

**Lecture 32 - The "Short"**  
**Metal-Oxide-Semiconductor Field-Effect**  
**Transistor** (*cont.*)

November 20, 2002

**Contents:**

1. MOSFET scaling

**Reading assignment:**

P. K. Ko, "*Approaches to Scaling.*"

## Key questions

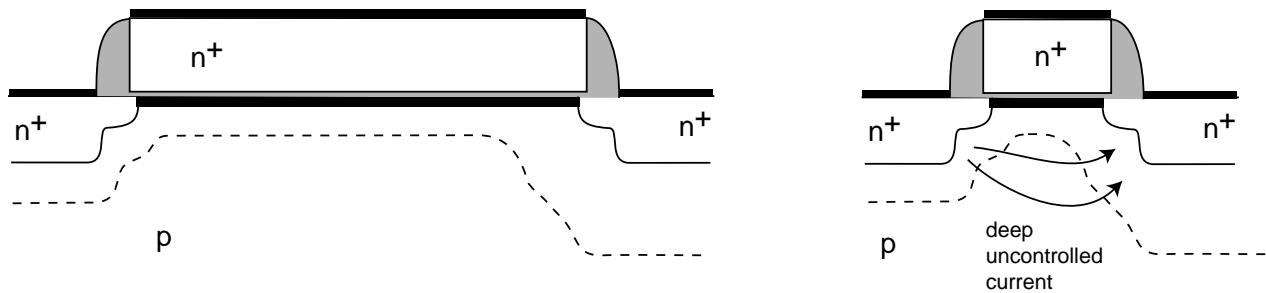
- What happens if a MOSFET gate length is simply shrunk in size without changing anything else?
- How should the MOSFET design change as it shrinks down in size?

## 1. MOSFET scaling

Several driving forces for scaling down size of MOSFET:

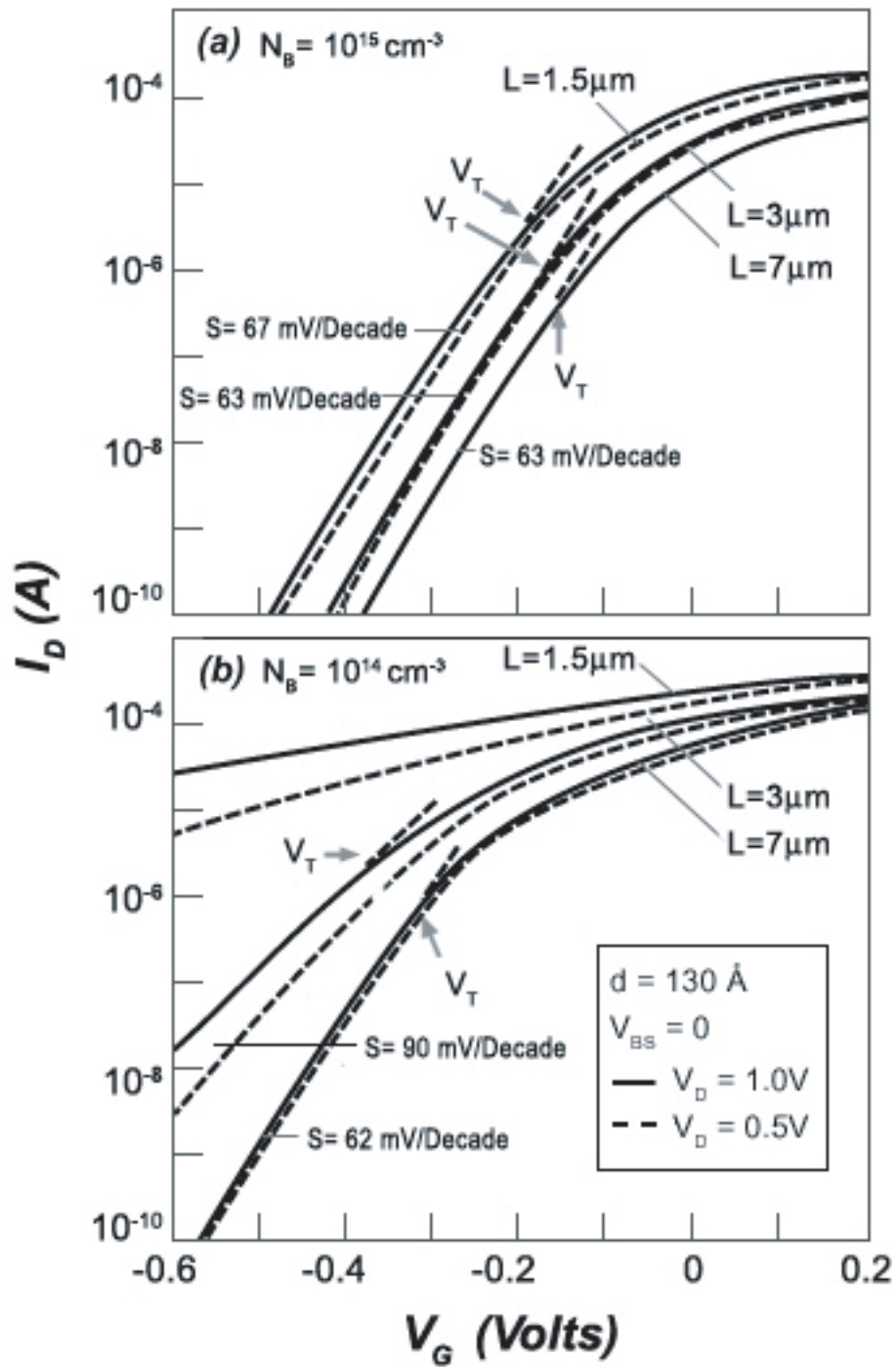
- higher density circuits: SSI, MSI, LSI, VLSI, ULSI, RLSI, ...
- higher performance:  $L \downarrow \Rightarrow I_D \uparrow \Rightarrow \tau_{switch} \downarrow$
- lower power consumption:  $L \downarrow \Rightarrow V_{DD} \downarrow$

