# NM6604: Semiconductor Process and Device Simulation

# Virtual Process Integration (VPI)

## *Dr Zhou Xing*

Office:   S1-B1c-95
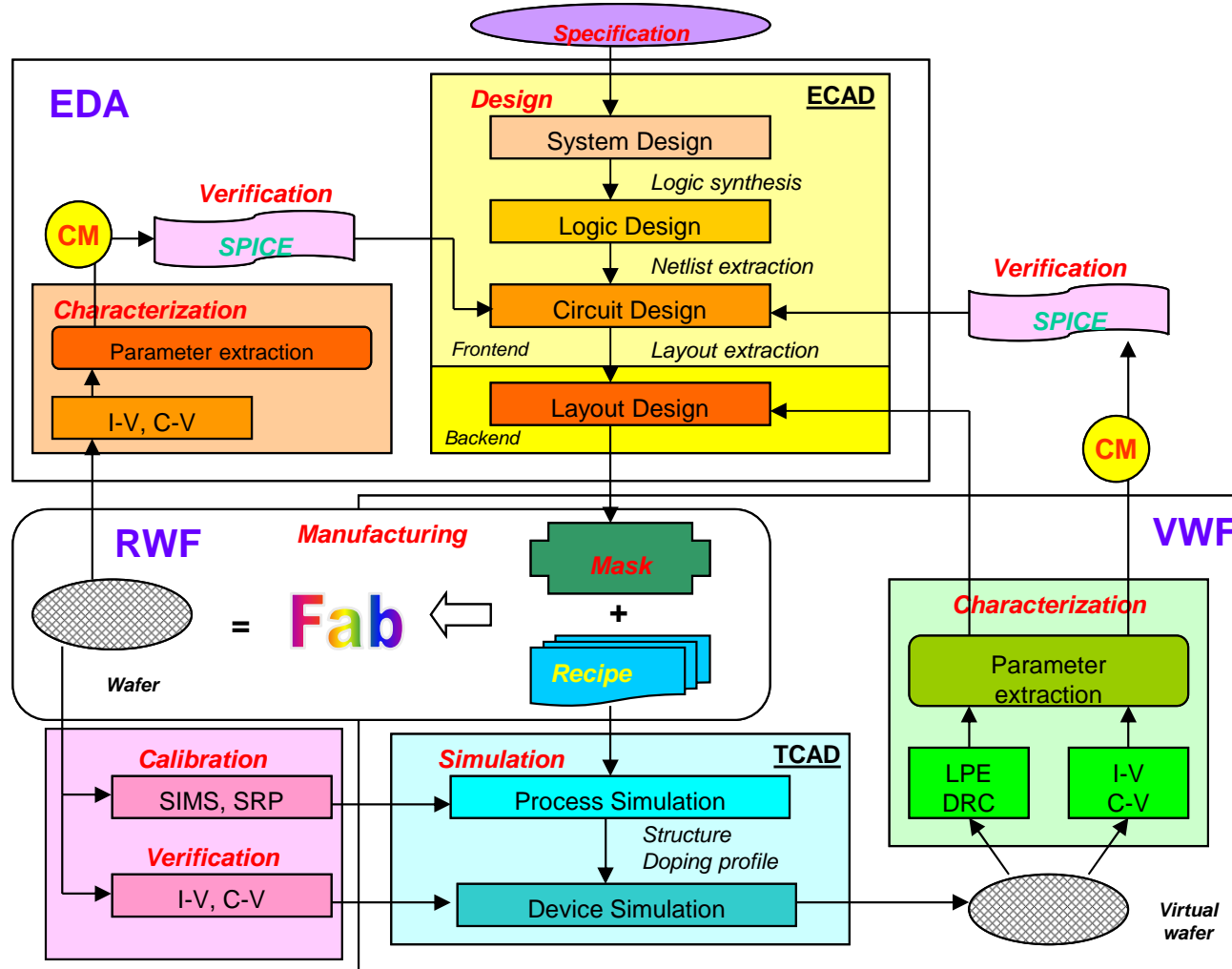Phone:  6790-4532
Email:   exzhou@ntu.edu.sg

