

Lecture 25 - CMOS Scaling Rules - Outline

- **Announcements**

Handouts - Lecture Outline and Summary; on web: CMOS scaling

Course evaluation - at the end of lecture today (approx. 11:45)

Final - Monday, Dec. 15, 9:00 am - noon, duPont Gymnasium

- **Review - Intrinsic high frequency limits for transistors: ω_T**

Short-circuit current gain: best can do from CE or CS

Unity gain frequency: BJT: $\omega_T = 2D_e/w_B^2 = 2\mu_e V_{th}/w_B^2$ (used Einstein rel.)

MOSFET: $\omega_T = 3\mu_e(V_{GS}-V_T)/2L^2$ (this needs discussion)

Revisiting the quasi-static assumption (^ same form; ain't it neat!!)

- **CMOS gate delay and power**

Review of Lecture 16 results: Gate Delay = $12 \text{ n} L_{\min}^2 V_{DD} / \mu_n (V_{DD} - V_T)^2$

$P_{\text{ave}} @ \text{max. } f \propto C_L V_{DD}^2 / \text{GD} = K_n V_{DD} (V_{DD} - V_T)^2 / 4$

Power density issue: have to add this as well

- **CMOS scaling rules**

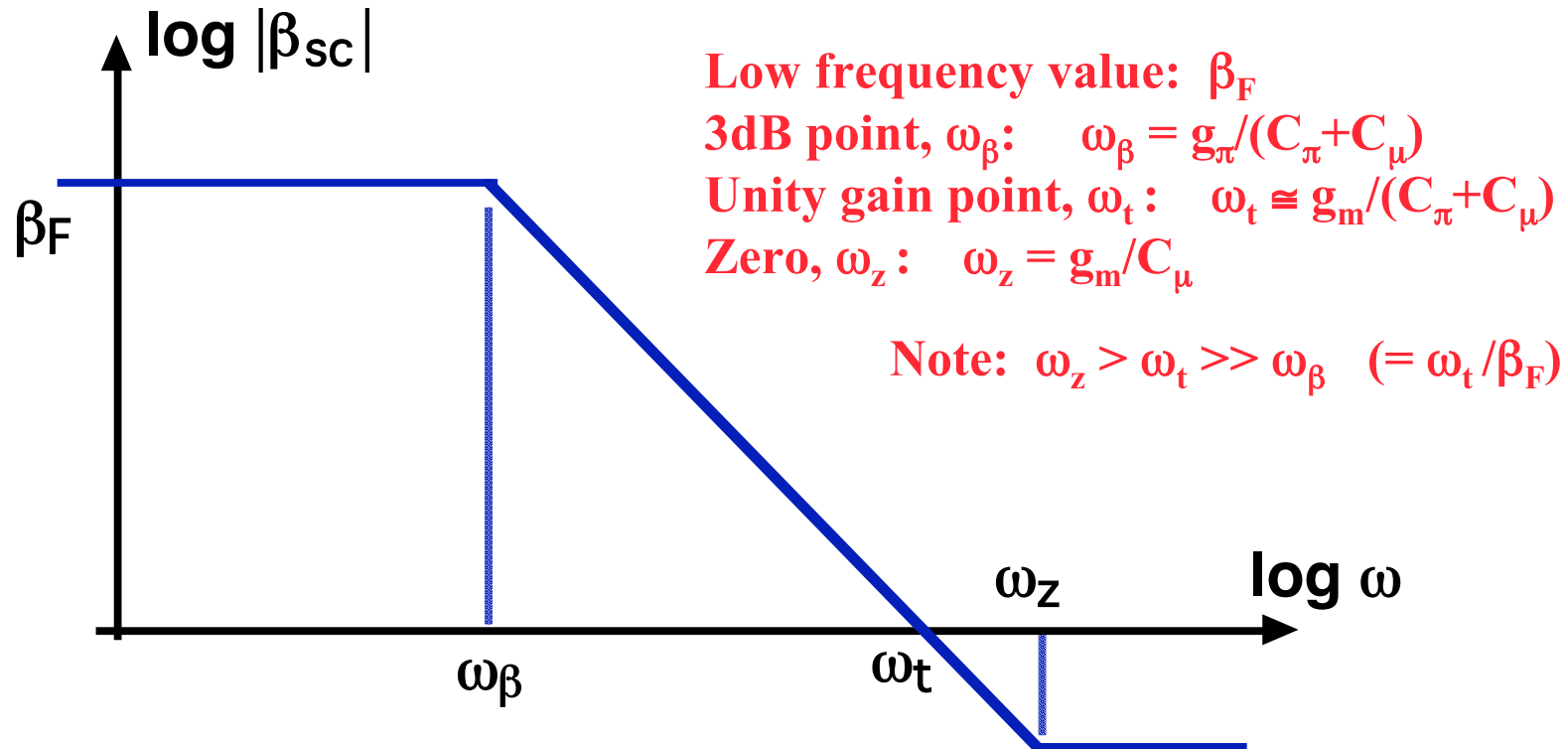
The issues and challenges

Approaches: Dimension scaling

Scaling voltages as well

Summary of rules

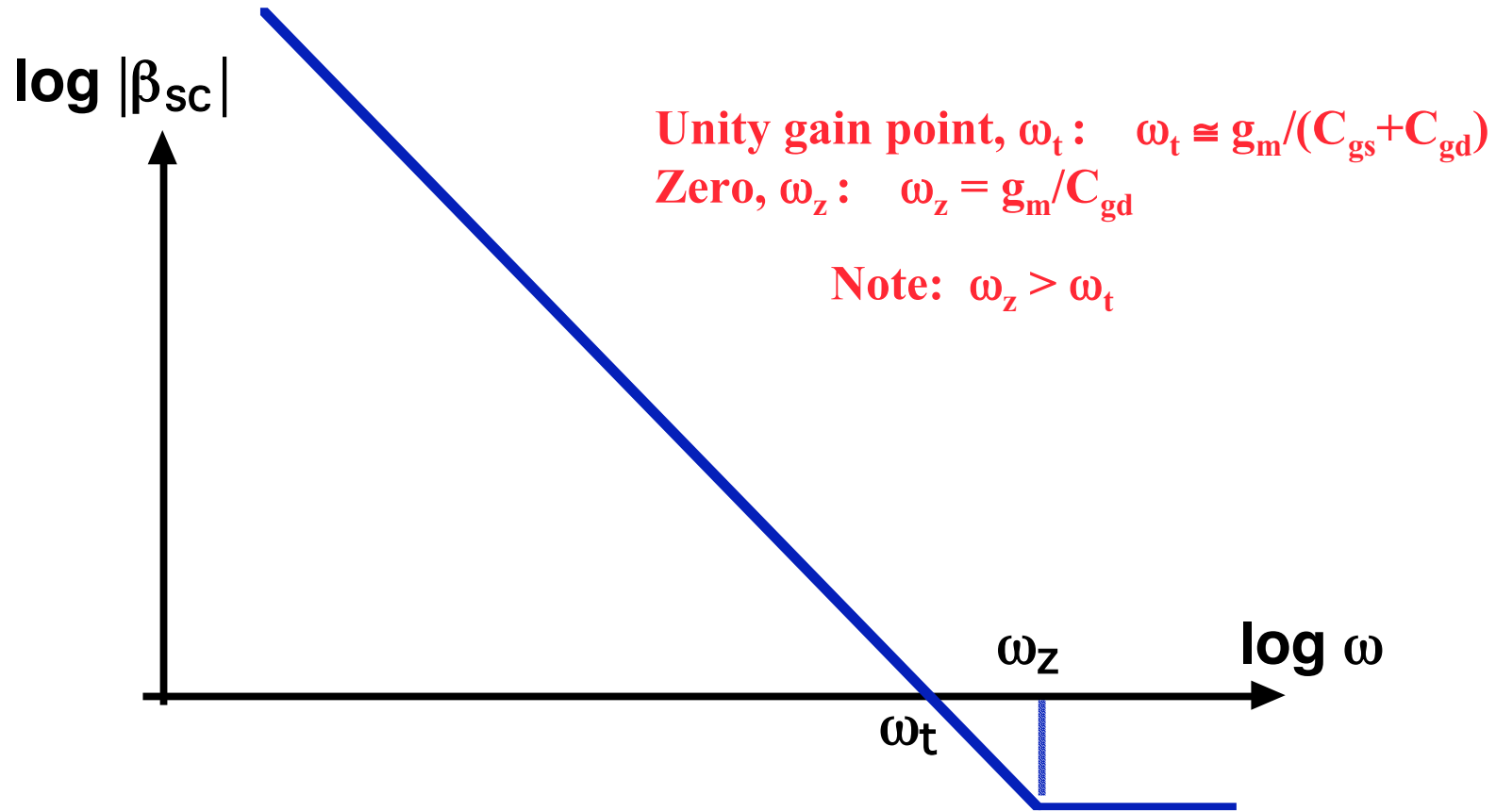
BJT short-circuit current gain, $\beta_{sc}(j\omega)$