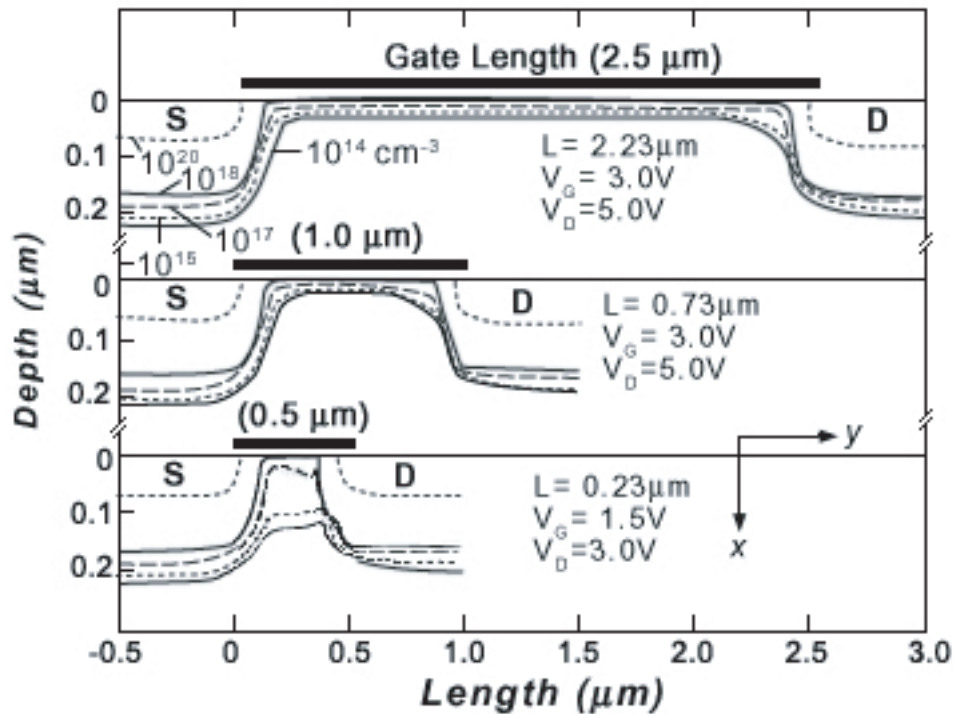
Simple  $L$  scaling compromises *electrostatic integrity* and produces *punchthrough* (extreme case of short-channel effects):



To avoid punchthrough:

- $N_A \uparrow \Rightarrow V_{th} \uparrow \Rightarrow I_D \downarrow$
- $V_{DD} \downarrow \Rightarrow I_D \downarrow$
- $x_{ox} \downarrow \Rightarrow V_{th} \downarrow \Rightarrow I_D \uparrow$





Constant electron density contours for 3 MOSFETs with channel lengths 2.23, 0.73, and 0.23  $\mu\text{m}$ .

Adapted from S. M. Sze,  
*Physics of Semiconductor Devices, 2nd ed., Wiley, 1981 (481).*

Need smart way of scaling:

- constant field scaling
- constant voltage scaling
- generalized scaling

## □ Constant field scaling

Scale keeping vertical and horizontal electric fields constant.

Define: *scaling factor*  $S > 1$

parameter	scaling factor
device dimensions ( $L, W, x_{ox}$ )	$1/S$
doping level ( $N_A$ )	$S$
supply voltage ( $V_{DD}$ )	$1/S$

Consequences (use simple long-channel theory):

- gate capacitance:

$$C'_{gs} = C'_{ox} L' W' = S C_{ox} \frac{L W}{S S} = \frac{C_{gs}}{S} \downarrow$$

- threshold voltage:

$$V'_{th} = V_{FB} + \phi_{sth} + \gamma \sqrt{\phi_{sth}} \simeq \frac{1}{C'_{ox}} \sqrt{2\epsilon_s q N'_A \phi_{sth}} \sim \frac{V_{th}}{\sqrt{S}} \downarrow$$

- drive current:

$$I'_D = \frac{W'}{2L'} \mu_e C'_{ox} (V'_{DD} - V'_{th})^2 = \frac{\frac{W}{S}}{2\frac{L}{S}} \mu_e S C_{ox} \left( \frac{V_{DD}}{S} - \frac{V_{th}}{\sqrt{S}} \right)^2 = \frac{I_D}{S} \downarrow$$

- gate delay:

$$\tau' = \frac{C'_{gs} V'_{DD}}{I'_D} = \frac{\frac{C_{gs}}{S} \frac{V_{DD}}{S}}{\frac{I_D}{S}} = \frac{\tau}{S} \downarrow$$

- power dissipation:

$$I'_D V'_{DD} = \frac{I_D V_{DD}}{S} \frac{V_{DD}}{S} = \frac{I_D V_{DD}}{S^2} \downarrow\downarrow$$

- power density:

$$\frac{I'_D V'_{DD}}{L'W'} = \frac{\frac{I_D V_{DD}}{S} \frac{V_{DD}}{S}}{\frac{LW}{S} \frac{S}{S}} = \frac{I_D V_{DD}}{LW} \text{ unchanged}$$

