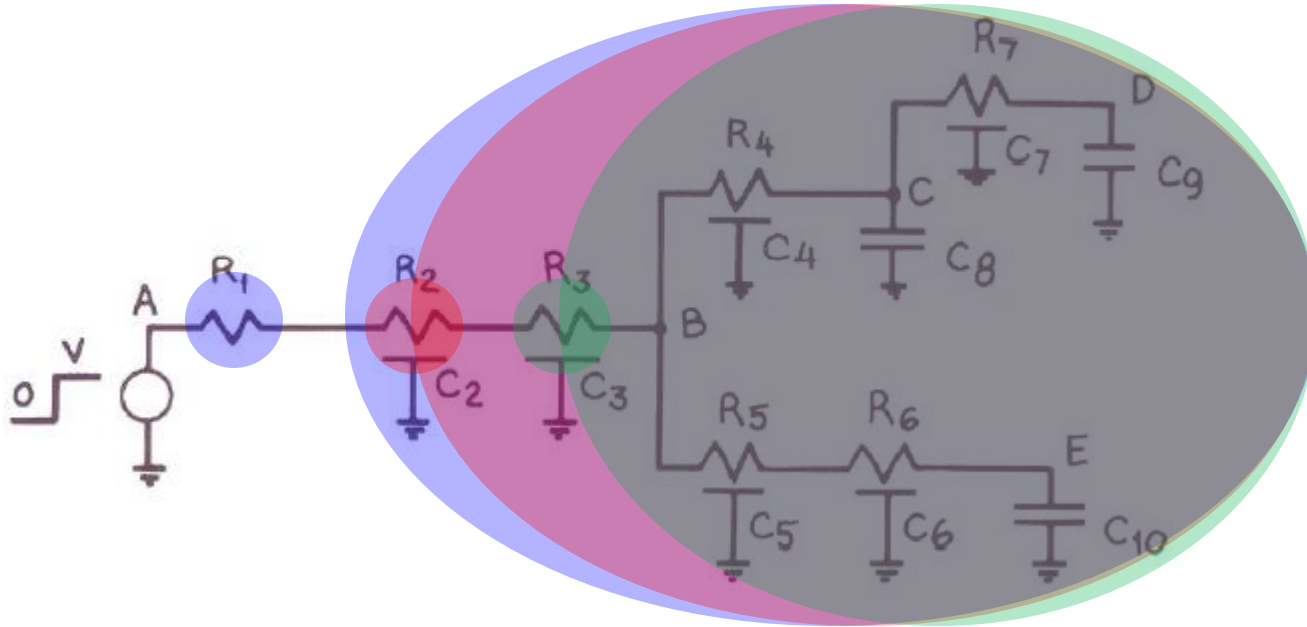
For example, the Elmore delay at node 7 is give by:

