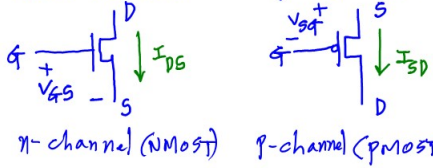
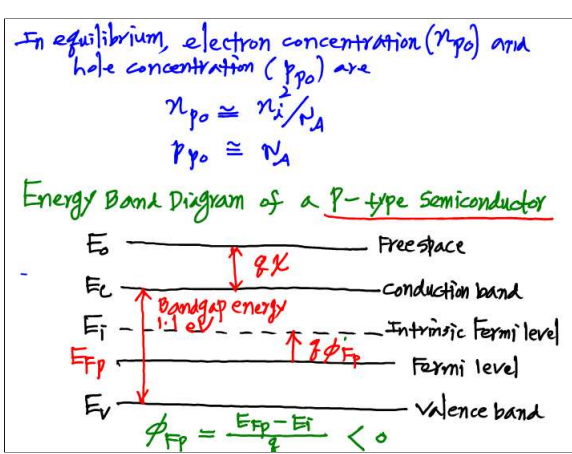
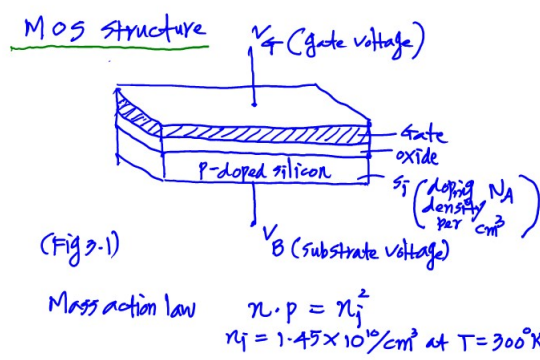
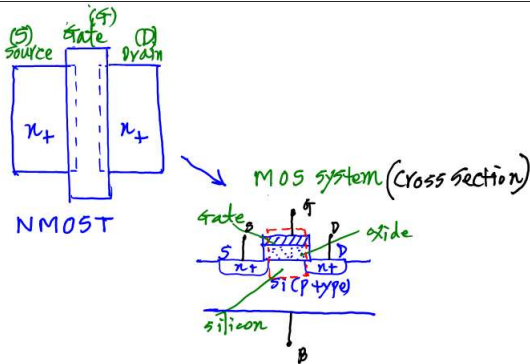
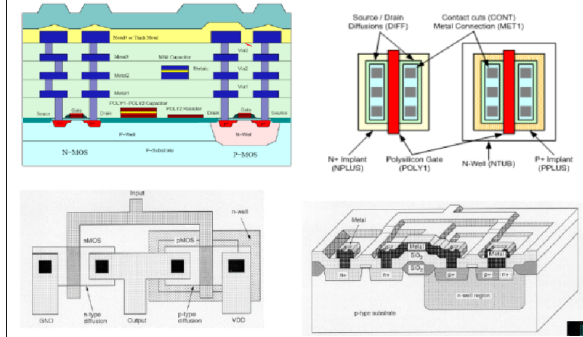
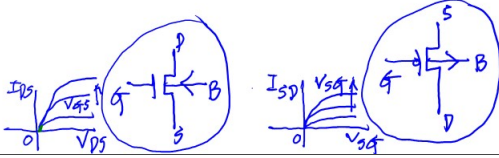


ECE222 Lecture 2 Oct 1, 2019

MOS Transistors chap 3 (pp 51-78; 79-104)



n-channel (NMOS) p-channel (PMOS)



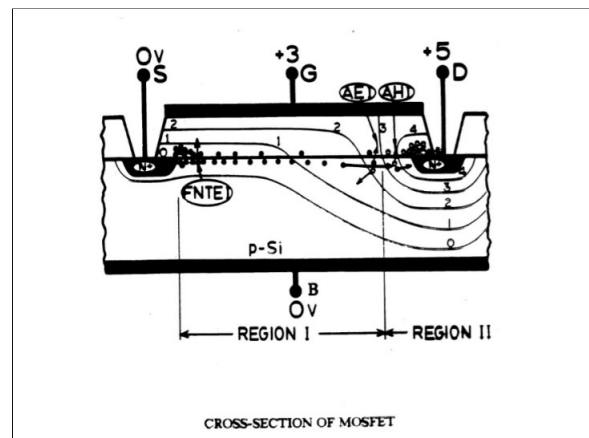
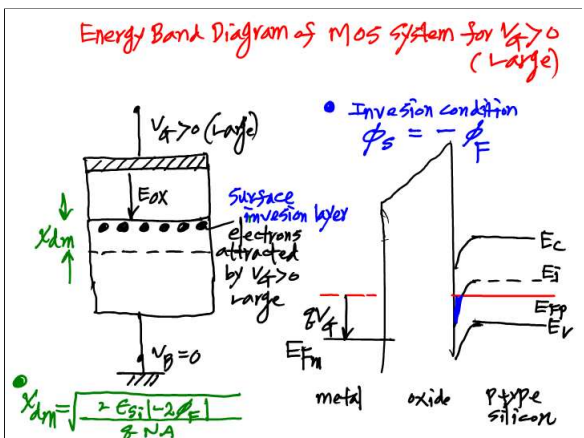
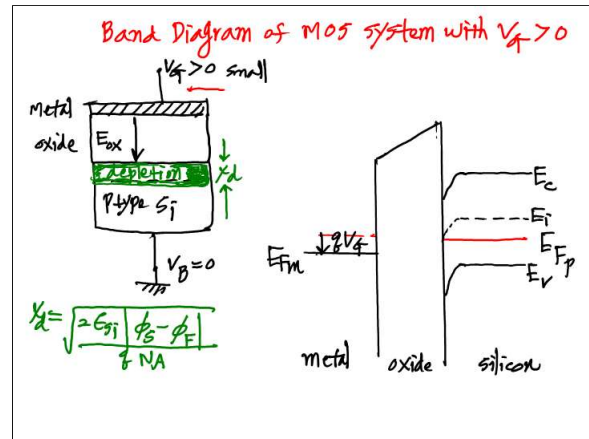
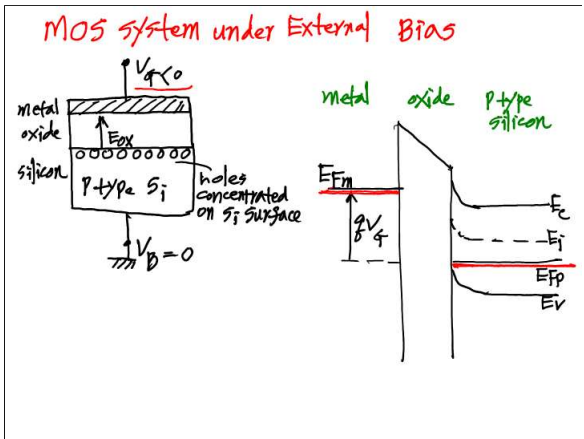
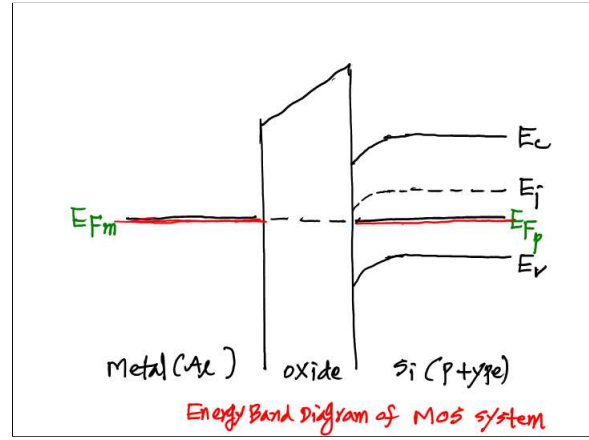
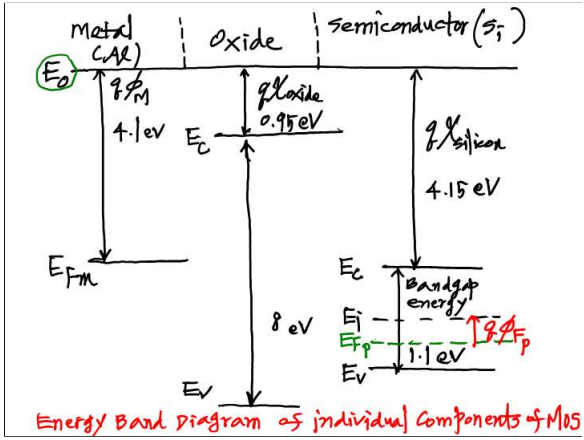
$$\phi_{Fp} = \frac{kT}{q} \ln \left(\frac{n_i}{n_A} \right) < 0 \quad (3.4)$$

