

# CS/EE 5710/6710

MOS Transistor Models  
Electrical Effects  
Propagation Delay

## Transistor Characteristics

### ▶ Three conduction characteristics

#### ▶ Cutoff Region

- ▶ No inversion layer in channel
- ▶  $I_{ds} = 0$

#### ▶ Nonsaturated, or linear region

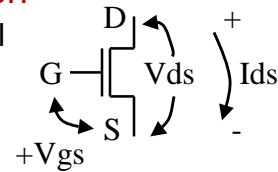
- ▶ Weak inversion of the channel
- ▶  $I_{ds}$  depends on  $V_{gs}$  and  $V_{ds}$

#### ▶ Saturated region

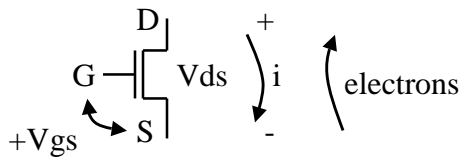
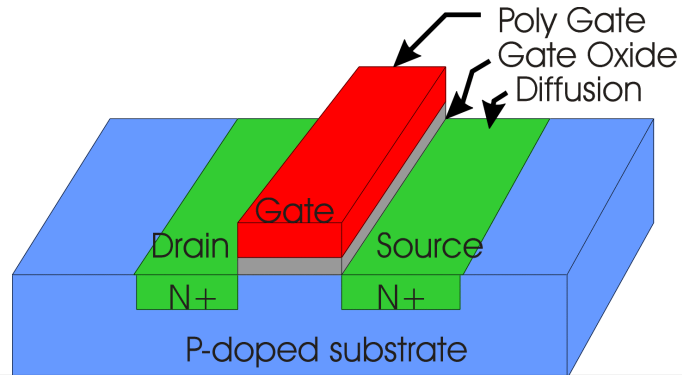
- ▶ Strong inversion of channel
- ▶  $I_{ds}$  is independent of  $V_{ds}$

#### ▶ As an aside, at very high drain voltages:

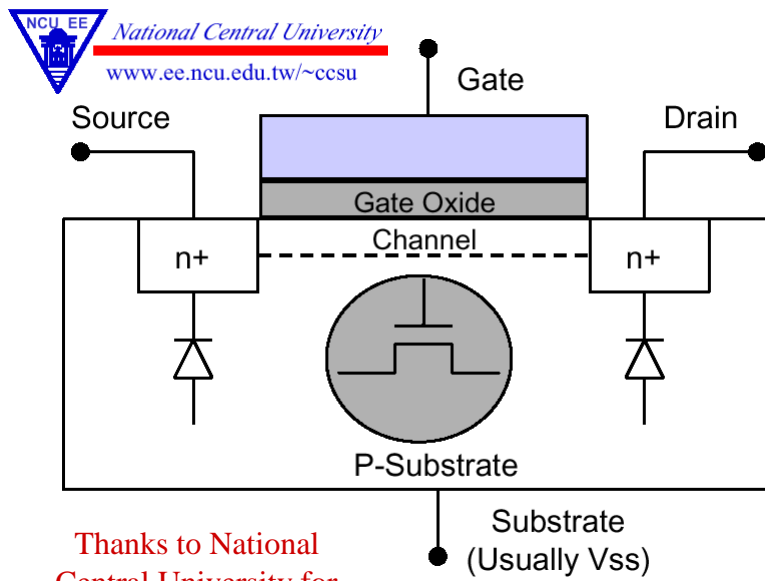
- ▶ “avalanche breakdown” or “punch through”
- ▶ Gate has no control of  $I_{ds}$ ...



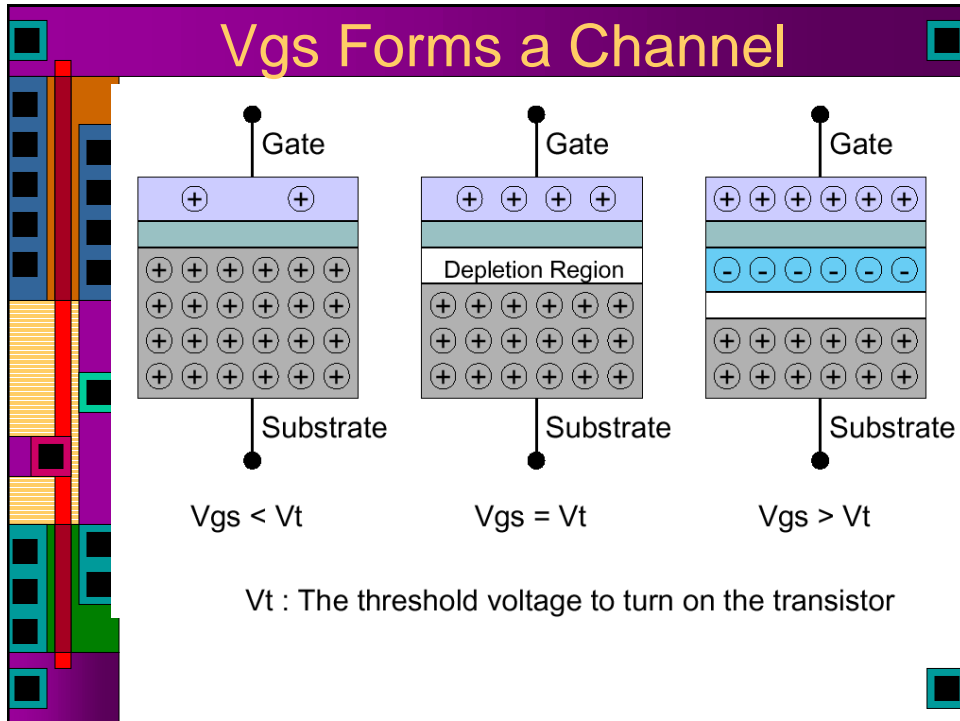
# N-type Transistor



# Another Cutaway View



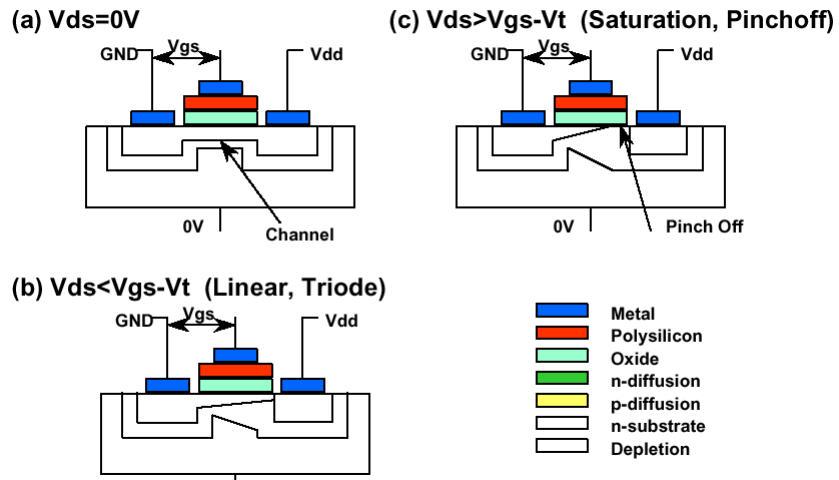
Thanks to National Central University for Some images



