

# Module 5

- ❖ MOSFET scaling
- ❖ Types Of scaling
  - Constant field scaling
  - Constant voltage scaling

# Moore's Law<sup>[1]</sup>

- No. of transistors on a chip doubled every 18 to 24 months.
- Semiconductor technology will double its effectiveness every 18 months.

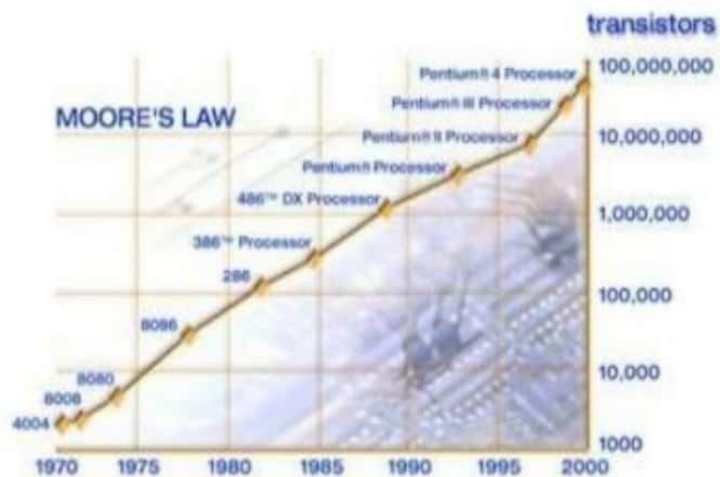
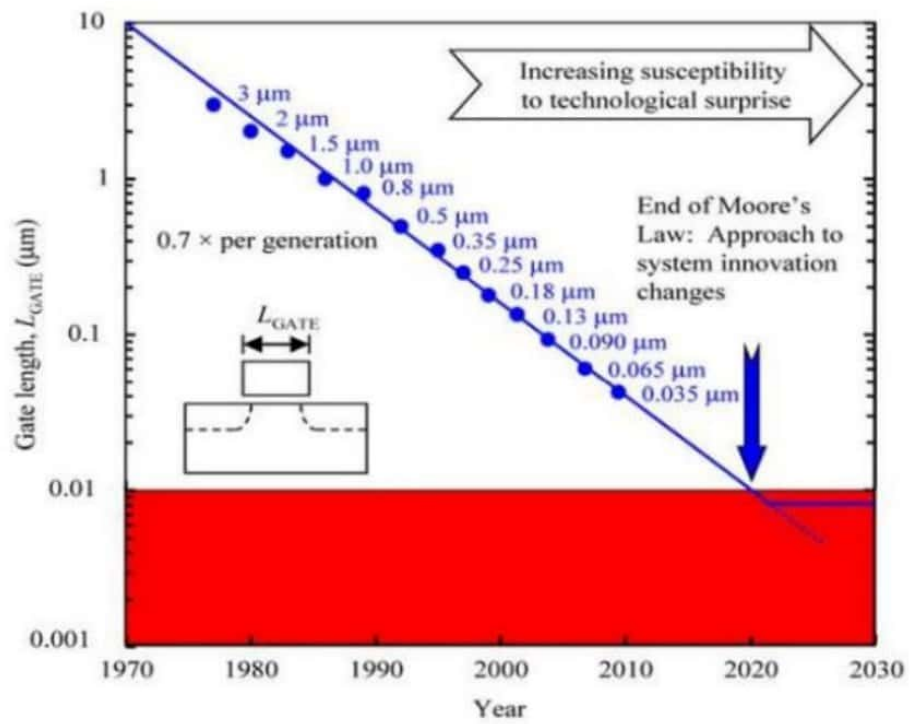


Fig. 1 : No. of transistors with years

# Need for Scaling

- The ultimate goal of **scaling** is to build an individual transistor that is smaller, faster, cheaper, and consuming low power
- **Scaling** of the **MOS transistor** improves its size, cost and performance. Today's fabricated integrated circuits are many times faster and occupy much less area, like today's microprocessors that contain nearly one billion transistors on a single chip

- The reduction of the size, i.e., the dimensions of **MOSFETs**, is commonly referred to as **scaling**. In order to meet the demand of high density chips in MOS technology, it is **required** that **MOSFET** are **scaled** down i.e. reduction in the size of transistor, so that high packaging density can be achieved.
- The reduction of the dimensions of a MOSFET has been dramatic during the last three decades. Starting at a minimum feature length of 10  $\mu\text{m}$  in 1970 the gate length was gradually reduced to 0.15  $\mu\text{m}$  minimum feature size in 2000, resulting in a 13% reduction per year. Proper scaling of MOSFET however requires not only a size reduction of the gate length and width but also requires a reduction of all other dimensions including the gate/source and gate/drain alignment, the oxide thickness and the depletion layer widths. Scaling of the depletion layer widths also implies scaling of the substrate doping density. In short, we will study simplified guidelines for shrinking device dimensions to increase transistor density & operating frequency and reduction in power dissipation & gate delays.



**Fig. 2 : End of Moore's Law**

## Statistics<sup>[3]</sup>

Processor name	Year of introduction	Transistors
4004	1971	2,250
8008	1972	2,500
8080	1974	5,000
8086	1978	29,000
286	1982	120,000
386	1985	275,000
486 DX	1989	1,180,000
Pentium	1993	3,100,000
Pentium II	1997	7,500,000
Pentium III	1999	24,000,000
Pentium 4	2000	42,000,000

**Fig. 4 : Data showing no. of transistors on processors of Intel Corporation with Years**

# Types of Scaling

- Two types of scaling are common:
- 1) constant field scaling and 2) constant voltage scaling
- Constant field scaling yields the largest reduction in the power-delay product of a single transistor. However, it requires a reduction in the power supply voltage as one decreases the minimum feature size. Constant voltage scaling does not have this problem and is therefore the preferred scaling method since it provides voltage compatibility with older circuit technologies. The disadvantage of constant voltage scaling is that the electric field increases as the minimum feature length is reduced. This leads to velocity saturation, mobility degradation, increased leakage currents and lower breakdown voltages.
- After scaling, the different Mosfet parameters will be converted as given by table below:

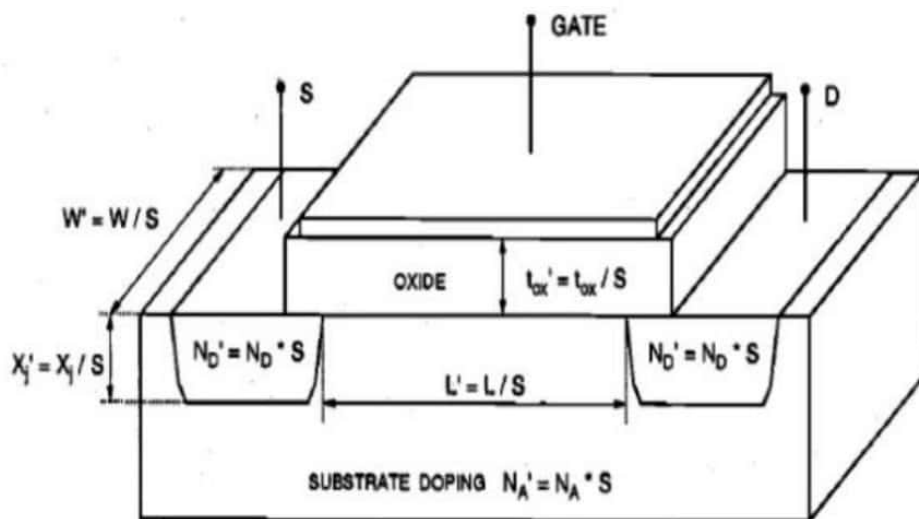


Fig. 5 : MOSFET Scaling by a factor  $S$



# Before Scaling After Constant Field Scaling After Constant Voltage Scaling

Before Scaling	After Constant Field Scaling	After Constant Voltage Scaling
$L$	$L' = L/s$	$L' = L/s$
$W$	$W' = W/s$	$W' = W/s$
$t$	$t_{ox}' = t_{ox}/s$	$t_{ox}' = t_{ox}/s$
$x_i$	$x_i' = x_i/s$	$x_i' = x_i/s$
$V_{DD}$	$V_{DD}' = V_{DD}/s$	$V_{DD}' = V_{DD}$
$V_{Th}$	$V_{Th}' = V_{Th}/s$	$V_{Th}' = V_{Th}$
$N_a$ or $N_d$	$N_a' = N_a * s$ or $N_d' = N_d * s$	$N_a' = N_a * s^2$ or $N_d' = N_d * s^2$
$C_{ox}$	$C_{ox}' = C_{ox} * s$	$C_{ox}' = C_{ox} * s$
$I_{DS}$	$I_{DS}' = I_{DS}/s$	$I_{DS}' = I_{DS} * s$
$P_D$	$P_D' = P_D/s^2$	$P_D' = P_D * s$

Where  $s$  = scaling parameter of MOS

## 1. Constant field scaling or full scaling :

- Magnitude of internal electric fields is kept constant.
- Only lateral dimensions are changed.
- Threshold voltage is also effected.

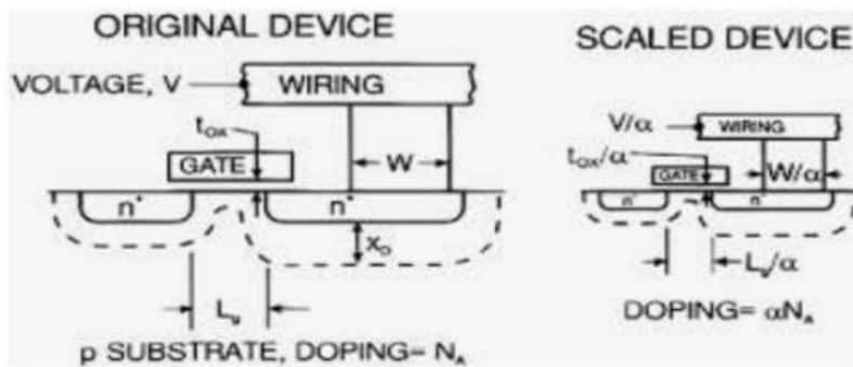


Fig. 6: Full – Scaling of MOSFET

## Consequences of Constant Field Scaling :

Quantity	Before Scaling	After Scaling
Channel length	$L$	$L' = L / S$
Channel width	$W$	$W' = W / S$
Gate oxide thickness	$t_{ox}$	$t_{ox}' = t_{ox} / S$
Junction depth	$x_j$	$x_j' = x_j / S$
Power supply voltage	$V_{DD}$	$V_{DD}' = V_{DD} / S$
Threshold voltage	$V_{T0}$	$V_{T0}' = V_{T0} / S$
Doping densities	$N_A$	$N_A' = S \cdot N_A$
	$N_D$	$N_D' = S \cdot N_D$

Quantity	Before Scaling	After Scaling
Oxide capacitance	$C_{ox}$	$C_{ox}' = S \cdot C_{ox}$
Drain current	$I_D$	$I_D' = I_D / S$
Power dissipation	$P$	$P' = P / S^2$
Power density	$P / \text{Area}$	$P' / \text{Area}' = P / \text{Area}$

Fig. 7 : Change in parameters due to full scaling

➤ **Most significant reduction :**

Power dissipation is reduced by a factor of  $S^2$  as  $P' = P/S^2$

- Power density remains unchanged.
- Gate oxide capacitance is scaled down as  $C_g' = C_g/S$
- Overall performance improvement.

## 2. Constant Voltage Scaling :

- More preferred.
- All dimensions are scaled down except power supply and terminal voltages.