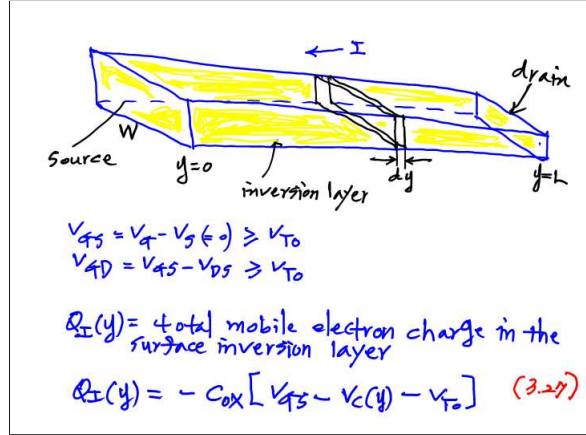
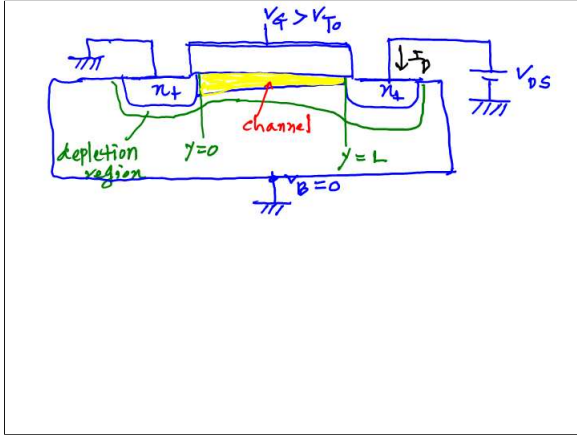
$$\phi_{Fn} = \frac{kT}{q} \ln \left(\frac{N_D}{n_i} \right) > 0 \quad (3.5)$$

$q\chi$ = electron affinity = potential difference between the conduction band and the vacuum (free space)

$q\phi_s$ = work function = the energy required for an electron to move from the Fermi level to free space

$$q\phi_s = q\chi + E_C - E_F$$





$$V_{gs} = V_g - V_s (\neq 0) \geq V_{to}$$

$$V_{gd} = V_{gs} - V_{ds} \geq V_{to}$$

$Q_I(y)$ = total mobile electron charge in the surface inversion layer

$$Q_I(y) = -C_{ox} [V_{gs} - V_c(y) - V_{to}] \quad (3.27)$$

dR = incremental resistance of the differential channel segment

$$dR = \frac{dy}{W \cdot \mu_n \cdot Q_I(y)} \quad (3.28)$$

where μ_n = surface electron mobility
 W = channel width (transistor size)

$$-dV_c = I_D dR = \frac{I_D}{W \cdot \mu_n \cdot Q_I(y)} dy$$

$$-W \cdot \mu_n \cdot Q_I(y) dV_c = I_D dy$$

$$-W \cdot \mu_n \int_0^{V_{ds}} Q_I(y) dV_c = \int_0^L I_D dy$$

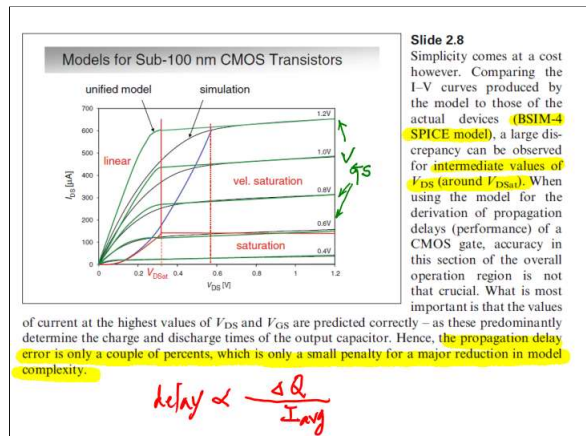
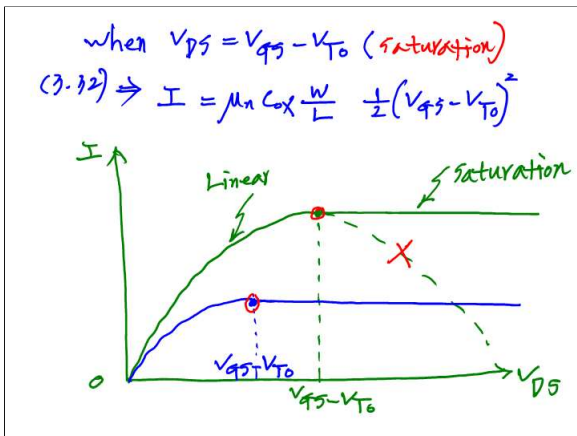
$$\Rightarrow I_D L = -W \mu_n \int_0^{V_{ds}} -C_{ox} (V_{gs} - V_c(y) - V_{to}) dV_c$$