## Basic N-Type MOS Transistor

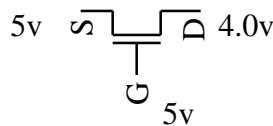
- ▶ Conditions for the regions of operation
  - ▶ **Cutoff:** If  $V_{gs} < V_t$ , then  $I_{ds}$  is essentially 0
    - ▶  $V_t$  is the “Threshold Voltage”
  - ▶ **Linear:** If  $V_{ds} < (V_{gs} - V_t)$  then  $I_{ds}$  depends on both  $V_{gs}$  and  $V_{ds}$ 
    - ▶ Channel becomes deeper as  $V_{gs}$  goes up
  - ▶ **Saturated:** If  $(V_{gs} - V_t) < V_{ds}$  then  $I_{ds}$  is essentially constant (Saturated)

## Visualizing the Channel

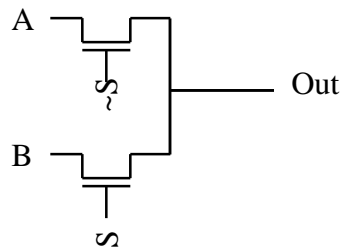


## "Saturated" Transistor

- ▶ In the  $0 < (V_{gs} - V_t) < V_{ds}$  case
  - ▶  $I_{ds}$  Current is effectively constant
  - ▶ Channel is "pinched off" and conduction is accomplished by drift of carriers
  - ▶ Voltage across pinched off channel (i.e.  $V_{ds}$ ) is fixed at  $V_{gs} - V_t$ 
    - ▶ This is why you don't use an N-type to pass 1's!
    - ▶ High voltage is degraded by  $V_t$
    - ▶  $V_{gs} - V_t$  is 1.0v, 5v in one side, 4.0v out the other

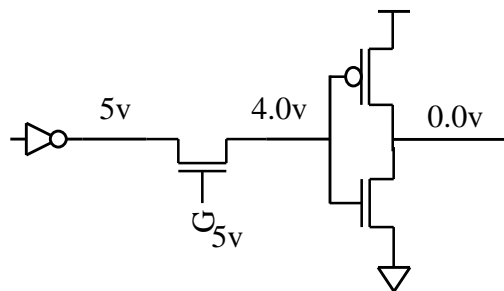


## Aside: N-type Pass Transistors



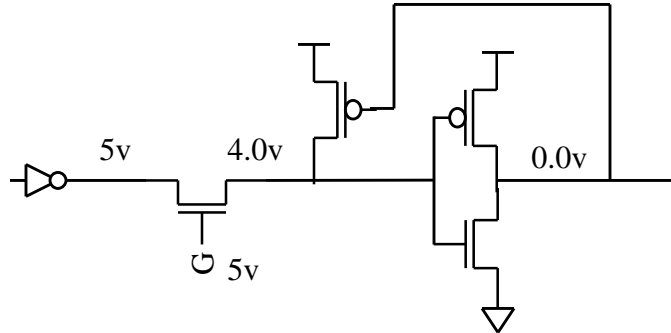
- ▶ If it weren't for the threshold drop, N-type pass transistors (without the P-type transmission gate) would be nice
  - ▶ 2-way Mux Example...

## N-type Pass Transistors



- ▶ On one hand, the degraded high voltage from the pass transistor will be restored by the inverter
- ▶ On the other hand, the P-device may not turn off completely resulting in extra power being used

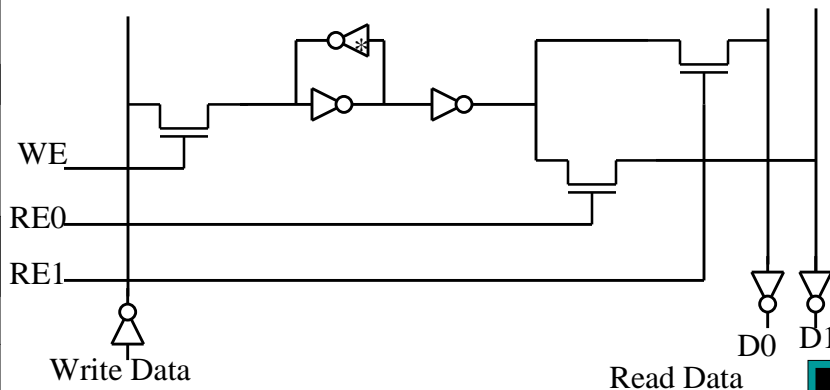
## N-type Pass Transistors



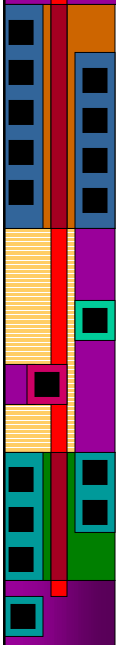
- ▶ Another option is a “keeper” transistor fed back from the output
  - ▶ This pulls the internal node high when the output is 0
  - ▶ But is disconnected when output is high
- ▶ Make sure the size is right...

## N-type Pass Transistors

- ▶ In practice, they are used fairly often, but be aware of what you’re doing
  - ▶ For example, read/write circuits in a **Register File**



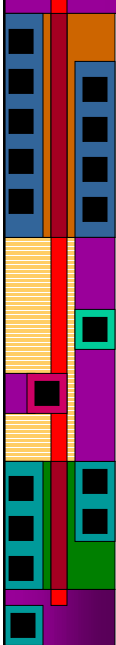
## Back to the Saturated Transistor



A vertical cross-sectional diagram of a MOSFET transistor. From top to bottom, it shows: a gate stack (orange and blue layers), a channel region (yellow and purple layers), a source/drain region (green and blue layers), and a substrate (purple layer). A red vertical line indicates the channel length. Small black squares are scattered throughout the diagram.

- ▶ What influences the constant  $I_{ds}$  in the saturated case?
  - ▶ Channel length
  - ▶ Channel width
  - ▶ Threshold voltage  $V_t$
  - ▶ Thickness of gate oxide
  - ▶ Dielectric constant of gate oxide
  - ▶ Carrier mobility  $\mu$

## Threshold Voltage



A vertical cross-sectional diagram of a MOSFET transistor, identical to the one in the first slide. It shows the gate stack, channel region, source/drain region, and substrate. A red vertical line indicates the channel length. Small black squares are scattered throughout the diagram.