$$R1 ( C1 + C2 + C3 + C4 + C5+ C6+ C7 + C8 ) +$$

$$R6 ( C6+ C7+ C8 )+$$

$$R7 ( C7 + C8)$$



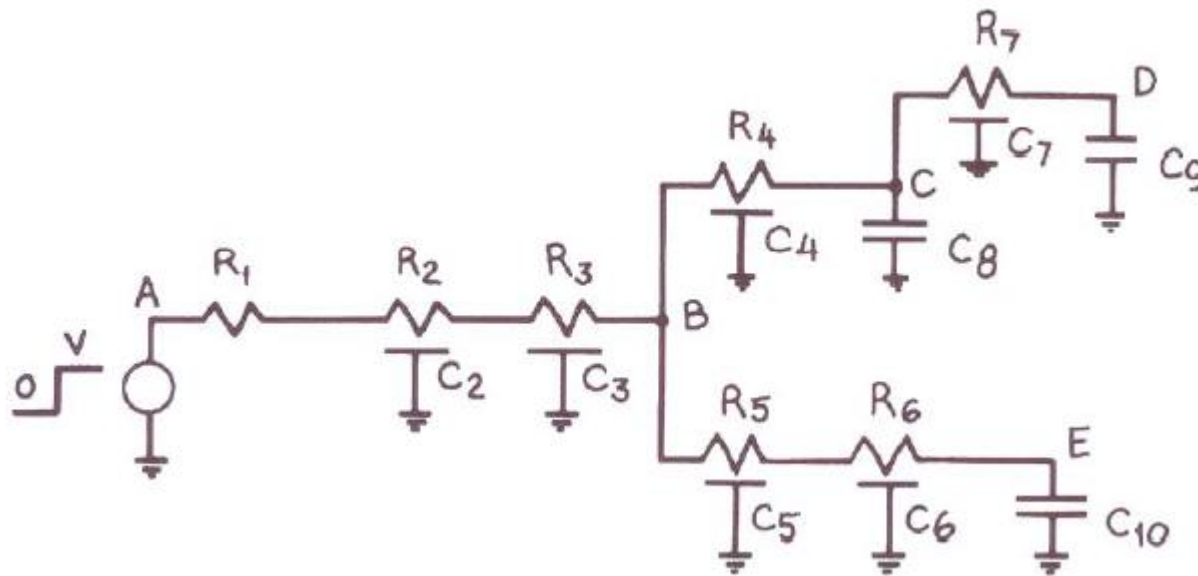


The 50 percent delay from  $A$  to  $B$  is

$$T_{AB} = R_1(C_2 + C_3 + C_4 + C_5 + C_6 + C_7 + C_8 + C_9 + C_{10})$$

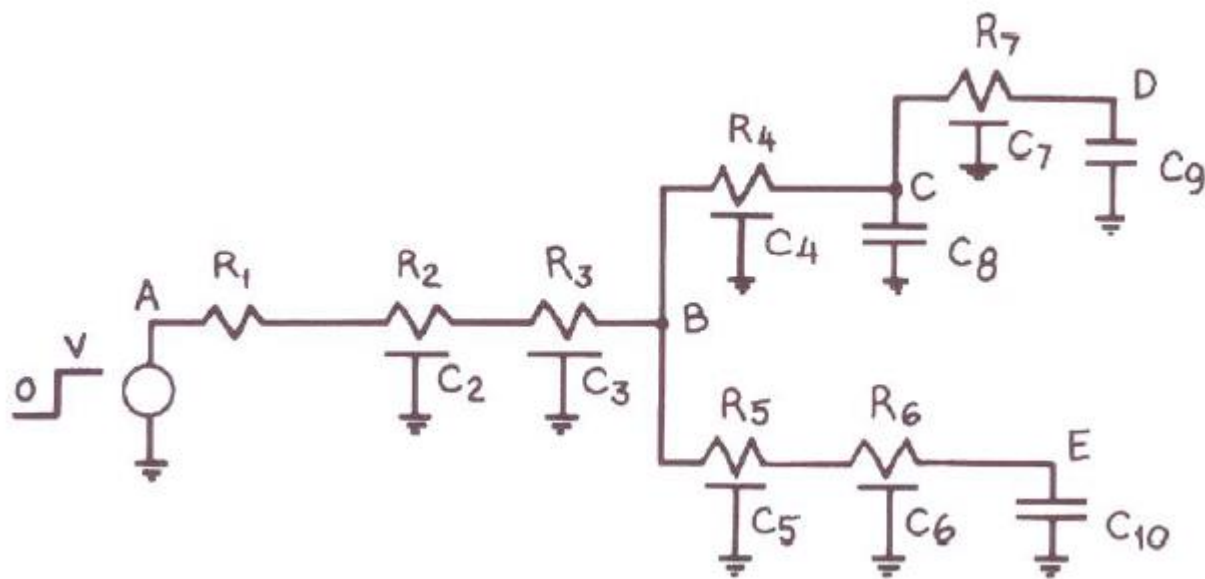
$$+ R_2 \left( \frac{C_2}{2} + C_3 + C_4 + C_5 + C_6 + C_7 + C_8 + C_9 + C_{10} \right)$$

$$+ R_3 \left( \frac{C_3}{2} + C_4 + C_5 + C_6 + C_7 + C_8 + C_9 + C_{10} \right).$$



The 50 percent delay from  $B$  to  $D$  is

$$T_{BD} = R_4 \left( \frac{C_4}{2} + C_7 + C_8 + C_9 \right) + R_7 \left( \frac{C_7}{2} + C_9 \right).$$



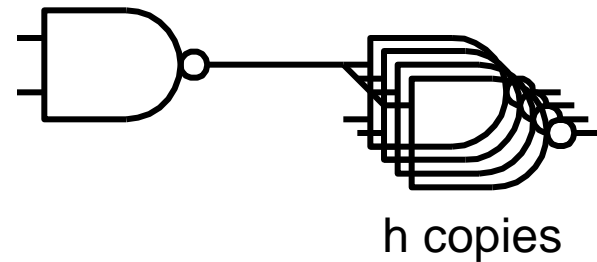
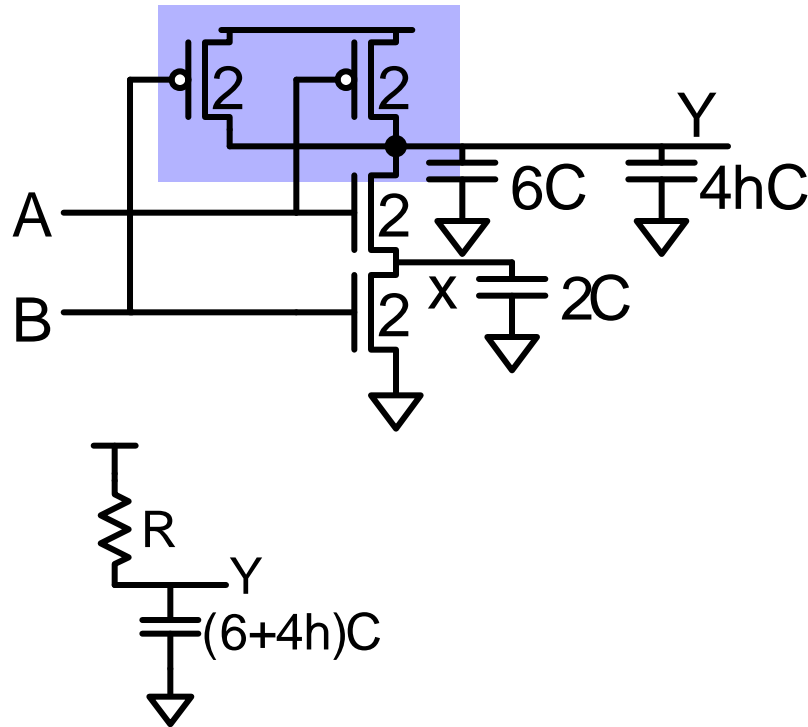
The 50 percent delay from  $B$  to  $E$  is

$$T_{BE} = R_5 \left( \frac{C_5}{2} + C_6 + C_{10} \right) + R_6 \left( \frac{C_6}{2} + C_{10} \right).$$



# Example: 2-input NAND

- Estimate **rising** and falling propagation delays of a 2-input NAND driving  $h$  identical gates.