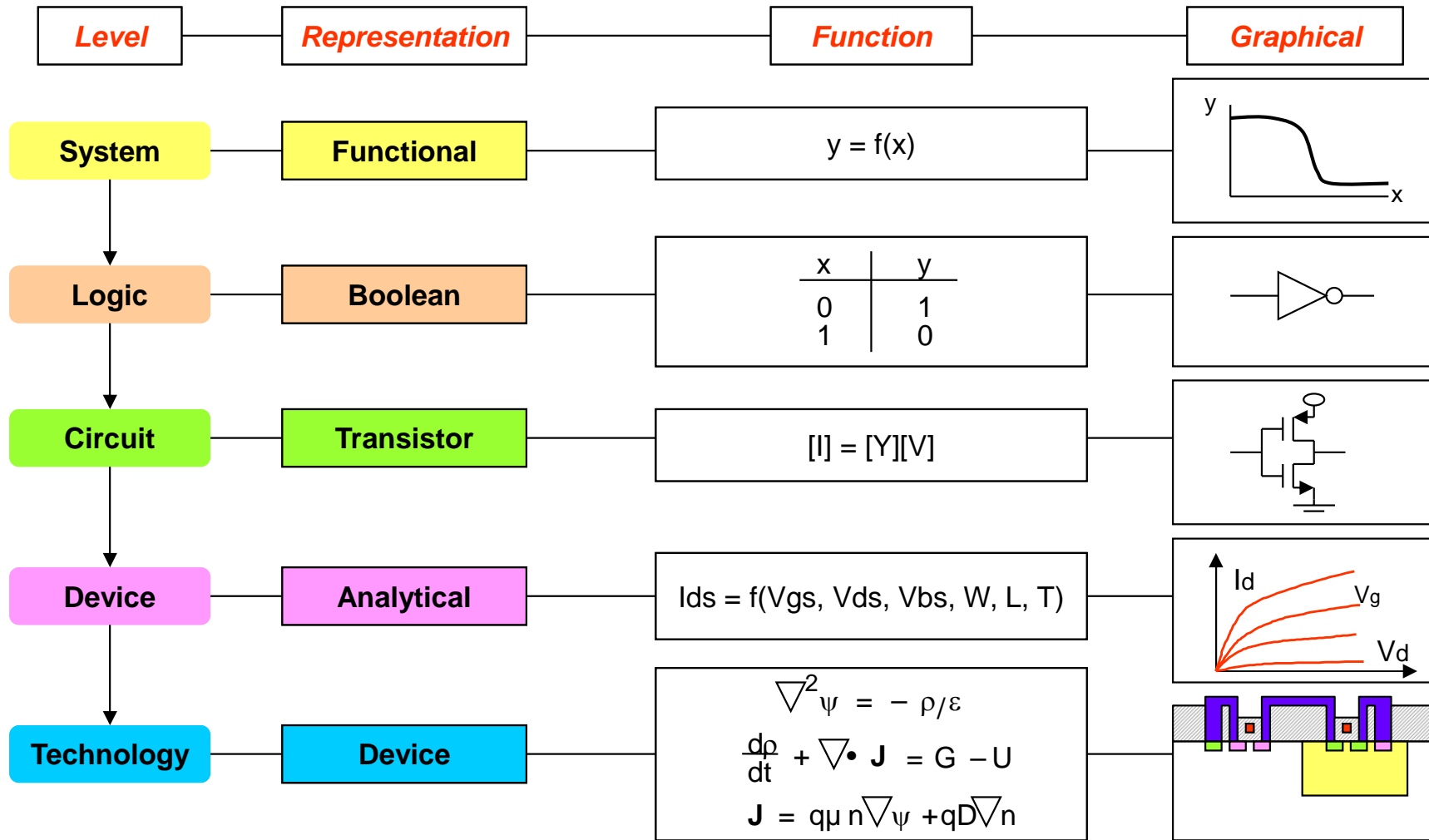Web:     https://www3.ntu.edu.sg/home/exzhou/Teaching/TUM-NM6604/
