

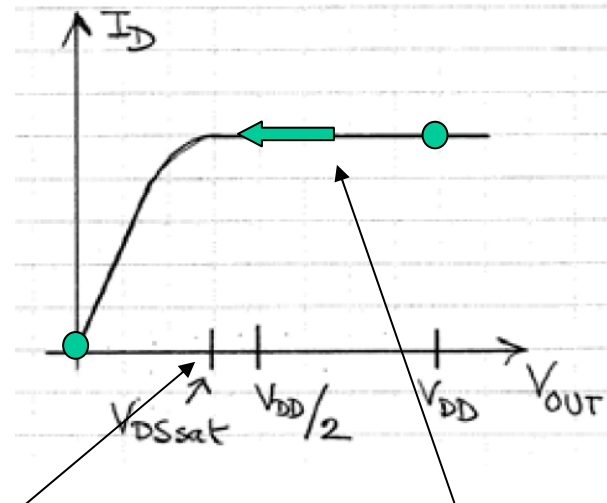
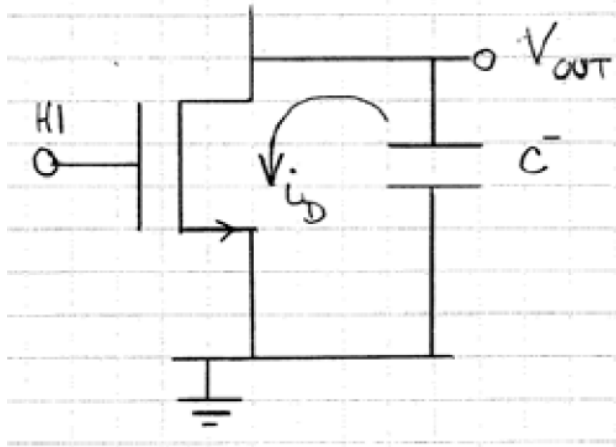
# High performance CMOS

## LECTURE 21

- MOSFET scaling
- short-channel effect
- shallow S & D, halos, SOI
- strain engineering and mobility

Sec.  
13.0

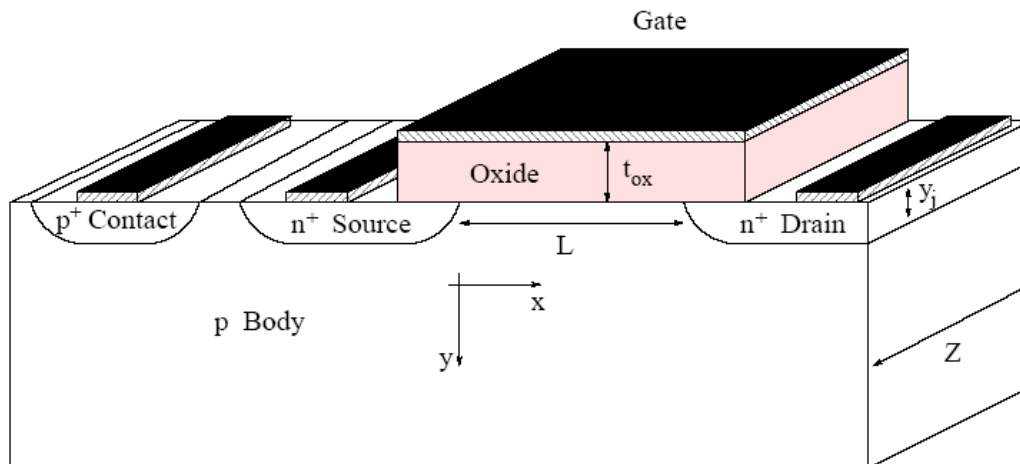
# Logic speed is about Q and I



$$\tau = \frac{C\Delta V}{I}$$

$$V_{DS} \equiv V_{DSsat} = \frac{V_{GS} - V_T}{m}$$

$$I_{Dsat} = \frac{Z}{L} C_{ox} \mu_{eff} \frac{(V_{GS} - V_T)^2}{2m}$$



Last 40 years:

Improve speed via  
 $L$ ,  $C_{ox}$ ,  $V_T$ .

What about  
 $\mu$  and  $y_j$ ?

# Co-ordinated scaling

		CMOS3 1987	CMOS180 2001	CMOS130 2002	CMOS90 2003	CMOS65 2005	CMOS45 2008	CMOS32 2011
L	nm	3000	180	130	90	65	45	32
YJ	nm	1000	160	39	28	19.6	14	11
TOX	nm	85	4.1	2.8	2.3	1.85*	1.75*	1.65*
NCH	cm <sup>-3</sup>	1.00E+16	3.90E+17	6.15E+17	8.37E+17	2.54E+18	3.24E+18	4.12E+18
VDD	V	5.0	1.8	1.2	1.0	1.0	1.0	1.0
VT**	V	0.95	0.47	0.35	0.24	0.42	0.47	0.51

\* based on relative permittivity of 3.9      \*\* “long-channel” threshold voltage

CMOS65/45/32 from <http://ptm.asu.edu/>

- VDD and L should scale together to keep
- Lower limit to  $V_T$  is set by
- TOX lowered to improve
- NCH raised to keep  acceptable
- YJ reduced to combat actual threshold voltage lowering by

Can these trends continue?

# CMOS: the Industrial drive

## CMOS Device Scaling Demonstration

90nm Node  
2003



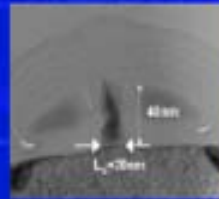
65nm Node  
2005



50nm Length  
(IEDM2002)

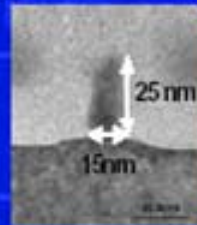
30nm  
Prototype  
(IEDM2000)

45nm Node  
2007



20nm Prototype  
(VLSI2001)

32nm Node  
2009



15nm Prototype  
(IEDM2001)

22nm Node  
2011



10nm Prototype  
(DRC 2003)

16 nm node  
2013



TBD

11nm node  
2015



TBD

8 nm node  
2017

Nodes relate to the DRAM half pitch, i.e., the width, and space in between, metal lines connecting DRAM bit cells

Intel research devices scale to 10nm (16nm node)  
Channel engineering solutions (Nanowires/Nanotubes)  
are being investigated to extend device scaling through  
end of next decade

