

# MOS Transistor (Scaling)

---

- A) Overview
- B) Short channel effects (use TCAD to explore)
- C) Idealistic (constant field) and realistic scaling
- D) Innovation in material and architecture for scaling
- E) Variability
- F) Future challenges

# Technology Scaling

---

**Scale transistor size → Faster switching, more transistors/area**

**Higher performance, more functionality, lower supply voltage (mobile applications)**

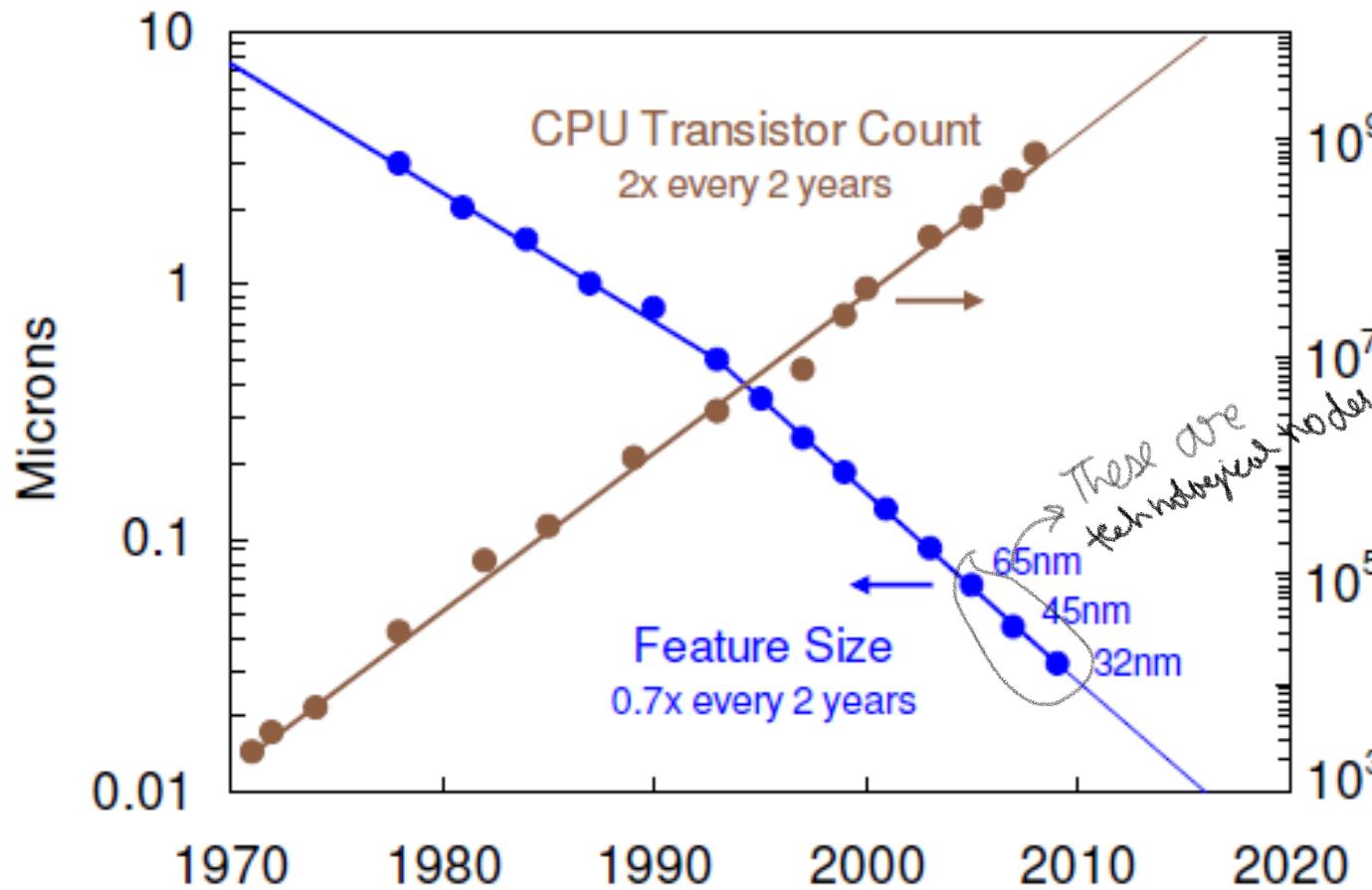
**Today's mobile phones can do much more than yesterday's PCs**