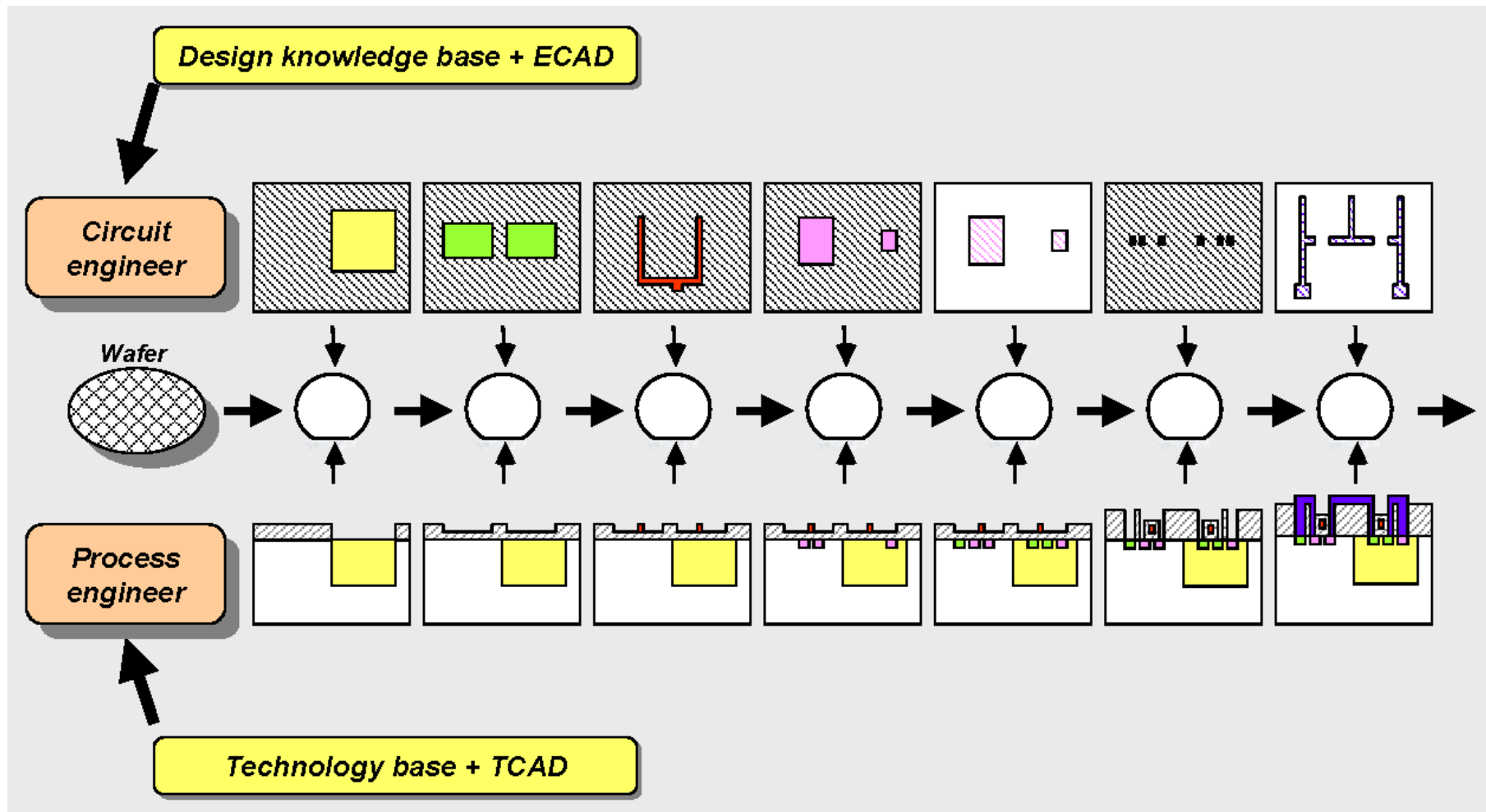
# Overall Picture: Chip Design and Wafer Fabrication