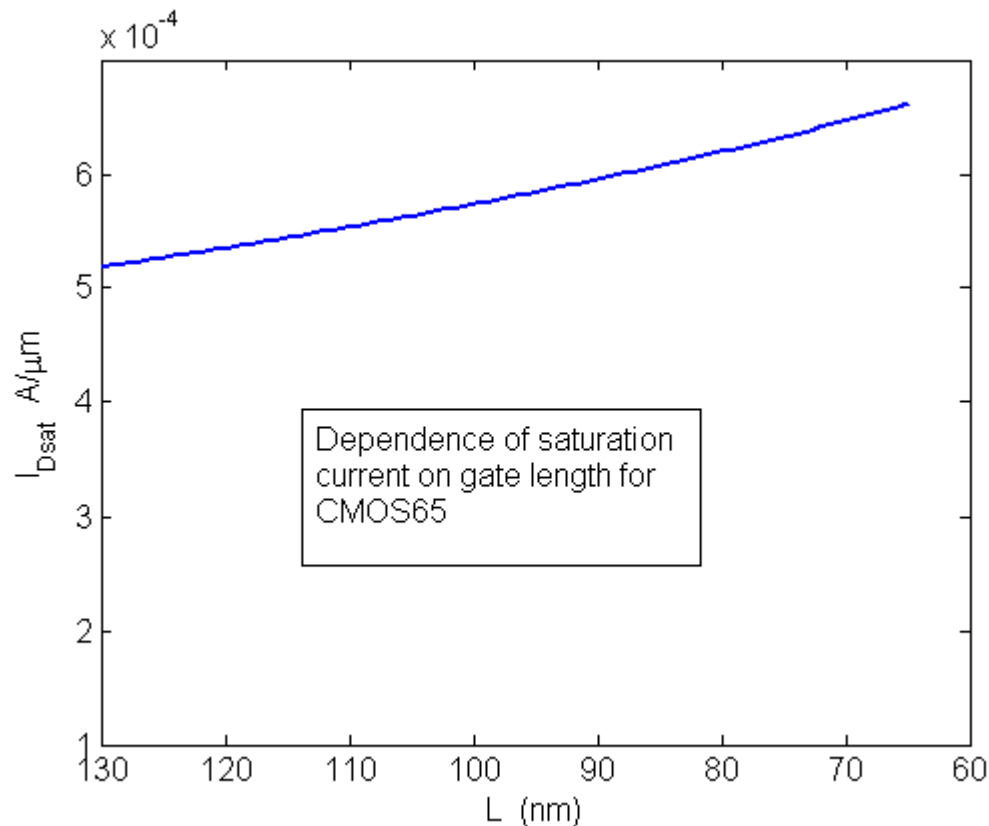
intel.

Source: Intel; Morales and  
Lieber  
Science, 279, 208, 1998

# Shrinking L no longer helps much for

## $I_{Dsat}$

$$I_{Dsat} = ZC_{ox}(V_{GS} - V_T) \cdot v_{sat} \frac{\sqrt{1 + 2\mu_{eff}(V_{GS} - V_T)/(mv_{sat}L)} - 1}{\sqrt{1 + 2\mu_{eff}(V_{GS} - V_T)/(mv_{sat}L)} + 1}$$



Why is this?

Why continue to shrink?

## 3 major concerns for digital CMOS

$$I_{Dsat} = ZC_{ox}(V_{GS} - V_T) \cdot v_{sat} \frac{\sqrt{1 + 2\mu_{eff}(V_{GS} - V_T)/(mv_{sat}L)} - 1}{\sqrt{1 + 2\mu_{eff}(V_{GS} - V_T)/(mv_{sat}L)} + 1}$$

### Concerns:

- L cannot be further reduced without adversely affecting  $V_T$  and  $I_{subt}$
- Some other way needs to be found to increase  $I_{ON}$
- TOX cannot be further reduced without causing excessive gate leakage current

### Solutions:

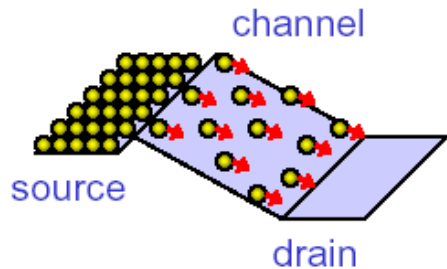
1. Controlling  $V_T$  and  $I_{subt}$  via suppression of the short-channel effect
2. Increasing  $I_{ON}$  via mobility improvement
3. Reducing gate leakage via thicker, high- $k$  dielectrics

1,2 started at  
CMOS90

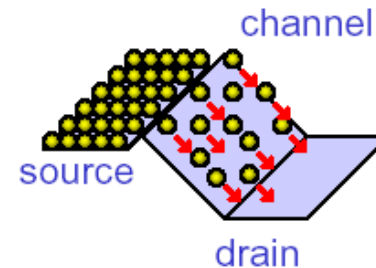
3 is new to  
CMOS45

# Moving More Charge in Less Time

$$I_{dsat} \approx \frac{1}{2} \mu C_{ox} \frac{W}{L} (V_{GS} - V_T)^2$$

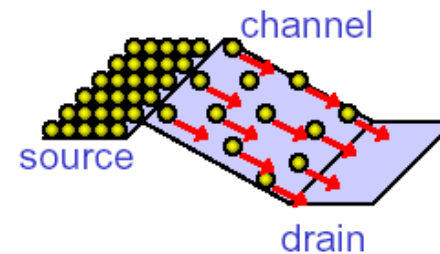


Decrease  $L$  – shorten distance to finish line



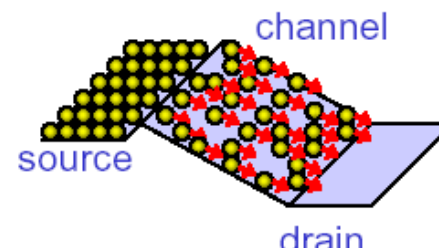
lithography  
scaling

Increase  $\mu$  – make electrons travel faster



strain engineering  
faster materials

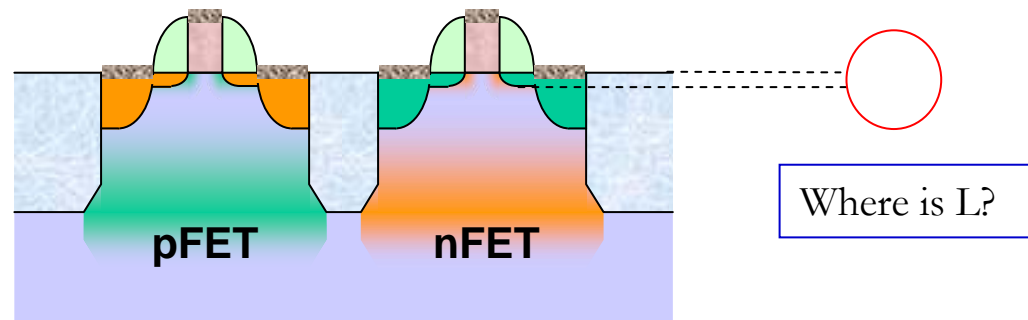
Increase  $C_{ox}$  – move more electrons