**Huge reduction in cost / functionality**

# Transistor Scaling

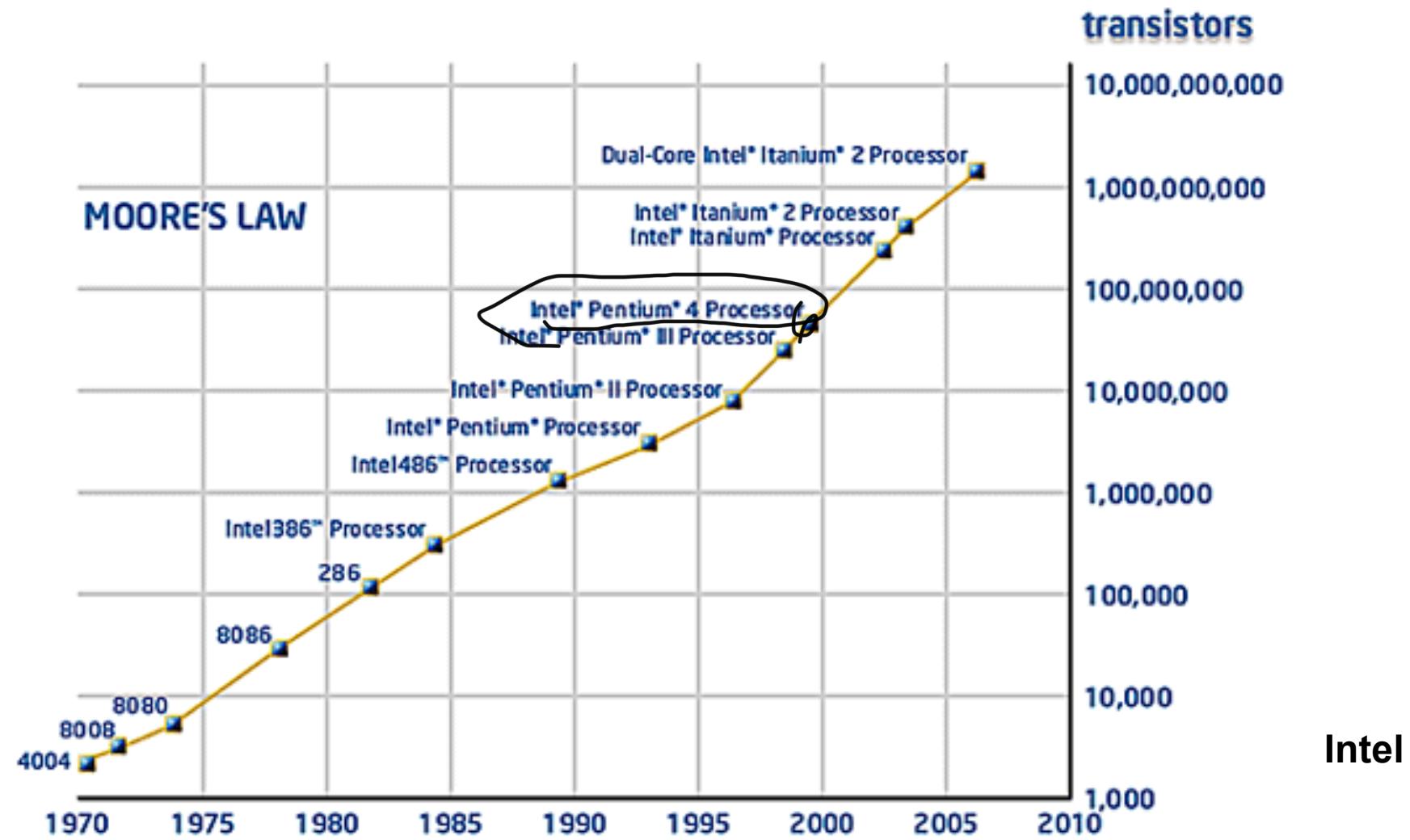
---

Exponential increase over the last ~40 years



# Transistor Count (Functionality)

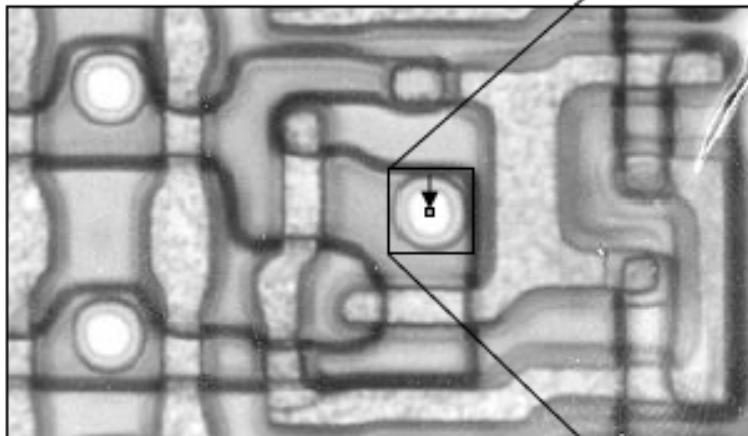
Exponential increase over the last ~40 years



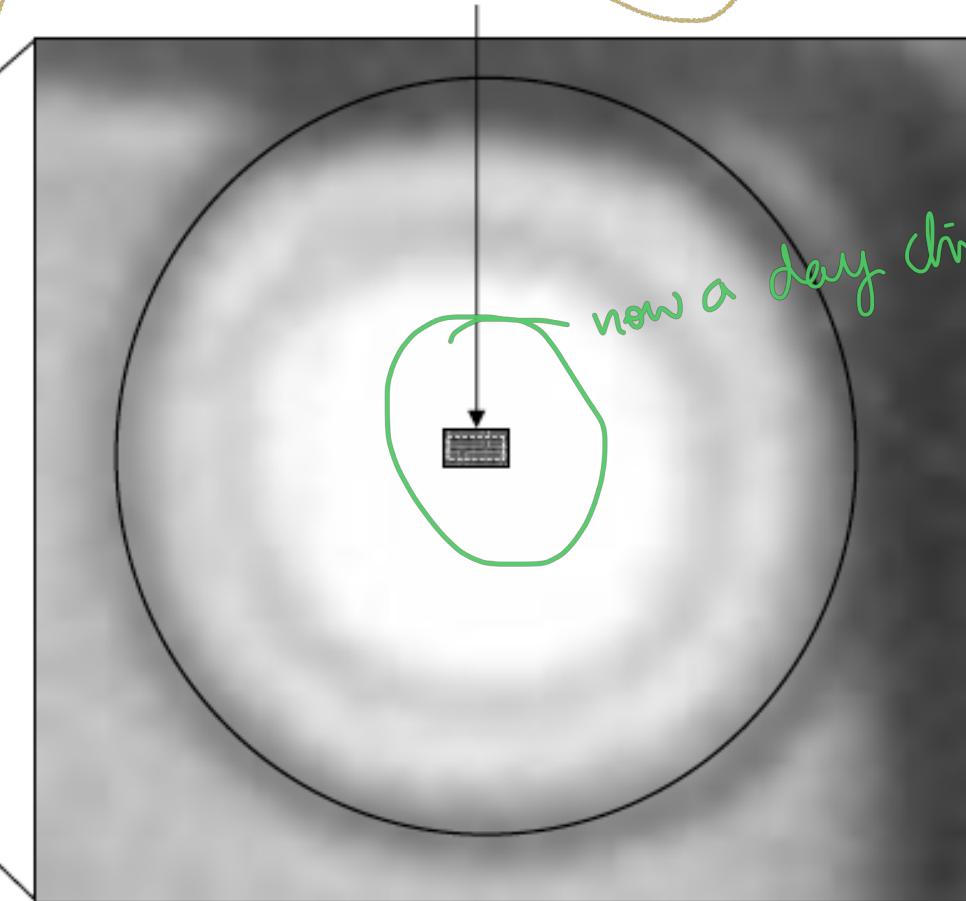
# Example of Scaling – 6T SRAM

1980 SRAM Cell:  $1700 \mu\text{m}^2$       22nm SRAM Cell:  $0.092 \mu\text{m}^2$

Intel



10000X

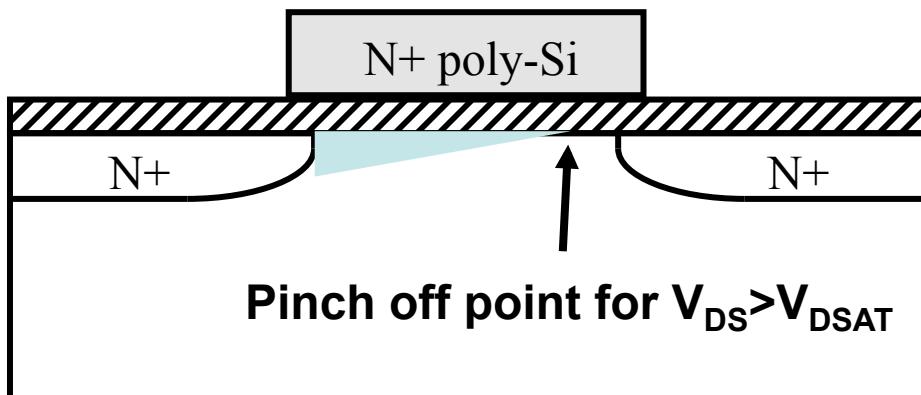
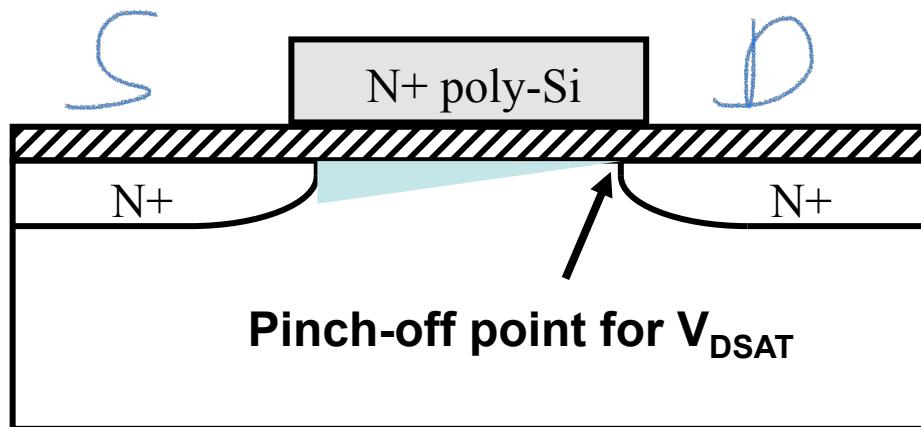