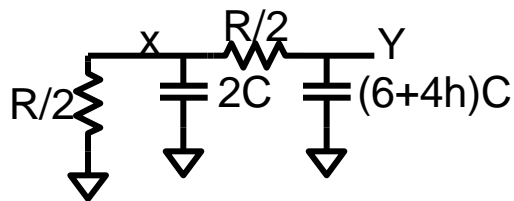
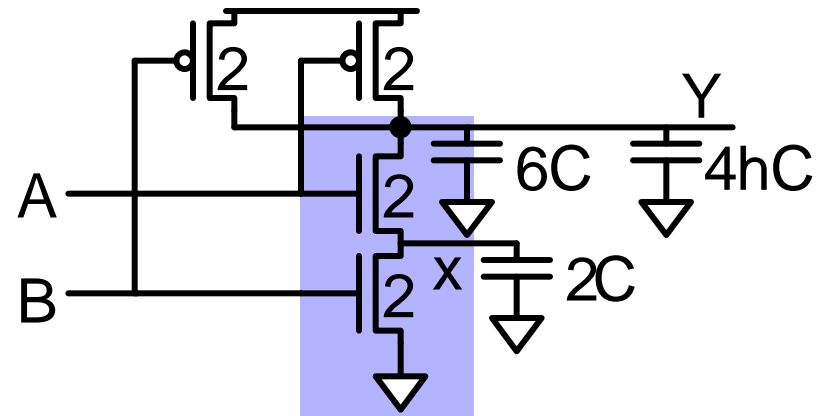
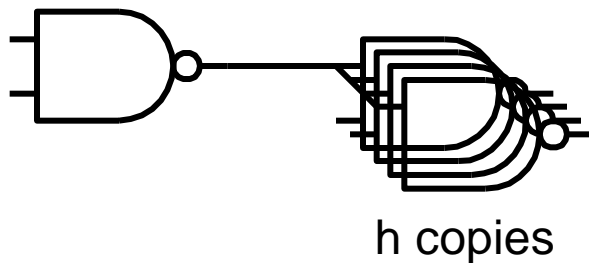


$$t_{pdr} = (6 + 4h)RC$$



# Example: 2-input NAND

- Estimate rising and **falling** propagation delays of a 2-input NAND driving  $h$  identical gates.



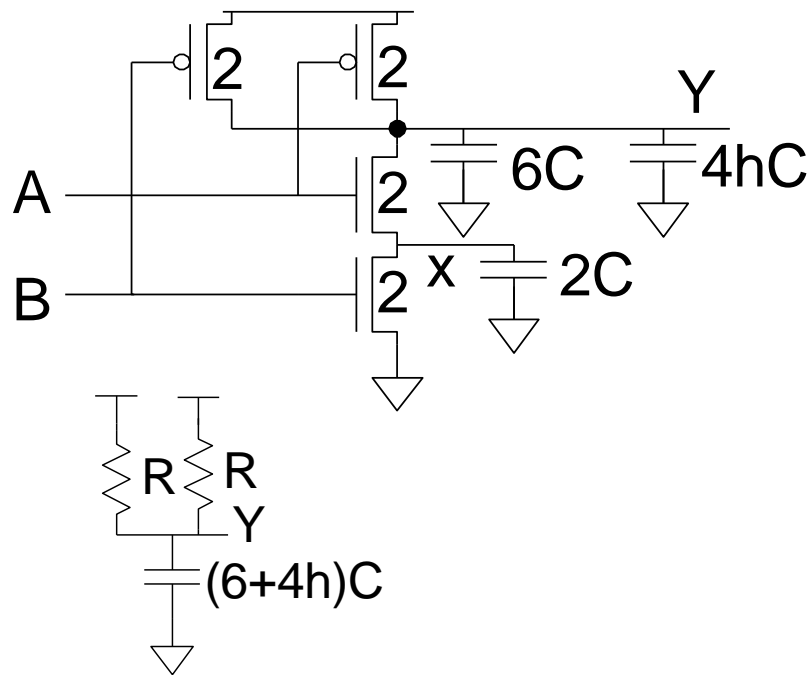
$$\begin{aligned}
 t_{pdf} &= (2C)\left(\frac{R}{2}\right) + \left[(6 + 4h)C\right]\left(\frac{R}{2} + \frac{R}{2}\right) \\
 &= (7 + 4h)RC
 \end{aligned}$$

# Delay Components

- ❑ Delay has two parts
  - *Parasitic delay*
    - 6 or 7 RC
    - Independent of load
  - *Effort delay*
    - 4h RC
    - Proportional to load capacitance

