

Appendix 1.

MOSFET Scaling Rule and Issues

Introduction

- Scaling of MOSFETs → mostly down-scaling
- Down-scaling → generally channel length (L) shrinks (shorter channel)
- But only reduction of L gives punch-through where the depletion regions at source and drain are merged → abnormal MOSFET operation, and many other issues (e.g. SCE like the DIBL)
- So, we need a rule for down-scaling MOSFET while maintaining the MOSFET's normal operation.
- **What and How?**

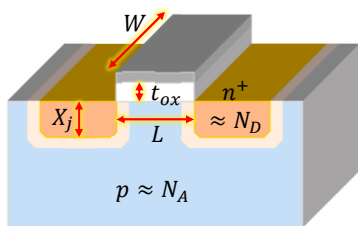
Part 1.

Scaling Rule

with constant power-density strategy and constant Transconductance

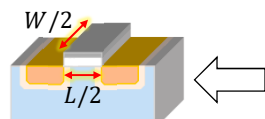
(= keeping the MOSFET operation the same)

Scaling Rule with the ratio of **S (>1)** (example: n-MOSFET)



Geometry	Doping density	Signal level
$L \rightarrow L/S$	$N_D \rightarrow N_D \times S$	$V_{DD} \rightarrow V_{DD}/S$
$W \rightarrow W/S$	$N_A \rightarrow N_A \times S$	$V_{th} \rightarrow V_{th}/S$
$X_j \rightarrow X_j/S$		
$t_{ox} \rightarrow t_{ox}/S \Rightarrow C_{ox} \rightarrow C_{ox} \times S$		

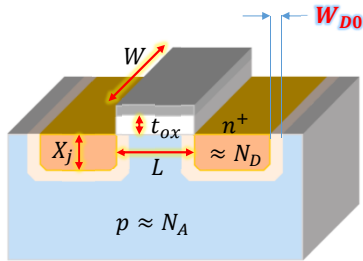
↓ **S = 2**



Example for S=2

Geometry	Doping density	Signal level
$L \rightarrow L/2$	$N_D \rightarrow N_D \times 2$	$V_{DD} \rightarrow V_{DD}/2$
$W \rightarrow W/2$	$N_A \rightarrow N_A \times 2$	$V_{th} \rightarrow V_{th}/2$
$X_j \rightarrow X_j/2$		
$t_{ox} \rightarrow t_{ox}/2 \Rightarrow C_{ox} \rightarrow C_{ox} \times 2$		

Why scaling-up of doping density (1)

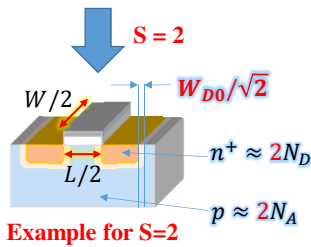


$$pn^+ \approx N_A^- N_D^+ = n_i^2 \exp\left(\frac{qV_{bi}}{kT}\right)$$

$$W_{D0} = \sqrt{\frac{2\epsilon}{q} \left(\frac{1}{N_A^-} + \frac{1}{N_D^+} \right) V_{bi}} \approx \sqrt{\frac{2\epsilon}{q} \left(\frac{1}{N_A^-} \right) V_{bi}}$$

$$W'_{D0} = \sqrt{\frac{2\epsilon}{q} \left(\frac{1}{sN_A^-} \right) V'_{bi}} = \frac{W_{D0}}{\sqrt{s}}$$

Assuming $V_{bi} \approx V'_{bi}$

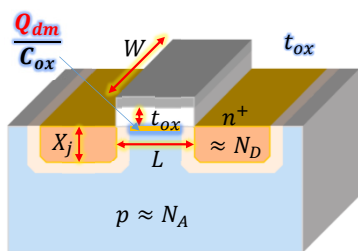


Example for $S=2$

← To reduce the depletion width of pn+ junctions.

← Why not S^2 to get $\frac{W_{D0}}{S}$ → Zener breakdown.