Small enough that a 2011 22nm SRAM cell is dwarfed by a 1980 SRAM cell CONTACT

# Channel Length Modulation

---

Beyond pinch-off ( $V_{DS} > V_{DSAT}$ ), pinch off point moves towards S

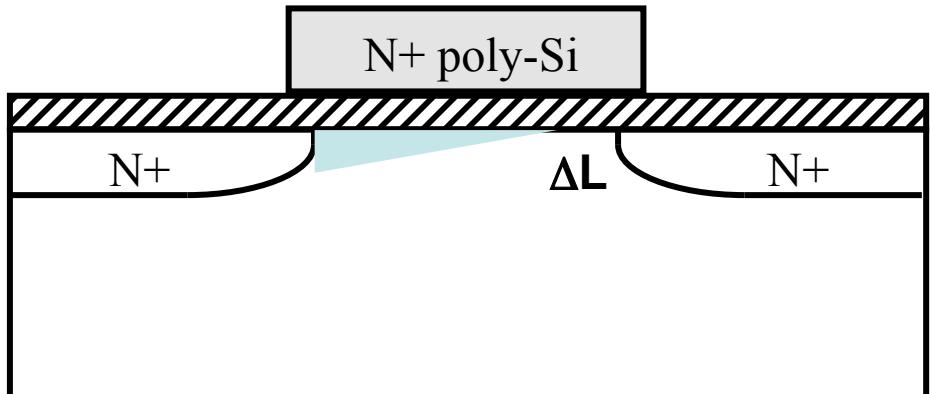
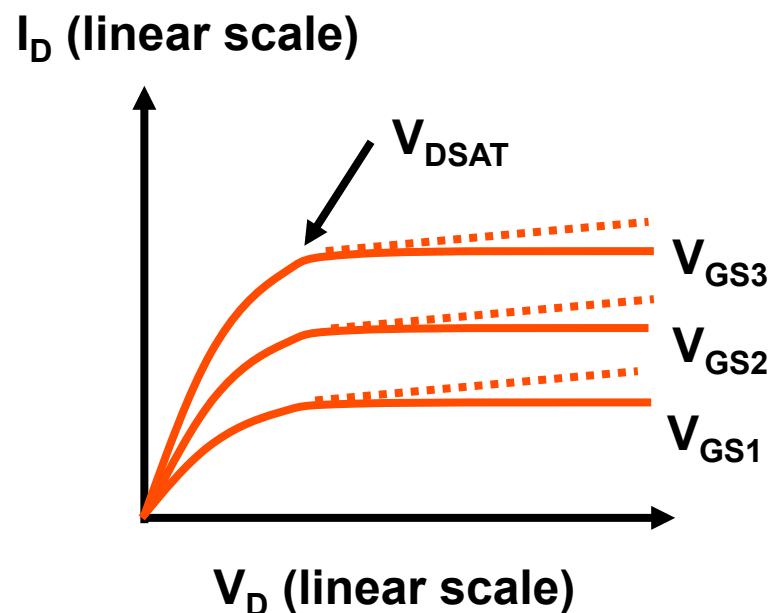


Channel length get modulated (shortened) by  $\Delta L$  (pure depletion layer), voltage drop in  $\Delta L$  is ( $V_{DS} - V_{DSAT}$ )

Current conduction similar to reverse-bias PN junction

# Channel Length Modulation

Reduction in  $L_{\text{eff}}$  to  $L_{\text{eff}} - \Delta L$



$$\Delta L = \sqrt{\frac{2\epsilon_{\text{si}}(V_D - V_{DSAT})}{qN_A}}$$

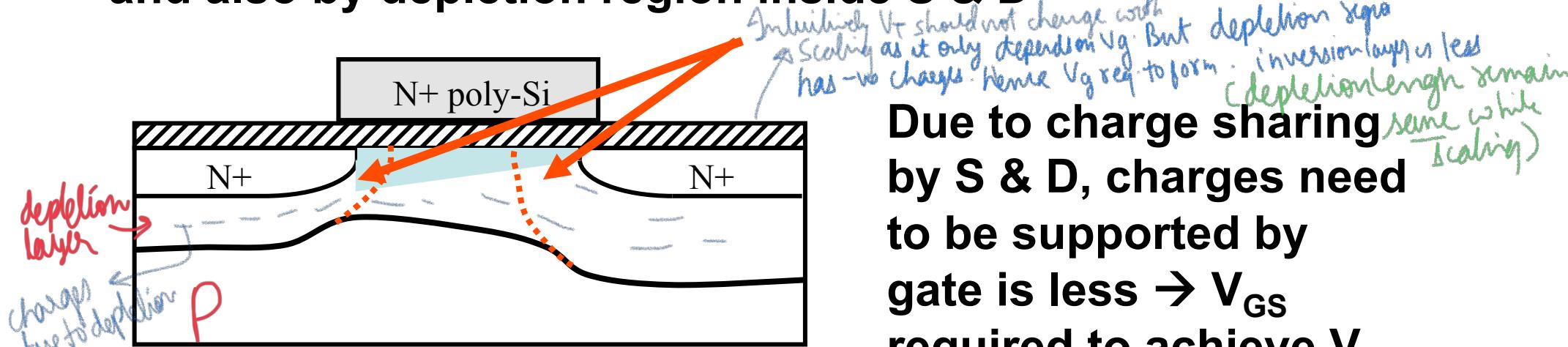
change in channel length due to pinch off

$$I_D = \mu_{\text{eff}} C_{OX} \frac{W}{L_{\text{eff}} - \Delta L} \frac{(V_{GS} - V_T)^2}{2m}$$

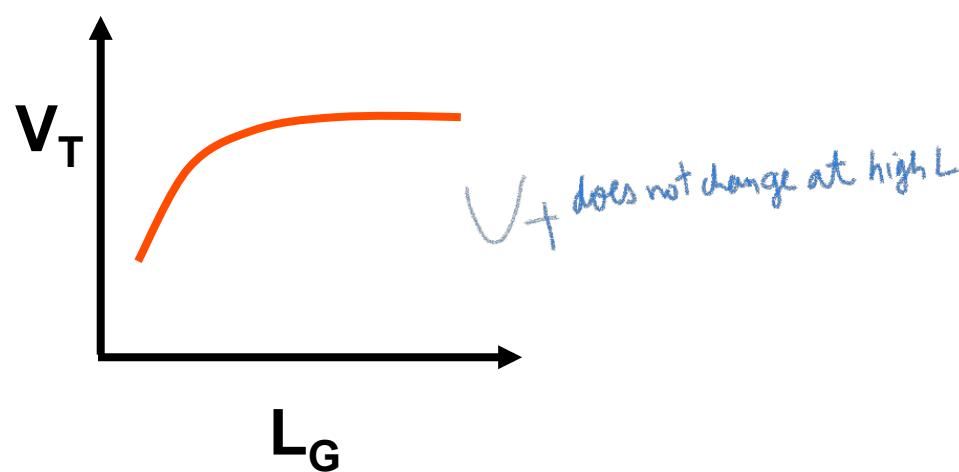
Slight increase in  $I_D$  for  $V_{DS} > V_{DSAT}$