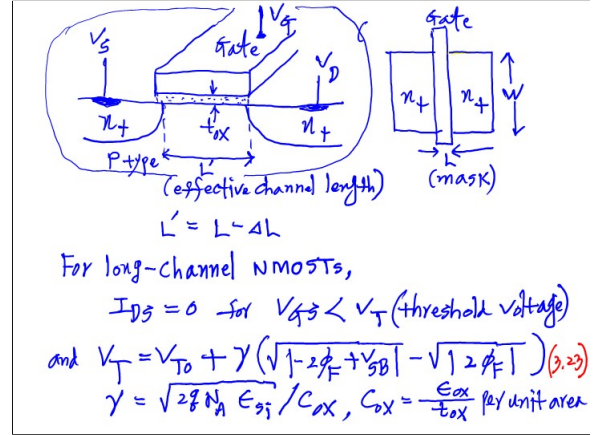
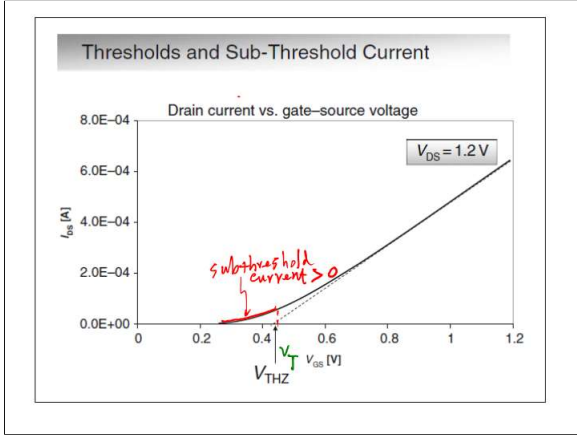
$$= \mu_n C_{ox} W \int_0^{V_{ds}} (V_{gs} - V_c(y) - V_{to}) dV_c$$

$$= \mu_n C_{ox} W \left[(V_{gs} - V_{to}) V_c - \frac{1}{2} V_c^2 \right]_0^{V_{ds}}$$

$$= \mu_n C_{ox} W \left[(V_{gs} - V_{to}) V_{ds} - \frac{1}{2} V_{ds}^2 \right]$$

$$\Rightarrow I_D = \mu_n C_{ox} \frac{W}{L} \left[(V_{gs} - V_{to}) V_{ds} - \frac{1}{2} V_{ds}^2 \right] \quad (3.32)$$





NMOS PMOS

ϕ_F substrate Fermi potential - +

γ substrate bias coefficient + -

V_{SB} substrate bias voltage + -

and $V_{T0} = \phi_{FC} - 2\phi_F - \frac{Q_{B0}}{C_{ox}} - \frac{Q_{ox}}{C_{ox}}$ Eq. (3.19)

$Q_{B0} = -\sqrt{2q N_A \epsilon_{Si} | -2\phi_F |}$

$Q_B = -\sqrt{2q N_A \epsilon_{Si} | -2\phi_F + V_{SB} |}$

$Q_{ox} = q N_{ox}$ (=oxide interface fixed charge density)

$\epsilon_{Si} = 11.7 \epsilon_0 = 11.7 \times (8.854 \times 10^{-14})$

N_A = acceptor concentration (typically Boron)

hole concentration in the p-type substrate (body)
 $p_{p0} \approx N_A$ (typically $10^{15} \sim 10^{16} / \text{cm}^3$, but can be much higher)

electron concentration
 $n_{p0} \approx \frac{n_i^2}{N_A}$, n_i = intrinsic carrier concentration in S_i
 (at $T = 300K$, $1.45 \times 10^{10} / \text{cm}^3$)

$q = 1.602 \times 10^{-19} C$

$\phi_F = \frac{E_F - E_i}{q}$
 ($\phi_{Fp} < 0$)

Energy Band of p-type substrate

E_0 free space

E_C conduction

E_i intrinsic

E_{Fp} Fermi level

E_V valence band

Band gap ($i.v.$) $\uparrow q\phi_{Fp}$

(e.g.) for p-type substrate $N_A = 4 \times 10^{18} / \text{cm}^3$
 polysilicon gate doping $N_D = 2 \times 10^{20} / \text{cm}^3$

$\phi_{F(\text{substrate})} = \frac{kT}{q} \ln\left(\frac{n_i}{N_A}\right)$

$\stackrel{\uparrow}{=} \frac{1.38 \times 10^{-23} [J/K] \times 300 [K]}{1.6 \times 10^{-19} [C]} \ln\left(\frac{1.45 \times 10^{10}}{4 \times 10^{18}}\right)$

$\stackrel{\uparrow}{=} \text{at room-temp}$

$= 0.026 [V] \times (-19.435) = -0.505 [V]$

$\phi_{F(\text{gate})} = \frac{kT}{q} \ln\left(\frac{N_D}{n_i}\right) = 0.026 [V] \ln\left(\frac{2 \times 10^{20}}{1.45 \times 10^{10}}\right)$

poly-silicon

$= 0.026 [V] \times (729.347) = +0.607 [V]$

work-function difference between gate and channel,
 $\phi_{FC} = \phi_{F(\text{substrate})} - \phi_{F(\text{gate})} = -0.505 - 0.607$
 $= -1.11 [V]$

Depletion region charge density at $V_{SB} = 0$

$Q_{B0} = -\sqrt{2q N_A \epsilon_{Si} | -2\phi_F(\text{substrate}) |}$

$= -\sqrt{2(1.6 \times 10^{-19})(4 \times 10^{18})(11.7 \times 8.85 \times 10^{-14}) | 2(0.505) |}$

$= -1.16 \times 10^{-6} [C] / \text{cm}^2$

oxide interface charge

$Q_{ox} = q N_{ox} = 1.6 \times 10^{-19} [C] \times 4 \times 10^{18} / \text{cm}^2$
 specified

$= 6.4 \times 10^{-9} [C] / \text{cm}^2$

$C_{ox} = \frac{\epsilon_{ox}}{t_{ox}} = \frac{3.97 \times 8.85 \times 10^{-14} F / \text{cm}}{16 \times 10^{-8} \text{cm}} = 2.20 \text{ nF} / \text{cm}^2$

$V_{T0} = \phi_{FC} - 2\phi_{F(\text{substrate})} - \frac{Q_{B0}}{C_{ox}} - \frac{Q_{ox}}{C_{ox}}$

$= -1.11 [V] - 2(-0.505) - \frac{-1.16 \times 10^{-6} [C]}{2.20 \times 10^{-6} [F]} - \frac{6.4 \times 10^{-9} [C]}{2.20 \times 10^{-6} [F]}$

$$= -1.11 + 1.01 + 0.527 - 2.91 \times 10^{-3}$$

$$= 0.427 \text{ [V]} = V_{T0}$$

$$\gamma = \frac{\sqrt{2qNA E_{Si}}}{C_{ox}}$$

$$= \frac{\sqrt{2 \times 1.6 \times 10^{-19} \text{ [C]} \times 4 \times 10^{19} \text{ [cm}^{-3}] \times 11.7 \times 8.85 \times 10^{-14} \text{ F/cm}}}{2.2 \times 10^{-6} \text{ F/cm}^2}$$

$$= \frac{\sqrt{1325.376 \times 10^{-15} \text{ C.F/cm}^2}}{2.2 \times 10^{-6} \text{ F/cm}^2} = \frac{11.51 \times 10^{-9} \text{ F/cm}^2 \sqrt{\text{V}}}{2.2 \times 10^{-6} \text{ F/cm}^2}$$

$$= 0.527 \text{ [V]}^{1/2}$$

$$V_T = V_{T0} + \gamma (\sqrt{|-2\phi_F + V_{SB}|} - \sqrt{|-2\phi_F|}) \quad (3.23)$$