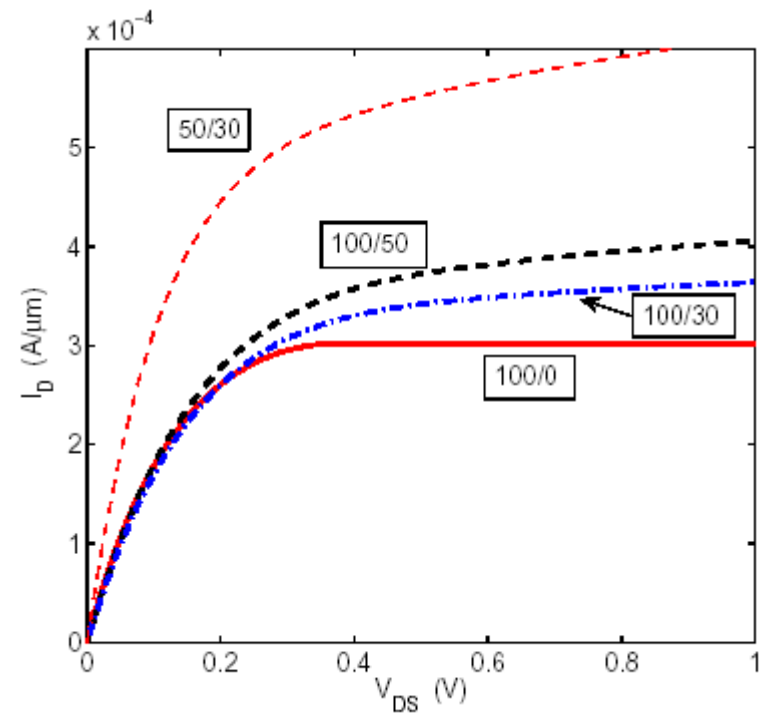
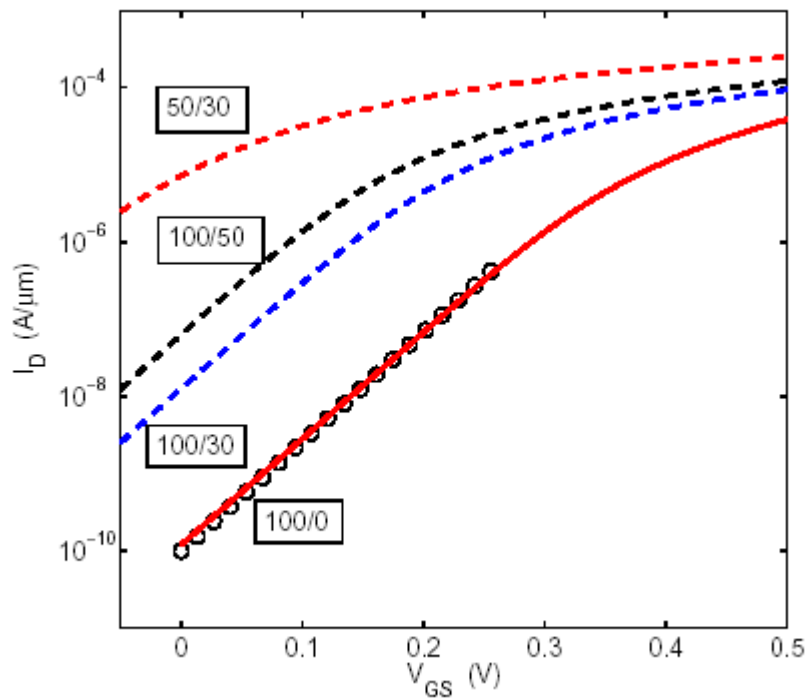
high-K dielectric  
metal gate

Sec.  
13.1.7

# Lithography scaling: reducing $L$ and $y_j$



The effect of changing  $L/y_j$

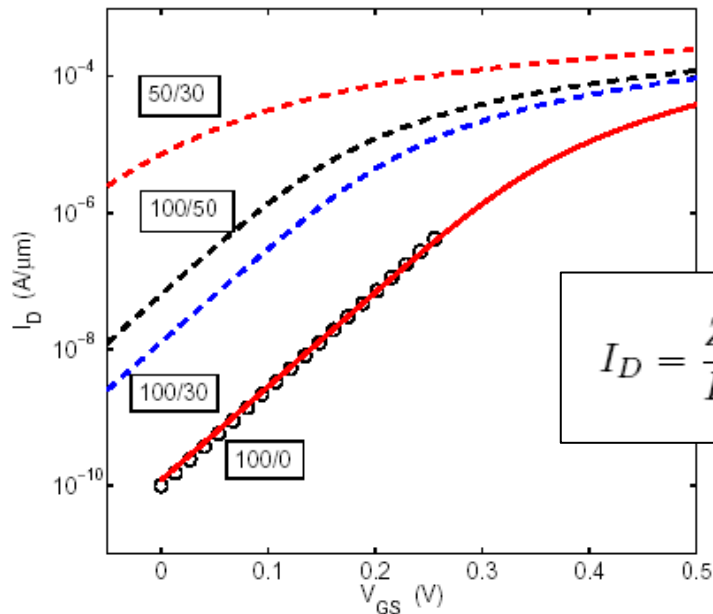


What is happening?

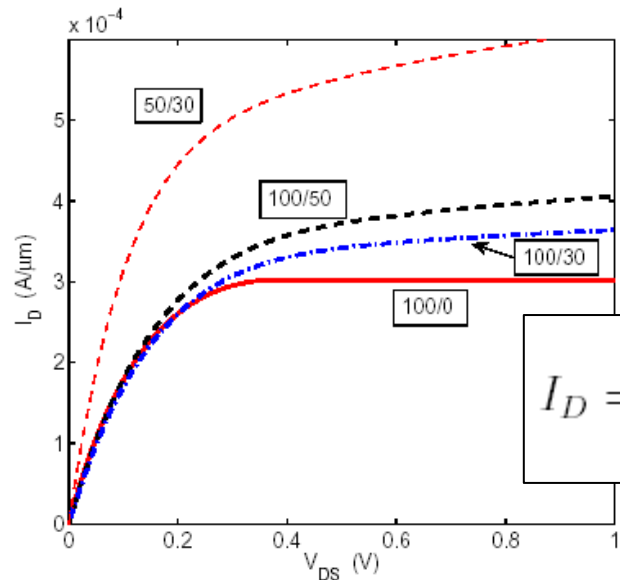


Sec.  
13.1.7

# Which model parameter is changing?



$$I_D = \frac{Z}{L} \mu_{\text{eff}} \left( \frac{k_B T}{q} \right)^2 C_{\text{ox}} (m - 1) e^{(V_{GS} - V_T)/mV_{\text{th}}} \left\{ 1 - e^{-V_{DS}/V_{\text{th}}} \right\}$$



$$I_D = Z C_{\text{ox}} \left[ V_{GS} - V_T - m \frac{V_{DS}}{2} \right] \cdot \mu_{\text{eff}} \frac{V_{DS}}{L + (\mu_{\text{eff}} V_{DS}/v_{\text{sat}})}$$