- ▶ The  $V_{gs}$  voltage at which  $I_{ds}$  is essentially 0
  - ▶ Tiny  $I_{ds}$  is exponentially related to  $V_{gs}$ ,  $V_{ds}$
  - ▶ Take 5720/6720 for “subthreshold” circuit ideas
- ▶  $V_t$  is affected by
  - ▶ Gate conductor material
  - ▶ Gate insulator material
  - ▶ Gate insulator thickness
  - ▶ Channel doping
  - ▶ Impurities at Si/insulator interface
  - ▶ Voltage between source and substrate ( $V_{sb}$ )

## Basic DC Equations for $I_{ds}$

### ▶ Cutoff Region

▶  $V_{gs} < V_t, I_{ds} = 0$

### ▶ Linear Region

▶  $0 < V_{ds} < (V_{gs} - V_t)$

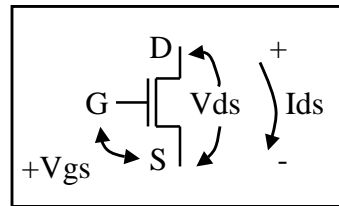
$$I_{ds} = \beta[(V_{gs} - V_t)V_{ds} - V_{ds}^2/2]$$

▶ Note that this is only “linear” if  $V_{ds}^2/2$  is very small, i.e.  $V_{ds} \ll V_{gs} - V_t$

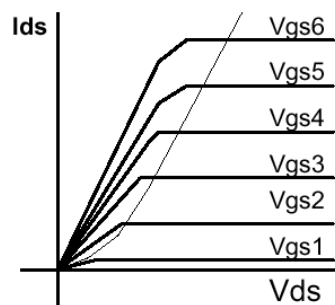
### ▶ Saturated Region

▶  $0 < (V_{gs} - V_t) < V_{ds}$

$$I_{ds} = \beta[(V_{gs} - V_t)^2/2]$$



## $I_{ds}$ Curves



$$\beta = \frac{\mu\epsilon}{t_{ox}} \left(\frac{W}{L}\right) = \mu C_{ox} \left(\frac{W}{L}\right)$$

### Cutoff Region

$$V_{gs} < V_t$$

$$I_{ds} = 0$$

### Triode (Linear) Region

$$V_{gs} - V_t > V_{ds} > 0$$

$$I_{ds} = \beta \left[ (V_{gs} - V_t)V_{ds} - \frac{V_{ds}^2}{2} \right]$$

### Saturation Region

$$V_{gs} - V_t > V_{ds} > 0$$

$$I_{ds} = \beta \frac{(V_{gs} - V_t)^2}{2}$$



## Transistor Gain

▶  $\beta$  is the MOS transistor gain factor

$$\beta = (\mu\epsilon/t_{ox})(W/L)$$

Process-dependent

Layout dependent

▶  $\mu$  = mobility of carriers

▶ Note that N-type is twice as good as P-type

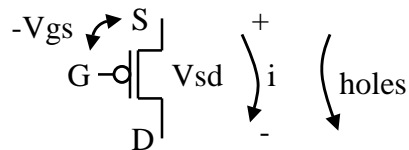
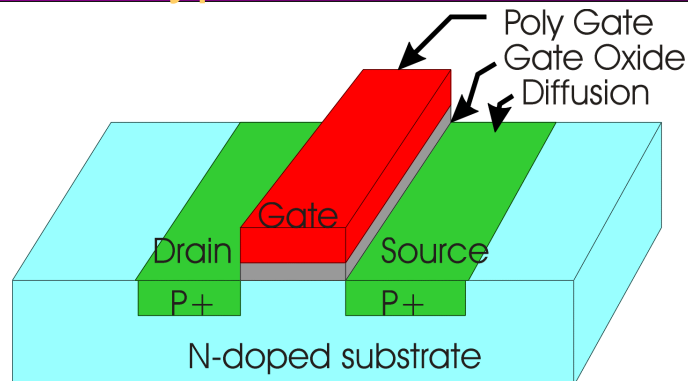
▶  $\epsilon$  = permittivity of gate insulator

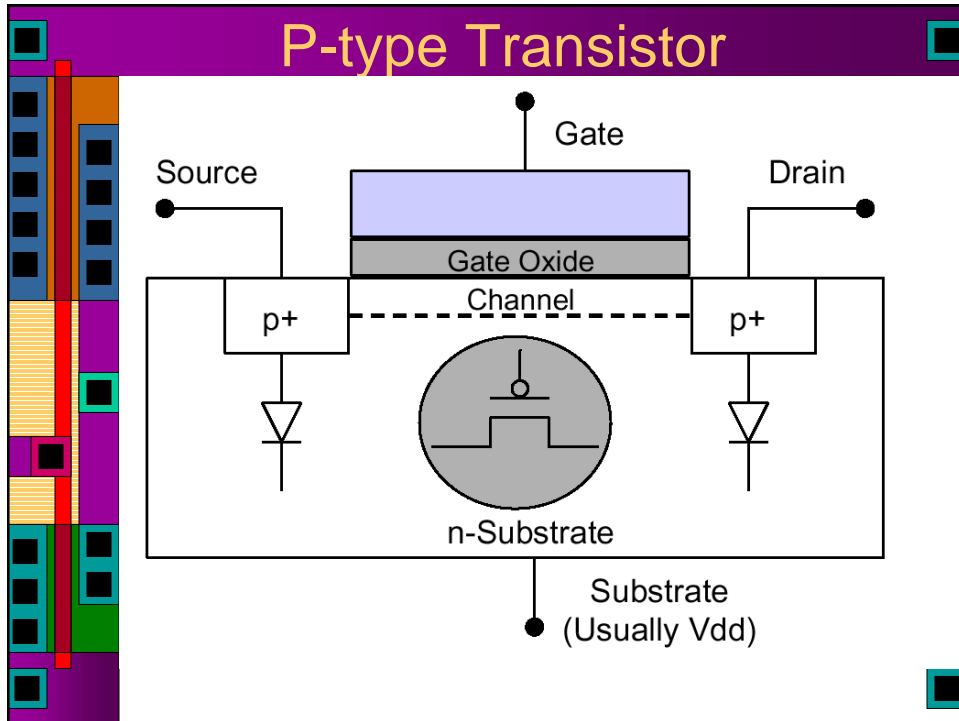
▶  $T_{ox}$  = thickness of gate oxide

▶ Book calls  $(\mu\epsilon/t_{ox}) = k'$

▶ Increase  $W/L$  to increase gain

## P-type Transistor



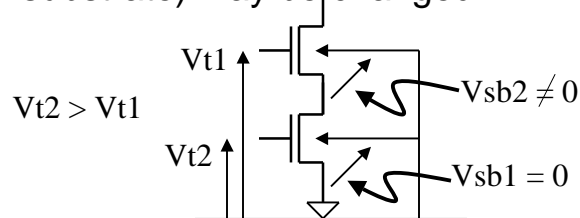


## P-type Transistors

- ▶ Source is Vdd instead of GND
  - ▶  $V_{sg} = (V_{dd} - V_{in})$ ,  $V_{sd} = (V_{dd} - V_{out})$ ,  $V_t$  is negative
- ▶ **Cutoff:**  $(V_{dd} - V_{in}) < -V_t$ ,  $I_{ds} = 0$
- ▶ **Linear Region**
  - ▶  $(V_{dd} - V_{out}) < (V_{dd} - V_{in} + V_t)$
  - $I_{ds} = \beta[(V_{dd} - V_{in} + V_t)(V_{dd} - V_{out}) - (V_{dd} - V_{out})^2/2]$
- ▶ **Saturated Region**
  - ▶  $((V_{dd} - V_{in}) + V_t) < (V_{dd} - V_{out})$
  - $I_{ds} = \beta[(V_{dd} - V_{in} + V_t)^2/2]$