# Charge Sharing by S / D (V<sub>T</sub> Roll-Off)

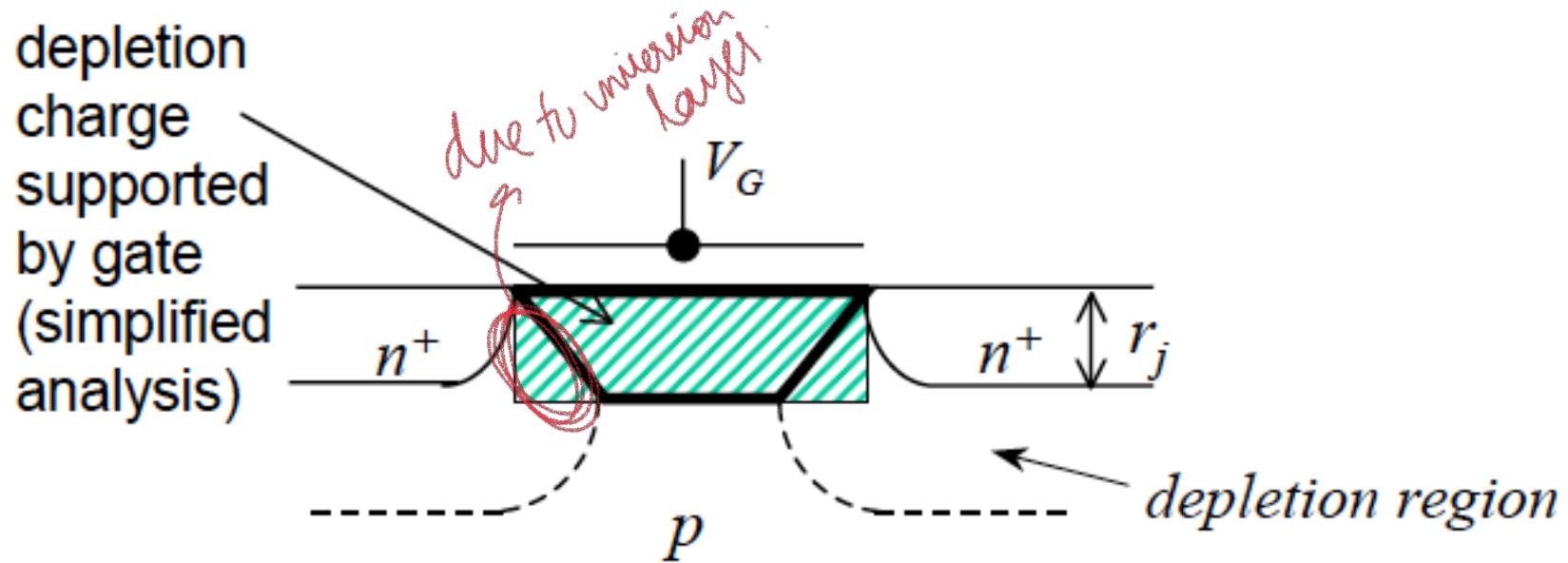
Charge in channel near S & D are supported by gate and also by depletion region inside S & D



Due to charge sharing by S & D, charges need to be supported by gate is less  $\rightarrow V_{GS}$  required to achieve  $V_T$  is less  $\rightarrow$  Significant effect as  $L_G$  is reduced  $\rightarrow V_T$  roll off



# Charge Sharing by S / D ( $V_T$ Roll-Off)



Large L:



Depletion charge supported by S/D

Small L:



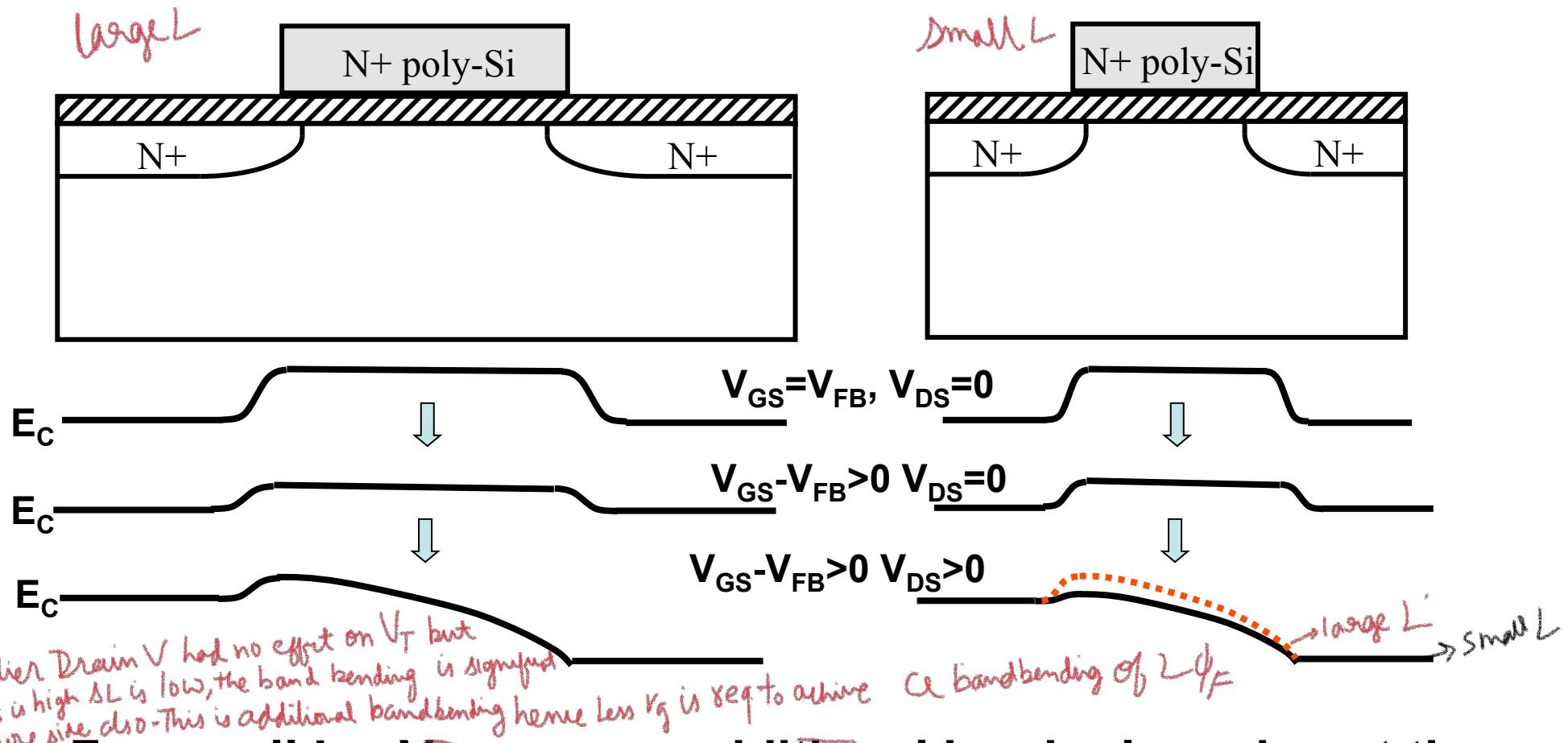
Depletion charge supported by S/D

Lower charge supported by gate

Can not neglect charge of depletion layer

$$\text{Total charge} = Q \text{ by gate} + Q \text{ by S/D}$$
$$\approx Q \text{ by gate}$$