Quantity	Before Scaling	After Scaling
Dimensions	$W, L, t_{ox}, x_j$	reduced by $S$ ( $W' = W / S, \dots$ )
Voltages	$V_{DD}, V_T$	remain unchanged
Doping densities	$N_A, N_D$	increased by $S^2$ ( $N_A' = S^2 \cdot N_A, \dots$ )

Quantity	Before Scaling	After Scaling
Oxide capacitance	$C_{ox}$	$C_{ox}' = S \cdot C_{ox}$
Drain current	$I_D$	$I_D' = S \cdot I_D$
Power dissipation	$P$	$P' = S \cdot P$
Power density	$P / Area$	$P' / Area' = S^3 \cdot (P / Area)$

Fig. 8 : Parameters effected due to Constant Voltage Scaling

## **Cons of Constant Voltage Scaling :**

- Increase in drain current density and power density by a factor of  $S^3$  adversely effecting device reliability.
- Causes problems like :
  - Electro Migration
  - Hot Carrier Degradation
  - Gate Oxide Breakdown
  - Electrical Over-stress

# SHORT CHANNEL EFFECTS

- Drain Induced Barrier Lowering(Refer Text book)
- Hot Carrier Effects
- Channel length modulation
- Velocity Saturation
- Threshold Voltage Variations

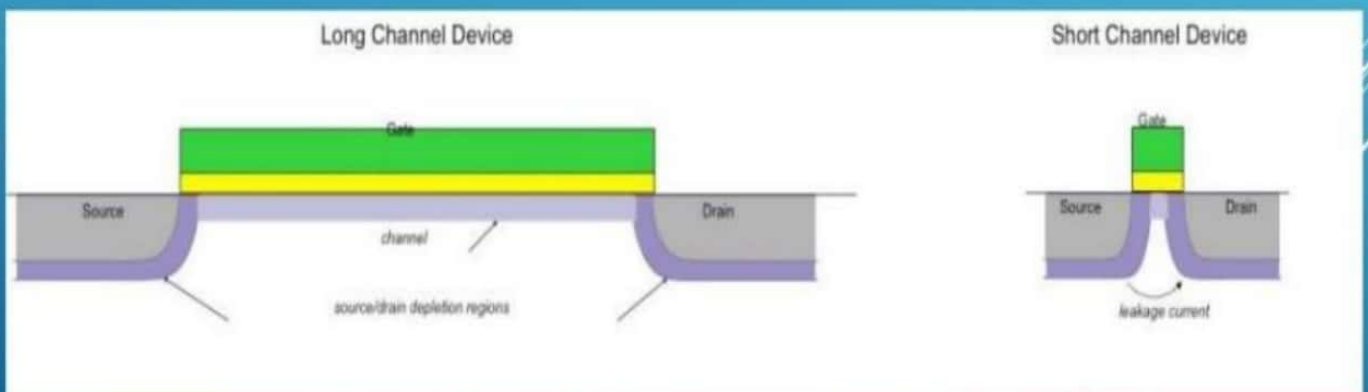
## PRE-REQUISITE

- Threshold voltage:- The minimum voltage across gate and source which allows current flow through the transistor.
- Sub-threshold current:- The current flowing through the transistor before threshold voltage is called sub-threshold current.
- Pinch-off voltage :- It is a voltage after which there is no effect of increase in drain to source voltage on the current flowing through transistor.



# What is short channel

- Channel length  $\approx$  depletion width of source and drain



## SHORT CHANNEL EFFECTS

- Short channel MOS has good processing speed, requires low operating potential and increases transistor density on the chip.
- Although the performance degrades with decrease in channel length.
- It faces some serious issues like DIBL, surface scattering, velocity saturation, impact ionisation, hot electron effect.

## MOSFET SCALING

One approach to size reduction is a scaling of the MOSFET that requires all device dimensions to reduce proportionally. The main device dimensions are the channel length, channel width, and oxide thickness. Lateral dimensions such as channel length and width are reduced by a factor of  $k$ , so should the vertical dimensions such as source/drain junction depths and gate insulator thickness.

Scaling of depletion width is achieved indirectly by scaling up doping concentrations. If we simply reduce the dimensions of the device and kept the power supply voltages same, the internal electric field in the device would increase.

Scaling improves

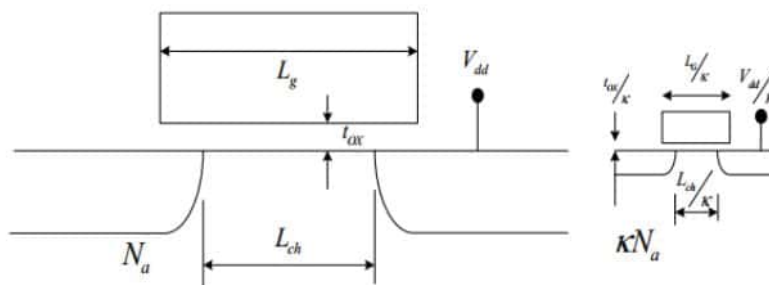
1. Packing density
2. Speed
3. Power dissipation

Two types of scaling are common:

- (i) constant field scaling
- (ii) constant voltage scaling.

**Full scaling (constant-field scaling) –**

- All dimensions are scaled by  $k$  and the supply voltage and other voltages are so scaled
- Magnitude of internal electric field is kept constant
- Only lateral dimensions are changed
- Threshold voltage is also affected



**Figure 2: 5** Illustration of MOSFET miniaturisation. The sketch on the right hand is the scaled device according to the constant field rule. (Reference [2.15])

Constant field scaling yields the largest reduction in the power-delay product of a single transistor. However, it requires a reduction in the power supply voltage as one decreases the minimum feature size.

- For ideal scaling, power supply voltages should be reduced to keep the internal electric field reasonably constant from one technology generation to the next. But power supply voltages are not scaled hand in hand with the device dimensions, partly because of other system related constraints. The longitudinal electric field in the pinch off region and the transverse electric field across the gate oxide increase with MOSFET scaling which causes **hot electron effects and short channel effects**.