# Contamination Delay

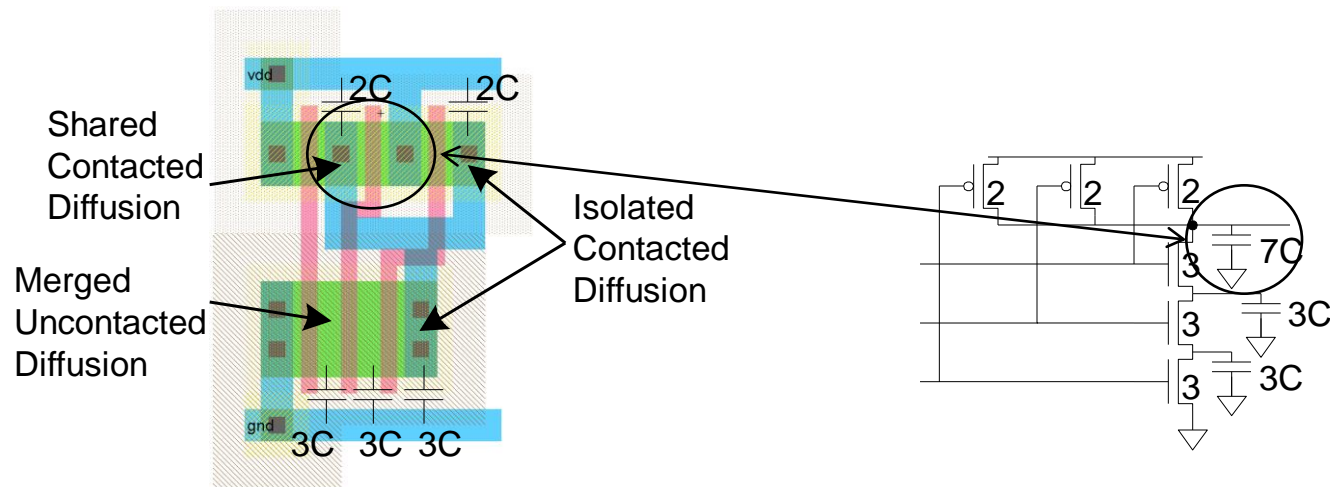
- ❑ Best-case (contamination) delay can be substantially less than propagation delay.
- ❑ Ex: If both inputs fall simultaneously



$$t_{cdr} = (3 + 2h)RC$$

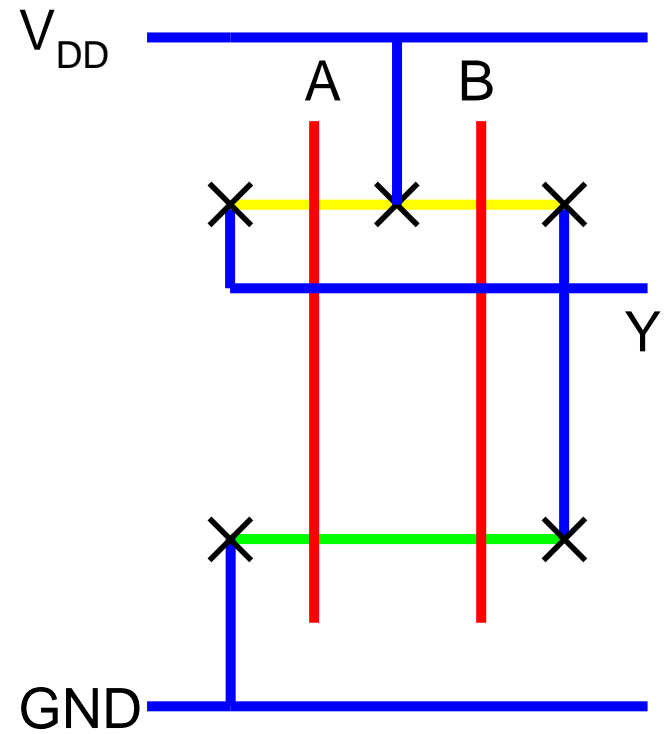
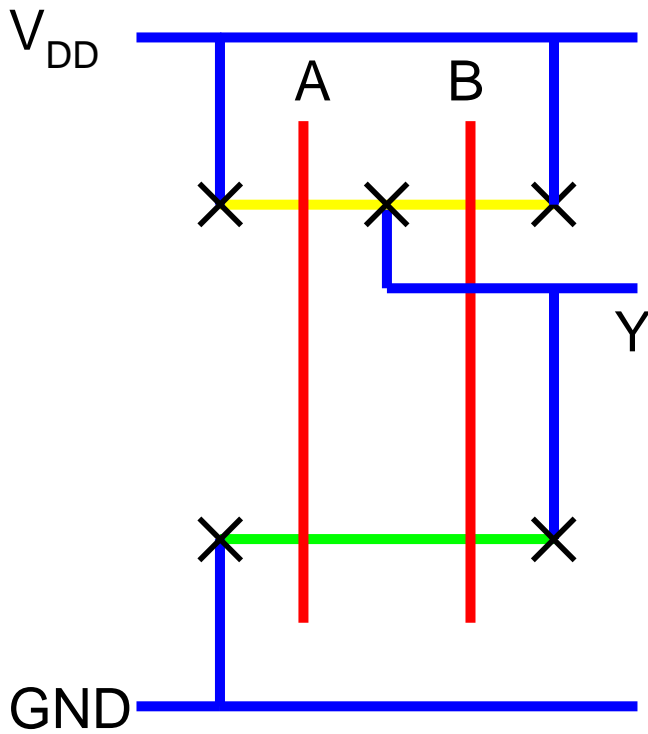
# Diffusion Capacitance

- ❑ we assumed contacted diffusion on every s / d
- ❑ Good layout minimizes diffusion area
- ❑ Ex: NAND3 layout shares one diffusion contact
  - Reduces output capacitance by  $2C$
  - Merged uncontacted diffusion might help too



# Layout Comparison

- ❑ Layout representation by *stick diagram*. What CKT?
- ❑ Which layout is better?



# Power and Energy

❑ Power is drawn from a voltage source attached to the  $V_{DD}$  pin(s) of a chip.