# Drain Induced Barrier Lowering (DIBL)



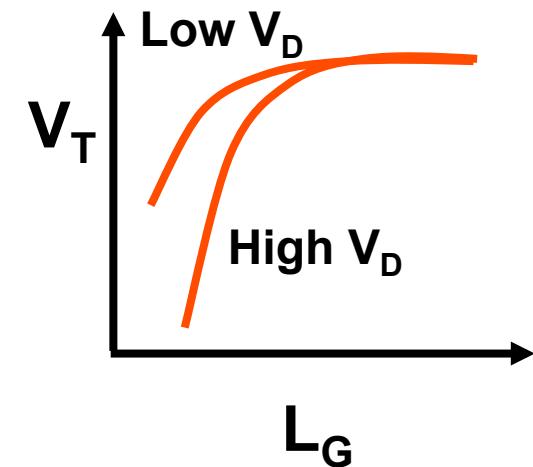
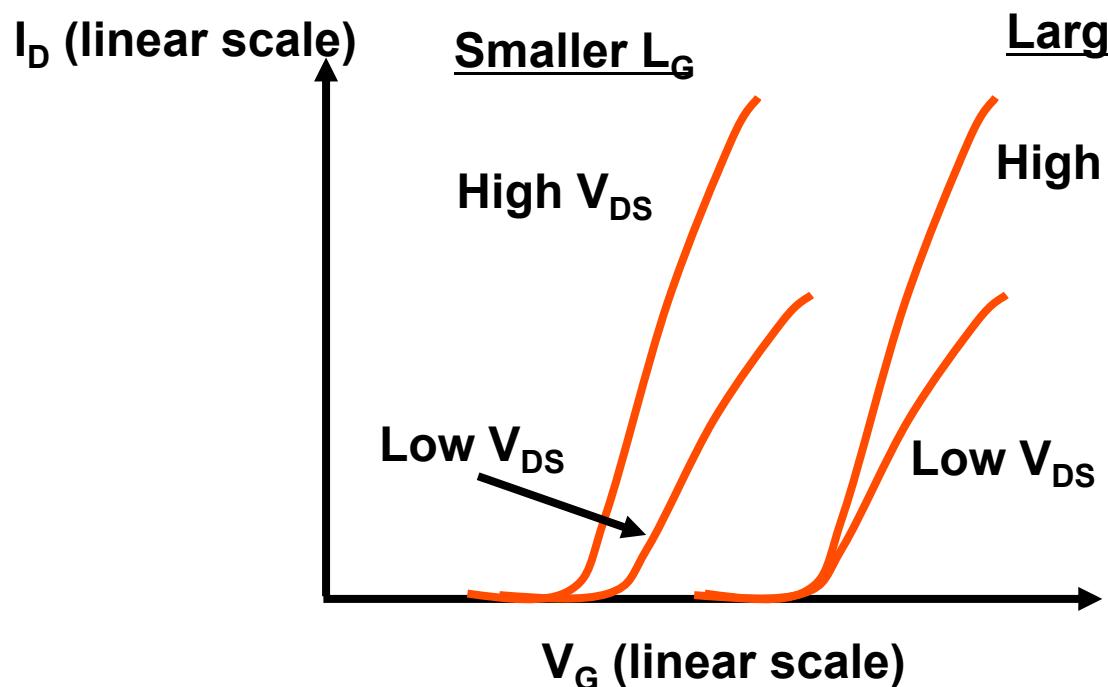
**For small  $L_G$ ,  $V_{DS}$  causes additional barrier lowering at the S end → Reduced amount of barrier to lowered by  $V_{GS}$  to achieve  $V_T \rightarrow V_T$  (saturation) is lower than  $V_T$  (linear)**

*Why*

# $V_T$ Roll-Off and DIBL

$V_T$  roll-off → reduction in  $V_T$  at smaller  $L_G$

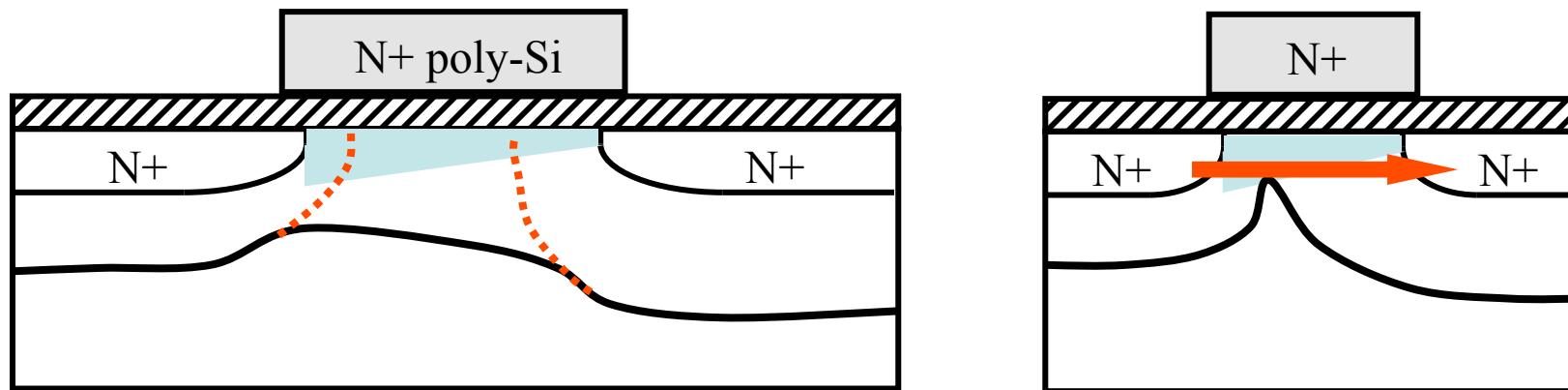
DIBL → reduction in  $V_T$  at higher VDS (especially for smaller  $L_G$ )