**Table 6-1** Scaling rules for MOSFETs according to a constant factor K. The horizontal and vertical dimensions are scaled by the same factor. The voltages are also scaled to keep the internal electric fields more or less constant, and the hot carrier effects manageable.

	Scaling factor
Surface dimensions (L,Z)	1/K
Vertical dimensions ( $d, x_i$ )	1/K
Impurity Concentrations	K
Current, Voltages	1/K
Current Density	K
Capacitance (per unit area)	K
Transconductance	1
Circuit Delay Time	1/K
Power Dissipation	1/K <sup>2</sup>
Power Density	1
Power-Delay Product	1/K <sup>3</sup>

Constant voltage scaling does not have this problem and is therefore the preferred scaling method since it provides voltage compatibility with older circuit technologies.

The disadvantage of constant voltage scaling is that the electric field increases as the minimum feature length is reduced.

**Constant-voltage scaling** The voltages are not scaled and, in some cases, dimensions associated with voltage are not scaled.

### Constant voltage scaling

Parameter	Scaled parameter
Channel length (L)	$1/\alpha$
Junction depth ( $x_j$ )	$1/\alpha$
Substrate doping ( $N_A$ )	$\alpha$
Depletion layer thickness (d)	$1/\alpha$
Transconductance ( $g_m$ )	$\alpha$
Static power dissipation ( $P_{\text{stat}}$ )	$\alpha$
Dynamic power dissipation ( $P_{\text{dyn}}$ )	$\alpha$
Current (I)	$\alpha$
Gate delay ( $\tau_p$ )	$1/\alpha^2$
Load capacitance ( $C_L$ )	$1/\alpha$
Channel width (W)	$1/\alpha$
Supply voltage (V)	1
Gate oxide thickness ( $t_{\text{ox}}$ )	$1/\alpha$
Current density (J)	$\alpha^3$

### Comparison

Quantity	Sensitivity	Constant Field	Constant Voltage
Scaling Parameters			
Length	$L$	$1/S$	$1/S$
Width	$W$	$1/S$	$1/S$
Gate Oxide Thickness	$t_{\text{ox}}$	$1/S$	$1/S$
Supply Voltage	$V_{dd}$	$1/S$	1

<b>Threshold Voltage</b>	$V_{T0}$	$1/S$	$1$
Doping Density	$N_A, N_D$	$S$	$S^2$
<b>Device Characteristics</b>			
Area (A)	$WL$	$1/S^2$	$1/S^2$
D-S Current ( $I_{DS}$ )	$\beta(V_{dd} - v_T)^2$	$1/S$	$S$
Gate Capacitance ( $C_g$ )	$WL/t_{ox}$	$1/S$	$1/S$
Power Dissipation ( $P$ )	$I_{DS}V_{dd}$	$1/S^2$	$S$
Power Dissipation Density ( $P/A$ )	$P/A$	$1$	$S^3$

### Subthreshold conduction / Subthreshold Characteristics

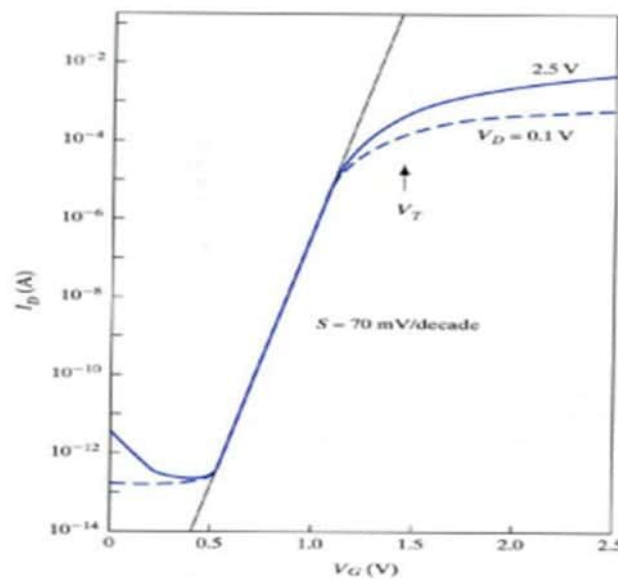
For the conduction to happen in the MOSFET; we need the  $V_G$  to be greater than the threshold voltage  $V_T$ . But, this threshold voltage is calculated at the point where the region below the oxides has entered into strong inversion.

From experimental results, one can observe that there is still some non-zero current flowing from drain to source even when we are operating at a region with  $V_G < V_T$  (sub-threshold region). This happens because, for the subthreshold region, the substrate near oxide-interface is in “Weak-Inversion”. At this point, if we apply a positive  $V_{DS}$ , there will be a small current  $I_D$  flowing. This effect is plotted in the transfer characteristics in figure below. We have,

$$I_D = \frac{W}{L} \cdot \mu_n C_{ox} \left[ (V_{GS} - V_T) V_{DS} - \frac{V_{DS}^2}{2} \right]$$

In this equation, current abruptly goes to zero as soon as  $V_G$  is reduced to  $V_T$ . In reality there is still some drain conduction below threshold, and is known as subthreshold conduction. This current is due to weak inversion in the channel between flat band and threshold which leads to a diffusion current from source to drain.