## 2<sup>nd</sup> Order Effect: Body Effect

- ▶ A second order effect that raises  $V_t$
- ▶ Recall that  $V_t$  is affected by  $V_{sb}$  (voltage between source and substrate)
  - ▶ Normally this is constant because of common substrate
  - ▶ But, when transistors are in series,  $V_{sb}$  ( $V_s - V_{\text{substrate}}$ ) may be changed



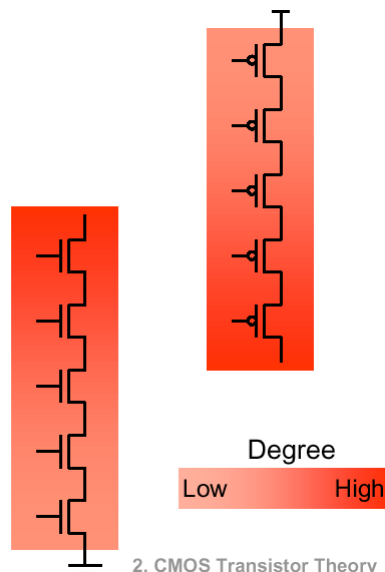
## 2<sup>nd</sup> Order Effect: Body Effect

- **Body Effect -**  
 $V_t$  is a function of voltage between source and substrate

$$V_t = V_{t0} + \gamma \sqrt{(2\phi_b + |V_{sb}|) + 2\sqrt{\phi_b}}$$

$$\phi_b = \frac{kT}{q} \ln\left(\frac{N_A}{N_i}\right)$$

$$\gamma = \frac{t_{ox}}{\epsilon_{ox}} \sqrt{2q\epsilon_{si}N_A} = \frac{1}{C_{ox}} \sqrt{2q\epsilon_{si}N_A}$$



## 2<sup>nd</sup> Order Effect

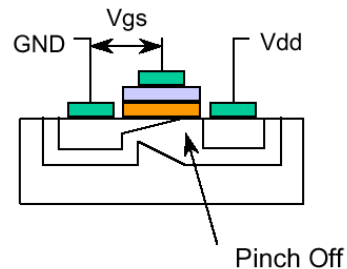
### • Channel Length Modulation -

Channel length is a function of  $V_{ds}$ . When  $V_{ds}$  increase, the depletion region of the pinch off at drain shorten the channel length.

$$L_{eff} = L = L_{short}$$

$$L_{short} = \sqrt{2 \frac{\epsilon_{si}}{qN_A} (V_{ds} - (V_{gs} - V_t))}$$

$$I_{ds} = \frac{kW}{2L} (V_{gs} - V_t)^2 (1 + \lambda V_{ds})$$



## 2<sup>nd</sup> Order Effect

### • Mobility Variation -

The mobility of the carrier decreases when the carrier density increases. Therefore, when  $V_{gs}$  is large. The density of the carrier in the channel increases. As a result, the mobility decreases.

$$\mu = \frac{\text{Average\_carrier\_drift\_velocity}(V)}{\text{Electrical\_Field}(E)}$$

$$\mu_n = 600 \text{ cm}^2 / V \cdot \text{sec}$$

$$\mu_p = 250 \text{ cm}^2 / V \cdot \text{sec}$$

## 2<sup>nd</sup> Order Effect

- **Fowler-Nordheim Tunneling**

When the gate oxide is very thin, a current can flow from gate to source by electron tunneling through the gate oxide.

$$I_{FN} = C_1 W L E_{ox}^2 e^{\frac{-E_o}{E_{ox}}}$$

$$E_{ox} = \frac{V_{gs}}{t_{ox}}$$

- **Drain Punchthrough**

When the drain voltage is high enough, the depletion region around the drain may extend to the source. Thus, causing current to flow irrespective of the gate voltage.

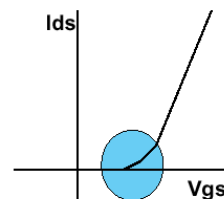
## 2<sup>nd</sup> Order Effect

- **Impact Ionization - Hot Electrons**

When the source-drain electric field is too large, the electron speed will be high enough to break the electron-hole pair. Moreover, the electrons will penetrate the gate oxide, causing a gate current.

- **Subthreshold Region**

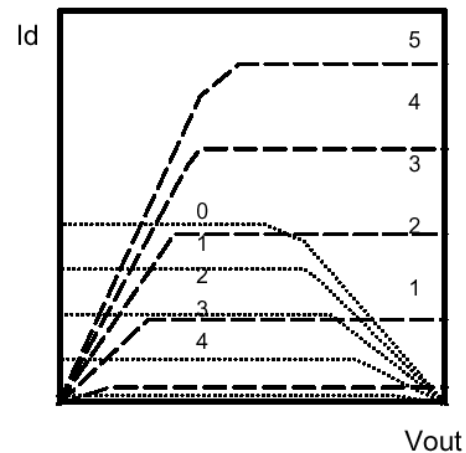
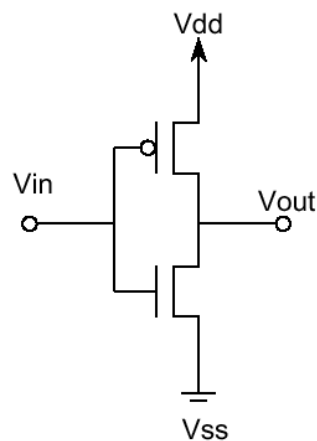
The cutoff region is also referred to as the subthreshold region, where  $I_{ds}$  increase exponentially with  $V_{ds}$  and  $V_{gs}$ .



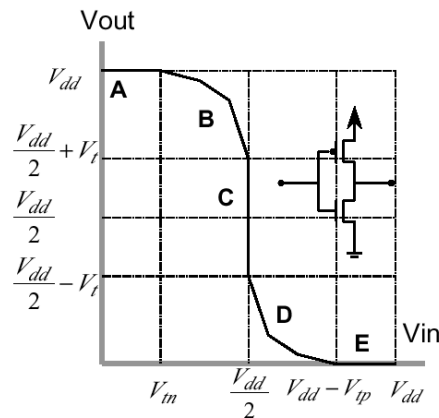
## Inverter Switching Point

- ▶ Inverter switching point is determined by ratio of  $\beta_n/\beta_p$ 
  - ▶ If  $\beta_n/\beta_p = 1$ , then switching point is  $V_{dd}/2$
- ▶ If W/L of both N and P transistors are equal
  - ▶ Then  $\beta_n/\beta_p = \mu_n/\mu_p =$  electron mobility / hole mobility
  - ▶ This ratio is usually between 2 and 3
  - ▶ Means ratio of  $W_{p\text{tree}}/W_{n\text{tree}}$  needs to be between 2 and 3 for  $\beta_n/\beta_p = 1$
  - ▶ For this class, we'll use  $W_{p\text{tree}}/W_{n\text{tree}} = 2$