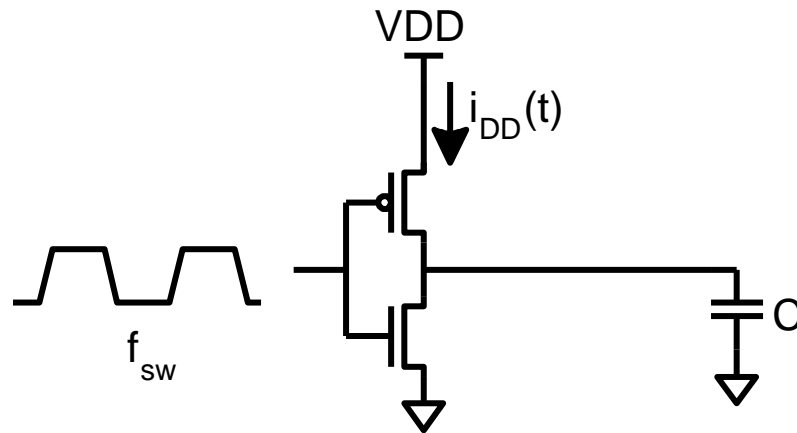
❑ Instantaneous Power:  $P(t) = i_{DD}(t)V_{DD}$

❑ Energy: 
$$E = \int_0^T P(t)dt = \int_0^T i_{DD}(t)V_{DD}dt$$

❑ Average Power: 
$$P_{\text{avg}} = \frac{E}{T} = \frac{1}{T} \int_0^T i_{DD}(t)V_{DD}dt$$

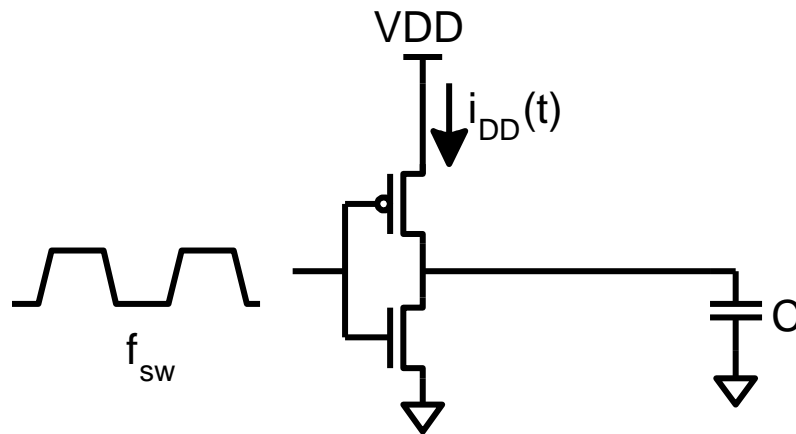
# Dynamic Power

- ❑ Dynamic power is required to charge and discharge load capacitances when transistors switch
- ❑ One cycle involves a rising and falling output
- ❑ On rising output, charge  $Q = CV_{DD}$  is required
- ❑ On falling output, charge is dumped to GND
- ❑ This repeats  $Tf_{sw}$  times over an interval of  $T$



$$P_{\text{dynamic}} = \frac{E}{T} = \frac{1}{T} \int_0^T i_{DD}(t) V_{DD} dt =$$

$$\frac{V_{DD}}{T} \int_0^T i_{DD}(t) dt = \frac{V_{DD}}{T} [T f_{\text{sw}} C V_{DD}] = C V_{DD}^2 f_{\text{sw}}$$





# Activity Factor

- ❑ Suppose the system clock frequency =  $f$
- ❑ Let  $f_{sw} = \alpha f$ , where  $\alpha$  = activity factor
  - If the signal is a clock,  $\alpha = 1$
  - If the signal switches once per cycle,  $\alpha = 1/2$
  - Static gates:
    - Depends on design, but typically  $\alpha = 0.1$
  - Dynamic gates:
    - Switch either 0 or 2 times per cycle,  $\alpha = 1/2$

❑ Dynamic power: 
$$P_{\text{dynamic}} = \alpha C V_{DD}^2 f$$

# Short Circuit Current

- ❑ When transistors switch, both nMOS and pMOS networks may be momentarily ON at once
- ❑ Leads to a blip of “short circuit” current.
- ❑  $< 10\%$  of dynamic power if rise/fall times are comparable for input and output

# Power Dissipation Sources

- ❑  $P_{\text{total}} = P_{\text{dynamic}} + P_{\text{static}}$
- ❑ Dynamic power:  $P_{\text{dynamic}} = P_{\text{switching}} + P_{\text{shortcircuit}}$ 
  - Switching load capacitances
  - Short-circuit current
- ❑ Static power:  $P_{\text{static}} = (I_{\text{sub}} + I_{\text{gate}} + I_{\text{junct}} + I_{\text{contention}})V_{\text{DD}}$ 
  - Sub-threshold leakage
  - Gate leakage
  - Junction leakage
  - Contention current (ratioed logic)

# Dynamic Power Example

- ❑ 1 billion transistor chip
  - 50M logic transistors
    - Average width:  $12 \lambda$
    - Activity factor = 0.1
  - 950M memory transistors
    - Average width:  $4 \lambda$
    - Activity factor = 0.02
  - 1.0 V 65 nm process
  - $C = 1 \text{ fF}/\mu\text{m}$  (gate) +  $0.8 \text{ fF}/\mu\text{m}$  (diffusion)
- ❑ Estimate dynamic power consumption @ 1 GHz.  
Neglect wire capacitance and short-circuit current.