$$\omega_t \approx \frac{g_m}{[C_\pi + C_\mu]} = \frac{g_m}{\{g_m \tau_{trB} + C_{eb,dp} + C_{cb,dp}\}} = \frac{1}{\{\tau_{trB} + kT[C_{eb,dp} + C_{cb,dp}] / qI_C\}}$$

In the limit of large I_C : $\omega_t(\text{BJT}) \approx \frac{1}{\tau_{trB}} = \frac{2D_{eB}}{W_B^2} = \frac{2\mu_e V_{thermal}}{W_B^2}$

MOSFET short-circuit current gain, $\beta_{sc}(j\omega)$



$$\omega_t(\text{MOSFET}) \approx \frac{g_m}{[C_{gs} + C_{gd}]} \approx \frac{g_m}{C_{gs}} = \frac{\frac{W}{L} \mu_e C_{ox}^* [V_{GS} - V_T]}{\frac{2}{3} W L C_{ox}^*} = \frac{3}{2} \frac{\mu_e [V_{GS} - V_T]}{L^2} \approx \frac{1}{\tau_{trCh}}$$

Looking more at ω_T for BJTs and MOSFETs:

For a MOSFET we have

$$\omega_t(\text{MOSFET}) \approx \frac{3}{2} \frac{\mu_e [V_{GS} - V_T]}{L^2}$$

**The average E-field in the channel, E_y , is
So we can also write ω_T as**

$$E_y \approx \frac{[V_{GS} - V_T]}{L}$$

$$\omega_t(\text{MOSFET}) \approx \frac{3}{2} \mu_e \frac{[V_{GS} - V_T]}{L} \frac{1}{L} \approx \frac{\bar{s}_y}{L} \approx \frac{1}{\tau_{trChannel}}$$

This is identical to the form we have for ω_T in a BJT

$$\omega_t(\text{BJT}) \approx \frac{1}{\tau_{trB}} = \frac{2D_{eB}}{w_B^2} = \frac{2\mu_e V_{thermal}}{w_B^2}$$

What happens when we have velocity saturation?

