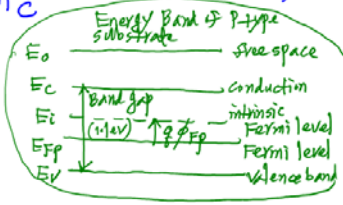
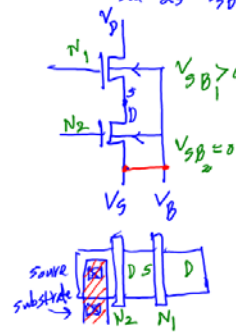


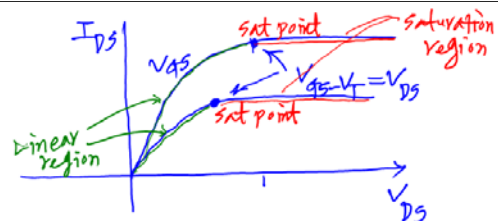
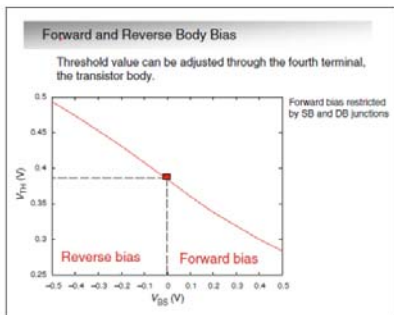
- $\epsilon_{Si} = 11.7 \epsilon_0 = 11.7 \times (8.854 \times 10^{-14})$
- N_A = acceptor concentration (typically Boron)
- hole concentration in the p-type substrate (body)
 $p_0 \approx N_A$ (typically $10^{15} \sim 10^{16}/\text{cm}^3$, but can be much higher)
- electron concentration
 $n_0 \approx \frac{n_i^2}{N_A}$, n_i = intrinsic carrier concentration in Si
(At $T=300\text{K}$, $1.45 \times 10^{10}/\text{cm}^3$)
- $q = 1.602 \times 10^{-19} \text{ C}$
- $\phi_F = \frac{E_F - E_i}{q}$
($\phi_{Fp} < 0$)



From $V_T = V_{T0} + \gamma \left(\sqrt{|-2\phi_F + V_{SB}|} - \sqrt{|-2\phi_F|} \right)$
 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 NMOS, I_{DS}

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

Saturation region $I_{DS} = \mu_n C_{ox} \frac{W}{L} \left[2(V_{GS} - V_T)(V_{GS} - V_T) - (V_{GS} - V_T)^2 \right]$
 $V_{DS sat} = V_{GS} - V_T$
 $= \mu_n C_{ox} \frac{W}{L} (V_{GS} - V_T)^2$ (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) $\leftarrow I_{DS sat} = \mu_n C_{ox} \frac{W}{L(1 - \lambda V_{DS})} (V_{GS} - V_T)^2$
 $= \mu_n C_{ox} (V_{GS} - V_T)^2 (1 + 2\lambda V_{DS})$ (3.49)
($\frac{1}{1-\epsilon} = 1 + \epsilon$)

In summary for NMOS

$I_{DS linear} = \mu_n C_{ox} \frac{W}{L} \left[(V_{GS} - V_T) V_{DS} - V_{DS}^2 \right]$ (3.4)
 $I_{DS sat} = \mu_n C_{ox} \frac{W}{L} (V_{GS} - V_T)^2 (1 + 2\lambda V_{DS})$ (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 $V_{GS} < V_T < 0$ & $V_{DS} > V_{GS} - V_T$

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 - λV_{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

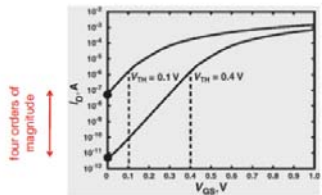
$$V_T = V_{T0} + K_1 \left(\sqrt{|-2\phi_F + V_{SB}|} - \sqrt{|-2\phi_F|} \right) + K_2 V_{SB}$$

$-\Delta V_{T, SCE} + \Delta V_{T, NWF} - \Delta V_{T, DIBL}$
short channel effect, V_{GS} and V_{DS} effect, V_{GS} and V_{DS} effect, V_{GS} and V_{DS} effect
 $+ \Delta V_{T, RSCE} - \Delta V_{T, DIOS}$
(reverse SCE), V_{GS} and V_{DS} effect, V_{GS} and V_{DS} effect
drain induced barrier lowering, V_{GS} and V_{DS} effect
drain induced threshold shift