# Power Estimate Ex (Cont'd)

$$C_{\text{logic}} = (50 \times 10^6)(12\lambda)(0.025 \mu m / \lambda)(1.8 fF / \mu m) = 27 nF$$

$$C_{\text{mem}} = (950 \times 10^6)(4\lambda)(0.025 \mu m / \lambda)(1.8 fF / \mu m) = 171 nF$$

$$P_{\text{dynamic}} = [0.1C_{\text{logic}} + 0.02C_{\text{mem}}](1.0)^2 (1.0 GHz) = 6.1 W$$

# Dynamic Power Reduction

□  $P_{\text{switching}} = \alpha C V_{DD}^2 f$

□ Try to minimize:

- Activity factor
- Capacitance
- Supply voltage
- Frequency

# Activity Factor Estimation

- ❑ Let  $P_i = \text{Prob}(\text{node } i = 1)$
- ❑  $\alpha_i = P_i * (1 - P_i)$
- ❑ Completely random data has  $P = 0.5$  and  $\alpha = 0.25$
- ❑ Data is often not completely random
  - e.g. MSBs of 64-bit words in memory address bus. MSBs of data representing measurements of physical phenomena.
- ❑ Data propagating through ANDs and ORs has lower activity factor
  - Depends on design, but typically  $\alpha \approx 0.1$

# Switching Probability

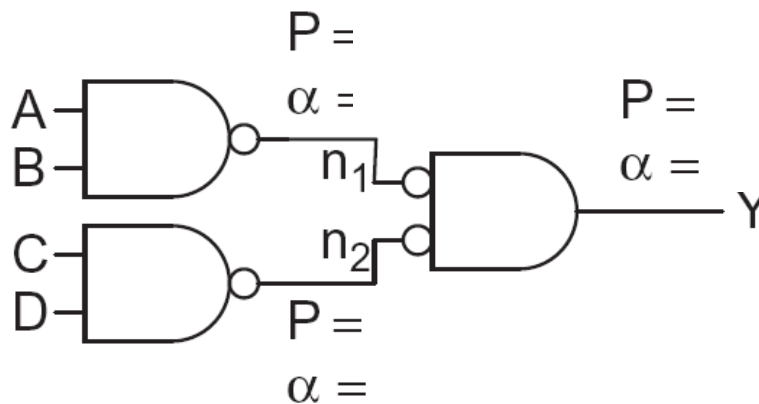
Gate	$P_Y$
AND2	$P_A P_B$
AND3	$P_A P_B P_C$
OR2	$1 - \bar{P}_A \bar{P}_B$
NAND2	$1 - P_A P_B$
NOR2	$\bar{P}_A \bar{P}_B$
XOR2	$P_A \bar{P}_B + \bar{P}_A P_B$

What is the switching probability?



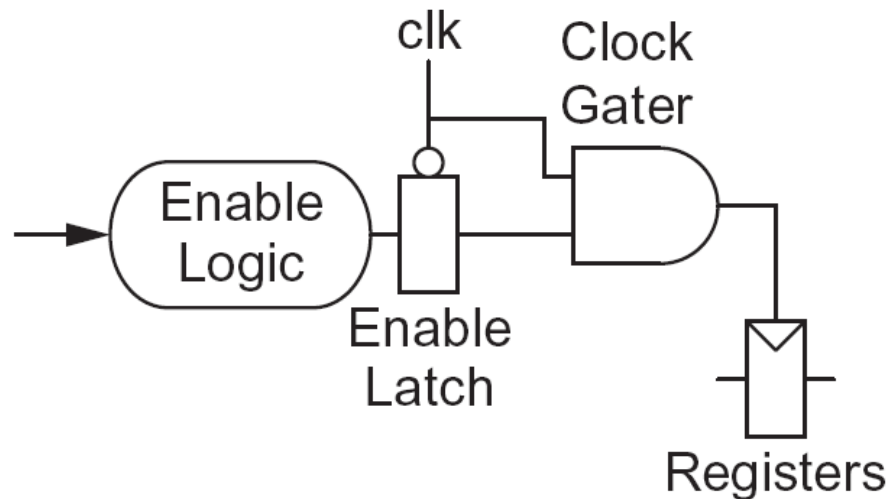
# Example

- ❑ A 4-input AND is built out of two levels of gates
- ❑ Estimate the activity factor at each node if the inputs have  $P = 0.5$



# Clock Gating

- ❑ The best way to reduce the activity is to turn off the clock to registers in unused blocks
  - Saves clock activity ( $\alpha = 1$ )
  - Eliminates all switching activity in the block
  - Requires determining if block will be used