Why scaling-up of doping density (2)



$$V_{th} = \phi_{ms} - \frac{Q_{ox}}{C_{ox}} - \frac{Q_{dm}}{C_{ox}} + 2\phi_F$$

$$Q_{dm} \approx -qpW_{dm} \quad p = n_i \exp\left(\frac{q\phi_F}{kT}\right)$$

$$W_{dm} = \sqrt{\frac{2\epsilon_s}{q} \frac{1}{p} 2\phi_F} \quad \Rightarrow \phi_F = \frac{kT}{q} \ln\left(\frac{p}{n_i}\right)$$

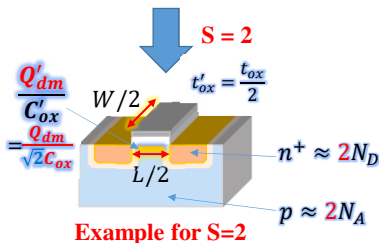
← saturated

$$W'_{dm} = \sqrt{\frac{2\epsilon_s}{q} \frac{1}{sp} 2\phi_F} = \frac{W_{dm}}{\sqrt{s}} \quad Q'_{dm} \approx -qspW'_{dm} \approx -qsp \frac{W_{dm}}{\sqrt{s}} = \sqrt{s} Q_{dm}$$

$$\frac{Q'_{dm}}{C'_{ox}} \approx \frac{Q_{dm}}{\sqrt{s} C_{ox}} \quad \rightarrow C'_{ox} = s C_{ox}$$

← To reduce the depletion voltage ($-\frac{Q_{dm}}{C_{ox}}$)

← Why not S^2 to get $\frac{W_{D0}}{S}$ → Zener breakdown.
→ A reason why V_{th} is tardy to be scaled down.



Example for $S=2$

Scaling rule by the factor of “s” for a constant DC-power density

Current

$$I_D = \frac{1}{2} \mu_n \frac{W}{L} C_{ox} (V_G - V_{th})^2 \Rightarrow \mu_n \frac{W}{L} \frac{\epsilon_{ox}}{t_{ox}} V_{DD}^2 \Rightarrow I'_D \Rightarrow \frac{1}{2} \mu_n \frac{W}{L} \frac{\epsilon_{ox}}{t_{ox}} \left(\frac{V_{DD}}{s} \right)^2 = \frac{I_D}{s}$$

$$s > 1$$

Voltage

$$V_{DD} \Rightarrow V'_{DD} = \frac{V_{DD}}{s}$$

$$s = s_i = s_v$$

DC-Power consumption

$$P_{dc} = V_{DD} I_D \Rightarrow P'_{dc} = \frac{V_{DD} I_D}{s^2} = \frac{P_{dc}}{s^2}$$

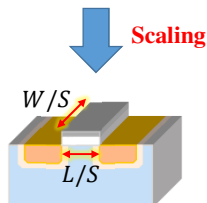
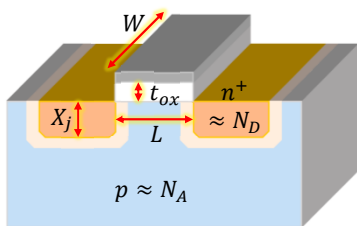
DC-Power density:

$$PD = \frac{P_{dc}}{A} = \frac{V_{DD} I_D}{LW} \Rightarrow PD' = \frac{P'_{dc}}{A'} = \frac{P_{dc}/s^2}{LW/s^2} = PD \leftarrow \text{Constant !}$$

Transconductance:

$$g_m = \mu_n \frac{W}{L} \frac{\epsilon_{ox}}{t_{ox}} V_{DD} \Rightarrow g'_m = \mu_n \frac{W/s}{L/s} \frac{\epsilon_{ox}}{t_{ox}/s} V_{DD}/s = g_m \leftarrow \text{Constant !}$$