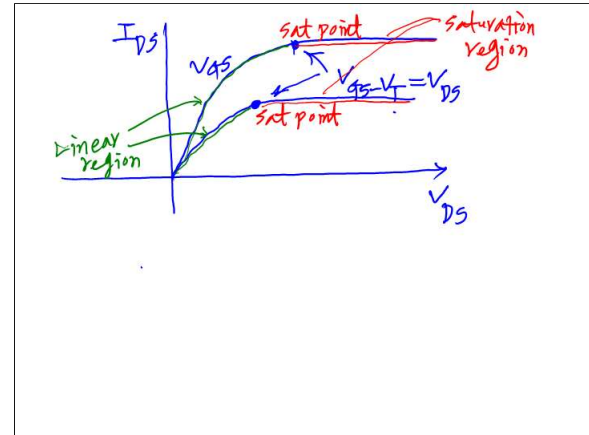
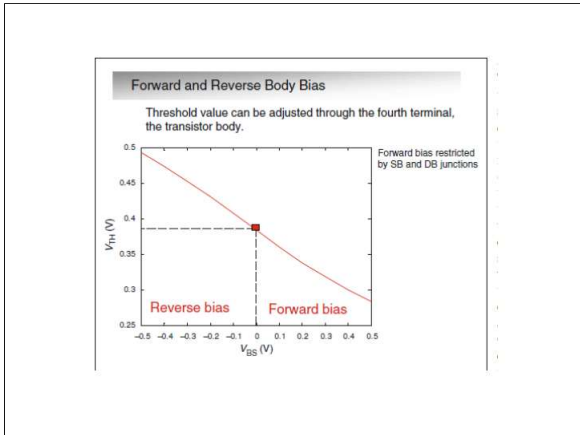
$$\begin{aligned} &= 0.427 \text{ [V]} + 0.527 \text{ [V]}^{1/2} (\sqrt{|-2(-0.509) + 1|} \text{ [V]} - \sqrt{|-2(-0.509)|} \text{ [V]}) \\ & \quad (\phi_{F(\text{substrate})} = -0.509) \\ &= 0.427 \text{ [V]} + 0.527 \text{ [V]}^{1/2} (\sqrt{2.018} \text{ [V]} - \sqrt{1.018} \text{ [V]}) \\ &= (0.427 + 0.527 \times 0.413) \text{ [V]} = 0.64 \text{ [V]} = V_T \end{aligned}$$

$V_{SB} = 1 \text{ V}$

From $V_T = V_{T0} + \gamma (\sqrt{|-2\phi_F + V_{SB}|} - \sqrt{|-2\phi_F|})$
 when $V_S = V_B$ $V_T = V_{T0}$
 But as $V_{SB} > 0$ increase $V_T \uparrow$

$V_{T2} > V_{T1} = V_{T0}$

note that due to a Body Effect $V_{T2} > V_{T1}$



For Long channel NMOSs via gradual channel approx.

Linear region $I_{DS} = \mu_n C_{ox} \frac{W}{L} \left[2(V_{GS} - V_T) V_{DS} - V_{DS}^2 \right] \quad (3.34)$
electron mobility for $V_{GS} > V_T > 0$

Saturation region $I_{DS} = \mu_n C_{ox} \frac{W}{L} \left[2(V_{GS} - V_T) V_{DS} - V_{DS}^2 \right]$
 $(V_{GS} - V_T = V_{DS})$
 $= \mu_n C_{ox} \frac{W}{L} (V_{GS} - V_T)^2 \quad (3.38)$

Channel length modulation parameter λ
 $L \leftarrow L' = L - \Delta L = L \left(1 - \frac{\Delta L}{L} \right)$
 $\approx L (1 - \lambda V_{DS})$
 $\lambda = \text{empirical parameter}$

With channel length modulation

(3.38) $\leftarrow I_{DS \text{ sat}} = \mu_n C_{ox} \frac{W}{L(1-\lambda V_{DS})} (V_{GS} - V_T)^2$
 $= \mu_n C_{ox} (V_{GS} - V_T)^2 (1 + 2\lambda V_{DS}) \quad (3.49)$
 $(\frac{1}{1-\epsilon} = 1 + \epsilon)$

In summary for NMOSs

$I_{DS \text{ linear}} = \mu_n C_{ox} \frac{W}{L} \left[2(V_{GS} - V_T) V_{DS} - V_{DS}^2 \right] \quad (3.51)$
 $I_{DS \text{ sat}} = \mu_n C_{ox} \frac{W}{L} (V_{GS} - V_T)^2 (1 + 2\lambda V_{DS}) \quad (3.52)$

Similarly for PMOSTA *not* μ_n (+Vps)

$$I_{SD \text{ linear}} = \frac{\mu_p C_{ox} W}{2 L} \left[2(V_{GS} - V_T) V_{DS} - V_{DS}^2 \right] \quad (3.58)$$

for $V_{GS} < V_T$ & $V_{DS} > V_{GS} - V_T$

which is same as

$$\frac{\mu_p C_{ox} W}{2 L} \left[2(V_{SG} + V_T) V_{SD} - V_{SD}^2 \right]$$

$$I_{SD \text{ sat}} = \frac{\mu_p C_{ox} W}{2 L} (V_{GS} - V_T)^2 (1 + \lambda V_{DS}) \quad (3.59)$$

$1 - \lambda V_{DS}, V_{DS} < 0$
(error in the book)

I-V Equations for short channel MOSTA

$$\mu_n(\text{eff}) = \frac{\mu_{n0}}{1 + \gamma(V_{GS} - V_T)} \quad (3.69)$$

γ = empirical coefficient

μ_{n0} = low-field electron mobility

For NMOSTA

$$I_{DS \text{ linear}} = \frac{\mu_n C_{ox} W}{2 L} \frac{1}{1 + \frac{V_{DS}}{E_c L}} \left[2(V_{GS} - V_T) V_{DS} - V_{DS}^2 \right]$$

for $V_{GS} > V_T > 0$ & $V_{DS} < \frac{(V_{GS} - V_T) E_c L}{(V_{GS} - V_T) + E_c L}$ (3.85)

where E_c = channel electric field

With V_{sat} (saturated drift velocity of electrons)

$$I_{DS \text{ sat}} = W V_{sat} C_{ox} \frac{(V_{GS} - V_T)^2}{(V_{GS} - V_T) + E_c L} (1 + \lambda V_{DS}) \quad (3.86)$$

for $V_{GS} \geq V_T, V_{DS} \geq \frac{(V_{GS} - V_T) E_c L}{(V_{GS} - V_T) + E_c L}$