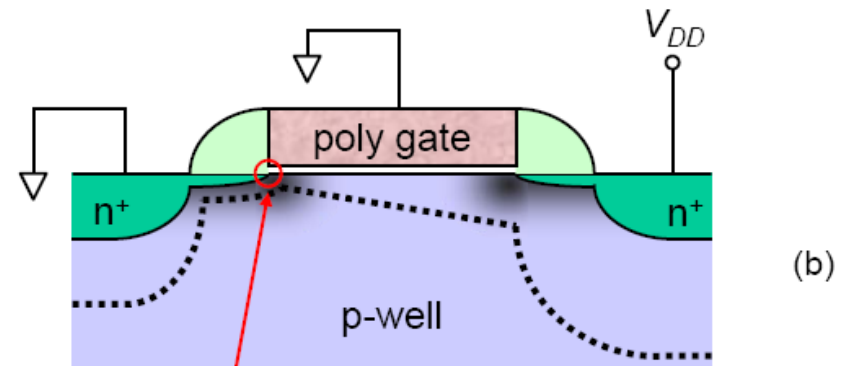
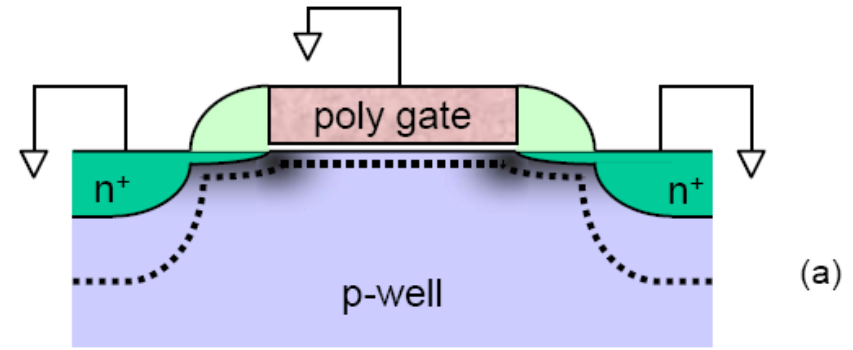
Sec.  
13.1.7

# The Short-Channel Effect

- The change in characteristics with  $y_j$  occurs at short  $L$ .
  - At short  $L$ , the characteristics also change with  $L$ .
  - These changes are known as the
- 
- They indicate a change in  $\psi_s(0)$  due to  $V_{DS}$ .

Why is this undesirable?

What can be done to avoid it?

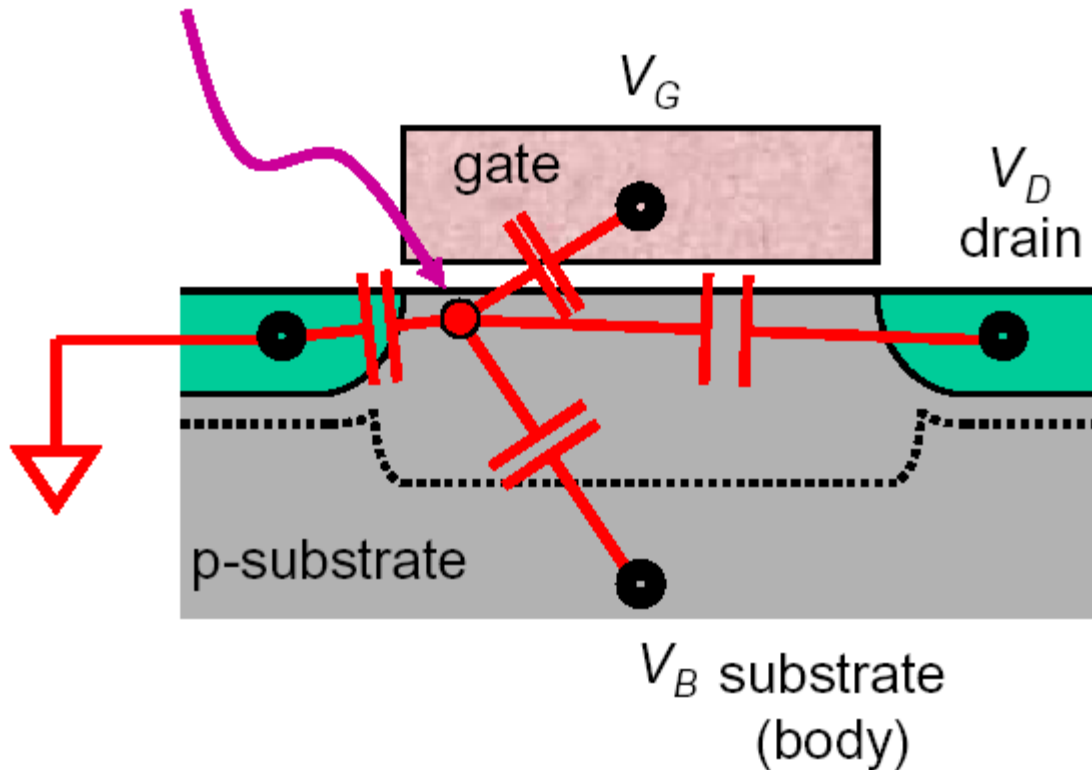


reduction of electron barrier  
height in conduction band (CB)  
at edge of source

# The Short-Channel Effect and Capacitance

$$\psi_s = f(L, y_j, V_{DS}) \quad \therefore V_T = f(L, y_j, V_{DS})$$

How is source-to-channel barrier height affected?



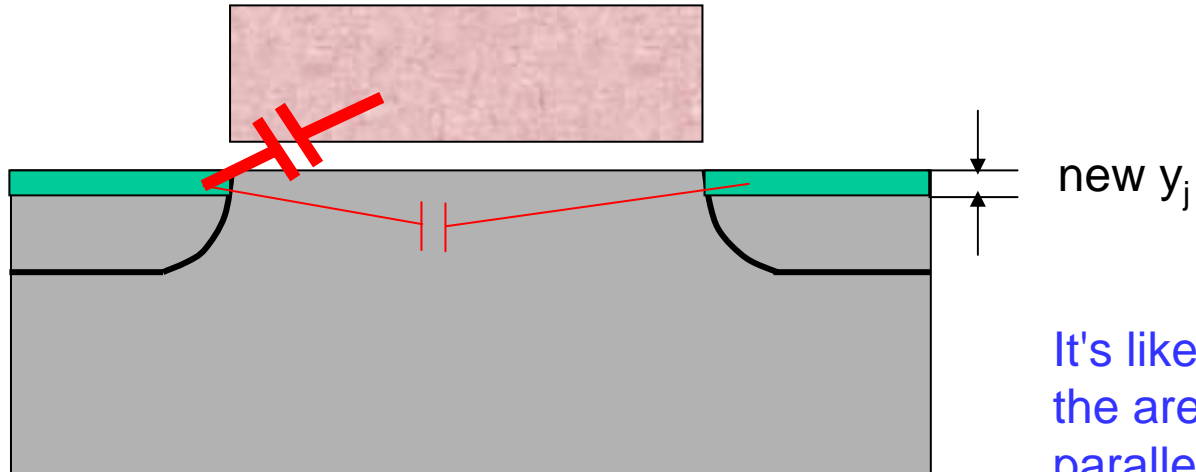
$\psi_s$  is determined by capacitive coupling via  $C_{ox}$  and  $C_{body}$ ,