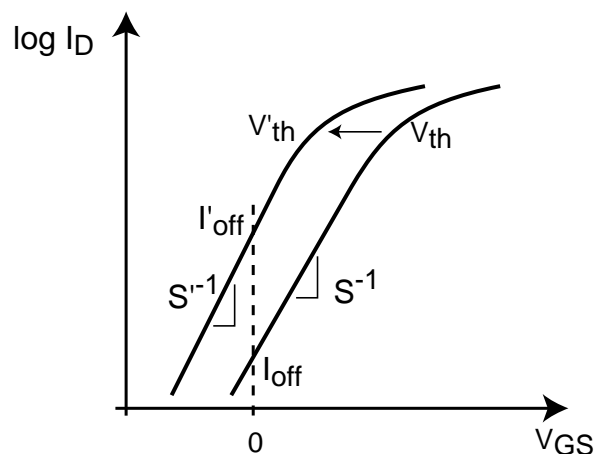
- power-delay product (or *switching energy*):

$$C'_{gs} V'_{DD}{}^2 = \frac{C_{gs}}{S} \left(\frac{V_{DD}}{S}\right)^2 = \frac{C_{gs} V_{DD}^2}{S^3} \downarrow\downarrow\downarrow$$

- inverse subthreshold slope:

$$n' = 1 + \frac{C'_{sth}}{C'_{ox}} = 1 + \frac{\sqrt{S} C_{sth}}{S C_{ox}} = 1 + \frac{C_{sth}}{\sqrt{S} C_{ox}} \downarrow$$

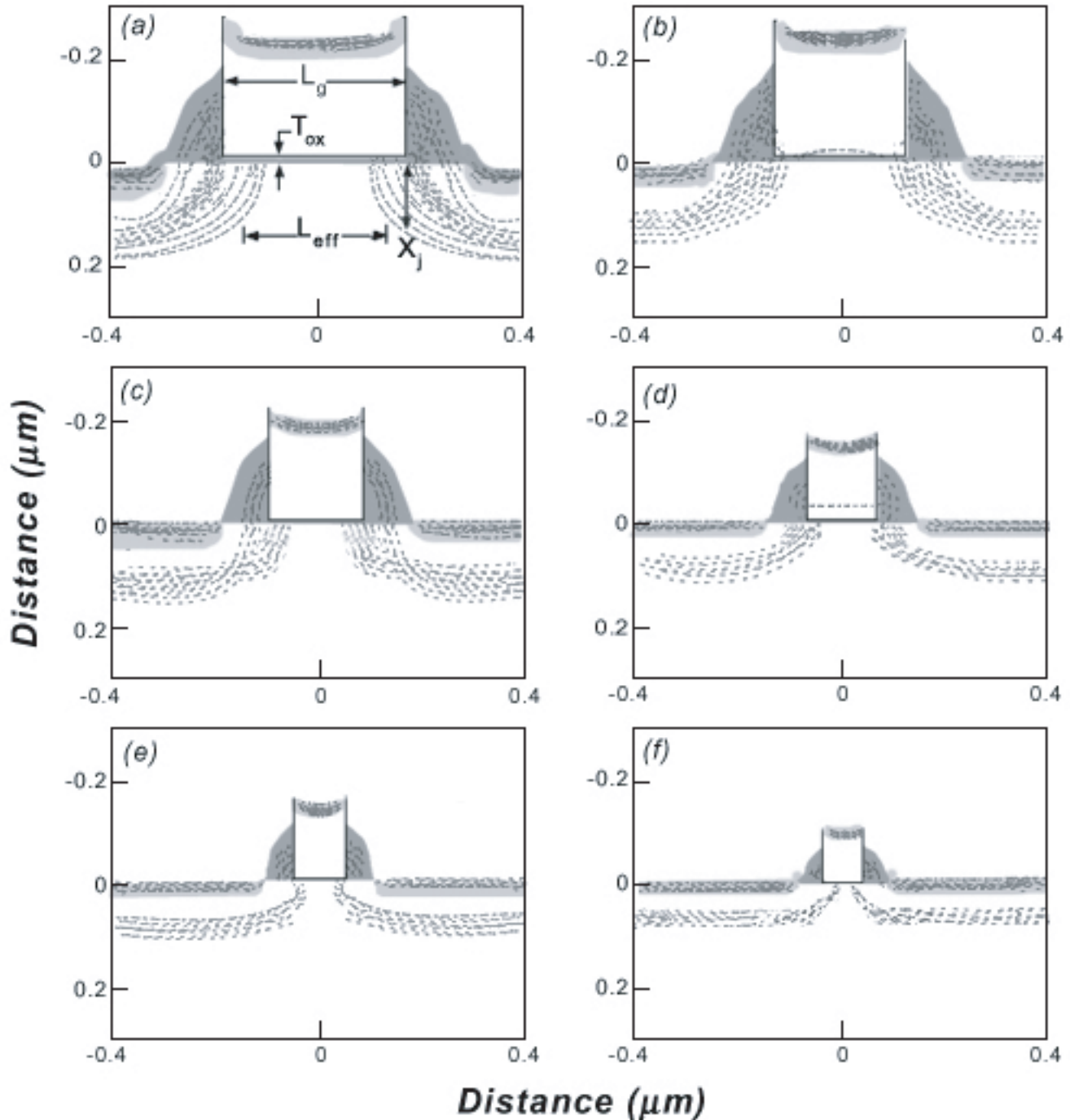
but since  $V_{th} \downarrow$ ,  $I_{off} \uparrow\uparrow$ .