## Inverter Switching Point



## Inverter Operating Regions



Region	NMOS	PMOS
A	Off	-
B	Sat	Linear
C	Sat	Sat
D	Linear	Sat
E	-	Off

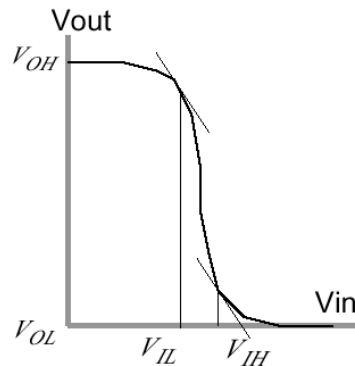
## Gate Sizes

- ▶ Assume minimum inverter is  $W_p/W_n = 2/1$  ( $L = L_{min}$ ,  $W_n = W_{min}$ ,  $W_p = 2W_n$ )
  - ▶ This becomes a 1x inverter
- ▶ To drive larger capacitive loads, you need more gain, more  $I_{ds}$ 
  - ▶ Increase widths to get 2x inverter
  - ▶  $W_p/W_n$  is still 2/1, but  $W_p$  and  $W_n$  are double the size
  - ▶ For most gates, diminishing returns after about 4x size





## Inverter Noise Margin



If  $\beta_n = \beta_p$  and  $V_{tn} = V_{tp}$

$$V_{IH} = \frac{1}{8}(5V_{dd} - 2V_t)$$

$$V_{OH} = V_{dd}$$

$$NM_H = \frac{1}{8}(3V_{dd} + 2V_t)$$

$$V_{IL} = \frac{1}{8}(3V_{dd} + 2V_t)$$

$$V_{OH} = 0$$

$$NM_H = \frac{1}{8}(3V_{dd} + 2V_t)$$

## Performance Estimation

- ▶ First we need to have a model for resistance and capacitance
  - ▶ Delays are caused (to first order) by RC delays charging and discharging capacitors
- ▶ All these layers on the chip have R and C associated with them
- ▶ Mostly this is handled in the Spectre simulator
  - ▶ But it's good to have an idea what's going on

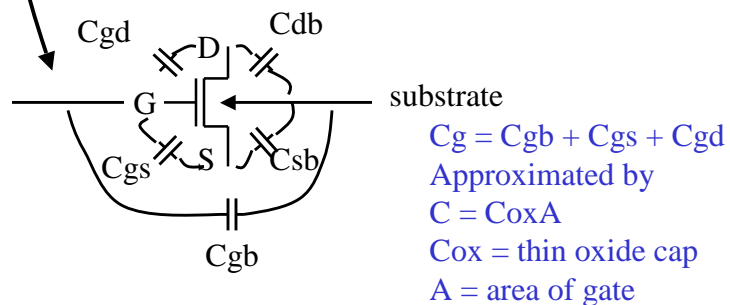
## Resistance

- ▶  $R = (\rho/t)(L/W) = R_s(L/W)$ 
  - ▶  $\rho$  = resistivity of the material
  - ▶  $t$  = thickness
  - ▶  $R_s$  = sheet resistance in  $\Omega/\text{square}$
- ▶ Typical values of  $R_s$

	Min	Typ	Max
M3	0.03	0.04	0.05
M1, M2	0.05	0.07	0.1
Poly	15	20	30
Silicide	2	3	6
Diffusion	10	25	100
Nwell	1k	2k	5k

## Capacitance

- ▶ Three main forms:
  - ▶ Gate capacitance (gate of transistor)
  - ▶ Diffusion capacitance (drain regions)
  - ▶ Routing capacitance (metal, etc.)

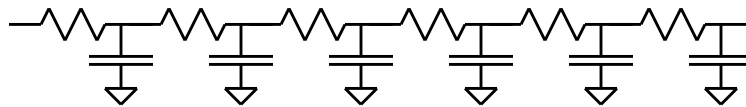


## Routing Capacitance

- ▶ First order effect is layer->substrate
  - ▶ Approximate using parallel plate model
  - ▶  $C = (\epsilon/t)A$ 
    - ▶  $\epsilon$  = permittivity of insulator
    - ▶  $t$  = thickness of insulator
    - ▶  $A$  = area
  - ▶ Fringing fields increase effective area
- ▶ Capacitance between layers becomes very complex!
  - ▶ Crosstalk issues...

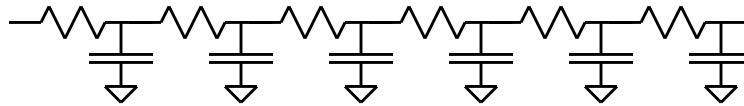
## Distributed RC on Wires

- ▶ Wires look like distributed RC delays
  - ▶ Long resistive wires can look like transmission lines
  - ▶ Inserting buffers can really help delay



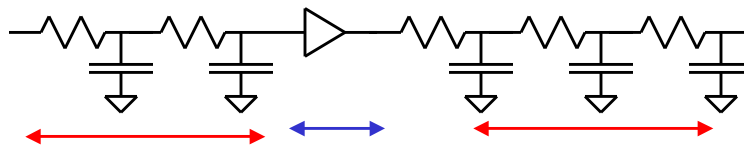
- ▶  $T_n = RCn(n+1)/2$
- ▶  $T = RCL^2/2$  as the number of segments becomes large
  - ▶  $R$  = resistance per unit length
  - ▶  $C$  = capacitance per unit length
  - ▶  $L$  = length of wire

## RC Wire Delay Example



- ▶  $R = 20\Omega/\text{sq}$
- ▶  $C = 4 \times 10^{-4} \text{ pF}/\mu\text{m}$
- ▶  $L = 2\text{mm}$
- ▶  $T = 4 \times 10^{-15} (2000)^2 \text{ s}$ 
  - ▶ delay = 16 ns

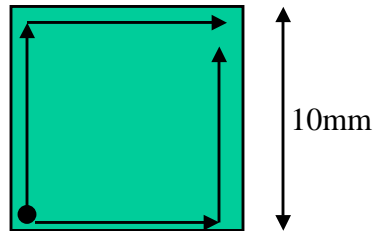
## RC Wire/Buffer Delay Example



- ▶ Now split into 2 1mm segments with a buffer
- ▶  $T = 2 \times (4 \times 10^{-15} (1000)^2) + T_{\text{buf}}$   
 $= 8\text{ns} + T_{\text{buf}}$
- ▶ Assuming  $T_{\text{buf}}$  is less than 8ns (which it will be), the split wire is a win

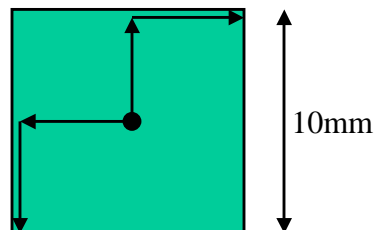
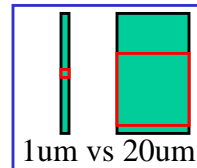
## Another Example: Clock

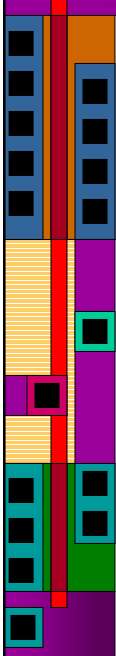
- ▶ 50pF clock load distributed across 10mm chip in 1um metal
  - ▶ Clock length = 20mm
  - ▶  $R = 0.05\Omega/\text{sq}$ ,  $C = 50\text{pF}/20\text{mm}$
  - ▶  $T = (RC/2)L^2 = (6.25 \times 10^{-17})(20,000)^2 = 25\text{ns}$