Similarly for PMOSTA

(3.87) & (3.88)

threshold voltage of small geometry Devices

$$V_T = V_{T0} + K_1 \left(\sqrt{|-2\phi_F + V_{SB}|} - \sqrt{|-2\phi_F|} \right) + K_2 V_{SB}$$

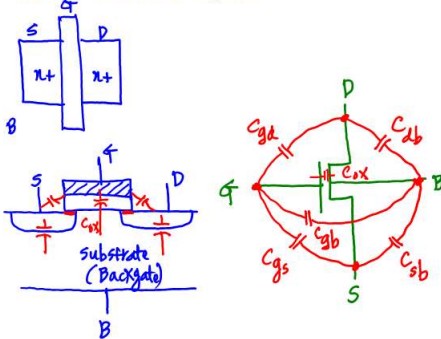
$$- \Delta V_{T, SCE} + \Delta V_{T, NWF} - \Delta V_{T, DEBL} + \Delta V_{T, RSCF} - \Delta V_{T, DL} + \Delta V_{T, LIS}$$

short channel effect, narrow width effect, drain induced barrier lowering, (reverse SCE), drain induced threshold shift

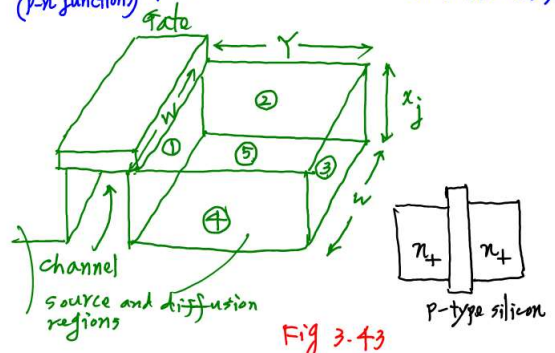
(3.116)

$$I_{DS \text{ subthreshold}} = \frac{q D_n W x_n n_0}{L B} e^{-\frac{q\phi}{kT}} e^{-\frac{q}{kT} (A V_{GS} + B V_{DS})} \quad (3.115)$$

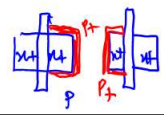
MOS Transistor Capacitances



Diffusion Capacitances in MOST (C_{sb}, C_{db}) (p-n junctions)



p-n Junction	Area (A)	Type
①	w · x _j	n ⁺ /p
②	Y · x _j	n ⁺ /p ← p ⁺
③	w · x _j	n ⁺ /p ← p ⁺
④	Y · x _j	n ⁺ /p ← p ⁺
⑤	w · Y	n ⁺ /p



Depletion capacitance of a reversely biased "abrupt" p-n junction: (approximation)

$$x_d = \sqrt{\frac{2 \epsilon_{si}}{q} \frac{N_A + N_D}{N_A \cdot N_D} (\phi_0 - V)} \quad (3.123)$$

where $\phi_0 = \frac{kT}{q} \ln \frac{N_A \cdot N_D}{n_i^2}$
and $V < 0$

$$Q_j = A \cdot q \cdot \frac{N_A N_D}{N_A + N_D} x_d \quad (3.124)$$

$$= A \sqrt{2 \epsilon_{si} q \frac{N_A N_D}{N_A + N_D} (\phi_0 - V)} \quad (3.125)$$

$$C_j(V) = \left| \frac{dQ_j}{dV} \right|$$

$$= A \sqrt{\frac{\epsilon_{si} q N_A N_D}{2(N_A + N_D)}} \frac{1}{\sqrt{\phi_0 - V}} \quad (3.127)$$

More generally, with junction grading

$$C_j(V) = \frac{A C_{j0}}{\left(1 - \frac{V}{\phi_0}\right)^m} \quad (3.128)$$

where $m =$ grading coefficient
 $= \frac{1}{2}$ abrupt
 $= \frac{1}{3}$ linearly graded profile
 $C_{j0} = C_j(V) \big|_{V=0}$ in (3.128)

$$C_{df} = \frac{\Delta Q}{\Delta V} = \frac{Q_j(V_2) - Q_j(V_1)}{V_2 - V_1}$$

$$= \frac{1}{V_2 - V_1} \int_{V_1}^{V_2} C_j(V) dV$$

$$= \frac{A C_{j0} \phi_0}{(V_2 - V_1)(1-m)} \left[\left(1 - \frac{V_2}{\phi_0}\right)^{1-m} - \left(1 - \frac{V_1}{\phi_0}\right)^{1-m} \right] \quad (3.131)$$

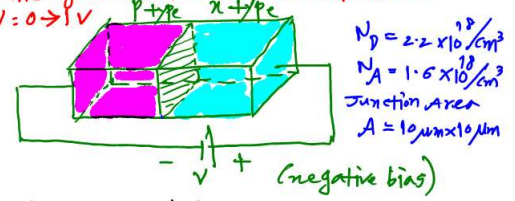
For $m = \frac{1}{2}$ (abrupt junction)

$$C_{df} = \frac{2 A C_{j0} \phi_0}{V_2 - V_1} \left[\sqrt{1 - \frac{V_2}{\phi_0}} - \sqrt{1 - \frac{V_1}{\phi_0}} \right]$$

$$= A C_{j0} K_{df} \quad (3.132)$$

where $K_{df} = \frac{2 \sqrt{\phi_0}}{V_2 - V_1} (\sqrt{\phi_0 - V_2} - \sqrt{\phi_0 - V_1})$
and $0 < K_{df} < 1$

Find the equivalent p-n junction capacitance for $V = 0 \rightarrow 1V$



$$\phi_0 = \frac{kT}{q} \ln \frac{N_A N_D}{n_i^2}$$

$$= 0.026 \ln \left(\frac{1.5 \times 10^{18} \cdot 2.2 \times 10^{18}}{(1.45 \times 10^{10})^2} \right)$$

$$= 0.026 \times 37.43 = \underline{0.97 [V]}$$

$$C_{j0} = \sqrt{\frac{\epsilon_{Si} q}{2} \left(\frac{N_A N_D}{N_A + N_D} \right) \frac{1}{\phi_0}}$$

$$= \sqrt{\frac{11.7 \times 1.6 \times 10^{-19} \text{ C} \times 1.6 \times 10^{19} \text{ cm}^{-3}}{2 \times 0.97} \left(\frac{1.6 \times 10^{19} \times 2.2 \times 10^{19}}{1.6 \times 10^{19} + 2.2 \times 10^{19}} \right) \times \left[\frac{1}{\text{cm}^2} \right]}$$