(3.116)

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

Impact of Reduced Threshold Voltages on Leakage



Leakage: sub-threshold current for $V_{DS} = 0$

Sub-threshold Current

• Sub-threshold behavior can be modeled physically

$$I_{DS} = 2n\mu C_{ox} \frac{W}{L} \left(\frac{kT}{q} \right)^2 e^{\frac{V_{GS} - V_{TH}}{n kT/q}} \left(1 - e^{-\frac{V_{DS}}{kT/q}} \right) = I_0 e^{\frac{V_{GS} - V_{TH}}{n kT/q}} \left(1 - e^{-\frac{V_{DS}}{kT/q}} \right)$$

where n is the slope factor (≥ 1 , typically around 1.5) and $I_0 = 2n\mu C_{ox} \frac{W}{L} \left(\frac{kT}{q} \right)^2$

• Very often expressed in base 10

$$I_{DS} = I_0 10^{\frac{V_{GS} - V_{TH}}{S}} \left(1 - 10^{-\frac{V_{DS}}{S}} \right) \quad = 1 \text{ for } V_{DS} > 100 \text{ mV}$$

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

Alpha Power Law Model

- Alternate approach, useful for hand analysis of propagation delay

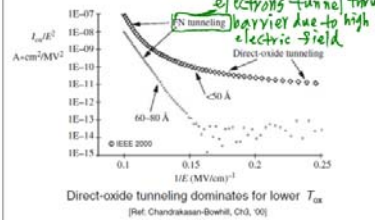
$$I_{DS} = \frac{W}{2L} \mu C_{ox} (V_{GS} - V_{TH})^\alpha$$

- 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]

Gate-Leakage Mechanisms



Slide 2.26

Gate leakage finds its source in two different mechanisms: Fowler-Nordheim (FN) tunneling, and direct-oxide tunneling. FN tunneling is an effect that has been effectively used in the design of non-volatile memories, and is already quite substantial for oxide thickness larger than 6 nm. Its onset requires high electric-field strengths, though. With reducing oxide thicknesses, tunneling starts to occur at far lower field strengths. The dominant effect under these conditions

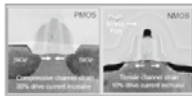
is direct-oxide tunneling.

Device and Technology Innovations

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



Strained Silicon

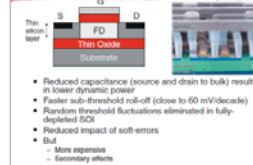


Improved ON-Current (10–25%) translates into:

- 64–67% leakage current reduction
- or 15% active power reduction

[Ref: R. Gauran, DAC'04]

Silicon-on-Insulator (SOI)



Slide 2.43

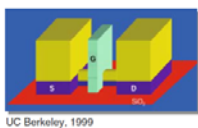
Silicon-on-insulator (SOI) is a technology that has been "on the horizon" for quite a long time, yet it never managed to really break ground, though with some exceptions here and there. An SOI MOS transistor differs from a "bulk" device in that the channel is formed in a thin layer of silicon deposited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

- 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
- 3d
- More expensive
- Secondary effects

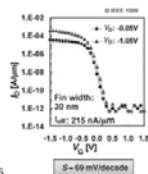
down to the insulator layer, their junction capacitances are substantially reduced, which translates directly into power savings. Another advantage is the higher sub-threshold slope factor (approaching the ideal 60 mV/decade), reducing leakage. Finally, the sensitivity to soil errors is reduced owing 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 (no 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 normal 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).

FinFETs – An Entirely New Device Architecture



UC Berkeley, 1999



- 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.

Illustration: Evolution of Microprocessors

Teraflops Research Chip
Introduced 2006
65nm Technology

80 Processor cores

- 3.16 GHz 62W 1.0 Tflops
- 5.1 GHz 175W 1.6 Tflops
- 5.7 GHz 265W 1.8 Tflops

For comparison: ASCI Red was the first supercomputer to reach Tflops in 1996. That system used nearly 10,000 Pentium® Pro processors running at 200MHz and consumed 500kW of power plus an additional 500kW just to cool the room that housed it.