Two key problems with constant field scaling:

- system designers don't want to scale  $V_{DD}$
- $I_{off} \uparrow\uparrow \Rightarrow$  more static power

- More rigorous study of constant field scaling using 2D simulations  
[P. Vande Voorde, HP Journal, 1997]

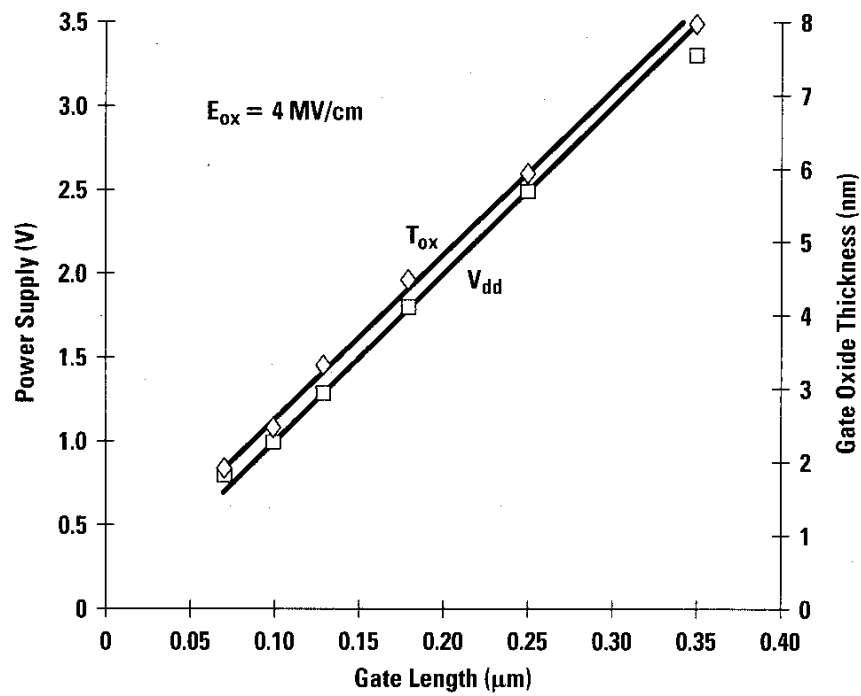


Simulated device structures. Dark shading is oxide. Lighter shading is silicide.  
Dashed lines are doping contours.

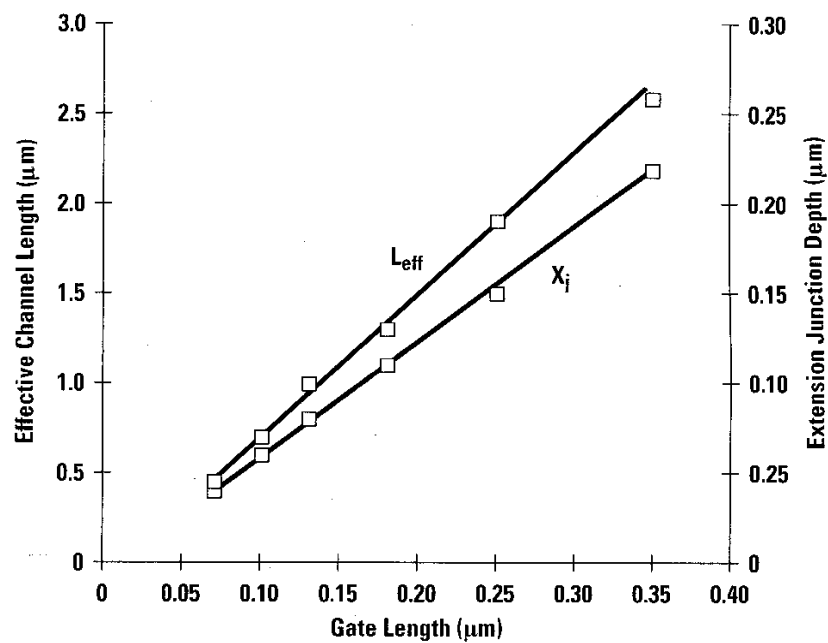
(a)  $L_g = 0.35 \mu\text{m}$ ,  $T_{ox} = 8.0 \mu\text{m}$ . (b)  $L_g = 0.25 \mu\text{m}$ ,  $T_{ox} = 6.0 \mu\text{m}$ . (c)  $L_g = 0.18 \mu\text{m}$ ,  $T_{ox} = 4.5 \mu\text{m}$ .  
(d)  $L_g = 0.13 \mu\text{m}$ ,  $T_{ox} = 3.4 \mu\text{m}$ . (e)  $L_g = 0.10 \mu\text{m}$ ,  $T_{ox} = 2.5 \mu\text{m}$ . (f)  $L_g = 0.07 \mu\text{m}$ ,  $T_{ox} = 1.9 \mu\text{m}$ .

Adapted from P. Vande Voorde, *HP Journal*, 1997.