## Design–Manufacturing–Characterization–Simulation–Verification

# Multi-Level Representation
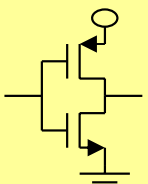
| Level | Representation | Function | Graphical |
|---|---|---|---|
| **System** | **Functional** | $y = f(x)$ |  |
| **Logic** | **Boolean** | $\begin{array}{c\|c} x & y \\ \hline 0 & 1 \\ 1 & 0 \end{array}$ |  |
| **Circuit** | **Transistor** | $[I] = [Y][V]$ |  |
| **Device** | **Analytical** | $Ids = f(Vgs, Vds, Vbs, W, L, T)$ |  |
| **Technology** | **Device** | $\nabla^2 \psi = -\rho/\varepsilon$ $\frac{d\rho}{dt} + \nabla \bullet \mathbf{J} = G - U$ $\mathbf{J} = q\mu\, n\nabla\psi + qD\nabla n$ |  |

# Layout + Process = Chip

# Target–Variable Relationship
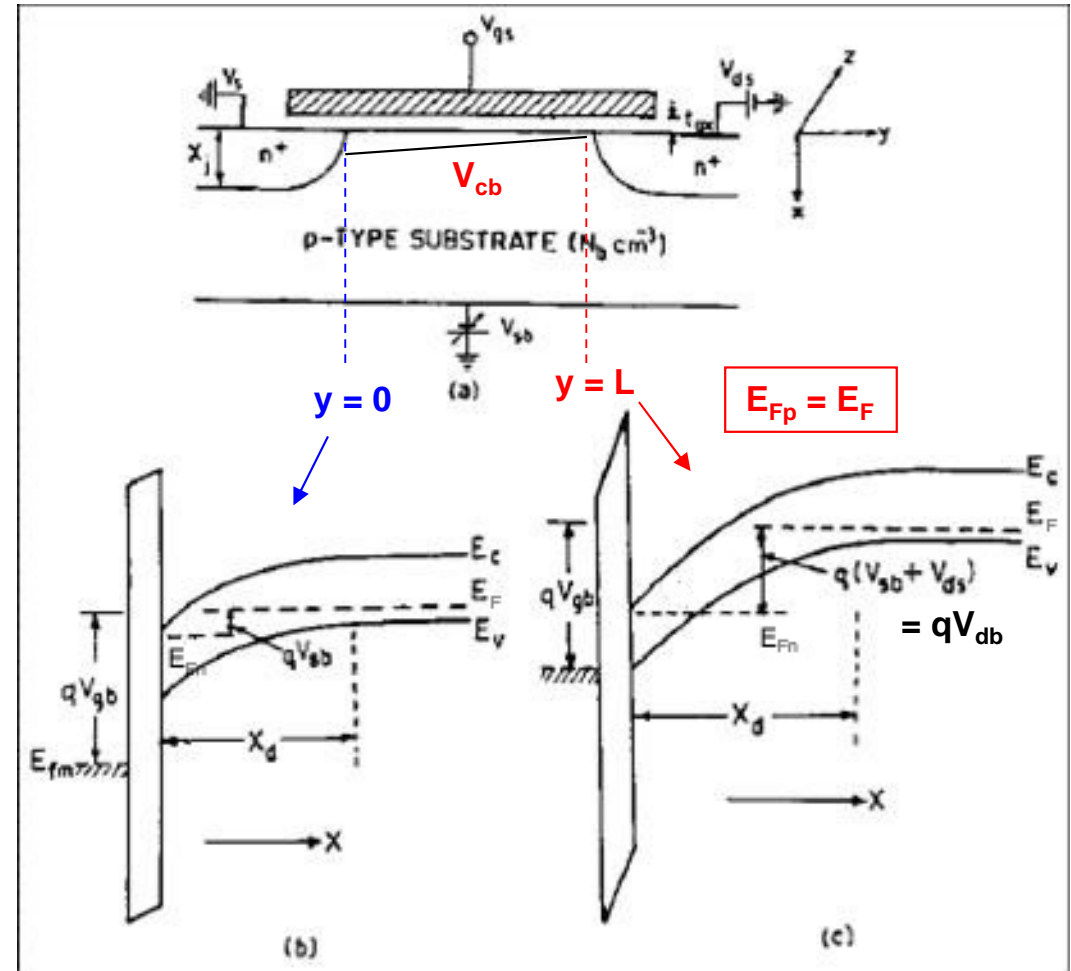
| Level | Variables | Targets |
|---|---|---|
| **Circuit**<br> | Spice: Model parameters, ...<br><br>Geometrical: Channel length, width, ...<br><br>Electrical: Supply voltage, substrate bias, ... | Digital: Delay, rise/fall time, drivability, off-state current, noise margin, ...<br><br>Analog: Voltage gain, cutoff frequency, slew rate, gain-bandwidth, ... |
| **Device**<br> | Structural: Oxide thickness, junction depth, sheet resistance, ...<br><br>Doping: Peak/surface concentration, ...<br><br>Electrical: Supply voltage, substrate bias, ... | Electrical: Threshold, transconductance, subthreshold swing, saturation current, punchthrough current, junction capacitance, lifetime, ...<br><br>Physical: Potential, field, charge, current, carriers, velocity, ... |
| **Process**<br> | Oxidation: Temperature, time, ambient, ...<br><br>Implantation: Dose, energy, tilt, damage, ...<br><br>Diffusion: Defect, stress, OED, TED, ... | Layer: Oxide thickness, junction depth, sheet resistance, ...<br><br>Profile: Peak/surface concentration, projected range/straggle, ... |

# MOSFET Operation: Due to Inversion Carrier Imref-Split

- MOSFET at equilibrium ($V_{ds} = 0$): no current flow even if channel is created at $V_{gs} = V_t$ ($\psi_s = 2\phi_F$)
- When $V_{sb} \neq 0$ ($V_{SB} > 0$ in NMOS), electron imref will "split" from hole imref with $qV_{sb} = E_{Fn} - E_{Fp}$, so $\psi_s = 2\phi_F + V_{sb}$.



- When $V_{ds} \neq 0$ ($V_{ds} > 0$ in NMOS), holes are still at quasi-equilibrium (since no 'source' nor 'drain'), so we can assume $E_{Fp} = E_F$. However, electron imref will change from $V_{sb}$ at source end to $V_{db} = V_{sb} + V_{ds}$ at drain end relative to $E_F$, and varying along the channel as $V_{cb}(y)$ ['c' stands for 'channel'].
- It is the <u>gradient</u> of $V_{cb}(y)$ that drives electrons *drifting/diffusing* from source to drain along y.



Key to understanding MOSFET operation: Band diagram in the x direction along a cutline at (b) source-end (y = 0) and (c) drain-end (y = L).

# MOSFET Source-Referenced Threshold Voltage

## MOSFET threshold voltage definition (source-referenced)

We define the **threshold voltage** ($V_t$) to be the gate-to-*source* voltage ($V_{gs}$) at which source-end surface potential is equal to twice of the bulk Fermi potential ($2\phi_F$) with reference to source–bulk voltage $V_{sb}$.