# Punch Through

---

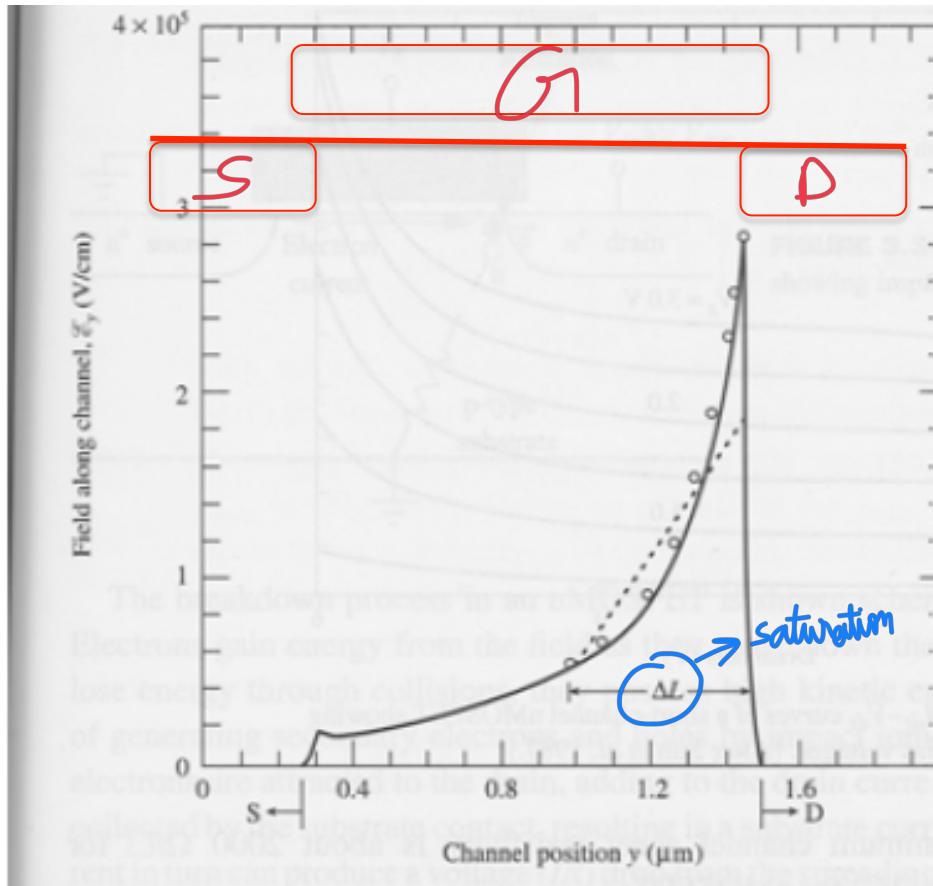
**For short  $L_G$  at high  $V_{DS}$ , S and D depletion layer can touch each other (away from interface) → direct flow of carriers from S to D with no control of G**



**Punch through is a sub-surface phenomenon (happens away from the interface) → Same impact on  $V_T$  as DIBL**

# Lateral Electric Field

---



Dimension scaling >  
Voltage Scaling (discussed  
later)

High lateral electric field  
near the drain junction

(1) Velocity saturation

(2) Impact ionization

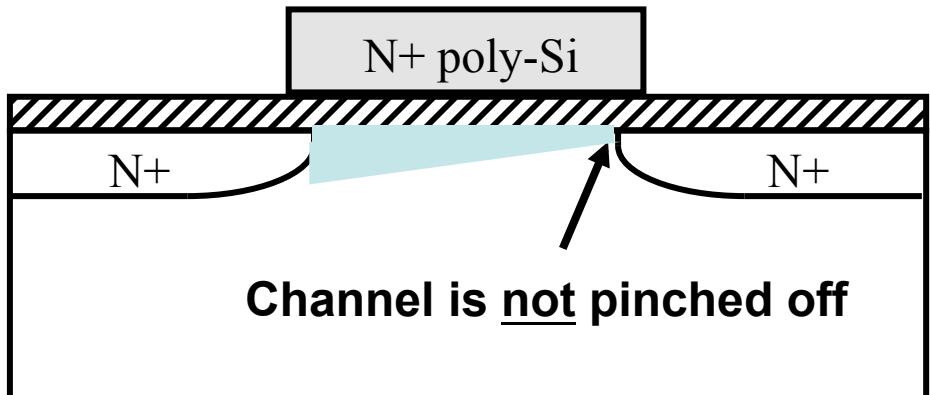
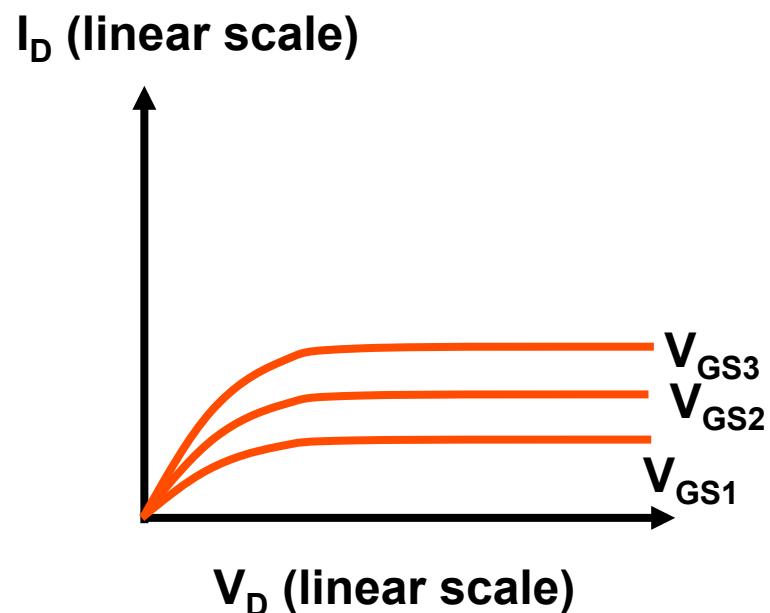
(3) Hot Carrier Injection

(4) Parasitic bipolar action

Ref: Taur and Ning

# Velocity Saturation

For short  $L_G$ ,  $E_y$  close to D junction becomes very high as  $V_{DS}$  is increased



Electron velocity saturates (does not increase with  $E_y$ )

$$I_{DSAT} \sim Q v_{sat} \sim C_{OX} (V_{GS} - V_T) v_{sat}$$

*saturation velocity*

$I_{DSAT}$  varies linearly with overdrive (not square law)

# Mobility at High Vertical and Lateral Field

---

**Reduction in mobility due to transverse field**

$$\mu_{eff} = \frac{\mu_0}{1 + \alpha(V_{GS} - V_T)}$$

**Reduction in mobility due to lateral field (reduction in carrier velocity)**

$$v = \frac{\mu_{eff} E}{[1 + (E / E_c)^n]^{1/n}}$$

n=2 for electrons

n=1 for holes

v<sub>sat</sub> =  $\mu_{eff} E_c$  *At  $E_c$*

$$\vec{V} = \frac{U_{eff}}{1 + \frac{\Sigma}{\Sigma_c}} \quad \Sigma_c = \frac{V_{sat}}{\mu_{eff}}$$

## Analytical Solution for n=1

---

$$I_D = -WQ_i(V) \frac{\mu_{eff} dV/dy}{[1 + (\mu_{eff} / v_{sat}) dV/dy]}$$

$$I = Wq\vec{V}$$

$$I_D = -[\mu_{eff} WQ_i(V) + \frac{\mu_{eff} I_D}{v_{sat}}] \frac{dV}{dy}$$

$$I_D = \frac{-\mu_{eff} (W/L) \int_0^{V_{DS}} Q_i(V) dV}{1 + (\mu_{eff} V_{DS} / v_{sat} L)}$$