# Capacitance

## ❑ Gate capacitance

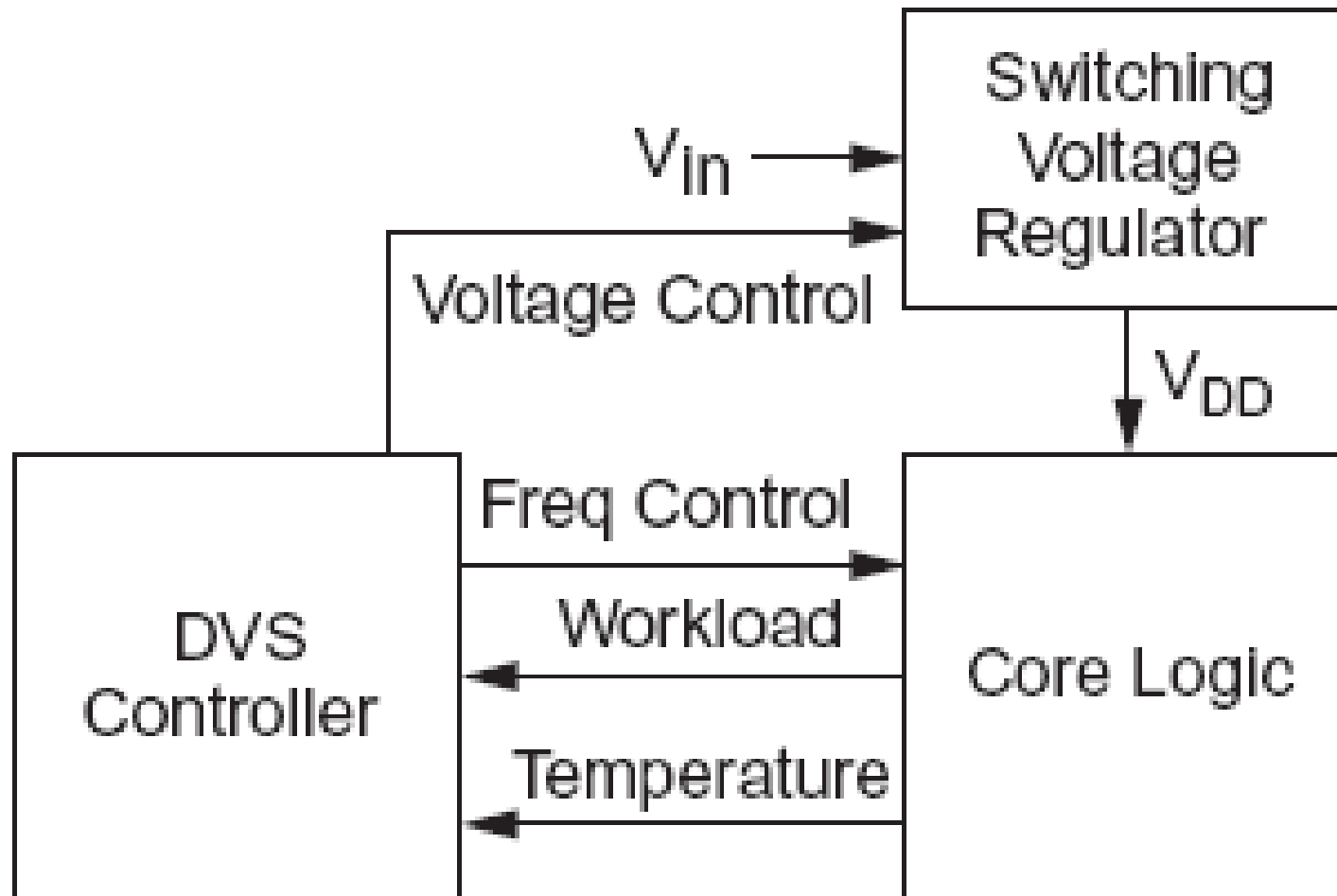
- Fewer stages of logic
- Small gate sizes

## ❑ Wire capacitance

- Good floorplanning to keep communicating blocks close to each other
- Drive long wires with inverters or buffers rather than complex gates

# Voltage / Frequency

- ❑ Run each block at the lowest possible voltage and frequency that meets performance requirements
- ❑ Voltage Domains
  - Provide separate supplies to different blocks
  - Level converters required when crossing from low to high  $V_{DD}$  domains
- ❑ Dynamic Voltage Scaling
  - Adjust  $V_{DD}$  and  $f$  according to workload



# Static Power

- ❑ Static power is consumed even when chip is quiescent
  - Ratioed circuits burn power in fight between ON transistors. Occurs when output is low (0).
  - Leakage draws power from nominally OFF devices

# Static Power Example

- ❑ Revisit power estimation for 1 billion transistor chip
- ❑ Estimate static power consumption
  - Subthreshold leakage
    - Normal  $V_t$ : 100 nA/ $\mu\text{m}$
    - High  $V_t$ : 10 nA/ $\mu\text{m}$
    - High  $V_t$  used in all memories and in 95% of logic gates
  - Gate leakage 5 nA/ $\mu\text{m}$
  - Junction leakage negligible

# Solution

$$W_{\text{normal-}V_t} = (50 \times 10^6)(12\lambda)(0.025 \mu\text{m} / \lambda)(0.05) = 0.75 \times 10^6 \mu\text{m}$$

$$W_{\text{high-}V_t} = \left[ (50 \times 10^6)(12\lambda)(0.95) + (950 \times 10^6)(4\lambda) \right] (0.025 \mu\text{m} / \lambda) = 109.25 \times 10^6 \mu\text{m}$$

$$I_{\text{sub}} = \left[ W_{\text{normal-}V_t} \times 100 \text{ nA}/\mu\text{m} + W_{\text{high-}V_t} \times 10 \text{ nA}/\mu\text{m} \right] / 2 = 584 \text{ mA}$$

$$I_{\text{gate}} = \left[ (W_{\text{normal-}V_t} + W_{\text{high-}V_t}) \times 5 \text{ nA}/\mu\text{m} \right] / 2 = 275 \text{ mA}$$

$$P_{\text{static}} = (584 \text{ mA} + 275 \text{ mA})(1.0 \text{ V}) = 859 \text{ mW}$$



# Leakage Control

$$I_{ds} = I_{ds0} e^{\frac{V_{gs}-V_t}{nv_T}} \left( 1 - e^{\frac{-V_{ds}}{v_T}} \right)$$