**Fig. 2.** Scaling of power supply voltage  $V_{dd}$  and oxide thickness  $T_{ox}$ .



**Fig. 3.** Scaling of effective channel length  $L_{eff}$  and extension junction depth  $X_j$ .

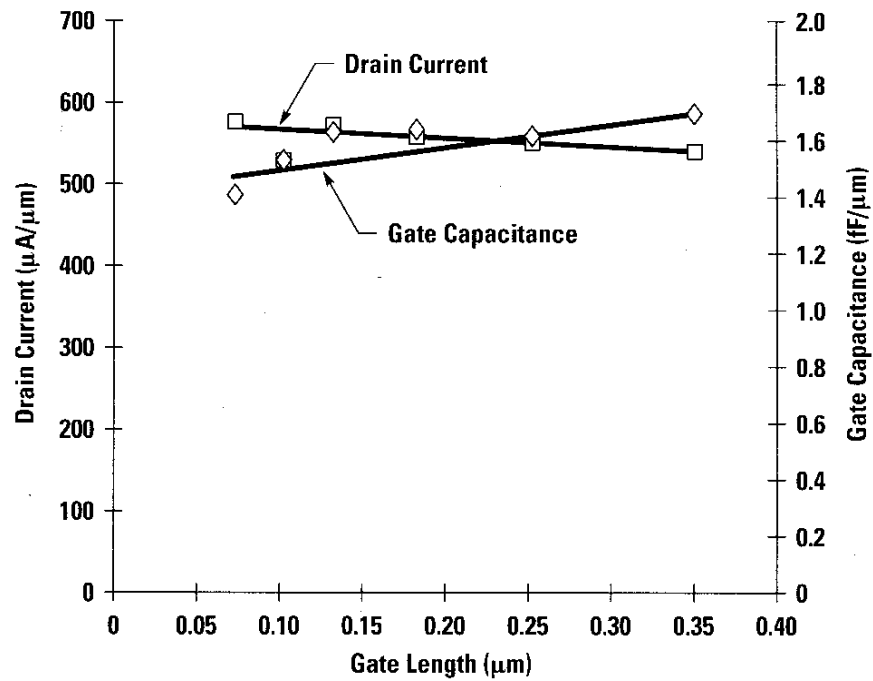


Fig. 5. Scaling of maximum drain current and total gate capacitance.

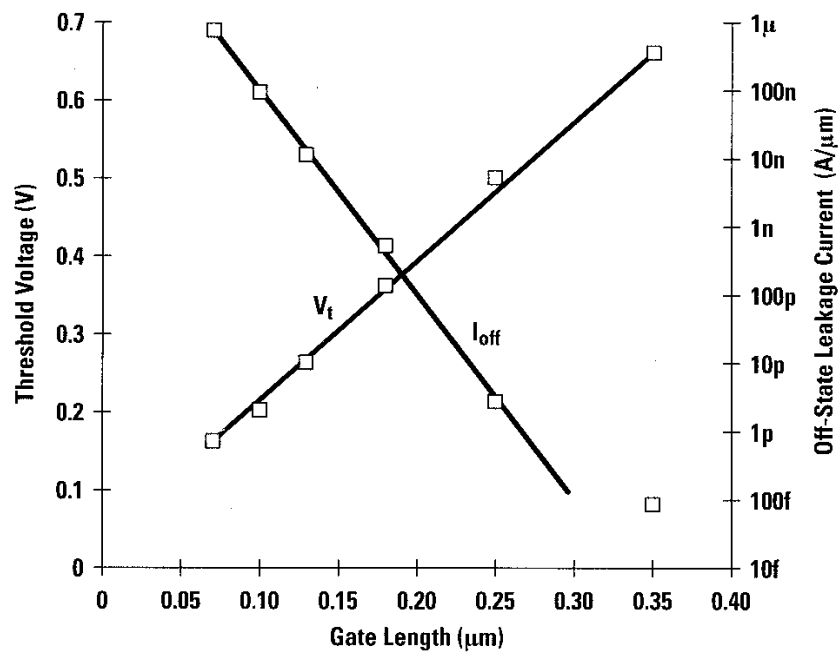


Fig. 4. Scaling of threshold voltage  $V_t$  and off-state leakage current  $I_{off}$ .

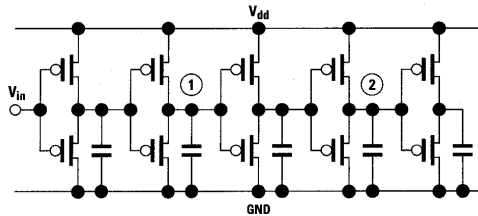


Fig. 8. Inverter chain.

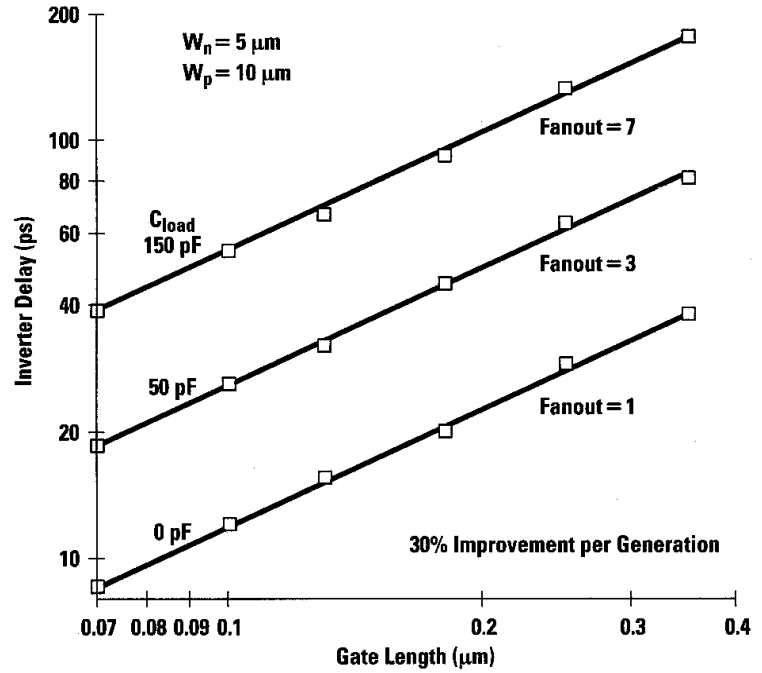


Fig. 10. Inverter delay versus gate length.