## Different Distribution Scheme

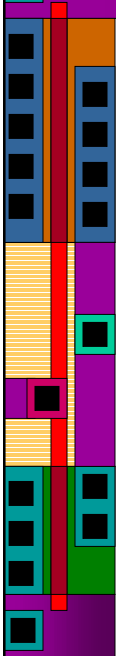
- ▶ Put clock driver in the middle of the chip
- ▶ Widen clock line to 20um wires
  - ▶ Clock length = 10mm
  - ▶  $R = 0.05\Omega/\text{sq}$ ,  $C = 50\text{pF}/20\text{mm}$
  - ▶  $T = (RC/2)L^2 = (0.31 \times 10^{-17})(10,000)^2 = 0.31\text{ns}$
  - ▶ Reduces  $R$  by a factor of 20,  $L$  by 2
  - ▶ Increases  $C$  a tiny bit





## Capacitance Design Guide

- ▶ Get a table of typical capacitances per unit square for each layer
  - ▶ Capacitance to ground
  - ▶ Capacitance to another layer
- ▶ Add them up...
- ▶ See, for example, Table 2-4 in your book

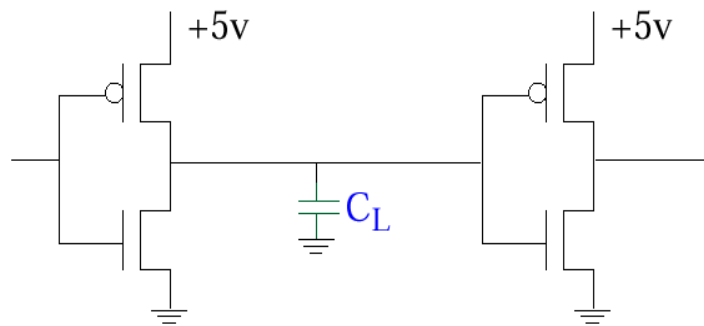


## Wire Length Design Guide

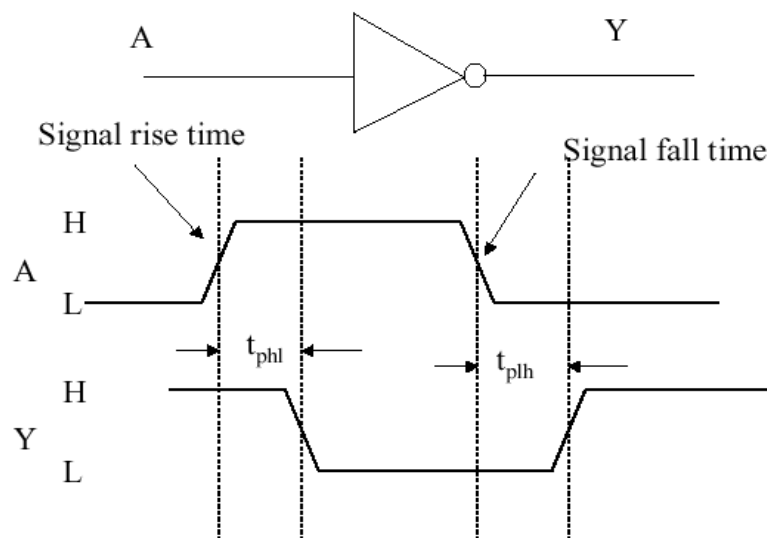
- ▶ How much wire can you use in a conducting layer before the RC delay approaches that of a unit inverter?
  - ▶ Metal3 = 2,500u
  - ▶ Metal2 = 2,000u
  - ▶ Metal1 = 1,250u
  - ▶ Silicide = 150u
  - ▶ Poly = 50u
  - ▶ Diffusion = 15u

## Propagation Delay

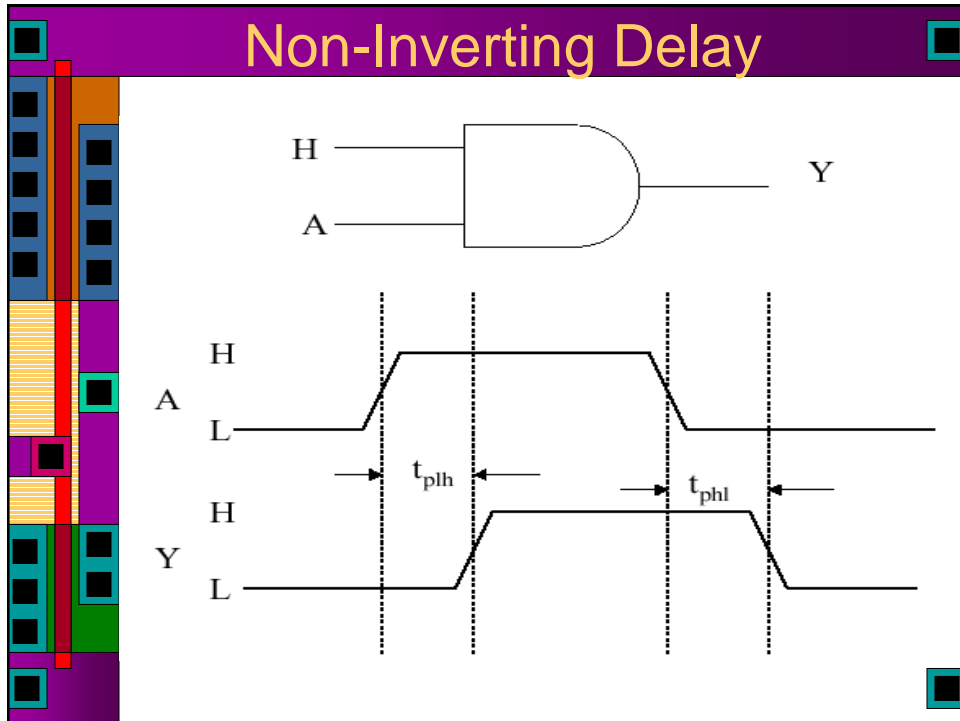
- Recall that it takes time to charge capacitors
- Recall that the gate of a transistor looks like a capacitor
- Wires have resistance and capacitance also!