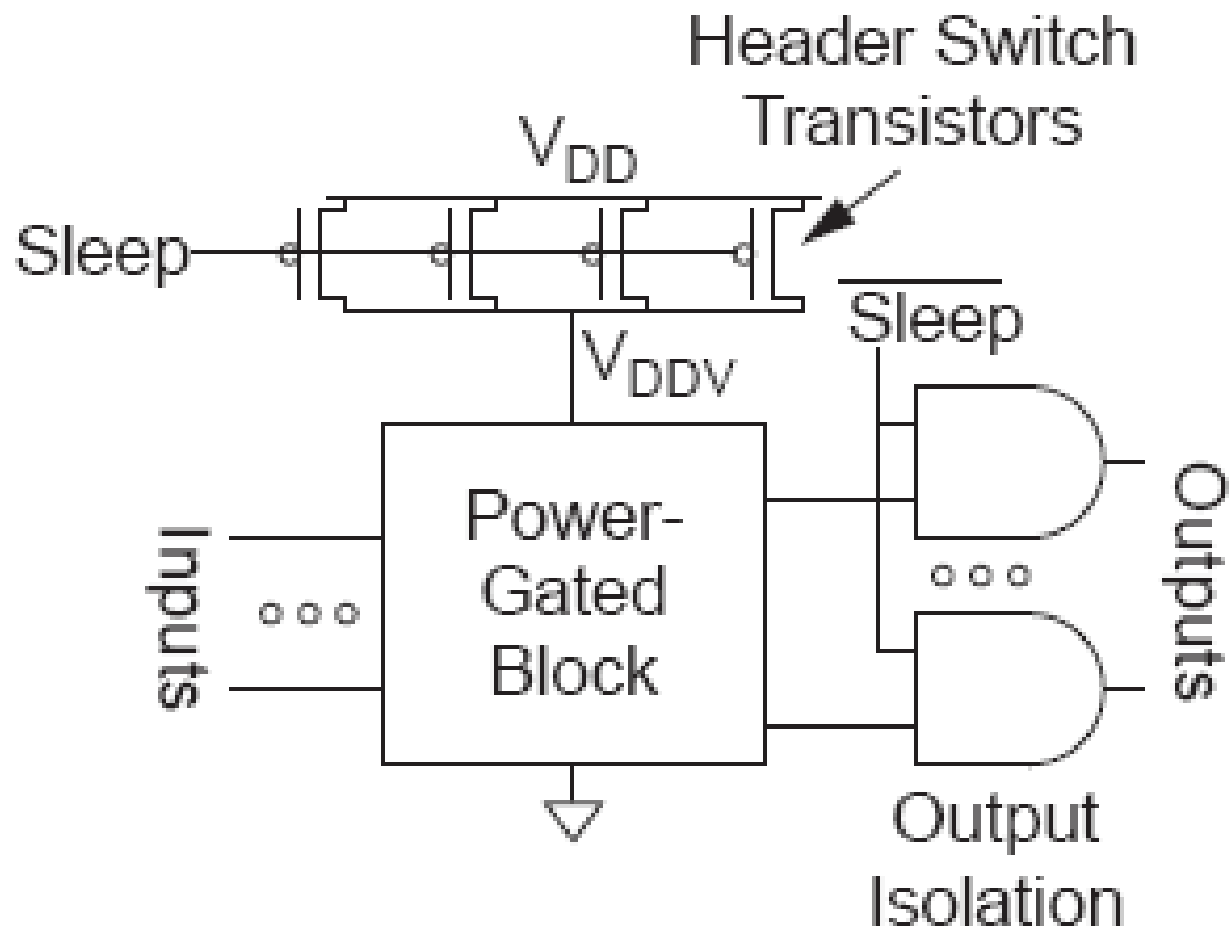
- ❑ Leakage and delay trade off
  - Aim for low leakage in sleep and low delay in active mode
- ❑ To reduce leakage:
  - Increase  $V_t$ : *multiple*  $V_t$ 
    - Use low  $V_t$  only in critical circuits
  - Increase  $V_s$ : *stack effect*
    - *Input vector control* in sleep

# Gate Leakage

- ❑ Extremely strong function of  $t_{ox}$  and  $V_{gs}$ 
  - Negligible for older processes
  - Approaches subthreshold leakage at 65 nm and below in some processes
- ❑ An order of magnitude less for pMOS than nMOS
- ❑ Control leakage in the process using  $t_{ox} > 10.5 \text{ \AA}$ 
  - High-k gate dielectrics help
  - Some processes provide multiple  $t_{ox}$ 
    - e.g. thicker oxide for 3.3 V I/O transistors
- ❑ Control leakage in circuits by limiting  $V_{DD}$

# Power Gating

- ❑ Turn OFF power to blocks when they are idle to save leakage
  - Use virtual  $V_{DD}$  ( $V_{DDV}$ )
  - Gate outputs to prevent invalid logic levels to next block
- ❑ Voltage drop across sleep transistor degrades performance during normal operation
  - Size the transistor wide enough to minimize impact
- ❑ Switching wide sleep transistor costs dynamic power
  - Only justified when circuit sleeps long enough



# Low Power Design

- ❑ Reduce dynamic power
  - ❑  $\alpha$ : clock gating, sleep mode
  - C: small transistors (esp. on clock), short wires
  - $V_{DD}$ : lowest suitable voltage
  - f: lowest suitable frequency
- ❑ Reduce static power
  - Selectively use ratioed circuits (minimize)
  - Selectively use low  $V_t$  devices (minimize)
  - Leakage reduction:  
stacked devices, body bias, low temperature