$$= 2.81 \times 10^{-9} \text{ F/cm}^2$$

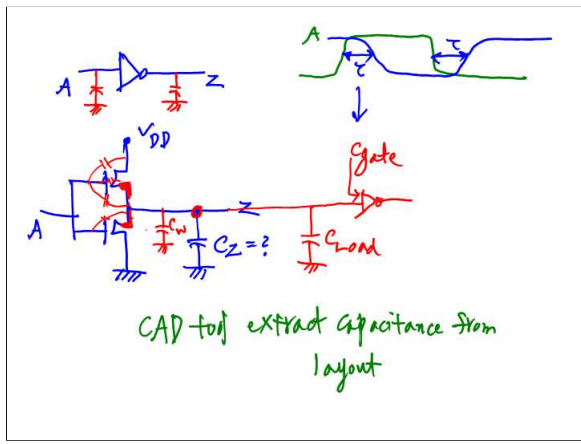
K_{eff} for V switching from 0 to 1V

$$K_{eff} = \frac{2\sqrt{\phi_0}}{V_2 - V_1} (\sqrt{\phi_0 - V_2} - \sqrt{\phi_0 - V_1}) \quad (3-194)$$

$$= \frac{2\sqrt{0.97}}{-1 - 0} (\sqrt{0.97 - (-1)} - \sqrt{0.97 - 0})$$

$$= 0.82$$

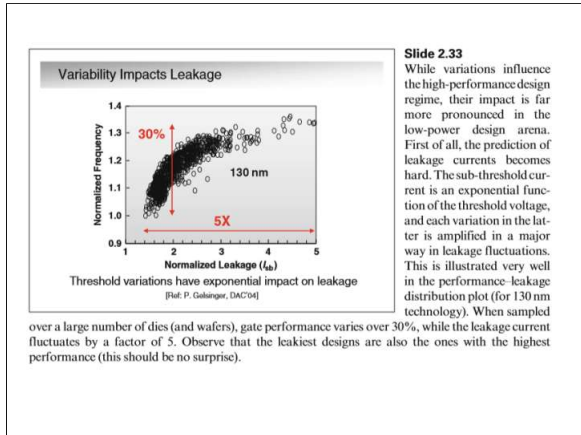
$$C_{eff} = A C_{j0} K_{eff} = 10^{-3} \text{ cm}^2 \times 2.81 \times 10^{-9} \text{ F/cm}^2 \times 0.82 = 230 \text{ fF}$$



Variability

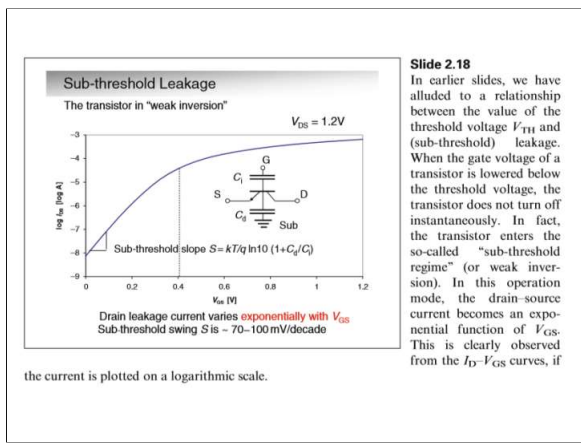
- Scaled device dimensions leading to increased impact of variations
 - Device physics
 - Manufacturing
 - Temporal and environmental
- Impacts performance, power (mostly leakage) and manufacturing yield
- More pronounced in low-power design due to reduced supply/threshold voltage ratios

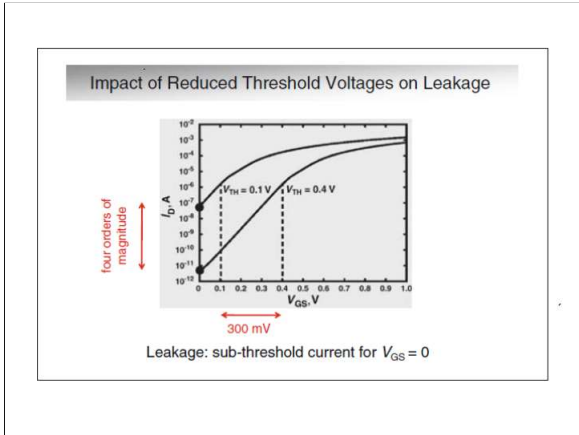
Slide 2.32
The topic of variability rounds out the discussion of the nanometer transistor and its properties. It has always been the case that transistor parameters such as the geometric dimensions or the threshold voltage are not deterministic. When sampled between wafers, within a wafer, or even over a die, each of these parameters exhibits a statistical nature. In the past, the projection of the parameter distributions onto the performance space yielded quite a narrow distribution. This is easily understandable. When the supply voltage is 3V and the threshold is at 0.5V, a 25mV variation in the threshold has only a small impact on the performance and leakage of the digital module. However, when the supply voltage is at 1V and the threshold at 0.3V, the same variation has a much larger impact. So, in past generation processors it was sufficient to evaluate a design over its worst-case corners (FF, SS, FS, SF) in addition to the nominal operation point to determine the yield distributions. Today, this is not sufficient, as the performance distributions have become much wider, and a pure worst-case analysis leads to wasteful design and does not give a good yield perspective either.



MOS Transistor Leakage Components

Slide 2.17
Quite a number of times in the introduction, we have alluded to the increasing effects of "leakage" currents in the nanometer MOS transistor. An ideal MOS transistor (at least from a digital perspective) should not have any currents flowing into the bulk (or well), should not conduct any current between drain and source when off, and should have an infinite gate resistance. As indicated in the accompanying slide, a number of effects are causing the contemporary devices to digress from this ideal model. Leakage currents, flowing through the reverse-biased source-bulk and drain-bulk pn junctions, have always been present. Yet, the levels are so small that their effects could generally be ignored, except in circuitry that relies on charge storage such as DRAMs and dynamic logic. The scaling of the minimum feature sizes has introduced some other leakage effects that are far more influential and exceed junction leakage currents by 3-5 orders of magnitude. Most important are the sub-threshold drain-source and the gate leakage effects, which we will discuss in more detail.





Sub-threshold Current

- Sub-threshold behavior can be modeled physically

$$I_{DS} = 2n\mu C_{ox} \frac{W}{L} \left(\frac{kT}{q} \right)^2 e^{\frac{V_{GS}-V_{TH}}{nV_T}} \left(1 - e^{-\frac{V_{DS}}{V_T}} \right) = I_s e^{\frac{V_{GS}-V_{TH}}{nV_T}} \left(1 - e^{-\frac{V_{DS}}{V_T}} \right)$$

where n is the slope factor (≥ 1 , typically around 1.5) and $I_s = 2n\mu C_{ox} \frac{W}{L} \left(\frac{kT}{q} \right)^2$

- Very often expressed in base 10

$$I_{DS} = I_s 10^{\frac{V_{GS}-V_{TH}}{S}} \left(1 - 10^{-\frac{V_{DS}}{S}} \right) \approx 1 \text{ for } V_{DS} > 100 \text{ mV}$$

where $S = n \left(\frac{kT}{q} \right) \ln(10)$, the sub-threshold swing, ranging between 60 mV and 100 mV

Models for Sub-100 nm CMOS Transistors

Slide 2.8
Simplicity comes at a cost however. Comparing the I-V curves produced by the model to those of the actual devices (BSIM4 SPICE model), a large discrepancy can be observed for intermediate values of V_{DS} (around V_{DSsat}). When using the model for the derivation of propagation delays (performance) of a CMOS gate, accuracy in this section of the overall operation region is not that crucial. What is most important is that the values of current at the highest values of V_{DS} and V_{GS} are predicted correctly – as these predominantly determine the charge and discharge times of the output capacitor. Hence, the propagation delay error is only a couple of percents, which is only a small penalty for a major reduction in model complexity.