**AND**

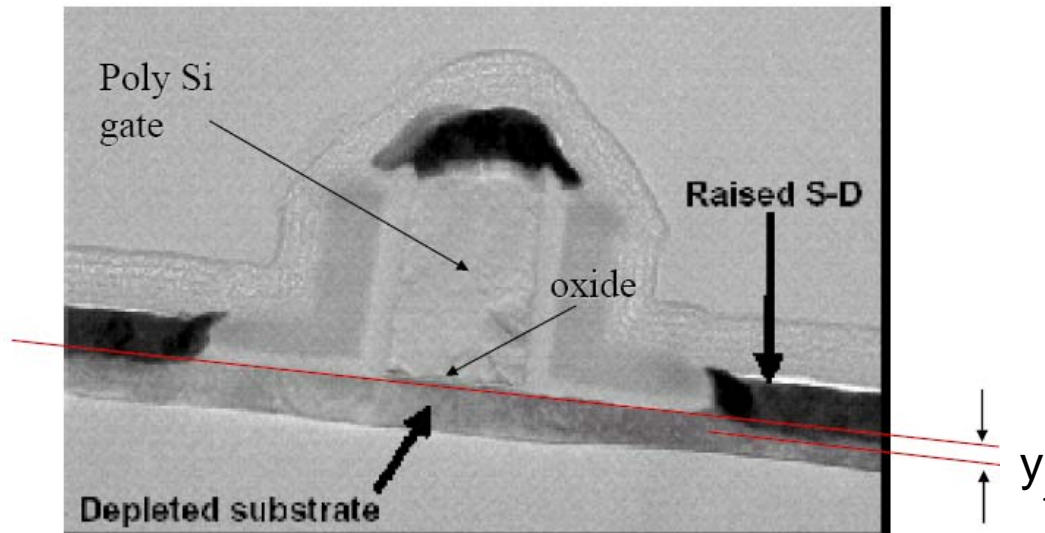
by capacitive coupling via  $C_{DS}$

Sec.  
13.1.7

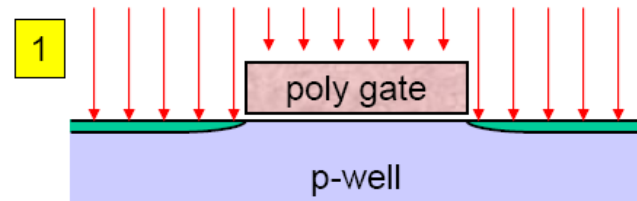
## Reduce $C_{DS}$ by shrinking $y_j$



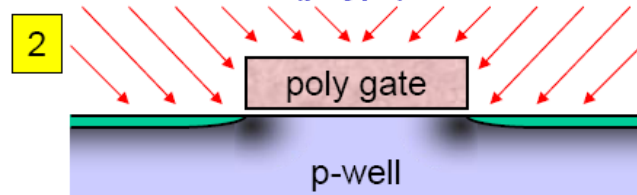
It's like reducing  
the area of a  
parallel plate  
capacitor



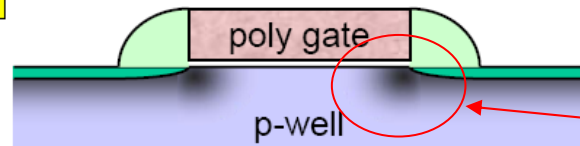
# Reduce $C_{DS}$ by screening $E_x$



self-aligned high-tilt halo/pocket implant  
(p-type)

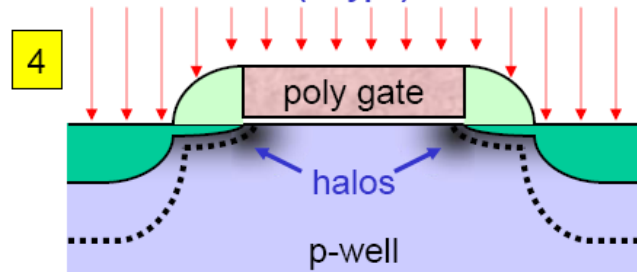


3 dielectric spacer formation



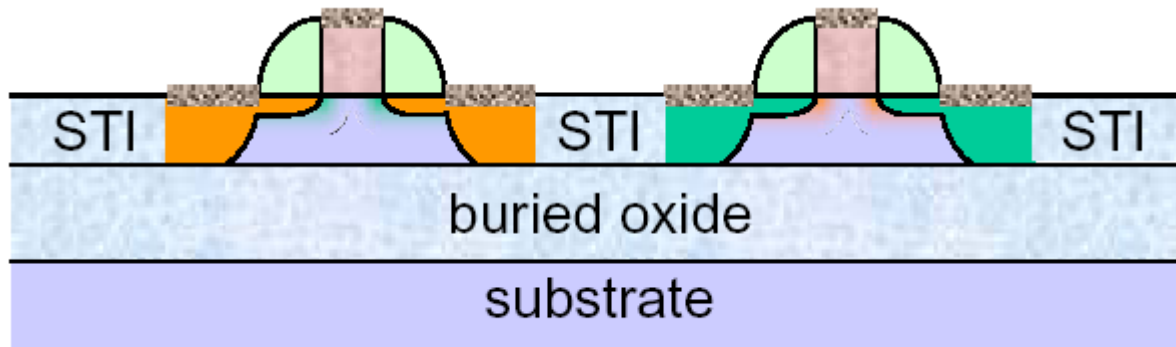
What is the  
doping density  
here?

self-aligned source/drain implant  
(n-type)

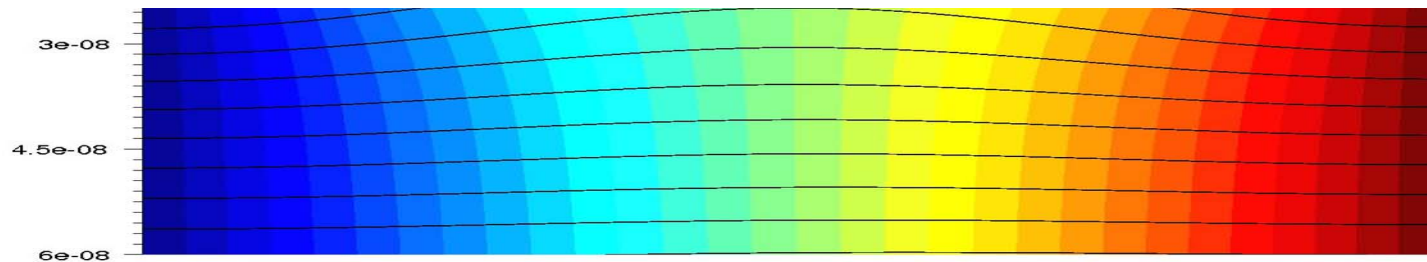
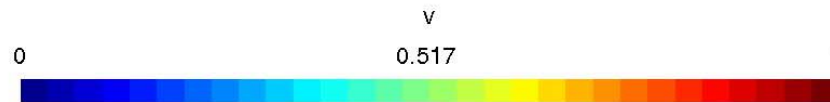
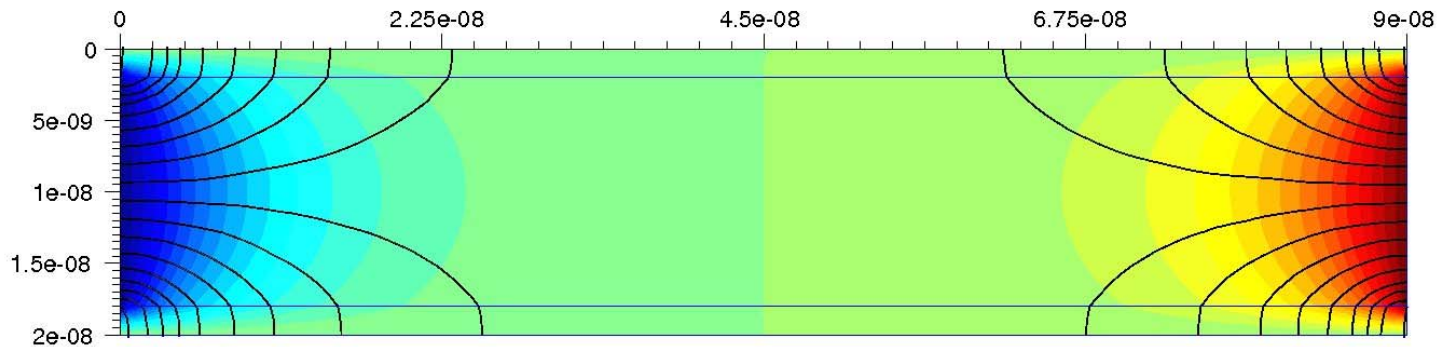