Subthreshold Swing,

$$S = \log \frac{d(V_G)}{d(\log I_D)}$$

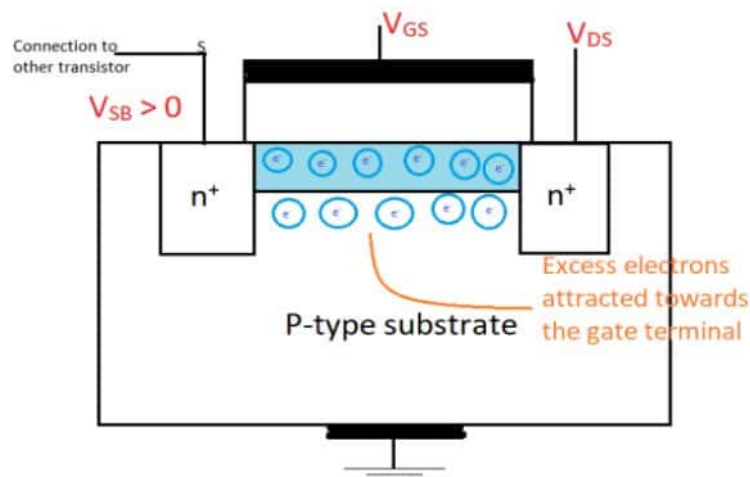
- Generally, in order to improve the performance and reduce the cost of production, one would prefer to scale down the size of the transistors.
- This scaling down also eliminates many stray capacitances that are present in the overall device. Ultimately increasing the speed of operation.
- But when the channel length is scaled down to the order of the depletion layer, a certain number of non-ideal effects come into play. These are called **second-order effects**

### **Substrate Bias Effect**

For the ideal IV characteristics, the biasing scheme we used had the source and the body both connected to ground.

But in practical design applications, Source is connected to substrate(body) so that there is a voltage  $V_{SB}$

In such scenarios, the difference in potential between the body and the source terminal causes a change in the threshold voltage of the MOSFET. This effect of change in threshold voltage is called the “Body Effect” or the “Back Gate Effect”.



When  $V_{SB}$  is positive, there is reverse bias between source and bulk. This causes depletion layer to widen.

The electrons in the bulk are repelled by the body terminal and are now attracted by the gate toward the oxide layer.

The threshold voltage of the MOS is also proportional to the density of electrons in the depletion layer.

Hence as we accumulate more and more electrons in the depletion layer below the oxide interface, there will be an increase in the value of threshold voltage.

As depletion region is widened, larger charge density is occupied. Therefore, the threshold required to achieve inversion increases.

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

$V_{T0}$  = zero - substrate - bias  $V_T$   
 $\gamma$  = body effect parameter  
 $|\phi_F|$  = surface potential parameter

### Short Channel Effects

Short-channel effects occur in MOSFETs in which the channel length is comparable to the depletion layer widths of the source and drain junctions.

A MOSFET device is considered to be short channel device when the channel length is the same order of magnitude as the depletion-layer widths ( $x_{dD}$ ,  $x_{dS}$ ) of the source and drain junction. (That is, the effective channel length  $L_{eff}$  is approximately equal to the source and drain junction depth  $x$ ).



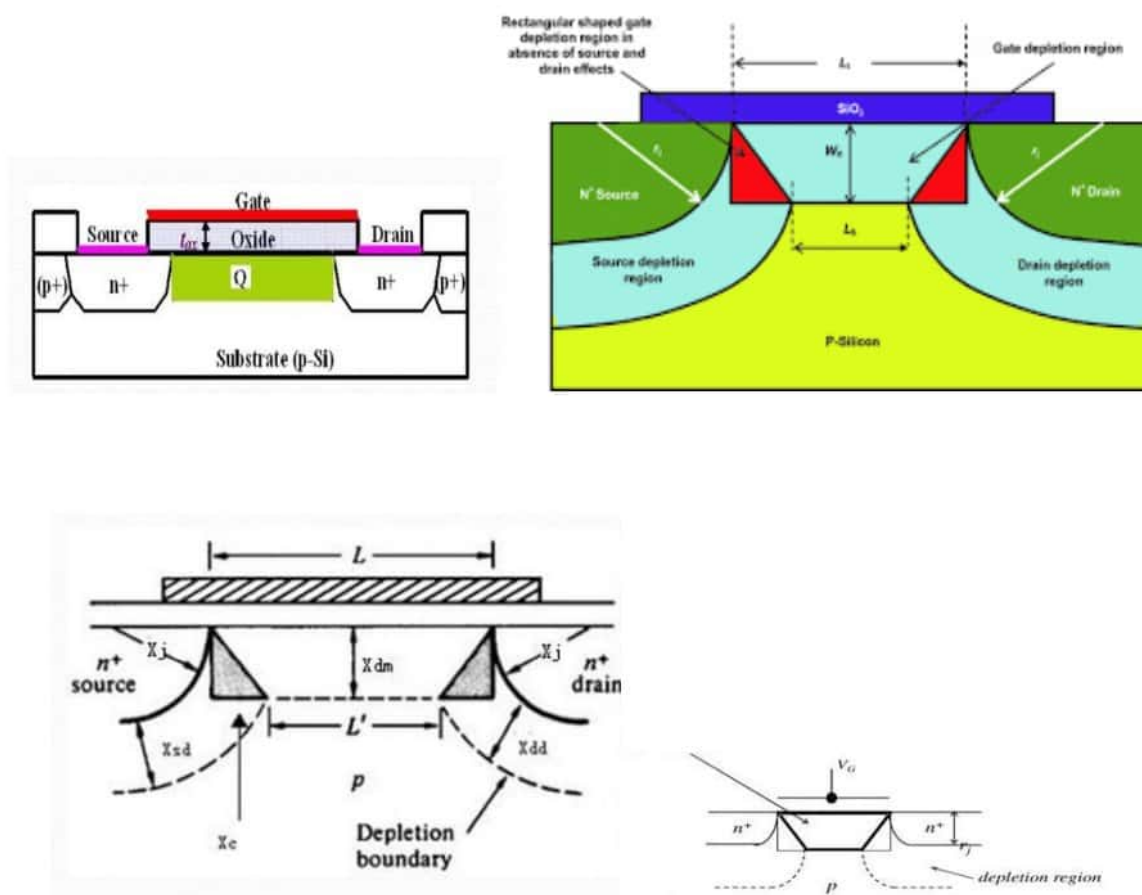
As the channel length  $L$  is reduced to increase both the operation speed and the number of components per chip, the so-called short-channel effects arise.

The short-channel effects are attributed to two physical phenomena:

1. The limitation imposed on electron drift characteristics in the channel
2. The modification of the threshold voltage due to the shortening channel length.