Alpha Power Law Model

- Alternate approach, useful for hand analysis of propagation delay

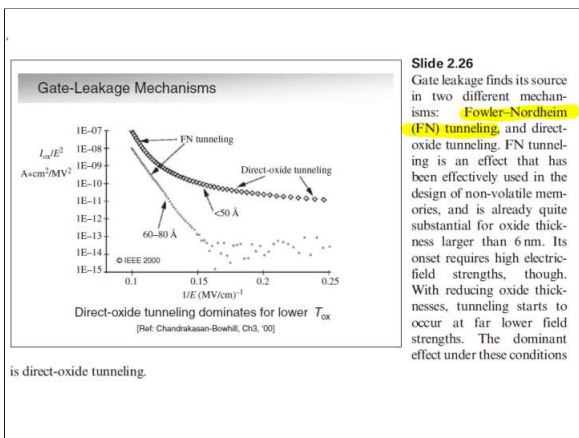
$$I_{DS} = \frac{W}{2L} \mu C_{ox} (V_{GS} - V_{TH})^\alpha$$

- Parameter α is between 1 and 2.
- In 65–180 nm CMOS technology $\alpha \sim 1.2-1.3$

- This is not a physical model
- Simply empirical:
 - Can fit (in minimum mean squares sense) to a variety of α 's, V_{TH}
 - Need to find one with minimum square error – fitted V_{TH} can be different from physical

[Ref: Sakurai, JSSC'90]

Slide 2.9
Even simpler is the alpha model, introduced by Sakurai and Newton in 1990, which does not even attempt to approximate the actual I-V curves. The values of α and V_{TH} are purely empirical, chosen such that the propagation delay of a digital gate, approximated by $t_p = \frac{C_{load} V_{DD}}{I_{DS}}$, best resembles the propagation delay curves obtained from simulation. Typically, curve-fitting techniques such as the minimum-mean square (MMS) are used. Be aware that these do not yield unique solutions and that it is up to the modeler to find the ones with the best fit. Owing to its simplicity, the alpha model is the corner stone of the optimization framework discussed in later chapters.



High-k Gate Dielectric

- Equivalent Oxide Thickness = $EOT = T_{ox} \cdot \epsilon_r \left(\frac{3.9}{\epsilon_r} \right)$, where 3.9 is relative permittivity of SiO_2 and ϵ_r is relative permittivity of high-k material
- Currently SiO_2/Ni ; Candidate materials: HfO_2 ($\epsilon_{ox} \sim 15-30$); $HfSiO_3$ ($\epsilon_{ox} \sim 12-16$)
- Often combined with metal gate

Reduced Gate Leakage for Similar Drive Current

Slide 2.28
The MOS transistor current is proportional to the process transconductance parameter $k' = \mu C_g = \mu \epsilon / t_g$. To increase k' through scaling, one must either find a way to increase the mobility of the carriers or increase the gate capacitance (per unit area). The former requires a fundamental change in the device structure (to be discussed later). With the traditional way of increasing the gate capacitance (i.e., scaling T_g) running out of steam, the only remaining option is to look for gate dielectrics with a higher permittivity ϵ – the so-called high-k

6 Executive Summary

This was the (First) Era of Geometrical (classical) Scaling. This type of scaling was the foundation of the National Technology Roadmap for Semiconductors (NTRS) initiated in 1991.

The ITRS laid out the foundations of the (Second) Era: Equivalent Scaling (e.g., strained silicon, high-K/metal gate, Multigate transistors and use of non-silicon semiconductors in general) between 1998 and 2000.

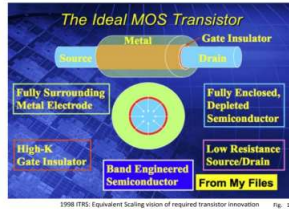


Fig. 3.1 The Ideal MOS Transistor

The implementation of these technologies successfully supported the growth of the semiconductor industry in the past decade and it will continue to do so until the end of the present decade and beyond.

In the next decade ITRS 2.0 predicts that the advent of the third phase of scaling "3D Power Scaling" will become the driver of the rejuvenated semiconductor industry and this answers the question posed before about the future of the semiconductor industry: **Yes the semiconductor industry will continue to be a key enabler of the IoT!**

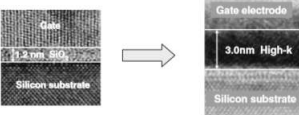
Moore's Law is now entering into a third phase characterized by vertical integration and performance specifications driven towards reduction of power in either the active or the stand-by modes.

34 Executive Summary

Electrical Properties of Transistors

YEAR OF PRODUCTION	2016	2017	2018	2019	2024	2027	2030
Logic device technology naming	P10M56	P4M24	P4M24	P3M20	P2M12G1	P2M12G2	P2M12G3
Logic industry "Node Range" Labeling (nm)	"16/14"	"11/10"	"7/7"	"6/5"	"4/3"	"3/2.8"	"2/1.6"
Logic device structure options	FinFET FDSON	FinFET FDSON	FinFET LGA	FinFET LGA VGAA MSD	VGAA MSD	VGAA MSD	VGAA MSD
DEVICE ELECTRICAL SPEEDS							
Power Supply Voltage - 1.0V (V)	0.80	0.75	0.70	0.65	0.55	0.45	0.40
Subthreshold slope - (mV/dec)	75	70	68	66	40	25	20
Inversion layer thickness - (nm)	1.10	1.00	0.90	0.85	0.80	0.80	0.80
V ₁₀₀ (m/s) at I _{off} =100pA/μm - HP Logic	129	129	133	136	84	52	52
V ₁₀₀ (m/s) at I _{off} =100pA/μm - LP Logic	351	336	333	326	201	125	125
Effective mobility (cm ² /V·s)	200	180	170	160	100	100	100
Resist (Ohm·μm) - HP Logic	290	238	202	172	146	124	106
Bulkivity, Injection velocity (cm/s)	1.20E-07	1.32E-07	1.45E-07	1.60E-07	1.76E-07	1.93E-07	2.13E-07
V ₁₀₀ (V) - HP Logic	0.115	0.127	0.136	0.128	0.141	0.155	0.170
V ₁₀₀ (V) - LP Logic	0.125	0.141	0.155	0.163	0.169	0.186	0.204
I _{on} (pA/μm) at I _{off} =100pA/μm - HP Logic after Rest	2311	2541	2782	2917	3001	2670	2408
I _{on} (pA/μm) at I _{off} =100pA/μm - HP Logic after Rest	1177	1287	1397	1476	1546	1456	1391
I _{on} (pA/μm) at I _{off} =100pA/μm - LP Logic after Rest	1455	1567	1614	1603	2008	1933	1552
I _{on} (pA/μm) at I _{off} =100pA/μm - LP Logic after Rest	896	637	637	629	890	955	821
C _{ch, total} (fF/μm ²) - HP LP Logic	31.38	34.52	36.35	40.61	43.14	43.14	43.14
C _{gate, total} (fF/μm ²) - HP Logic	1.81	1.49	1.29	0.97	1.04	1.04	1.04
C _{gate, total} (fF/μm ²) - LP Logic	1.96	1.66	1.47	1.17	1.24	1.24	1.24
C _{V100} (pF) - HP Logic	2.69	2.61	1.94	1.29	1.11	0.96	0.89
TCV3 (1/ps) - FD1 load, HP Logic	0.27	0.38	0.52	0.78	0.90	1.04	1.12
Energy per switching (fJ) - (Switching) - FD1 load, HP Logic	3.47	2.52	1.89	1.24	0.94	0.63	0.50