Sec.  
13.1.7

# Using SOI to beat SCE



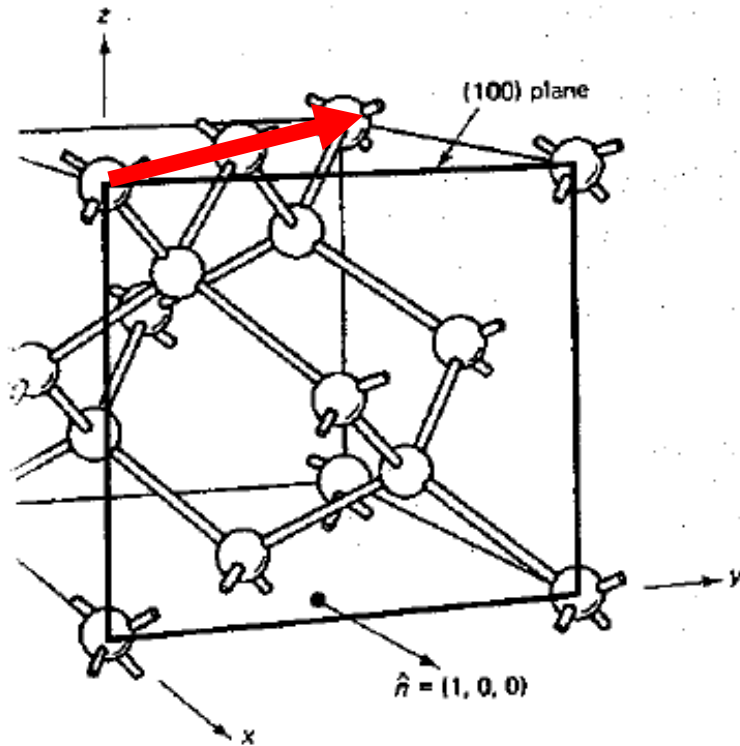
Alvin  
Loke



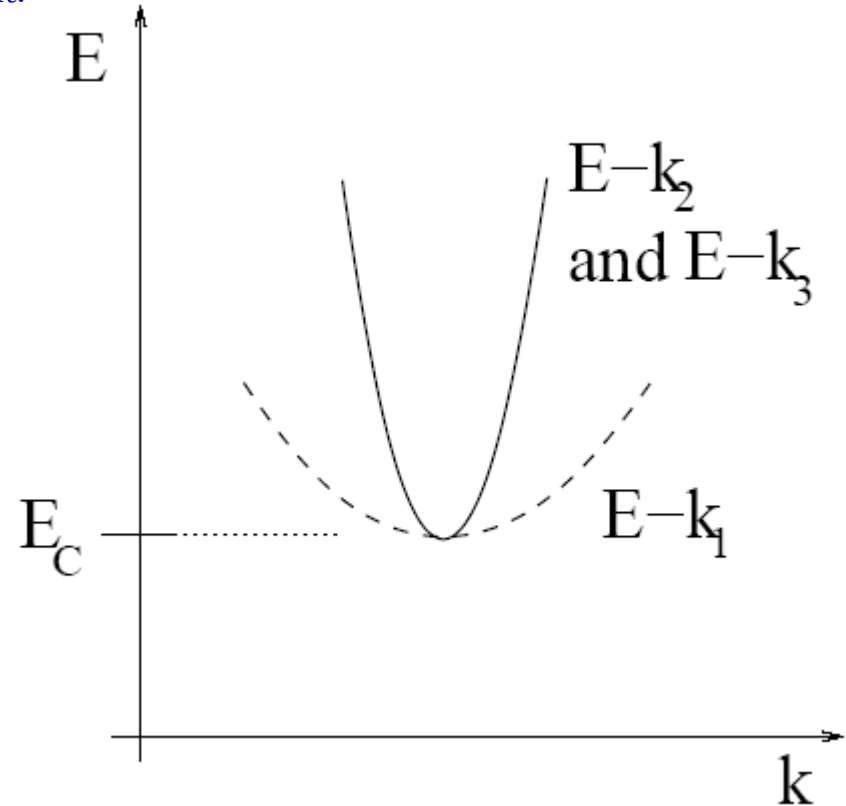
Daryl Van  
Vorst

# Strain engineering: improving $\mu$

Exploit the fact that  $m^*$  is direction-dependent.



Apply stress in  $\langle 110 \rangle$  to a (100) surface.