This occurs due to the charge sharing between source/drain and gate. A triangle region forms at both ends

Hence the rectangular area under the gate becomes Trapezoid



Different short-channel effects include

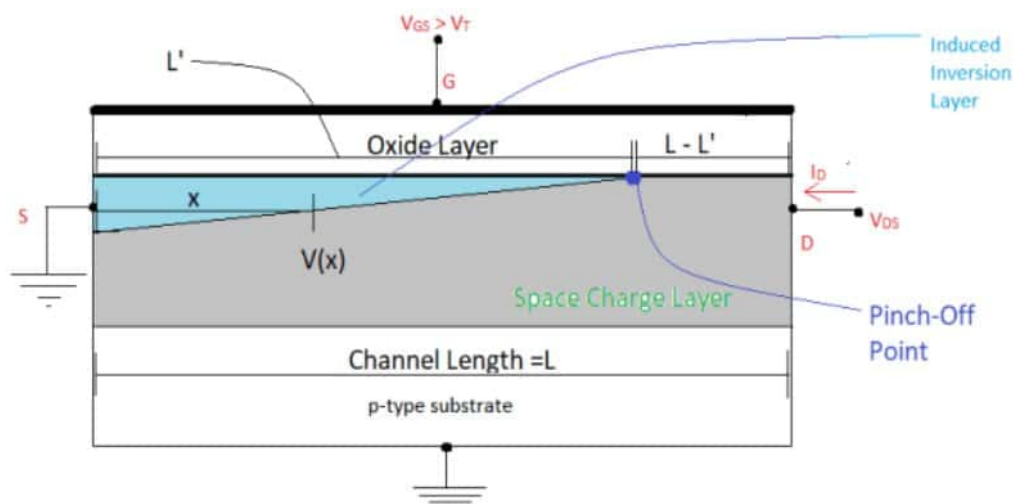
1. Channel Length Modulation
2. Drain-induced barrier lowering and “Punch through”
3. Velocity saturation
4. Threshold voltage variations
5. Hot carrier effects

### **Channel Length Modulation(CLM)**

As we keep on increasing  $V_{DS}$ , the region for which the inversion charge is zero keeps on increasing for a constant value of  $V_{GS}$  maintained. Thus channel length keeps on decreasing. This phenomenon is called Channel Length Modulation.

This is similar to “Base Width Modulation” Thus we get a  $V_{DS}$  term in the expression for  $I_D$  even when we are operating in the saturation region.

Generally, the fabrication of the MOSFET devices is done in a way such that the change in length given by  $\Delta L = L - L'$  is low with a change in  $V_{DS}$ .



### Drain Induced Barrier Lowering (DIBL)

When the depletion regions surrounding the drain extends to the source, so that the two-depletion layer merge (i.e., when  $x_{ds} + x_{dD} = L$ ), punch through occurs.

Punch through can be minimized with thinner oxides, larger substrate doping, shallower junctions, and obviously with longer channels.

The current flow in the channel depends on creating and sustaining an inversion layer on the surface.

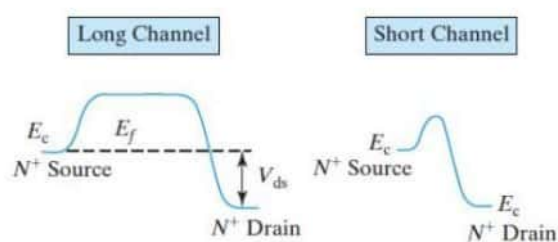
If the gate bias voltage is not sufficient to invert the surface ( $V_{GS} < V_T$ ), the carriers (electrons) in the channel face a potential barrier that blocks the flow. Increasing the gate voltage reduces this potential barrier and, eventually, allows the flow of carriers under the influence of the channel electric field.

In small-geometry MOSFETs, the potential barrier is controlled by both the gate-to-source voltage  $V_{GS}$  and the drain-to-source voltage  $V_{DS}$ .

If the drain voltage is increased, the potential barrier in the channel decreases, leading to drain-induced barrier lowering (DIBL).

The reduction of the potential barrier eventually allows electron flow between the source and the drain, even if the gate-to-source voltage is lower than the threshold voltage.

The channel current that flows under this condition ( $V_{GS} < V_T$ ) is called the sub-threshold current



## Velocity Saturation

The velocity of charge carriers, such as electrons or holes, is proportional to the electric field that drives them, but that is only valid for small fields.

As the field gets stronger, their velocity tends to saturate. That means that above a critical electric field, they tend to stabilize their speed and eventually cannot move faster.

Velocity saturation is specially seen in short-channel MOSFET transistors, because they have higher electric fields

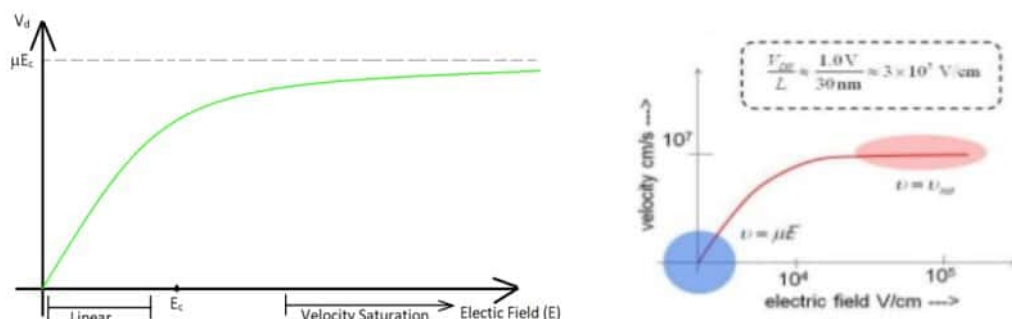
The drift velocity of the electrons in the inversion layer to be proportional to the lateral electric field applied. The proportionality constant was given by  $\mu_n$ .

The key point to understand the effect of velocity saturation is that the linearity of the drift velocity only holds true for low values of the applied electric field. The actual variation of drift velocity with respect to the applied electric field is shown in figure 6.