High-k Dielectrics



	High-k vs SiO ₂	Benefits
Gate capacitance	60% greater	Faster transistors
Gate dielectric leakage	>100% reduction	Lower power

Buys a few generations of technology scaling
[Courtesy: Intel]

Slide 2.29

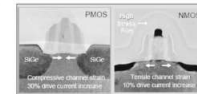
The advantages offered by high-k gate dielectrics are quite clear: faster transistors and/or reduced gate leakage.

Device and Technology Innovations

- Strained silicon
- Silicon-on-Insulator
- Dual-gated devices
- Very high mobility devices
- MEMS - transistors

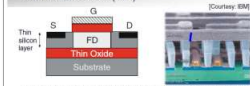


Strained Silicon



Improved ON-current (10-25%) translates into:
• 64-67% leakage current reduction
• of 15% active power reduction
[Ref: R. Gaurang, DAC'04]

Silicon-on-insulator (SOI)



- Reduced capacitance (source and drain to bulk) results in lower dynamic power
- Faster sub-threshold roll-off (close to 60 mV/decade)
- Random threshold fluctuations eliminated in fully-depleted SOI
- Reduced impact of soft-errors
- But:
 - More expensive
 - Secondary effects

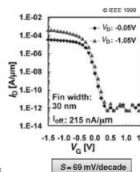
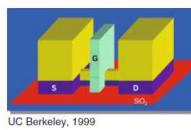
down to the insulator layer, their junction capacitances are substantially reduced, which translates directly into power savings. Another advantage is the higher sub-threshold slope factor (approaching the ideal 60 mV/decade), reducing leakage. Finally, the sensitivity to soft errors is reduced owing to the smaller collection efficiency, leading to a more reliable transistor. There are some important negatives as well. The addition of the SiO₂ layer and the thin silicon layer increases the cost of the substrate material, and may impact the yield as well. In addition, some secondary effects should be noted. The SOI transistor is essentially a three-terminal device without a bulk (or body) contact, and a "body" that is floating. This effectively eliminates body biasing as a threshold-control technique. The floating transistor body also introduces some interesting (ironically speaking...) features such as hysteresis and state-dependency.

Device engineers differentiate between two types of SOI transistors: partially-depleted (PD-SOI) and fully-depleted (FD-SOI). In the latter, the silicon layer is so thin that it is completely depleted under nominal transistor operation, which means that the depletion/inversion layer under the gate extends all the way to the insulator. This has the advantage of suppressing some of the floating-body effects, and an ideal sub-threshold slope is theoretically achievable. From a variation perspective, the threshold voltage becomes independent of the doping in the channel, effectively eliminating a source of random variations (as discussed in Slide 2.37). FD-SOI requires the depositing of extremely thin silicon layers (3-5 times thinner than the gate length).

Slide 2.43

Silicon-on-insulator (SOI) is a technology that has been "on the horizon" for quite a long time, yet it never managed to really break ground, though with some exceptions here and there. An SOI MOS transistor differs from a "bulk" device in that the channel is formed in a thin layer of silicon deposited above an electrical insulator, typically silicon dioxide. Doing so offers some attractive features. First, as drain and source diffusions extend all the way

FinFETs - An Entirely New Device Architecture



- Suppressed short-channel effects
- Higher on-current for reduced leakage
- Undoped channel - No random dopant fluctuations

[Ref: X. Huang, IEDM'99]

Slide 2.45

The FinFET (called a tri-gate transistor by Intel) is an entirely different transistor structure that actually offers some properties similar to the ones offered by the device presented in the previous slide. The term FinFET was coined by researchers at the University of California at Berkeley to describe a non-planar, double-gated transistor built on an SOI substrate. The distinguishing characteristic of the FinFET is that the controlling gate is wrapped around a thin silicon "fin",

which forms the body of the device. The dimensions of the fin determine the effective channel length of the device. The device structure has shown the potential to scale the channel length to values that are hard, if not impossible, to accomplish in traditional planar devices. In fact, operational transistors with channel lengths down to 7 nm have been demonstrated.

In addition to a suppression of deep submicron effects, a crucial advantage of the device is again increased control, as the gate wraps (almost) completely around the channel.

BackGated FinFET

Slide 2.46

This increased two-dimensional control can be exploited in a number of ways. In the dual-gated device, the fact that the gate is controlling the channel from both sides (as well as the top) leads to increased process transconductance. Another option is to remove the top part of the gate, leading to the back-gated transistor. In this structure, one of the gates acts as the standard control gate, whereas the other is used to manipulate the threshold voltage. In a sense, this device offers similar functionality as the buried-gate FD-SOI transistor discussed earlier. Controlling the work functions of the two gates through the selection of appropriate type and quantity of the dopants helps to maximize the range and sensitivity of the control knobs.

Independent front and back gates
One switching gate and V_{TH} control gate
Increased threshold control

New Transistors: FinFETs

Slide 2.47

The fact that the FinFET and its cousins are dramatically different devices compared to your standard bulk MOS transistor is best-illustrated with these pictures from Berkeley and Intel. The process steps that set and control the physical dimensions are entirely different. Although this creates new opportunities, it also brings challenges, as the process steps involved are vastly different. The ultimate success of the FinFET depends greatly upon how these changes can be translated into a scalable, low-cost and high-yield process – some formidable question, indeed! Also unclear at this time is how the adoption of such a different structure impacts variability, as critical dimensions and device parameters are dependent upon entirely different process steps.

Manufacturability still an issue – may even cause more variations

[Courtesy: T.J. King, UCB, Intel]

New Transistors: FinFETs

Manufacturability still an issue – may even cause more variations

[Courtesy: T.J. King, UCB, Intel]