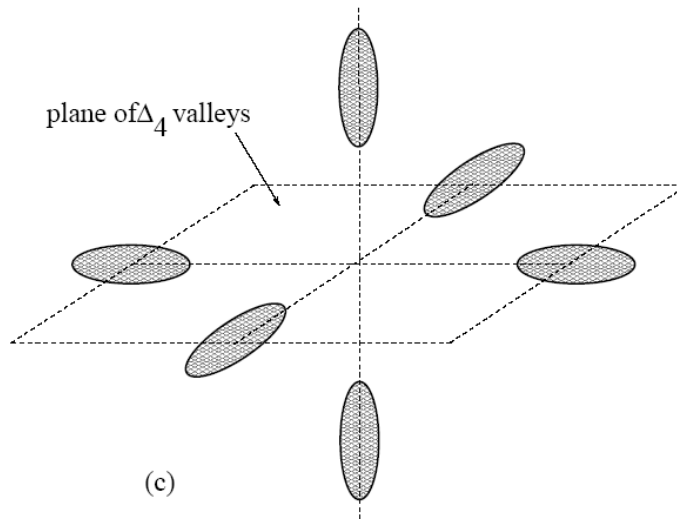
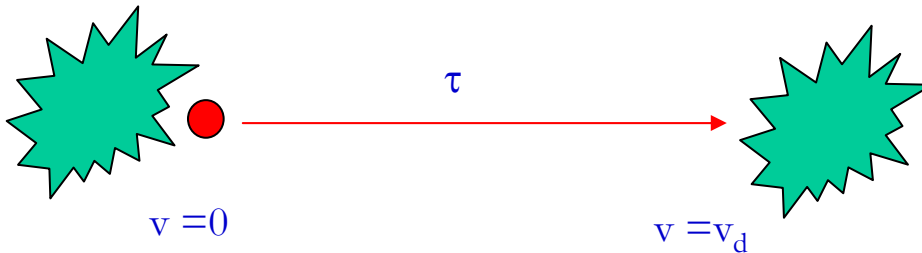
- $k_1$  is a  $\langle 110 \rangle$  direction
- $k_2$  and  $k_3$  are orthogonal at the point of the energy minimum  $E_C$

Which direction has the higher effective mass?

Sec.  
13.1.3

# Conductivity effective mass $m_c^*$

Electron accelerates in field  $E$  and reaches  $v_d$  on next collision  
after time  $\tau$



$$F = ma \approx \frac{mv_d}{\tau} = qE$$

$$\mu = \frac{v_d}{E} \equiv \frac{q\tau}{m_c^*}$$

$$J = \sigma E = qn \mu E = \frac{q^2 n \tau}{m_c^*} E$$

$$\sigma = q^2 \tau \frac{n}{6} \left[ \frac{2}{m_l^*} + \frac{4}{m_t^*} \right]$$

$$\sigma = q^2 n \tau \left\{ \frac{1}{3} \left[ \frac{1}{m_l^*} + \frac{2}{m_t^*} \right] \right\}$$

What happens when Si  
is tensioned?

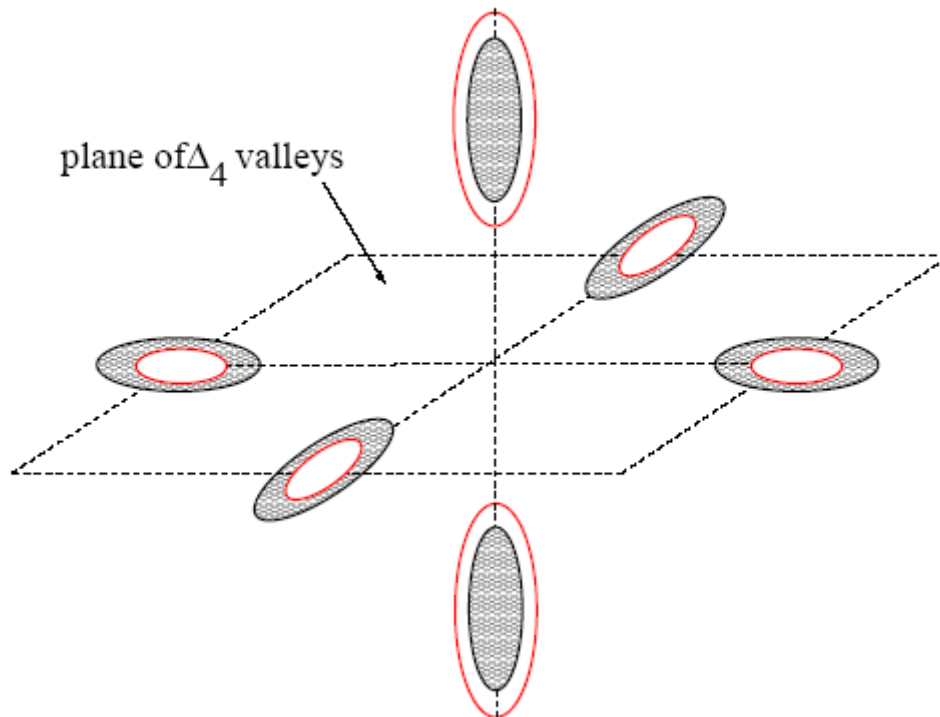
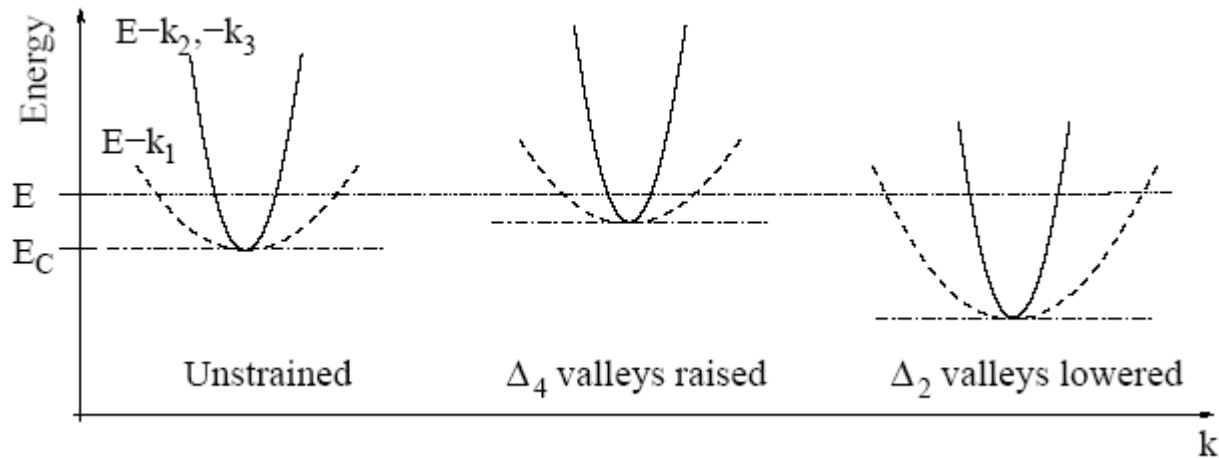
For unstrained  $\{001\}$  Si:  $m_c^* = 0.26m_0$

What is this mass called?