Summary of Scaling Rule and Its Electrical Aspects



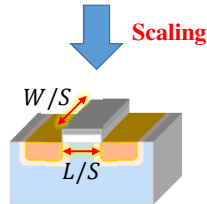
Scaling Rule:

Geometry	Doping density	Signal level
$L \rightarrow L/S$	$N_D \rightarrow N_D \times S$	$V_{DD} \rightarrow V_{DD}/S$
$W \rightarrow W/S$	$N_A \rightarrow N_A \times S$	$V_{th} \rightarrow V_{th}/S$
$X_j \rightarrow X_j/S$		
$t_{ox} \rightarrow t_{ox}/S \Rightarrow C_{ox} \rightarrow C_{ox} \times S$		

Electrical aspects:

Current (I) $\rightarrow 1/S$	Circuit delay (= $g_m C L W$) $\rightarrow 1/S$
Voltage (V) $\rightarrow 1/S$	
DC-Power consumption (P_{DC}) $\rightarrow 1/S^2$	$g_m \rightarrow \text{kept}$ $C \rightarrow S$ $LW \rightarrow 1/S^2$
DC-Power density (PD) and $g_m \rightarrow \text{kept}$	

Basic purpose of Scaling → Increase of Frequency capability of MOSFETs



Electrical aspects:

- Current (I) → 1/S
 - Voltage (V) → 1/S
 - DC-Power consumption (P_{DC}) → 1/S²
 - DC-Power density (PD) and g_m → kept
- Circuit delay ($= g_m C LW$) → 1/S
- $\left\{ \begin{array}{l} g_m \rightarrow \text{kept} \\ C \rightarrow S \\ LW \rightarrow 1/S_2 \end{array} \right.$

$$f_T = \frac{g_m}{2\pi(C_{gs} + C_{gd})} = \frac{g_m}{2\pi(c_{gs} + c_{gd})LW}$$

Operating frequency (BW) is increased by s-times

$$f'_T = \frac{g'_m}{2\pi(s c_{gs} + s c_{gd})LW/s^2} = \frac{s g_m}{2\pi(c_{gs} + c_{gd})LW}$$

However, there are several issues on this rule of scaling

- **Threshold voltage (V_{th})** is tardy to be scaled-down, so, **the supply voltage (V_{DD})** is not scaled-down (tardy). So the DC-power consumption is not scaled-down while the area is scaled-down, thus higher power density → hot plate
- And this is added with the AC-power consumption (dynamic power consumption Which is proportional to the operating frequency), leading to further higher power density → very hot plate → very high frequency operation is to be avoided.
- This is one of difficulties for a RF circuit design for the 5-G application (28 GHz).
- **Sub-threshold slope (SS)** based on diffusion current mechanism is not scalable.
- **Off-current ($I_{OFF} = I_{DS}(V_G=0)$)** is increased, so higher stand-by power consumption.
- **Gate oxide thickness (t_{ox})** is not easy to be very thin due to a dielectric breakdown.
- **And many other issues**, such as a limitation of lithography for a geometrical scaling while getting a failure of meeting the Moore's Law.

Example: Tardy Scaling-Down of Voltage Level

Current

$$I_D = \frac{1}{2} \mu_n \frac{W}{L} C_{ox} (V_G - V_{th})^2 \Rightarrow \mu_n \frac{W}{L} \frac{\epsilon_{ox}}{t_{ox}} V_{DD}^2 \Rightarrow I'_D \Rightarrow \frac{1}{2} \mu_n \frac{W}{L} \frac{\epsilon_{ox}}{t_{ox}} \left(\frac{V_{DD}}{s_v} \right)^2 = \frac{s}{s_v^2} I_D$$

Voltage

$$V_{DD} \Rightarrow V'_{DD} = \frac{V_{DD}}{s_v} \quad \text{Issue: } s > s_v > 1$$

DC-Power consumption

$$P_{dc} = V_{DD} I_D \Rightarrow P'_{dc} = \frac{s V_{DD} I_D}{s_v^3} = \frac{s}{s_v^3} P_{dc} \gg P_{dc}$$