**Potential balance**      Gauss law      Charge balance      Charge-sheet approximation

$$\left(V_{gs} + V_{sb}\right) - V_{FB} = V_{gb} - V_{FB} \equiv \boxed{V_{gf} = V_{ox} + \psi_s} = Q_g / C_{ox} + \psi_s = -Q_{sc} / C_{ox} + \psi_s \approx \underline{-Q_b / C_{ox} + \psi_s}$$

$$= -\left(-qN_A X_d\right) / C_{ox} + \psi_s = +\sqrt{2q\varepsilon_{Si} N_A \psi_s} / C_{ox} + \psi_s$$

Full-depletion approximation        $X_d = \sqrt{2\varepsilon_{Si} \psi_s / qN_A}$

$$V_{gs} = V_{FB} - V_{sb} - Q_b / C_{ox} + \psi_s$$

$$V_t \equiv V_{gs}\Big|_{\psi_s = 2\phi_F + V_{sb}} = V_{FB} - V_{sb} + \left[-Q_b(\psi_s)/C_{ox} + \psi_s\right]\Big|_{\psi_s = 2\phi_F + V_{sb}} = V_{FB} - V_{sb} - Q_b\left(\psi_s = 2\phi_F + V_{sb}\right)/C_{ox} + \left(2\phi_F + V_{sb}\right)$$

$$\boxed{\therefore \quad V_t \equiv V_{gs}\Big|_{\psi_s = 2\phi_F + V_{sb}} = V_{FB} + \Upsilon \sqrt{2\phi_F + V_{sb}} + 2\phi_F}$$      where $\Upsilon = \sqrt{2q\varepsilon_{Si} N_A} / C_{ox}$ is the **body factor**.

## Body effect — Threshold-voltage shift due to non-zero $V_{sb}$

For NMOS, $V_{sb} > 0$ so that source/drain-to bulk diodes always reverse biased.

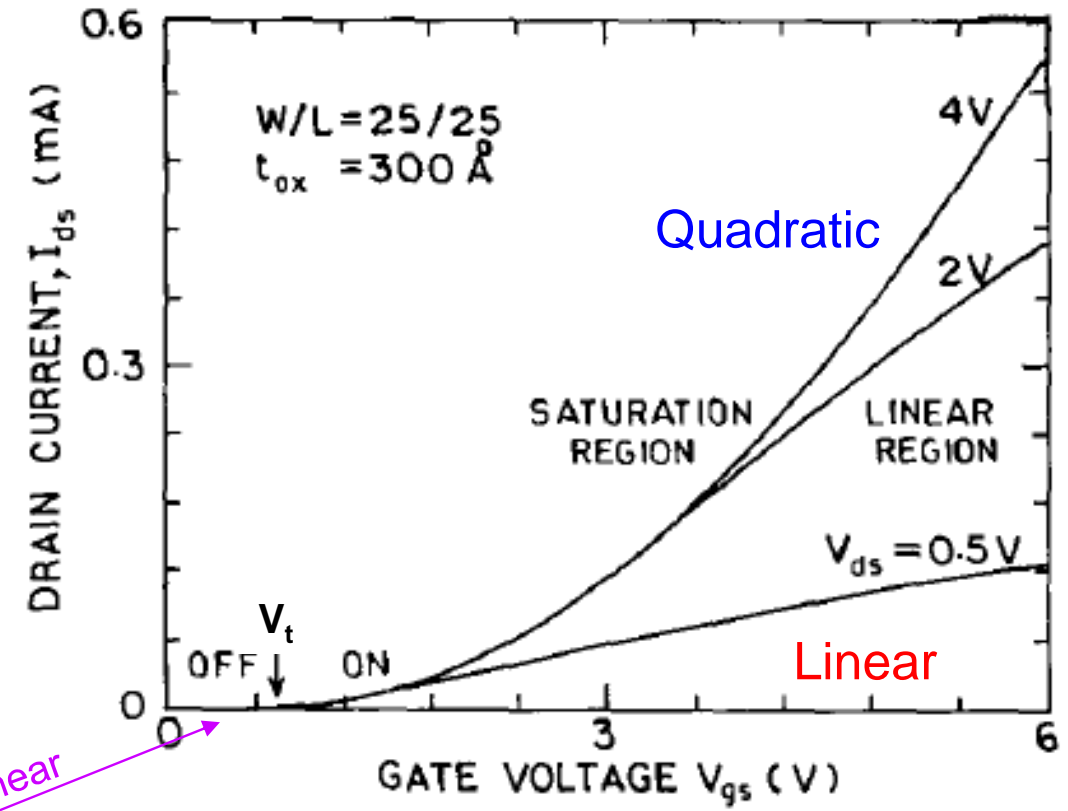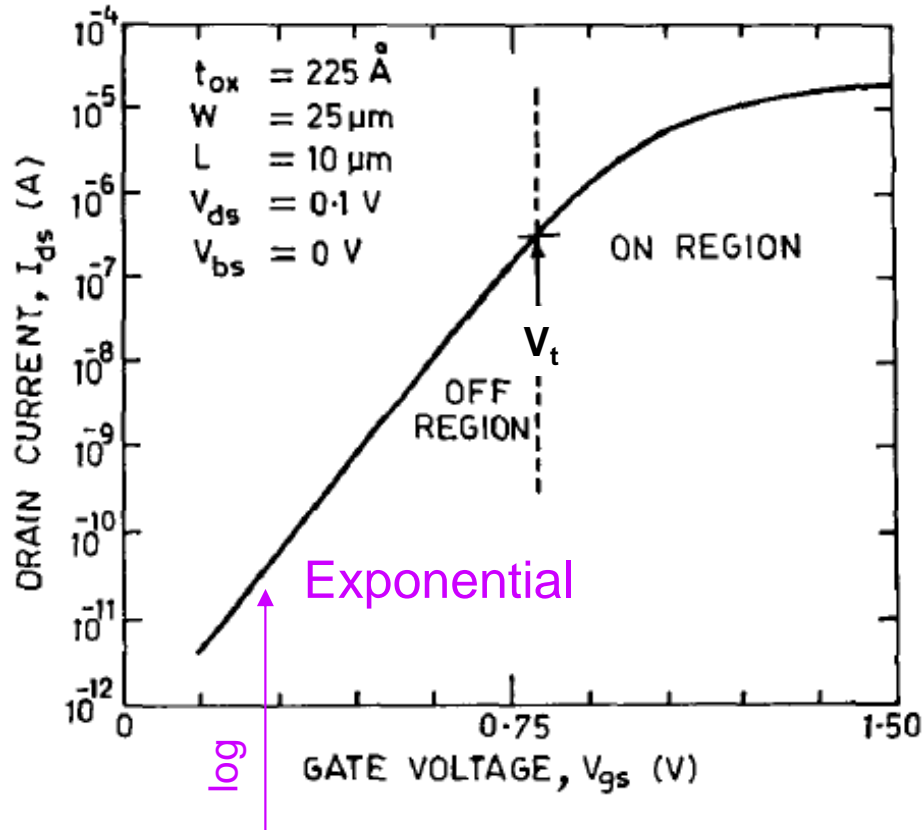$$\boxed{V_t(V_{sb}) = V_{t0} + \Upsilon\left(\sqrt{2\phi_F + V_{sb}} - \sqrt{2\phi_F}\right)}$$      $$\boxed{V_{t0} \equiv V_t\Big|_{V_{sb}=0} = V_{FB} + \Upsilon\sqrt{2\phi_F} + 2\phi_F}$$