The exact formula for the drift velocity can be given as:

$$v_d = \frac{\mu E}{1 + E/E_c}$$

The term  $E_c$  is called the critical electric field. Here the electric field  $E$  is equal to  $\frac{V_{DS}}{L}$ , i.e. the lateral voltage applied across the channel divided by the effective channel length. We can see that for large channel devices, the drift velocity formula simplifies to  $v_d = \mu E$ . Hence this is also a short channel effect because the lateral electric field is higher in case of short channel devices for similar range of drain-to-source voltage applied.



**Figure 6: Variation of drift velocity of electron w.r.t. applied electric field**

### **Threshold Variations**

The threshold voltage is only a function of the manufacturing technology and the applied body bias  $V_{SB}$ .

The threshold can therefore be considered as a constant over all NMOS (PMOS) transistors in a design. As the device dimensions are reduced, this model becomes inaccurate, and the threshold potential becomes a function of  $L$ ,  $W$ , and  $V_{DS}$ .

Two-dimensional second-order effects that were ignorable for long-channel devices suddenly become significant.

In the traditional derivation of the  $V_{T0}$ , for instance, it is assumed that the channel depletion region is solely due to the applied gate voltage and that all depletion charge beneath the gate originates from the MOS field effects.

This ignores the depletion regions of the source and reverse-biased drain junction, which become relatively more important with shrinking channel lengths.

Since a part of the region below the gate is already depleted (by the source and drain fields), a smaller threshold voltage suffices to cause strong inversion.

In other words,  $V_{T0}$  decreases with  $L$  for short-channel devices (Figure 3.35a).

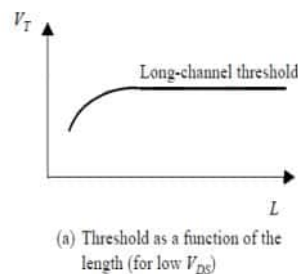


Figure 3.35 Threshold variations.

### **Hot-Carrier Effects**

Another problem, related to high electric fields, is caused by so-called hot electrons. This high energy electrons can enter the oxide, where they can be trapped, giving rise to oxide charging that can accumulate with time and degrade the device performance by increasing  $V_T$  and affect adversely the gate's control on the drain current.

Besides varying over a design, threshold voltages in short-channel devices also have the tendency to drift over time. This is the result of the hot-carrier effect



Over the last decades, device dimensions have been scaled down continuously, while the power supply and the operating voltages were kept constant. The resulting increase in the electrical field strength causes an increasing velocity of the electrons, which can leave the silicon and tunnel into the gate oxide upon reaching a high-enough energy level.

Electrons trapped in the oxide change the threshold voltage, typically increasing the thresholds of NMOS devices, while decreasing the  $V_T$  of PMOS transistors.

For an electron to become hot, an electrical field of at least  $10^4$  V/cm is necessary. This condition is easily met in devices with channel lengths around or below 1  $\mu$ m.

The hot-electron phenomenon can lead to a long-term reliability problem, where a circuit might degrade or fail after being in use for a while. This is illustrated in Figure 3.36, which shows the degradation in the I-V characteristics of an NMOS transistor after it has been subjected to extensive operation.

MOSFET technologies, therefore use specially-engineered drain and source regions to ensure that the peaks in the electrical fields are bounded, hence preventing carriers to reach the critical values necessary to become hot.

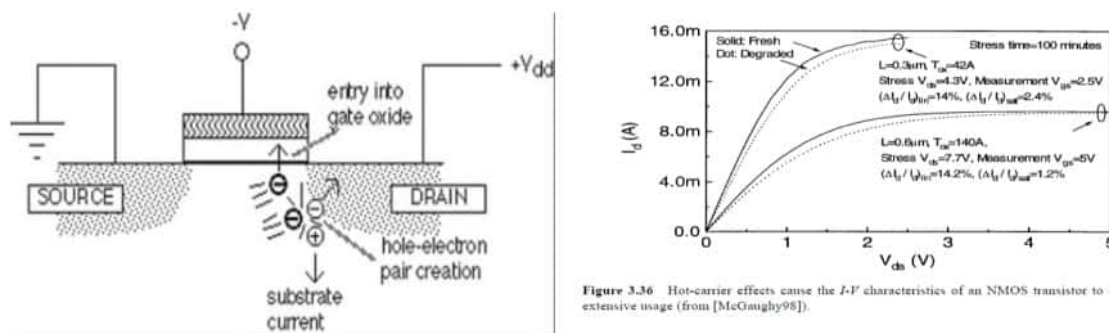


Figure 3.36 Hot-carrier effects cause the I-V characteristics of an NMOS transistor to degrade from extensive usage (from [McGaughey98]).

## FinFET

- FinFET, also known as Fin Field Effect Transistor, is a type of **non-planar or "3D" transistor** used in the design of modern processors
- FinFETs are new generation transistors which utilize **tri-gate structure**. In contrast to planar transistors where the Gate electrode was (usually) above the channel, the Gate electrode **"wraps" the channel**
- The distinguishing characteristic of the FinFET is that the conducting channel is wrapped by a thin silicon "fin", which forms the body of the device.

- The conducting channel is greatly controlled by the gate.
- The thickness of the fin (measured in the direction from source to drain) determines the effective channel length of the device.
- These effects make it harder for the voltage on a gate electrode to deplete the channel underneath and stop the flow of carriers through the channel – in other words, to turn the transistor Off. By raising the channel above the surface of the wafer instead of creating the channel just below the surface, it is possible to wrap the gate around up to three of its sides, providing much greater electrostatic control over the carriers within it.

## Advantages

**Reduced** short-channel effects (SCEs) and **leakage current** .

To overcome the worst types of short-channel effect encountered by deep submicron transistors, such as drain-induced barrier lowering (DIBL).

This technique provides increased operating speed by low-threshold MOSFET and **reduced leakage** by high-threshold voltage.