## Inverting Propagation Delay



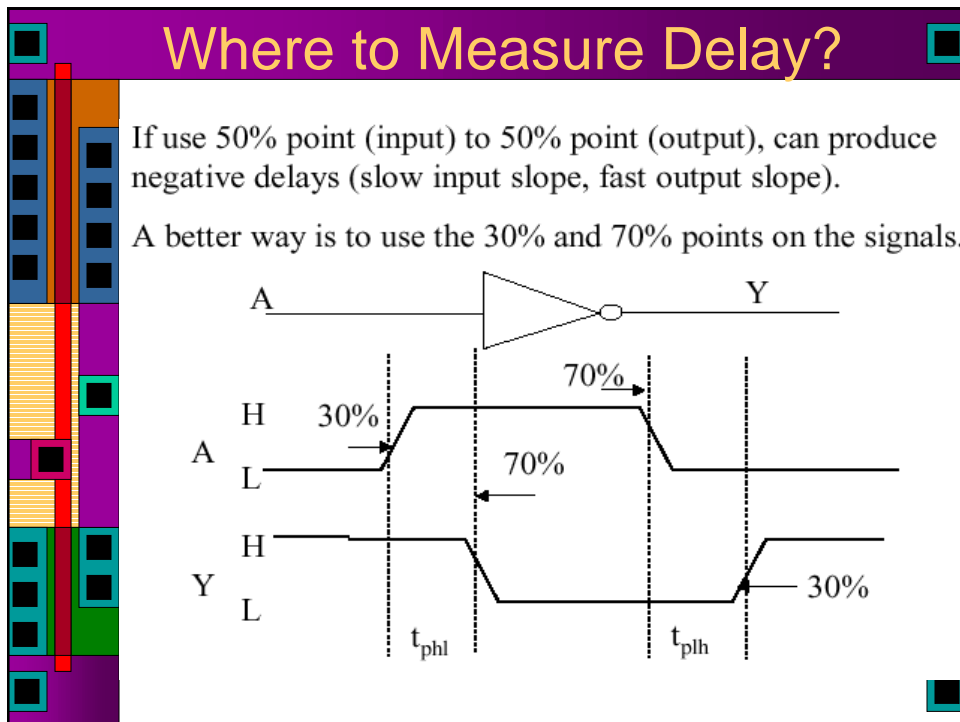
## Non-Inverting Delay



## Where to Measure Delay?

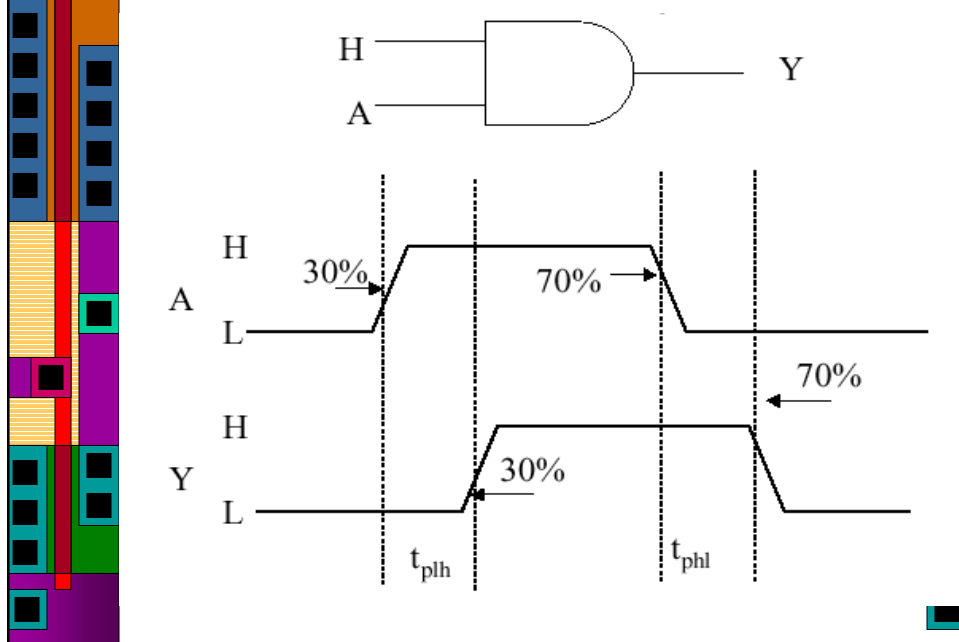
If use 50% point (input) to 50% point (output), can produce negative delays (slow input slope, fast output slope).

A better way is to use the 30% and 70% points on the signals.





## Example Non-Inverting Gate

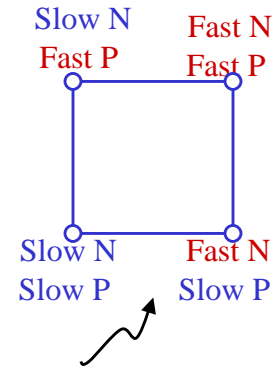


## What Affects Gate Delay?

- ▶ Environment
  - ▶ Increasing  $V_{dd}$  decreases delay
  - ▶ Decreasing temperature decreases delay
  - ▶ Fabrication effects, fast/slow devices
- ▶ Usually measure delay for at least three cases:
  - ▶ Best - high  $V_{dd}$ , low temp, fast N, Fast P
  - ▶ Worst - low  $V_{dd}$ , high temp, slow N, Slow P
  - ▶ Typical - typ  $V_{dd}$ , room temp (25C), typ N, typ P

## Process Corners

- ▶ When parts are specified, under what operating conditions?
- ▶ **Temp:** three ranges
  - ▶ Commercial: 0 C to 70 C
  - ▶ Industrial: -40 C to 85 C
  - ▶ Military: -55 C to 125 C
- ▶ **Vdd:** Should vary  $\pm 10\%$ 
  - ▶ 4.5 to 5.5v for example
- ▶ **Process variation:**
  - ▶ Each transistor type can be slow or fast



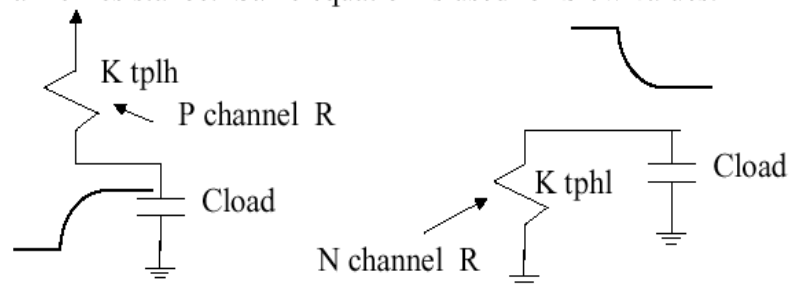
## What Else Affects Gate Delay?

Input slew and output load both effect timing. For a FIXED input slope, FIXED environment, a simple timing model is:

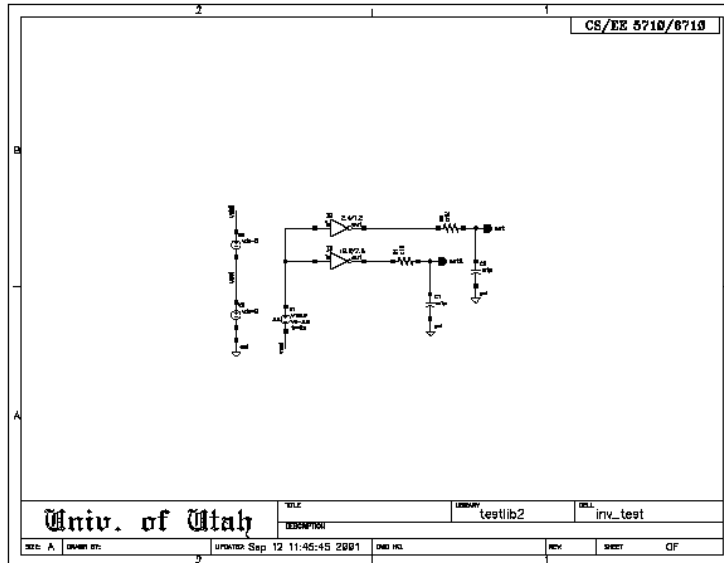
$$\text{delay} = T_{\text{noload}} + K * C_{\text{load}}$$

$T_{\text{noload}}$  is the delay of the gate with no external load.