DC-Power density:

$$PD = \frac{P_{dc}}{A} = \frac{V_{DD} I_D}{LW} \Rightarrow PD' = \frac{P'_{dc}}{A'} = \frac{\frac{s}{s_v^3} P_{dc}}{LW/s^2} = \frac{s^3}{s_v^3} PD \gg PD$$

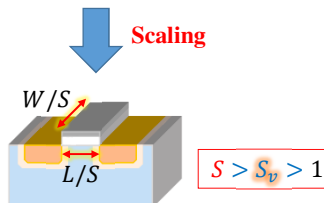
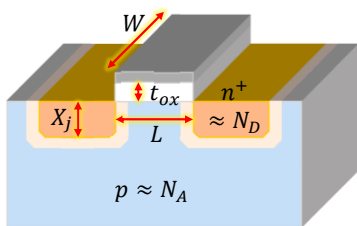
Transconductance:

$$g_m = \mu_n \frac{W}{L} \frac{\epsilon_{ox}}{t_{ox}} V_{DD} \Rightarrow g'_m = \mu_n \frac{W}{L} \frac{\epsilon_{ox}}{t_{ox}} V_{DD}/s_v = \frac{s}{s_v} g_m \gg g_m$$

← Increased !
← Increased !

PD is not kept but increased → Hot! → Burned Out !

Example: Summary of Tardy Scaling-Down of Voltage Level and its effect on Electrical Aspects



Scaling Rule:

Geometry	Doping density	Signal level
$L \rightarrow L/S$	$N_D \rightarrow N_D \times S$	$V_{DD} \rightarrow V_{DD}/s_v$
$W \rightarrow W/S$	$N_A \rightarrow N_A \times S$	$V_{th} \rightarrow V_{th}/s_v$
$X_j \rightarrow X_j/S$		
$t_{ox} \rightarrow t_{ox}/S \Rightarrow C_{ox} \rightarrow C_{ox} \times S$		

Electrical aspects:

Current (I) $\Rightarrow S/S_v^2$	Circuit delay (= $g_m C L W$) $\Rightarrow 1/S_v$
Voltage (V) $\Rightarrow 1/s_v$	$\left\{ \begin{array}{l} g_m \rightarrow \frac{s}{s_v} \\ C \rightarrow S \\ LW \rightarrow 1/S^2 \end{array} \right.$
DC-Power consumption (P_{DC}) $\Rightarrow \frac{s}{s_v^3}$	
DC-Power density (PD) $\Rightarrow \frac{s^3}{s_v^3} > 1$ and $g_m \rightarrow \frac{s}{s_v} > 1$	

Part 2.

Issues and Possible Solutions

Issue 1. Difficulty of Down-Scaling of V_{th} and Special MOSFETs to overcome: SOI-MOSFET

Typical MOSFET (bulk)

WF design is limited by nature to be negative value for cancelling $-\frac{Q_{dm}}{C_{ox}}$

$-\frac{Q_{ox}}{C_{ox}}$ is constant

$-\frac{Q_{dm}}{C_{ox}}$ is tardy to be scaled down

$$W_{dm} = \sqrt{\frac{2\epsilon_s}{q} \frac{1}{p} 2\phi_F}$$

$$\phi_F = \frac{kT}{q} \ln\left(\frac{p}{n_i}\right)$$

is saturated (\rightarrow constant)

$$V_{th} = \phi_{ms} - \frac{Q_{ox}}{C_{ox}} - \frac{Q_{dm}}{C_{ox}} + 2\phi_F$$

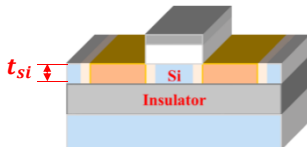
SOI-MOSFET (Fully depleted=FD)

$$Q_{dm}^{si} < Q_{dm}$$

$$t_{si} < W_{dm} \Rightarrow -qpt_{si} < -qpW_{dm}$$

Depletion width and charge is limited by the thickness of Silicon on Insulator (SOI).