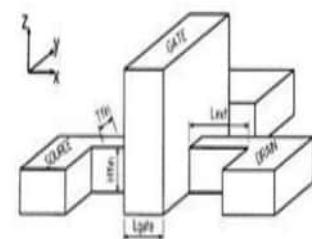
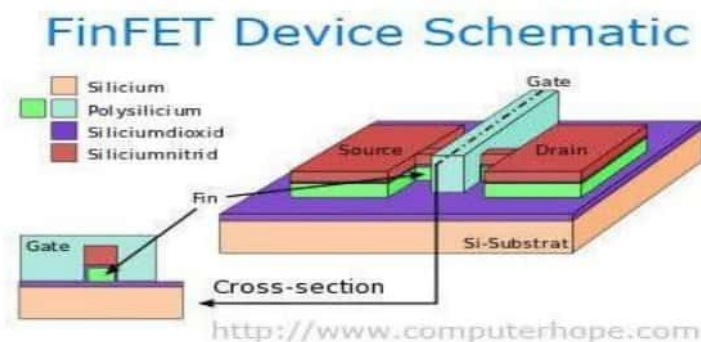
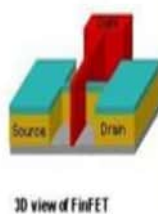
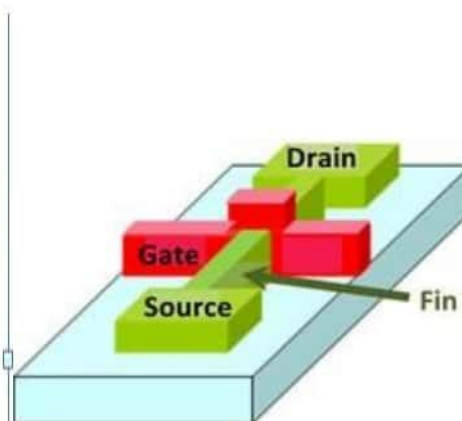


Fig. 1. Schematic of a FinFET structure



3D view of FinFET



- The very first finFETs were manufactured on top of insulating layer.
- The fact that the current can't flow "underneath" the gate when the transistor is in OFF state reduces the leakage current.
- Alternative techniques for stopping leakage current from flowing in the bulk were introduced later, which allowed for manufacturing of Bulk finFETs.
- This technique utilizes very high doping gradients along the height of the fin in order to prevent the current from flowing in the bulk.

## Output Characteristics

### Drain Characteristics

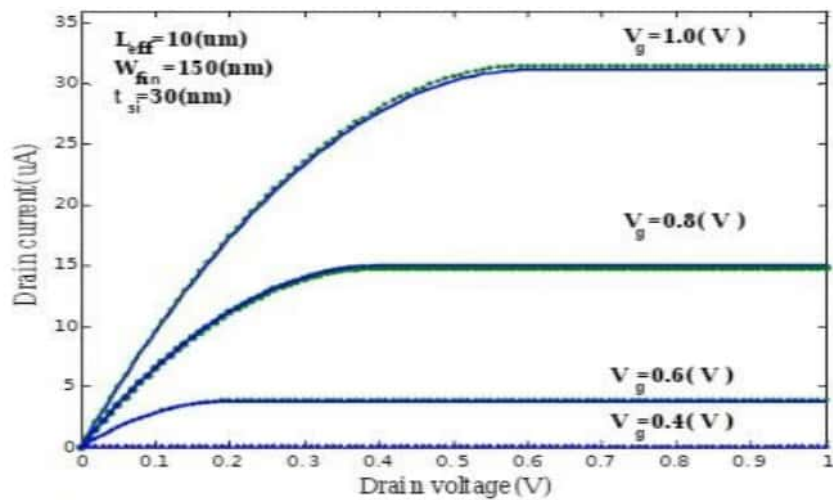
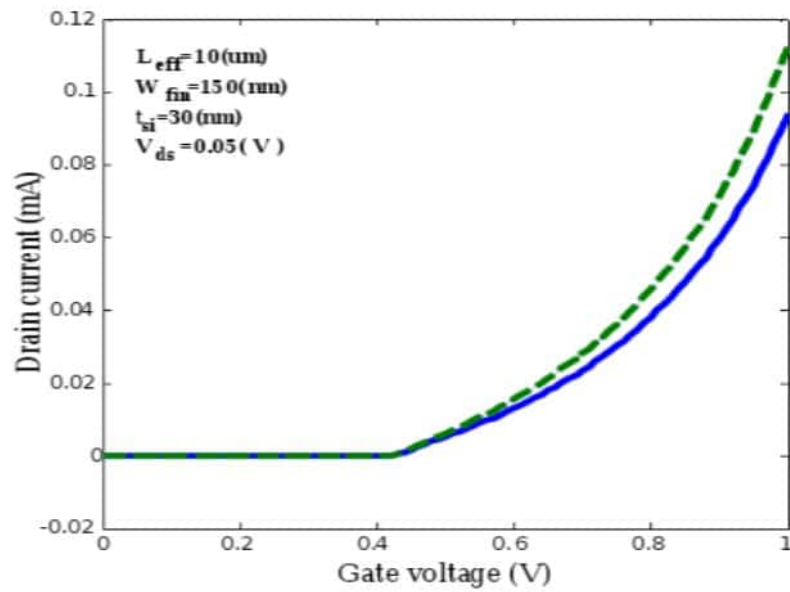


Fig. 2 : show the output characteristics of FinFET of  $L_{\text{ch}} = 10\mu\text{m}$ ,  $W_{\text{fin}} = 150\text{nm}$ ,  $t_{\text{si}} = 30\text{nm}$  for various gate voltage. Symbols are for experimental data and solid line for simulation result of this work.



## Transfer Characteristics



*Fig.3* : Transfer characteristics of FinFET of  $L_{ch} = 10 \mu\text{m}$ ,  $W_{fin} = 150 \text{ nm}$  and  $t_{si} = 30 \text{ nm}$ . Dashed line for the experimental data and solid line for simulation result of this work.