$K$  is different for TPLH, TPHL since it represents the channel resistance. Same equation is used for Slew values.

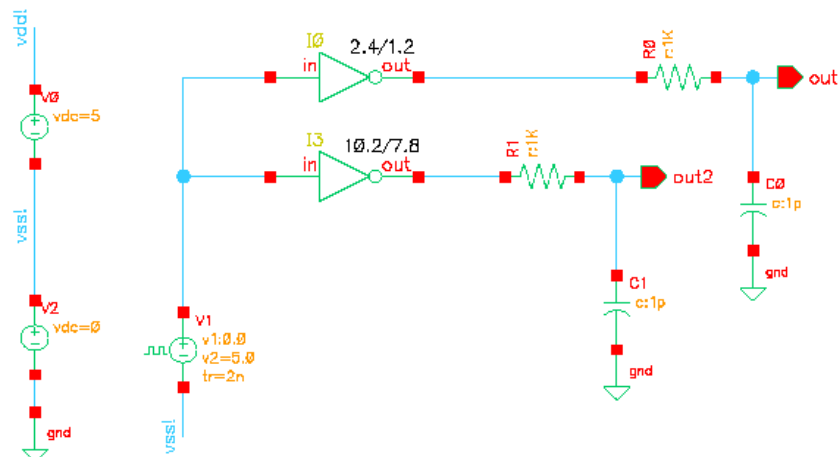


## Inv\_Test Schematic

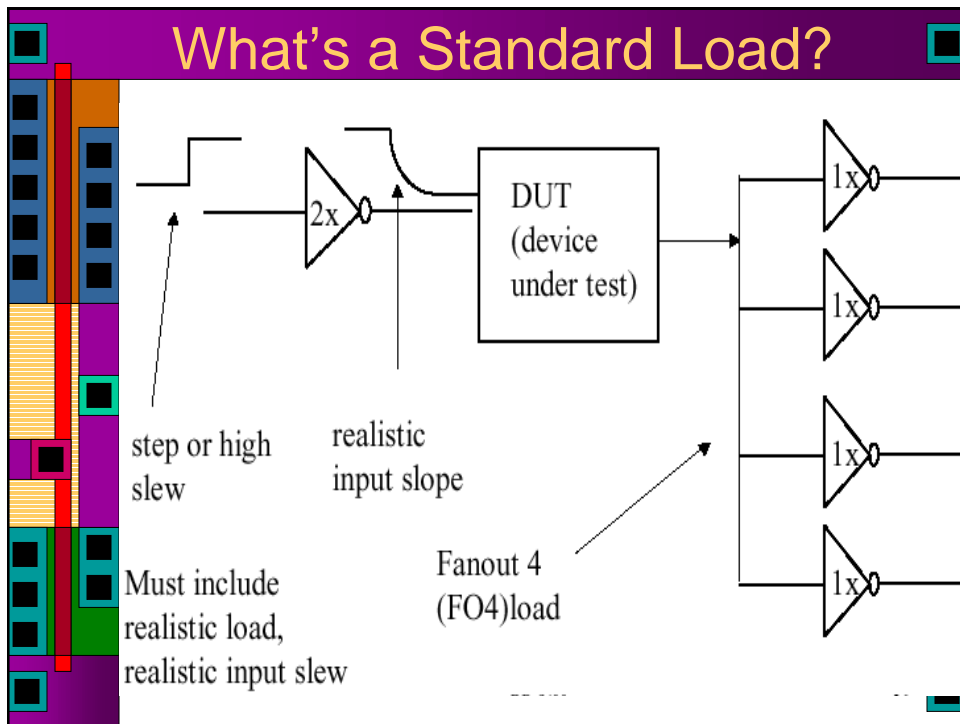
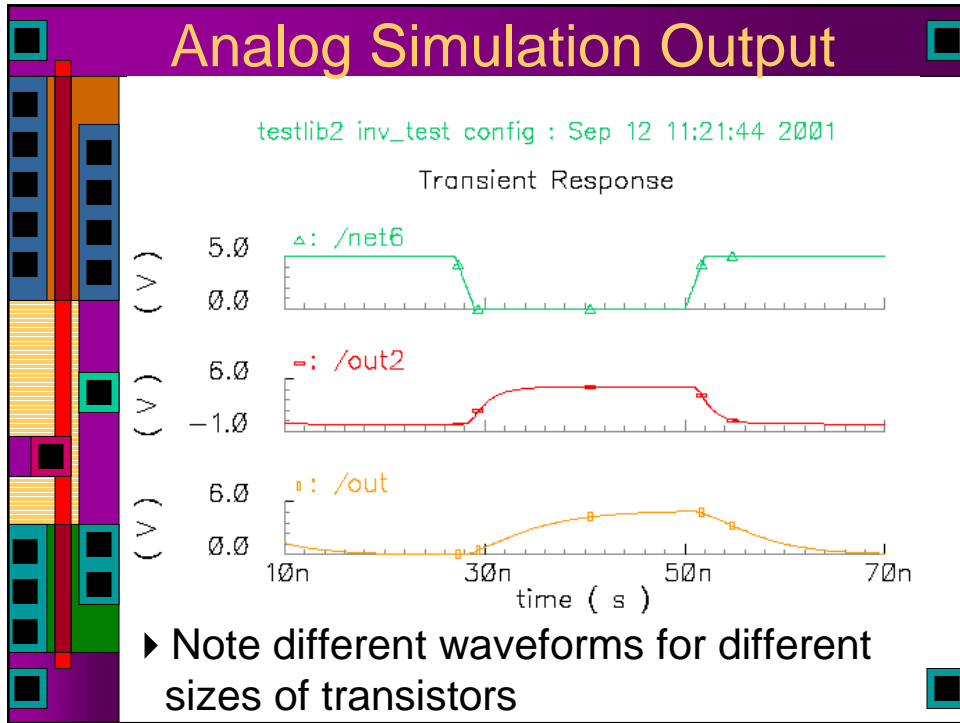


► From CAD #2

## Closeup of Inv-Test



► Note the sizes I used for this example...



## What About Gates in Series

- ▶ Basically we want every gate to have the delay of a “standard inverter”
  - ▶ Standard inverter starts with 2/1 P/N ratio
- ▶ Gates in series? Sum the conductance to get the series conductance
- ▶  $\beta_{n\text{-eff}} = 1/(1/\beta_1 + 1/\beta_2 + 1/\beta_3)$ 
  - ▶  $\beta_{n\text{-eff}} = \beta_n/3$
- ▶ Effect is like increasing L by 3
  - ▶ Compensate by increasing W by 3