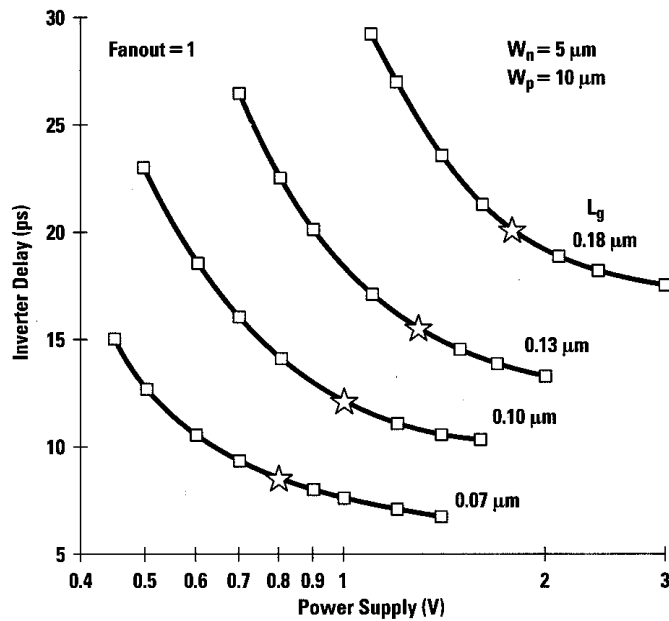


Fig. 11. Inverter delay versus power supply voltage. The stars show the expected operating points.  $W_n$  and  $W_p$  are the widths of the n-channel and p-channel transistors.

## □ Constant voltage scaling

Scale all device dimensions but do not scale  $V_{DD}$ .

<b>parameter</b>	<b>scaling factor</b>
device dimensions ( $L, W, x_{ox}$ )	$1/S$
doping level ( $N_A$ )	$S$
supply voltage ( $V_{DD}$ )	1

Consequences (using long-channel theory):

<b>figure of merit</b>	<b>scaling factor</b>
$C_{gs}$	$1/S$
$V_{th}$	$1/\sqrt{S}$
$I_D$	$S$
$\tau$	$1/S^2$
$I_D V_{DD}$	$S$
$I_D V_{DD} / LW$	$S^3$
$C_{gs} V_{DD}^2$	$1/S$

Features of constant voltage scaling:

- Performance  $\uparrow\uparrow$
- But:
  - It does not address  $I_{off}$  problem.
  - Electric field across oxide  $\uparrow$ :

$$\mathcal{E}_{ox} = \frac{V_{DD}}{x_{ox}} \propto S \uparrow$$

Reliability problems when  $\mathcal{E}_{ox} \simeq 4 \text{ MV/cm}$ .

- Electric field in semiconductor (at drain end of channel)  $\uparrow$ :

$$\mathcal{E}_m = \sqrt{\frac{V_{DS} - V_{DSsat}}{l^2} + \mathcal{E}_{sat}} \propto S \uparrow$$

with

$$l^2 = \frac{\epsilon_s}{\epsilon_{ox}} x_{ox} x_j \propto S^{-2}$$

Reliability problems when  $\mathcal{E}_m \simeq 0.5 \text{ MV/cm}$ .

- Power density  $\uparrow \Rightarrow$  system power  $\uparrow$

## □ Generalized scaling

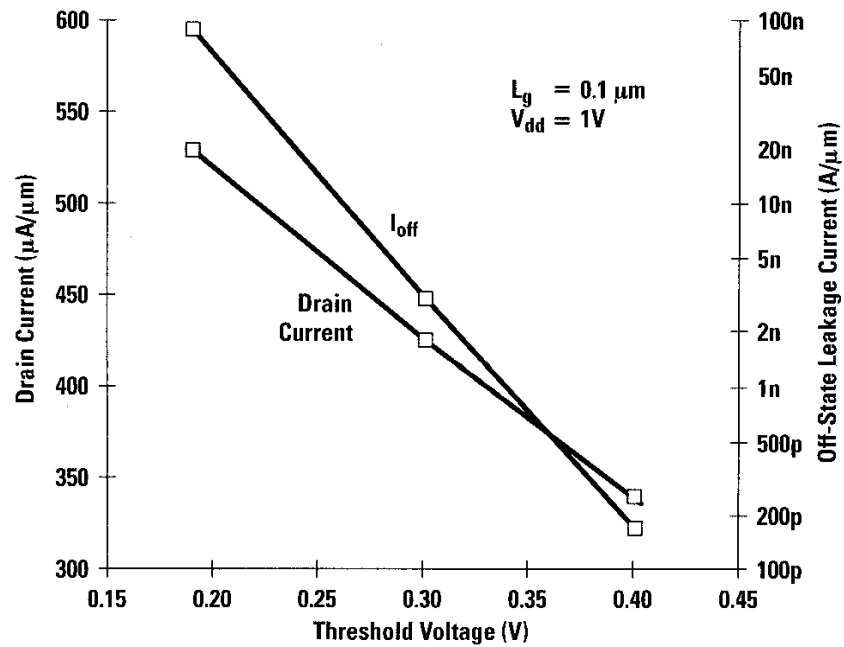
- scale oxide thickness more slowly than other device dimensions
- scale  $V_{DD}$  keeping  $\mathcal{E}_{ox}$  constant

parameter	scaling factor
$L, W$	$1/S$
$x_{ox}$	$1/R$
$N_A$	$S$
$V_{DD}$	$1/R$

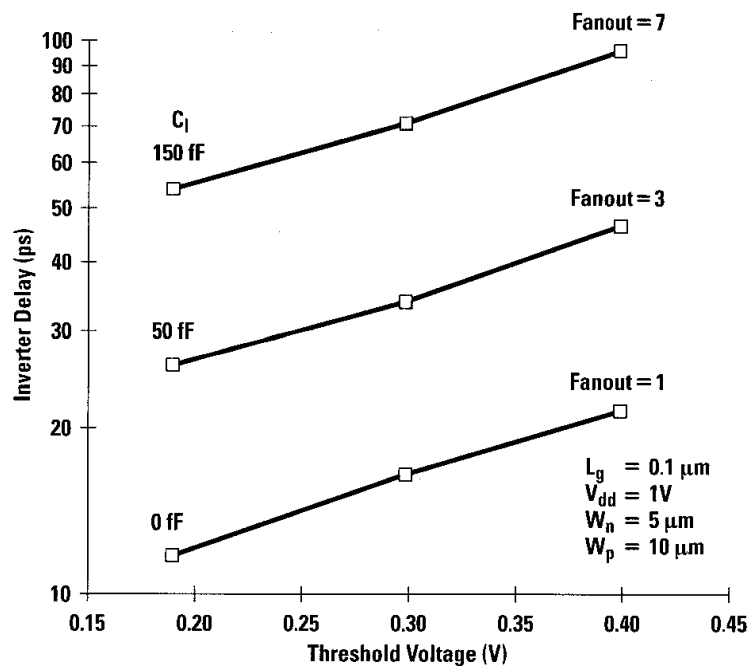
with  $1 < R < S$ .