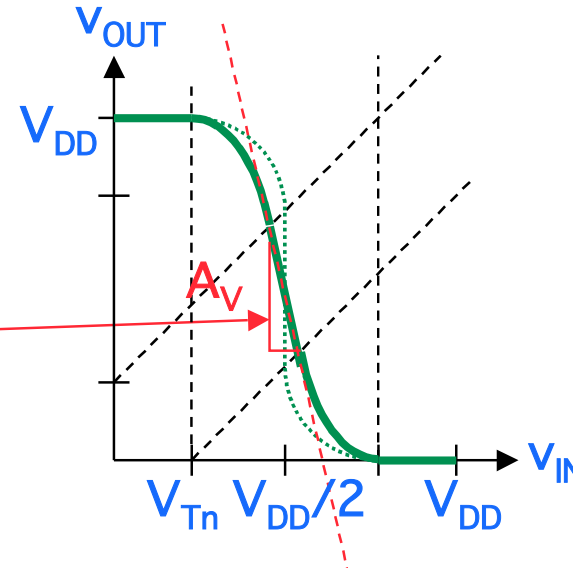
$$\omega_t(\text{MOSFET}) \approx \frac{s_{sat}}{L}, \quad \omega_t(\text{BJT}) \approx \frac{s_{sat}}{w_B}$$

ω_T still decreases with L and w_B , but not as quickly!

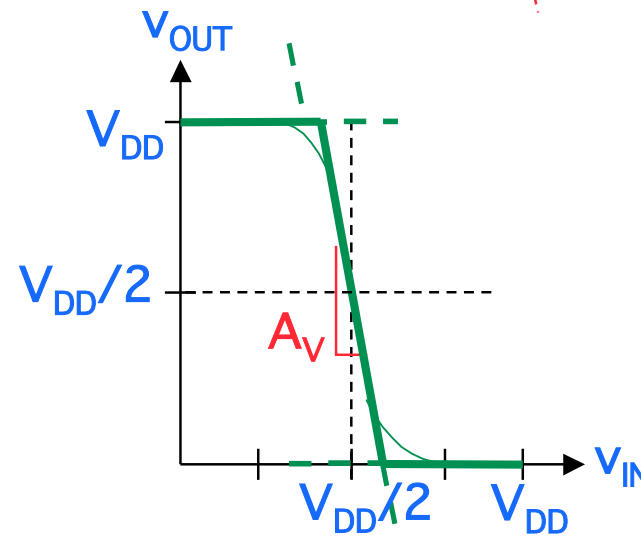
CMOS: transfer characteristic calculation, cont.

Returning to the transfer characteristic, we see that the slope in Region III is not infinite, but is instead:

$$A_v \equiv \left. \frac{\partial v_{OUT}}{\partial v_{IN}} \right|_Q = - \frac{2\sqrt{2K_n}}{[\lambda_n + \lambda_p]\sqrt{I_{Dn}}}$$



Final comment: A quick and dirty way to approximate the transfer curve of a CMOS gate is to simply draw the three straight line portions in Regions I, III, and V:



CMOS: switching speed; minimum cycle time

The load capacitance, C_L

- Assume to be linear
- Is proportional to MOSFET gate area
- In channel: $\mu_e = 2 \mu_h$ so to have $K_n = K_p$ we must have $W_p/L_p = 2W_n/L_n$
Typically $L_n = L_p = L_{\min}$, and $W_n = W_{\min}$, so we also have $W_p = 2 W_{\min}$.

$$C_L \approx n[W_n L_n + W_p L_p] C_{ox}^* = n[W_{\min} L_{\min} + 2W_{\min} L_{\min}] C_{ox}^* = 3nW_{\min} L_{\min} C_{ox}^*$$

Charging cycle

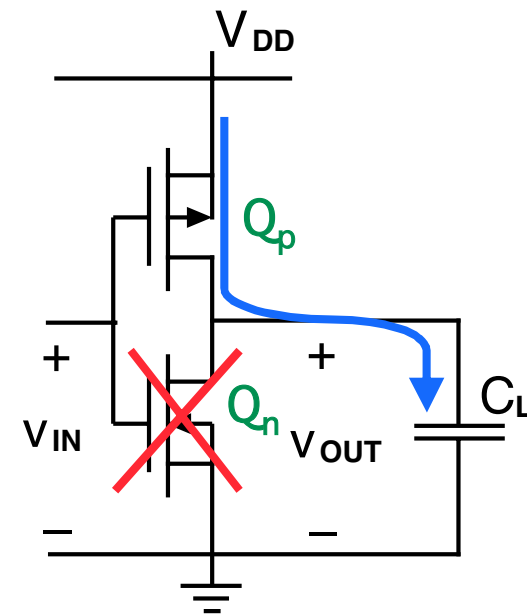
v_{IN} : Hi to Lo; Q_n off, Q_p on; v_{OUT} : Lo to Hi