V_{th} can be reduced slightly ...



Issue 1. Difficulty of Down-Scaling of V_{th} and Special MOSFETs to overcome: Fin-FET

Typical MOSFET (bulk)

WF design is limited by nature to be negative value for cancelling $-\frac{Q_{dm}}{C_{ox}}$

$-\frac{Q_{ox}}{C_{ox}}$ is constant

$-\frac{Q_{dm}}{C_{ox}}$ is tardy to be scaled down

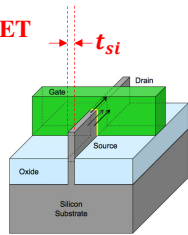
$$W_{dm} = \sqrt{\frac{2\epsilon_s}{q} \frac{1}{p} 2\phi_F}$$

$$\phi_F = \frac{kT}{q} \ln\left(\frac{p}{n_i}\right)$$

is saturated (\rightarrow constant)

$$V_{th} = \phi_{ms} - \frac{Q_{ox}}{C_{ox}} - \frac{Q_{dm}}{C_{ox}} + 2\phi_F$$

Fin-FET



22nm Bulk Fin-FET

$$|Q_{dm}^{si}| < |Q_{dm}|$$

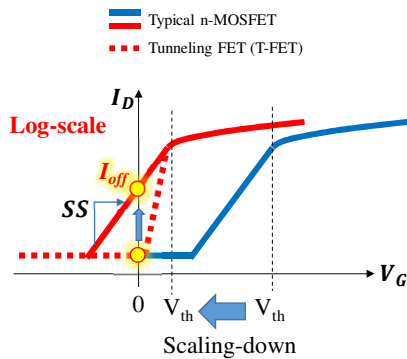
$$t_{si} < W_{dm} \Rightarrow |-qpt_{si}| < |-qpW_{dm}|$$

Depletion width and charge is limited by the thickness of Silicon Fin (지느러미).

V_{th} can be reduced slightly ...



Issue 2. Difficulty of Down-Scaling of SS and I_{off} and Special MOSFETs to overcome: Tunneling FET (T-FET)

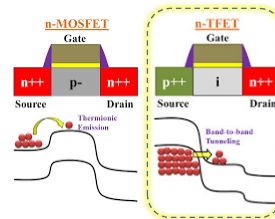


Although V_{th} is scaled down successfully, Off-current (I_{DS} at $V_G=0$) is increased due to SS which is not scaled-down below 60 mV/dec

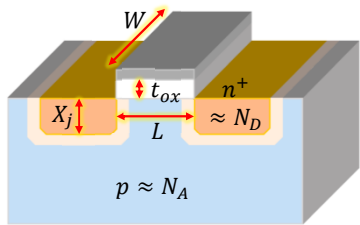
$$SS_{min} = 60 \text{ mV/dec at } T = 300K$$

This is limited by the diffusion mechanism in the Sub-threshold regime.

So, we need to use other mechanism rather than the diffusion. \rightarrow It is **Tunneling**



Issue 3. Difficulty of Down-Scaling of t_{ox} and Special MOSFETs to overcome: High-K insulator



Scaling by $S=2$

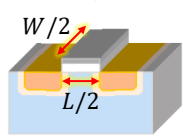


$$t_{ox} \rightarrow t_{ox} \\ \Rightarrow C_{ox} \rightarrow C_{ox} \times 2$$

$$C_{ox} = \frac{\epsilon_{ox}}{t_{ox}} = \frac{\epsilon_r \epsilon_0}{t_{ox}}$$

$$\Rightarrow \epsilon_{ox} \rightarrow \epsilon_{ox} \times 2$$

Scaling by $S=2$



$$t_{ox} \rightarrow t_{ox}/2 \\ \Rightarrow C_{ox} \rightarrow C_{ox} \times 2$$

Thinning the dielectric gives a breakdown
 → Limitation of scaling down
 → With the main purpose of scaling up C_{ox} ,
 To find the high ϵ_r insulator is better.