Integrate from 0 to L and 0 to  $V_{DS}$

$$Q_i(V) = -C_{ox}(V_G - V_T - mV)$$

# Analytical Solution for n=1

---

$$I_D = \frac{\mu_{eff} C_{OX} (W/L) [V_G - V_T - \frac{m}{2} V_{DS}] V_{DS}}{1 + (\frac{\mu_{eff} V_{DS}}{v_{sat} L})}$$

$$V_{DSAT} = \frac{2(V_G - V_T)/m}{1 + \sqrt{1 + 2\mu_{eff}(V_G - V_T)/(mv_{sat}L)}}$$

**Solve for  
 $dI_D/dV_{DS} = 0$**

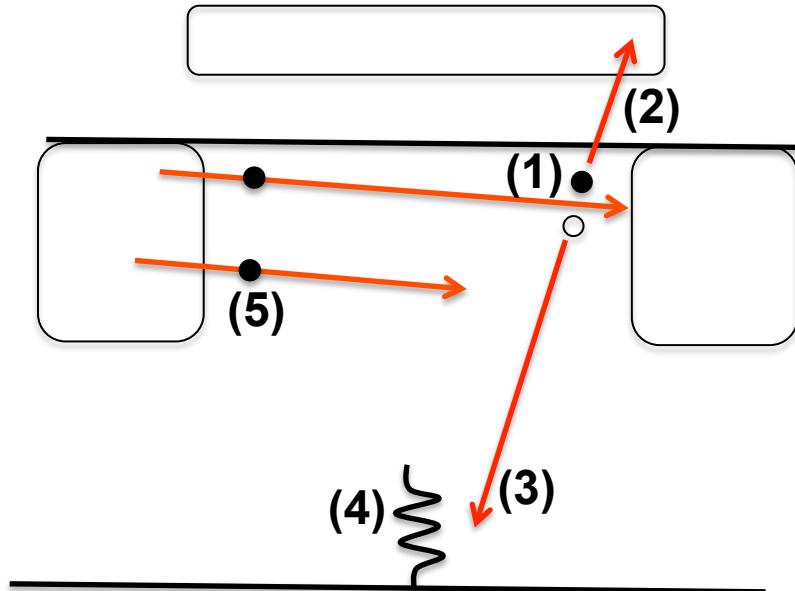
$$I_{DSAT} = C_{OX} W v_{sat} (V_G - V_T)$$

$$* \frac{\sqrt{1 + 2\mu_{eff}(V_G - V_T)/(mv_{sat}L)} - 1}{\sqrt{1 + 2\mu_{eff}(V_G - V_T)/(mv_{sat}L)} + 1}$$

**Substitute  $V_{DSAT}$**

# High Lateral Field Effects

---



- (1) Impact ionization, creation of electrons and holes**
- (2) Electron injection into gate – gate current, oxide damage**
- (3) Hole injection towards substrate – substrate current**
- (4) Forward bias of source substrate junction**
- (5) Additional electron injection from source**

# Scaling → Necessities and Challenges

---

**Reduce  $L_G$  → Higher speed, density**

**S & D junctions come closer → short channel effects  
(channel length modulation, velocity saturation, DIBL,  $V_T$   
roll off, punch through)**

**Solution:**

**Scale vertical dimensions →  $T_{ox}$  &  $X_j$**

**Increase doping, reduce operating voltage**

# Scaling → Necessities and Challenges

---

**On current → keep it high**

$$I_{DSAT} = \mu_{eff} C_{OX} \frac{W}{L_{eff}} \frac{(V_{GS} - V_T)^2}{2m}$$

**Classical**

$$I_{DSAT} = \mu_{eff} C_{OX} v_{sat} (V_{GS} - V_T)$$

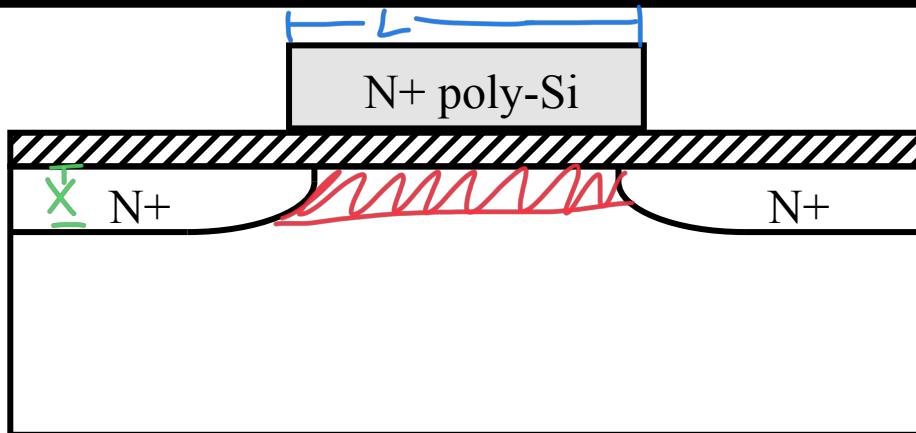
**Velocity  
saturation**

**Off current → keep it low**

$$I_D = \mu_{eff} C_{OX} \frac{W}{L_{eff}} (m-1) \left( \frac{kT}{q} \right)^2 e^{q(V_G - V_T)/mkT} [1 - e^{-qV_{DS}/kT}]$$

$$m = 1 + C_D / C_{OX} = 1 + 3T_{OX} / W_{DMAX}$$

# Constant Electric Field Scaling Rules



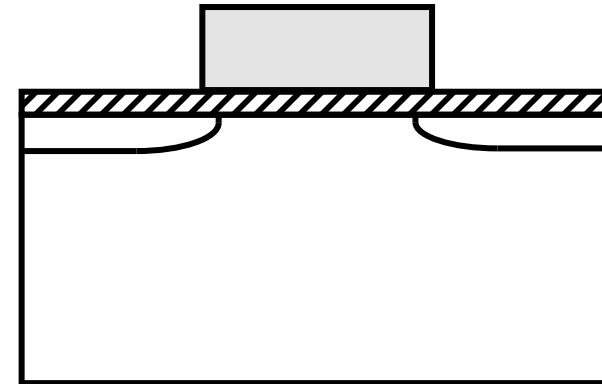
Long channel JUNCTION DEPTH

$W, L_G, T_{ox}, X_j$  Dimensions

$N_A, N_D$  Doping

$V_{DD}$  Supply voltage

Dennard, IBM



Short channel

$W/k, L_G/k, T_{ox}/k, X_j/k$

$N_A \cdot k, N_D \cdot k$

$V_{DD}/k$

# Scaling of Device Parameters ( $k > 1$ )

---

Electric field  $\rightarrow 1$  CONSTANT ELECTRIC FIELD SCALING

Carrier velocity  $\rightarrow 1$

Depletion layer width  $\rightarrow 1/k$   $\left( \frac{4\epsilon_s \phi}{w N_a} \right)$  should be  $\frac{1}{\sqrt{k}}$

Capacitance  $\rightarrow 1/k$  total  $C = \frac{\sum w L}{T_{ox}}$

Inversion charge density  $\rightarrow 1$

Drift current (on-state, per W)  $\rightarrow 1$

Diffusion current (off-state, per W)  $\rightarrow k$

?