- Assume charged by constant $i_{D,sat}$

$$i_{Charge} = i_{Dp} \approx \frac{K_p}{2} [V_{DD} - |V_{Tp}|]^2 = \frac{K_n}{2} [V_{DD} - V_{Tn}]^2$$

$$q_{Charge} = C_L V_{DD}$$

$$\begin{aligned} \tau_{Charge} &= \frac{q_{Charge}}{i_{Charge}} = \frac{2C_L V_{DD}}{K_n [V_{DD} - V_{Tn}]^2} \\ &= \frac{6nW_{\min} L_{\min} C_{ox}^* V_{DD}}{\frac{W_{\min}}{L_{\min}} \mu_e C_{ox}^* [V_{DD} - V_{Tn}]^2} = \frac{6nL_{\min}^2 V_{DD}}{\mu_e [V_{DD} - V_{Tn}]^2} \end{aligned}$$



CMOS: switching speed; minimum cycle time, cont.

Discharging cycle

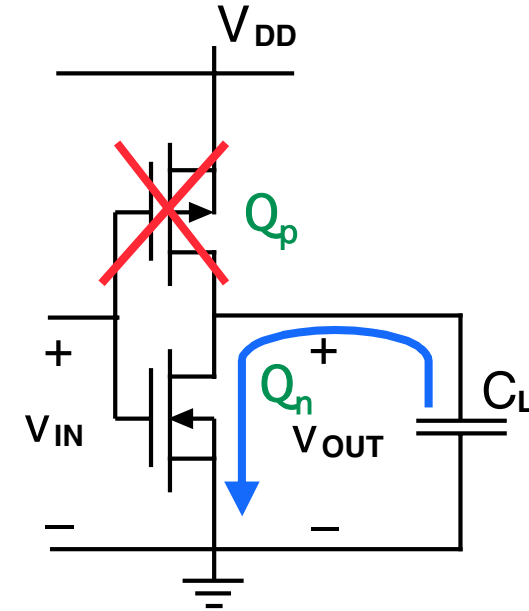
v_{IN} : Lo to Hi; Q_n on, Q_p off; v_{OUT} : Hi to Lo

- Assume discharged by constant $i_{D,sat}$

$$i_{Discharge} = i_{Dn} \approx \frac{K_n}{2} [V_{DD} - V_{Tn}]^2$$

$$q_{Discharge} = C_L V_{DD}$$

$$\begin{aligned} \tau_{Discharge} &= \frac{q_{Discharge}}{i_{Discharge}} = \frac{2C_L V_{DD}}{K_n [V_{DD} - V_{Tn}]^2} \\ &= \frac{6n W_{min} L_{min} C_{ox}^* V_{DD}}{\frac{W_{min}}{L_{min}} \mu_e C_{ox}^* [V_{DD} - V_{Tn}]^2} = \frac{6n L_{min}^2 V_{DD}}{\mu_e [V_{DD} - V_{Tn}]^2} \end{aligned}$$



Minimum cycle time

v_{IN} : Lo to Hi to Lo; v_{OUT} : Hi to Lo to Hi

$$\tau_{Min.Cycle} = \tau_{Charge} + \tau_{Discharge} \approx \frac{12n L_{min}^2 V_{DD}}{\mu_e [V_{DD} - V_{Tn}]^2}$$