High-K == High Dielectric Constant (ϵ_r)

For example, HfO_x , Al_2O_3 , etc,
 which have a much higher value of ϵ_r
 than SiO_2 ($\epsilon_r = 3.9$).

Issue 4. AC Power Density at Higher Frequency

← Solution = Stop increasing frequency → but Circuit optimization needed

DC-Power consumption

$$P_{dc} = V_{DD} I_D \quad \Rightarrow \quad P'_{dc} = \frac{s V_{DD} I_D}{s^3} = \frac{s}{s^3} P_{dc} \gg P_{dc}$$

← Increased !

DC-Power density:

$$PD_{dc} = \frac{P_{dc}}{A} = \frac{V_{DD} I_D}{LW} \quad \Rightarrow \quad PD'_{dc} = \frac{P'_{dc}}{A'} = \frac{\frac{s}{s^3} P_{dc}}{LW/s^2} = \frac{s^3}{s^3} PD_{dc} \gg PD_{dc}$$

← Increased !

AC-Power consumption (Dynamic Power Consumption)

$$P_{ac} = \frac{1}{2} C_{tot} V_{DD}^2 f_{CLK} \quad \Rightarrow \quad P'_{ac} = \frac{1}{2} \frac{C'_{tot} V_{DD}^2}{s} f_{CLK} = \frac{1}{s^2} P_{ac} < P_{ac}$$

← Decreased !

AC-Power density:

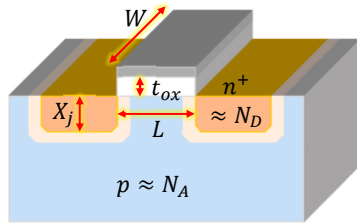
$$PD_{ac} = \frac{P_{ac}}{A} = \frac{V_{DD}^2 I_D}{LW} \quad \Rightarrow \quad PD'_{ac} = \frac{P'_{ac}}{A'} = \frac{\frac{1}{s^2} P_{ac}}{LW/s^2} = \frac{s^2}{s^2} PD_{ac} \gg \gg PD_{ac}$$

← Increased !

← Very High Frequency Circuit would be burned out !

Issue 4. AC Power Density at Higher Frequency

← Solution = Revising of Scaling Rule → But Poor Performance



Scaling Rule:

Geometry	Doping density	Signal level
$L \rightarrow L/S$	$N_D \rightarrow N_D \times S$	$V_{DD} \rightarrow V_{DD}/S_v$
$W \rightarrow W/S$	$N_A \rightarrow N_A \times S$	$V_{th} \rightarrow V_{th}/S_v$
$X_j \rightarrow X_j/S$		
$t_{ox} \rightarrow t_{ox}$	$\Rightarrow C_{ox} \rightarrow C_{ox}$	$S > S_v > 1$
Non-scaling C_{ox}		

Current → $I'_D \Rightarrow \frac{1}{2} \mu_n \frac{W}{S} \frac{\epsilon_{ox}}{L} \frac{V_{DD}}{S_v} \left(\frac{V_{DD}}{S_v} \right)^2 = \frac{1}{S_v^2} I_D$

Voltage → $V'_{DD} = \frac{V_{DD}}{S_v}$

AC-Power density (now OK):

→ $PD'_{ac} = \frac{P'_{ac}}{A'} = \frac{\frac{1}{S_v^3} P_{ac}}{LW/S^2} = \frac{S^2}{S_v^3} PD_{ac} < PD_{ac}$

Transconductance (not OK now):

$$g'_m = \mu_n \frac{W}{S} \frac{\epsilon_{ox}}{L} \frac{V_{DD}}{S_v} = \frac{1}{S_v} g_m < g_m$$

g_m decreased

→ Poor performance of circuits