In generalized scaling:

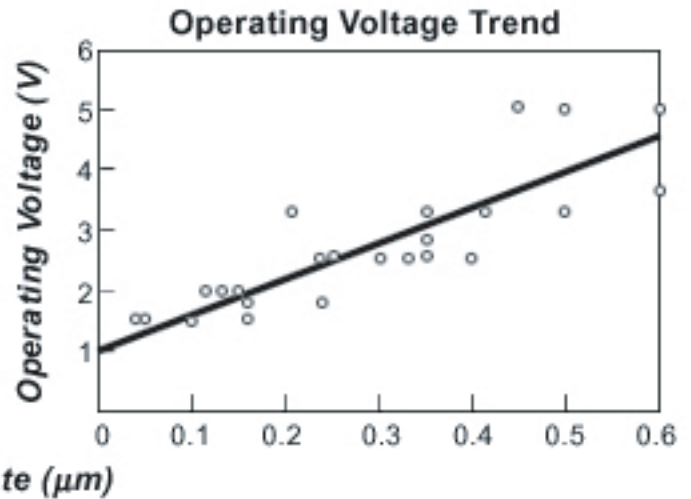
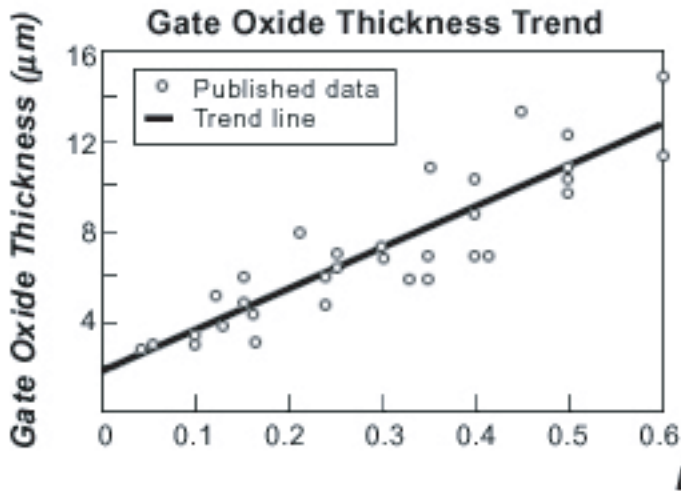
- $I_{off}$  problem alleviated by not scaling  $V_{th}$  so aggressively;  
*trade-off*: performance
- $V_{DD}$  scales;  
*trade-off*: performance



**Fig. 12.** Drain current and off-state leakage current  $I_{\text{off}}$  versus threshold voltage  $V_t$  for the 0.1- $\mu\text{m}$  generation.

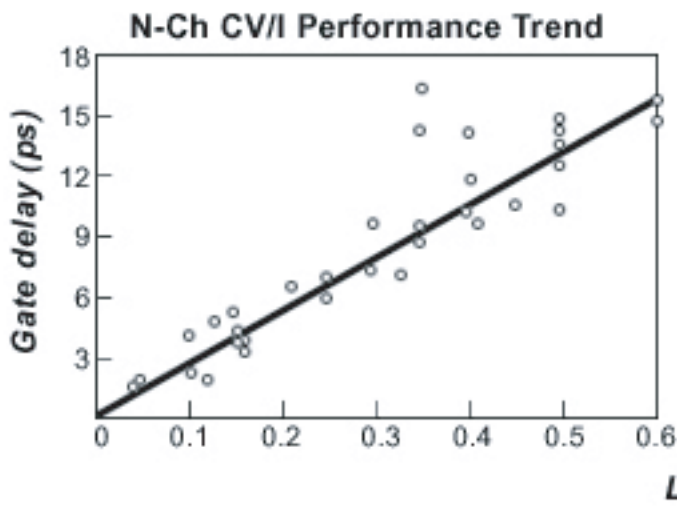


**Fig. 13.** Inverter delay versus threshold voltage  $V_t$  for the 0.1- $\mu\text{m}$  generation.

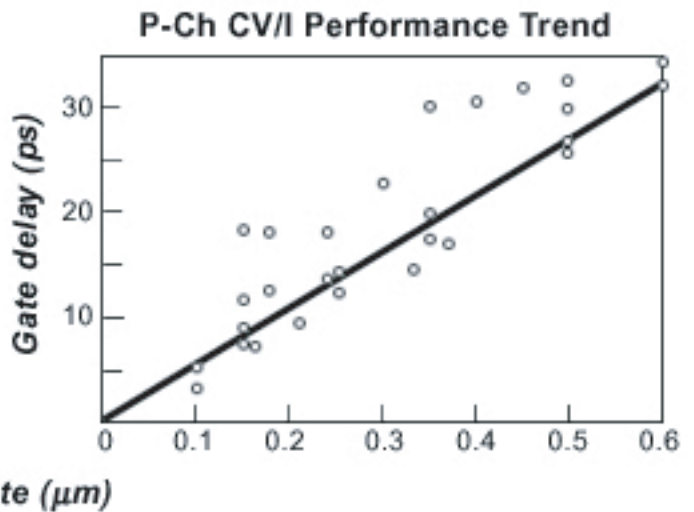


Oxide thickness will be in the range of 6-10 nm for 0.35  $\mu\text{m}$  generation technologies and will scale to less than 4 nm for 0.10  $\mu\text{m}$  generation technologies.