CMOS: power dissipation - total and per unit area

Average power dissipation

All dynamic

$$P_{ave} = C_L V_{DD}^2 f = 3nW_{\min}L_{\min}C_{ox}^*V_{DD}^2 f$$

Average power at maximum data rate

Maximum f will be $1/\tau_{\text{Min Cycle}}$

$$\begin{aligned} P_{ave@Max.f} &= 3nW_{\min}L_{\min}C_{ox}^*V_{DD}^2 \frac{\mu_e [V_{DD} - V_{Tn}]^2}{12nL_{\min}^2 V_{DD}} \\ &= \frac{1}{4} \frac{W_{\min}}{L_{\min}} \mu_e C_{ox}^* V_{DD} [V_{DD} - V_{Tn}]^2 = \frac{1}{4} K_n V_{DD} [V_{DD} - V_{Tn}]^2 \end{aligned}$$

Average power density at maximum data rate

Assume that the area per inverter will be proportional to $W_{\min}L_{\min}$

$$PD_{ave@Max.f} = \frac{P_{ave@Max.f}}{\text{Inverter area}} \propto \frac{P_{ave@Max.f}}{W_{\min}L_{\min}} = \frac{\mu_e C_{ox}^* V_{DD} [V_{DD} - V_{Tn}]^2}{4 L_{\min}^2}$$

CMOS: design for high speed

Maximum data rate

Proportional to $1/\tau_{\text{Min Cycle}}$

$$f_{\text{max}} \propto 1/\tau_{\text{Min.Cycle}} = \frac{\mu_e [V_{DD} - V_{Tn}]^2}{12 n L_{\text{min}}^2 V_{DD}}$$

Teaches us to make L_{min} small and/or V_{DD} large

Note: As we reduce L_{min} we must also reduce t_{ox} , but t_{ox} doesn't enter directly in f_{max} so it doesn't impact us here.

Average power density at maximum data rate

Assumes area per inverter is proportional to $W_{\text{min}} L_{\text{min}}$

$$P_{\text{ave}} = C_L V_{DD}^2 f \quad \text{and} \quad f_{\text{max}} \propto 1/\tau_{\text{Min.Cycle}} = \frac{\mu_e [V_{DD} - V_{Tn}]^2}{12 n L_{\text{min}}^2 V_{DD}}$$

$$PD_{\text{ave@Max.f}} = \frac{P_{\text{ave@Max.f}}}{\text{Inverter area}} \propto \frac{C_L V_{DD}^2 f_{\text{max}}}{W_{\text{min}} L_{\text{min}}} \propto \frac{\mu_e \epsilon_{\text{ox}} V_{DD} [V_{DD} - V_{Tn}]^2}{t_{\text{ox}} L_{\text{min}}^2}$$

Teaches us PD increases very quickly as we reduce L_{min} (and t_{ox}) unless we also reduce V_{DD} (which reduces f_{max}).

How do we make f_{max} larger without melting the silicon?
Through CMOS scaling rules - the topic of today's lecture.

- **Summary:**

Transfer characteristic: symmetric

$$V_{LO} = 0, \quad V_{HI} = V_{DD}, \quad I_{ON} = 0, \quad I_{OFF} = 0$$

$$N_{ML} = N_{MH} \Rightarrow K_n = K_p \text{ and } |V_{Tp}| = V_{Tn}$$

$$L_n = L_p = L_{\min}, \quad W_p = (\mu_n / \mu_p) W_n, \quad W_n = W_{\min}$$

Gate delay expressions

$$\tau_{Ch\ arg\ e} = \tau_{Disch\ arg\ e} \approx \frac{2V_{DD}C_L}{K_n [V_{DD} - V_{Tn}]^2}$$

$$C_L = n(W_n L_n + W_p L_p) C_{ox}^* = 3nW_{\min} L_{\min} C_{ox}^* \quad (\text{Assumes } \mu_n = 2\mu_p)$$

$$\tau_{Min.Cycle} = \tau_{Ch\ arg\ e} + \tau_{Disch\ arg\ e} \approx \frac{12nL_{\min}^2 V_{DD}}{\mu_e [V_{DD} - V_{Tn}]^2}$$