# Regions of Operation: Transfer Characteristics

Left graph labels: $t_{ox} = 225$ Å, $W = 25 \mu m$, $L = 10 \mu m$, $V_{ds} = 0.1$ V, $V_{bs} = 0$ V; ON REGION; OFF REGION; $V_t$; Exponential; log; DRAIN CURRENT, $I_{ds}$ (A); GATE VOLTAGE, $V_{gs}$ (V)

Right graph labels: $W/L = 25/25$, $t_{ox} = 300$ Å; Quadratic; 4V; 2V; SATURATION REGION; LINEAR REGION; $V_{ds} = 0.5$ V; $V_t$; OFF; ON; Linear; linear; DRAIN CURRENT, $I_{ds}$ (mA); GATE VOLTAGE $V_{gs}$ (V)

Low gate–source bias ($V_{gs} < V_t$): No inversion layer; diffusion dominant. MOS behaves like a wide-base (long-channel) BJT with $I_{ds} \propto \exp[q(V_{gs} - V_t)/nkT]$, **n = 1 + $C_d/C_{ox}$**.

High drain–source bias ($V_{ds} > V_{dsat}$): Drain side "pinched-off". MOS behaves like a current source.

Low drain–source bias ($V_{ds} < V_{dsat}$): Full channel. MOS behaves like a voltage-controlled resistor.