The 0.35  $\mu\text{m}$  generation marks a transition point between 3.3 V and 2.5 V operating voltage.



N-channel transistor performance trends based on data derived from papers published over the last few years, using the CV/I metric.



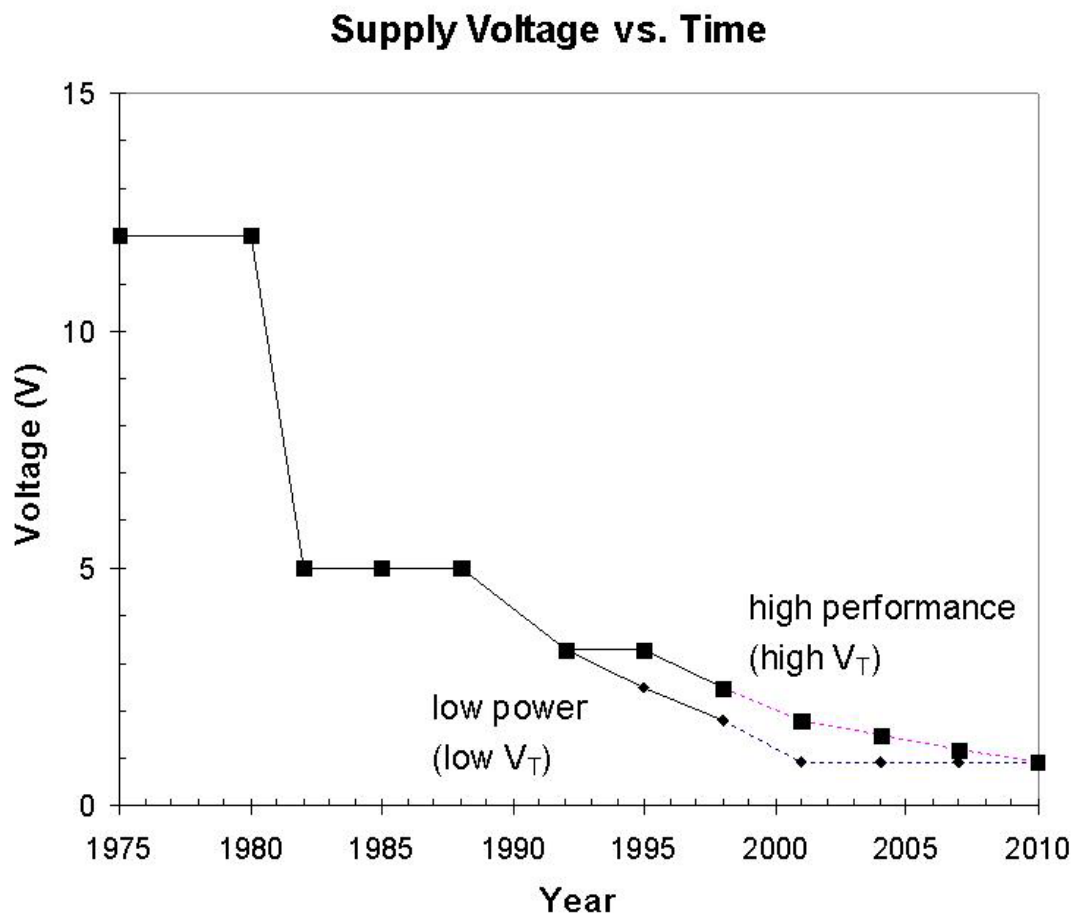
P-channel transistor performance is plotted using the CV/I metric.

Adapted from M. Bohr, *Semiconductor International*, July 1995 (75).



## □ Modern generalized scaling

- Concept of *generation*: every 2 years, new technology is deployed with 30% reduced transistor delay (microprocessor performance doubling every 2 years).
- Everything scales:  $L$  ( $\downarrow$ ),  $W$  ( $\downarrow$ ),  $x_{ox}$  ( $\downarrow$ ),  $N_A$  ( $\uparrow$ ),  $x_j$  ( $\downarrow$ ), and  $V_{DD}$  ( $\downarrow$ ).
- Scaling goal: *extract maximum performance from each generation* (maximize  $I_{on}$ ), for a given amount of:
  - short-channel effects (DIBL), *and*
  - off-current
- Currently two technology flavors:
  - *high-performance*: high  $V_{DD}$  (high  $I_D$ , low  $\tau$ ), high  $V_{th}$  (low  $I_{off}$ );
  - *low-power*: low  $V_{DD}$  (low  $I_D$ , high  $\tau$ ), low  $V_{th}$  (high  $I_{off}$ ).



## Key conclusions

- *Constant field scaling*: scale all device dimensions keeping vertical and horizontal electric fields constant.

Consequences:

- $I_{off} \uparrow$
- system designers don't want to scale  $V_{DD}$

- *Constant voltage scaling*: scale all device dimensions keeping voltage constant.

Consequences:

- $I_{off} \uparrow$
- fields everywhere  $\uparrow \Rightarrow$  reliability compromised

- For a long time scaling proceeded through constant  $V_{DD}$  path with abrupt drops in  $V_{DD}$ .

- Scaling goal: *extract maximum performance from each generation* (maximize  $I_{on}$ ), for a given amount of:

- short-channel effects (DIBL), *and*
- off-current

- *Generalized scaling* demands simultaneous scaling of  $L_g$ ,  $x_{ox}$ ,  $x_j$ ,  $N_A$ , and  $V_{DD}$ .