Average power at f_{\max} , and Power Density (dissipation per unit area)

$$P_{ave@Max.f} = C_L V_{DD}^2 f_{\max} \propto C_L V_{DD}^2 / \tau_{Min.Cycle} = \frac{\mu_e W_{\min} C_{ox}^* V_{DD} [V_{DD} - V_{Tn}]^2}{L_{\min}}$$

$$PD_{ave@Max.f} = \frac{P_{ave@Max.f}}{\text{Inverter area}} \propto \frac{P_{ave@Max.f}}{W_{\min} L_{\min}} = \frac{\mu_e \epsilon_{ox} V_{DD} [V_{DD} - V_{Tn}]^2}{t_{ox} L_{\min}^2}$$

Scaling Rules - making CMOS faster without melting Si

- **General idea:**

Reduce dimensions and/or voltages by factor $1/s$: $s > 1$
Evaluate impact on speed, power, power density

- **Scaling dimensions alone:**

$L_{\min} \rightarrow L_{\min}/s$ $w \rightarrow w/s$ $t_{\text{ox}} \rightarrow t_{\text{ox}}/s$
this yields
 $K \rightarrow sK$ $C_{\text{ox}}^* \rightarrow sC_{\text{ox}}^*$ (Also stay in grad. channel regime)
and ultimately
 $\tau \rightarrow \tau/s^2$ $P_{\text{ave}} \rightarrow s P_{\text{ave}}$ $P_{\text{density}} \rightarrow s^3 P_{\text{density}}$!!!

Scaling dimensions alone can yield melted silicon!!

- **Scaling dimensions and voltages in concert:**

still have **Add:** $V_{\text{DD}} \rightarrow V_{\text{DD}}/s$ $V_{\text{T}} \rightarrow V_{\text{T}}/s$
 $K \rightarrow sK$ $C_{\text{ox}}^* \rightarrow sC_{\text{ox}}^*$
but now
 $\tau \rightarrow \tau/s$ $P_{\text{ave}} \rightarrow P_{\text{ave}}/s^2$ $P_{\text{density}} \rightarrow P_{\text{density}}$

When scale both get: faster, lower power, same power density!!
Note: scaling the voltages is not as easy as scaling the dimensions.

Lecture 25 - CMOS Scaling Rules - Summary

- **CMOS gate delay and power**

Three key performance metrics: (We want to make them all smaller)

$$\tau_{Min.Cycle} = 12 n L_{min}^2 V_{DD} / \mu_e [V_{DD} - V_{Tn}]^2$$

$$P_{ave@Max.f} = C_L V_{DD}^2 f_{max} \propto \mu_e C_{ox}^* V_{DD} [V_{DD} - V_{Tn}]^2 [W_{min} / L_{min}]$$

$$PD_{ave@Max.f} = \frac{P_{ave@Max.f}}{\text{Inverter area}} \propto \frac{P_{ave@Max.f}}{W_{min} L_{min}} = \frac{\mu_e \epsilon_{ox} V_{DD} [V_{DD} - V_{Tn}]^2}{t_{ox} L_{min}^2}$$

- **CMOS scaling rules**

Summary of rules: best to reduce all dimensions and all voltages by 1/s

Scaling as:

$$\begin{aligned} L_{min} &\rightarrow L_{min} / s \\ W_{min} &\rightarrow W_{min} / s \\ t_{ox} &\rightarrow t_{ox} / s \\ V_{DD} &\rightarrow V_{DD} / s \\ V_T &\rightarrow V_T / s \end{aligned}$$

Results in:

$$\begin{aligned} K &\rightarrow sK \\ C_{ox}^* &\rightarrow sC_{ox}^* \\ \tau_{Min.cycle} &\rightarrow \tau_{Min.cycle} / s \\ P_{ave} &\rightarrow P_{ave} / s^2 \\ PD_{ave} &\rightarrow PD_{ave} \end{aligned}$$