# Current–Voltage in Linear (Triode) Region

- ❑ **MOSFET analysis** — major assumptions (NMOS as example)
  - ➢ *"GCA" – Gradual Channel Approximation*: $dE_y/dy << dE_x/dx$
  - ➢ *"Unipolar" – hole current can be neglected* ($E_{Fp} \approx E_F$) in normal region (excluding breakdown)
  - ➢ *Built-in voltages for the source/drain diodes can be ignored* (long channel)
  - ➢ *No recombination/generation and constant mobility*
  - ➢ *Current flows in the y direction only*
- ❑ **First-order equation derivation**
  - ➢ *Charge-sheet approximation (CSA)*
  - ➢ *"Pinned" surface potential at strong inversion ($2\phi_F$)*
  - ➢ *Constant bulk charge along channel*
  - ➢ *Drift-current only in linear region*

$$\boxed{\psi_s(y) = \psi_s(0) + V_{cb}(y) = 2\phi_F + V_{sb} + V(y)} \quad (0 \leq V \leq V_{ds})$$

$$V_{gb} - V_{FB} = V_{ox} + \psi_s = Q_g/C_{ox} + \psi_s = -(Q_b + Q_i)/C_{ox} + \psi_s$$

$$Q_i = -C_{ox}(V_{gb} - V_{FB} - \psi_s) - Q_b \qquad \boxed{Q_b \approx -\Upsilon C_{ox}\sqrt{2\phi_F + V_{sb}}}$$

$$= -C_{ox}\left[V_{gb} - V_{FB} - 2\phi_F - V_{sb} - V(y) - \Upsilon\sqrt{2\phi_F + V_{sb}}\right]$$

$$= -C_{ox}\left[V_{gs} - V_t - V(y)\right] \qquad \boxed{V_t \equiv V_{FB} + \Upsilon\sqrt{2\phi_F + V_{sb}} + 2\phi_F}$$

$$I_{ds}(y) \approx W\int_0^\infty J_{n,drift}(y)\,dx = W\int_0^\infty qn(x,y)\mu_n(-d\psi_s/dy)\,dx$$

$$= -W\mu_n Q_i(y)\,dV/dy \qquad \left[Q_i(y) \equiv \int_0^\infty qn(x,y)\,dx\right]$$

$$I_{ds} = \frac{W}{L}\mu_n\int_0^{V_{ds}} -Q_i(y)\,dV \quad \left(\int_0^L dy \sim \int_{\psi_s(0)}^{\psi_s(L)} d\psi_s = \int_{V_{sb}}^{V_{db}} dV_{cb} = \int_0^{V_{ds}} dV\right)$$

**Linear law** ($I_{ds}$ is a *linear* function of $V_{gs}$) ["Sah equation"]:

$$\boxed{I_{ds} = \mu_n C_{ox}\frac{W}{L}\left(V_{gs} - V_t - \frac{1}{2}V_{ds}\right)V_{ds}} \quad (V_{gs} > V_t, V_{gd} > V_t)$$

# Current–Voltage in Saturation (Pinch-off) Region
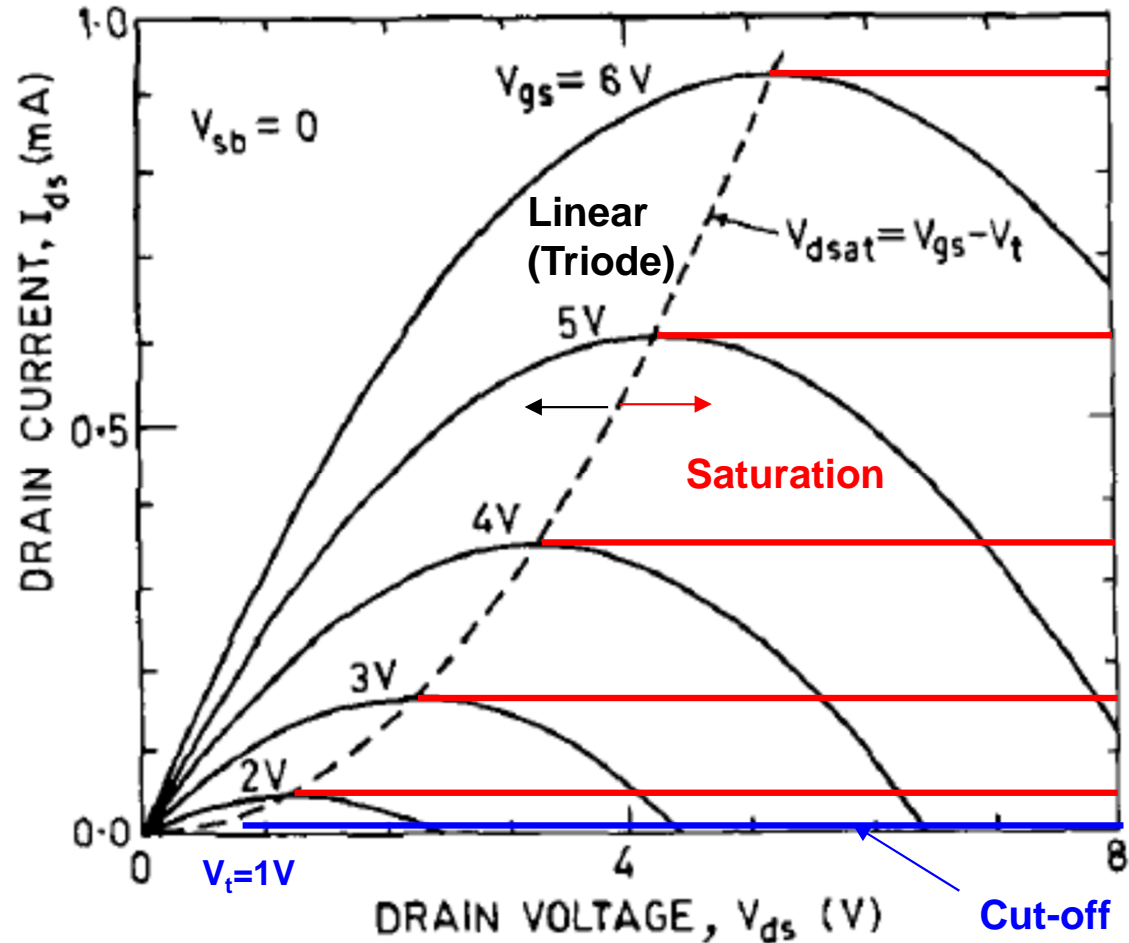
The peak of the linear current is reached when $dI_{ds}/dV_{ds} = \mu_n C_{ox} W/L(V_{gs} - V_t - V_{ds}) = 0$