Sec.  
13.1.3

# Effect of tensile strain on $E_C$



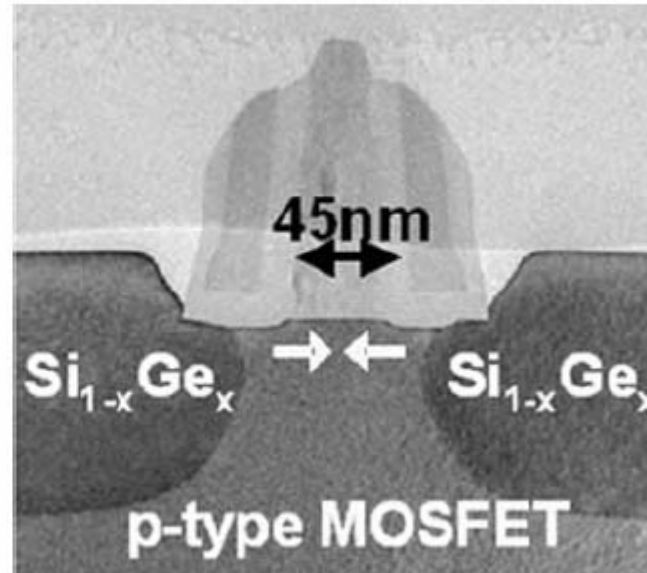
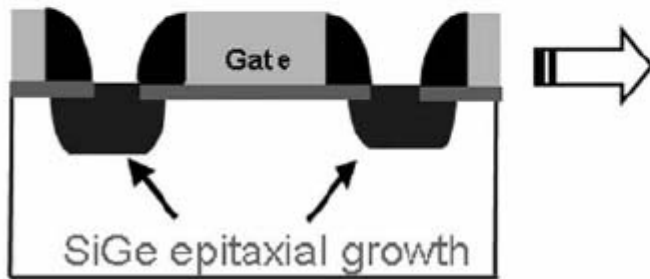
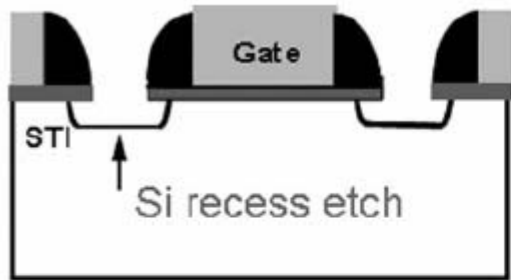
$$\sigma = q^2 \tau \frac{n}{6} \left[ \frac{2}{m_l^*} + \frac{4}{m_t^*} \right]$$

$$\frac{1}{m_C^*} = \left\{ \frac{1}{3} \left[ \frac{1}{m_l^*} + \frac{2}{m_t^*} \right] \right\}$$

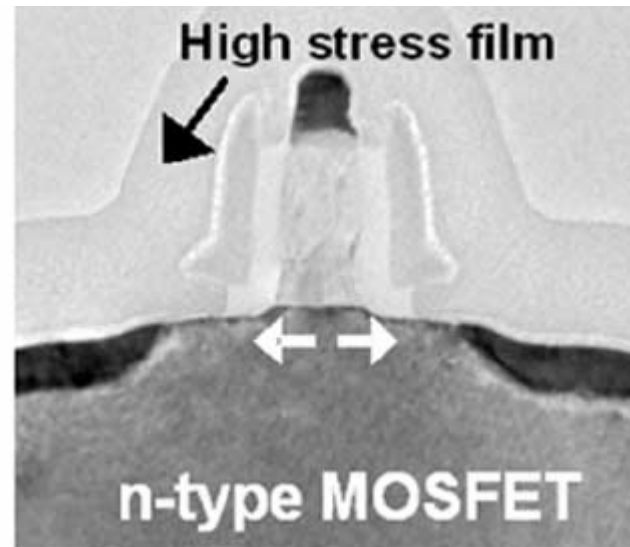
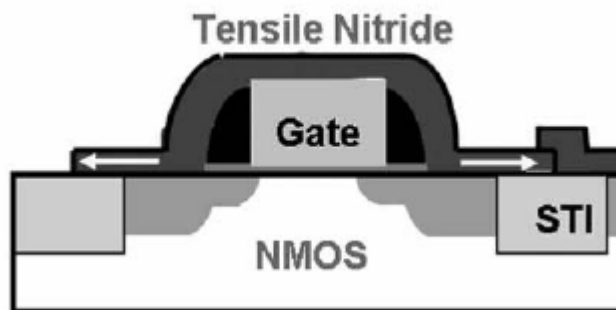
$$\Rightarrow \text{(ideally) to } \left\{ \frac{1}{2} \left[ \frac{2}{m_t^*} \right] \right\}$$

Sec.  
13.1.3

## Strained Si at the 45nm node



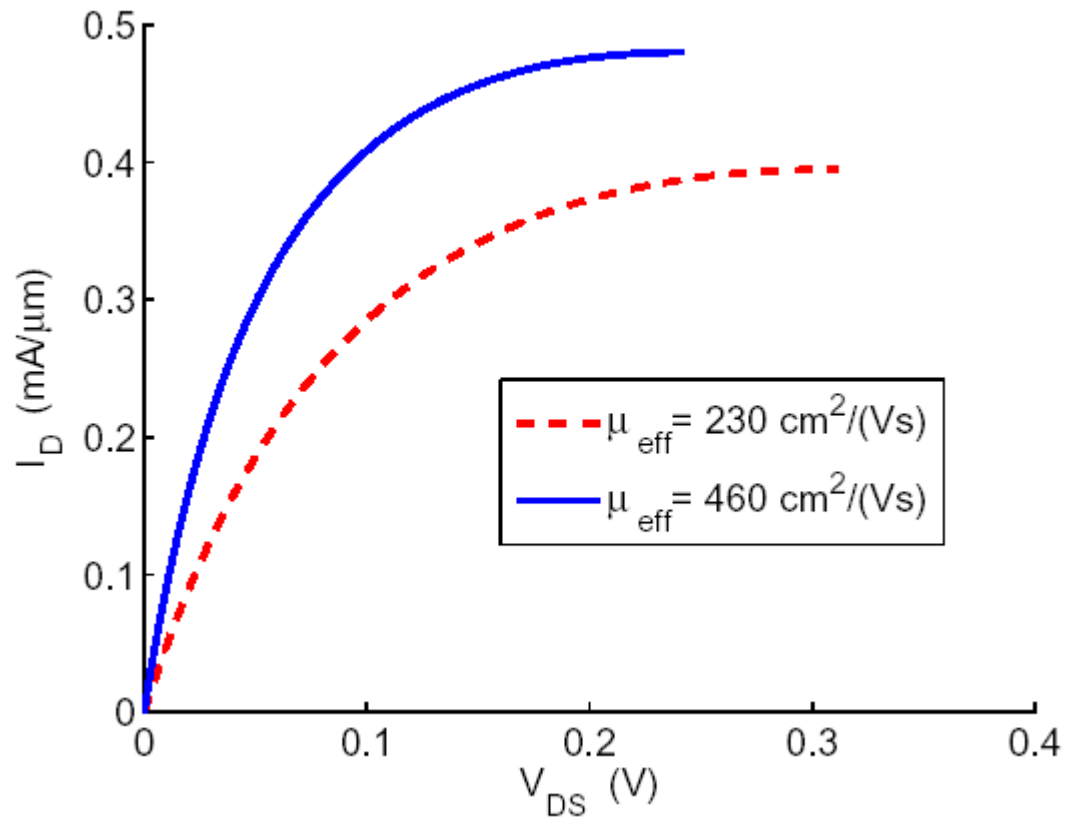
Compressive for  
P-FETs



Tensile for N-  
FETs

How much stress  
is involved?

# What a factor of 2 in $\mu$ brings



This is a 50 nm FET.

Why is  $I_{D\text{sat}}$  not directly proportional to  $\mu$  ?