# Scaling of Circuit Parameters

---

**Delay time ( $C \cdot V/I$ )  $\rightarrow 1/k$**

**Power dissipation ( $V \cdot I$ )  $\rightarrow 1/k^2$**

**Power-delay product  $\rightarrow 1/k^3$**

**Circuit density (1/Area)  $\rightarrow k^2$**

**Power density ( $P/A$ )  $\rightarrow 1$**

# Trouble: Non-Scaling Factors

---

## Threshold voltage

$$V_T = V_{FB} + 2\phi_F + \frac{\sqrt{2q\varepsilon_{si}N_A(2\phi_F)}}{C_{OX}}; \quad V_{FB} = -(E_G / 2q + \phi_F)$$

## Subthreshold current

$$I_D = \mu_{eff} \cdot C_{OX} \cdot \frac{W}{L_{eff}} (m-1) \left( \frac{kT}{q} \right)^2 e^{q(V_{GS}-V_T)/mkT} (1 - e^{-qV_{DS}/kT})$$

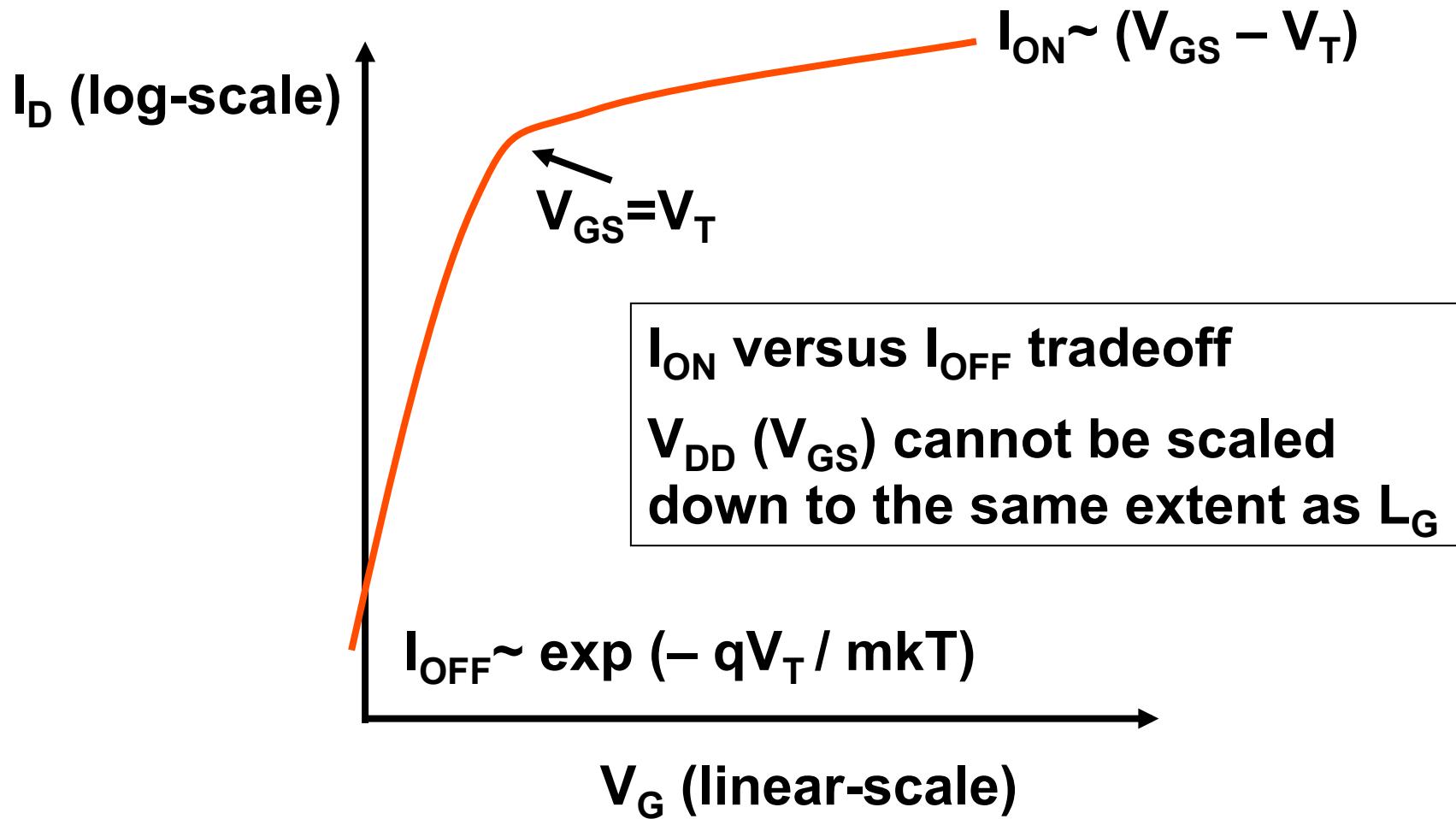
$$m = 1 + \frac{\sqrt{\varepsilon_{si} q N_A / 4\phi_F}}{C_{OX}}$$

**E<sub>G</sub>, kT/q does not scale**

# The $V_T$ scaling dilemma

---

Choice of proper  $V_T$



# Generalized scaling rules ( $k>1$ , $\alpha>1$ )

---

Long channel

$W, L_G, T_{ox}, X_j$

Dimensions

$N_A, N_D$

Doping

$V_{DD}$

Supply voltage

Short channel

$W/k, L_G/k, T_{ox}/k, X_j/k$

$N_A \cdot (\alpha k), N_D \cdot (\alpha k)$

$V_{DD} \cdot (\alpha/k)$

“Extra” increase in doping

“Smaller” reduction in supply voltage

# Scaling of device parameters ( $k>1$ , $\alpha>1$ )

---

**Electric field  $\rightarrow \alpha$**

**Carrier velocity  $\rightarrow \alpha$  (or 1 for velocity saturation)**

**Depletion layer width  $\rightarrow 1/k$**

**Capacitance  $\rightarrow 1/k$**

**Inversion charge density  $\rightarrow \alpha$**

**Drift current (on-state, per W)  $\rightarrow \alpha^2$  (or  $\alpha$  for velocity saturation)**

**Diffusion current (off-state, per W)  $\rightarrow \alpha.k$**

# Scaling of circuit parameters

---

conventional	(velocity saturation)
<b>Delay time (C.V/I) <math>\rightarrow 1/\alpha k</math></b>	<b>(or <math>1/k</math>)</b>
<b>Power dissipation (V.I) <math>\rightarrow \alpha^3/k^2</math></b>	<b>(or <math>\alpha^2/k^2</math>)</b>
<b>Power-delay product <math>\rightarrow \alpha^2/k^3</math></b>	
<b>Circuit density (1/Area) <math>\rightarrow k^2</math></b>	
<b>Power density (P/A) <math>\rightarrow \alpha^3</math></b>	<b>(or <math>\alpha^2</math>)</b>