For $V_{ds} \geq V_{gs} - V_t \equiv V_{dsat}$ , GCA is not valid. Also, $Q_i(V = V_{dsat}) \approx 0$ , channel is said to be "pinched-off." $V_{dsat}$ is called *saturation* or *pinch-off voltage*, and the corresponding current is the *saturation current*.

**Square law** ($I_{ds}$ is a *quadratic* function of $V_{gs}$):

$$\boxed{I_{ds} = \frac{1}{2}\mu_n C_{ox}\frac{W}{L}(V_{gs} - V_t)^2 = I_{dsat}} \quad \begin{pmatrix} V_{gs} > V_t, \\ V_{gd} < V_t \end{pmatrix}$$

The "pinch-off" picture ($Q_i = 0$ assumption) is not physically correct since it requires the field to be infinite $E_y(y) = J_{ds}(y)/\mu_n Q_i(y)$ at pinch-off and carriers travel with infinite drift velocity. A more correct picture is that $Q_i$ at pinch-off is very small but finite, with carriers drift under the large field in the pinch-off region at a saturated velocity.

MOSFET first-order piece-wise linear/square-law model.
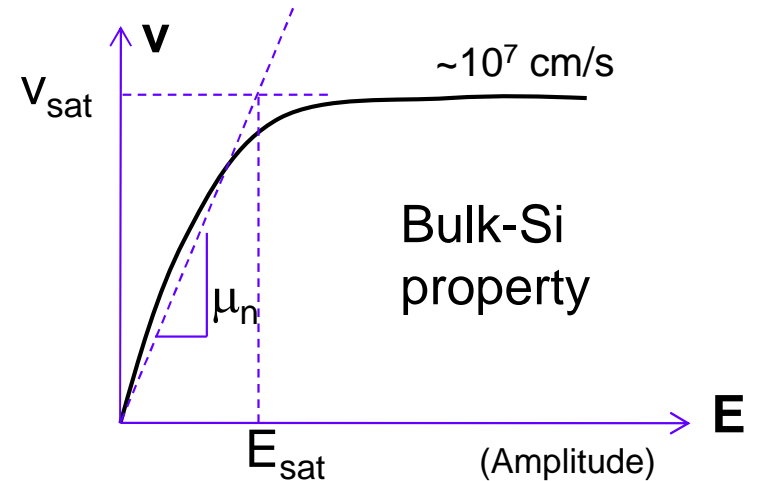
# Velocity Saturation and Saturation Current

❑ **Velocity–field relation** — piecewise model

$$v = \begin{cases} \dfrac{\mu_n E}{1 + E/E_{sat}} & E < E_{sat} \\[2ex] v_{sat} & E \geq E_{sat} \end{cases}$$

$$I_{ds}(y) \approx -WQ_i(y)\frac{\mu_n E}{1 + E/E_{sat}}$$

$$I_{ds}(y)\left(1 + \frac{1}{E_{sat}}\frac{dV}{dy}\right) = -WQ_i(y)\mu_n\frac{dV}{dy}$$

$$Q_i = -C_{ox}\left(V_{gb} - V_{FB} - (2\phi_F + V_{sb} + V)\right) - Q_b$$

V

$V_{sat}$

~$10^7$ cm/s

Bulk-Si property

$\mu_n$

$E_{sat}$

E

(Amplitude)

➢ **Saturation field**

$$v_{sat} = \frac{\mu_n E_{sat}}{1 + E_{sat}/E_{sat}} \rightarrow E_{sat} = \frac{2v_{sat}}{\mu_n}$$

➢ **Lateral-field mobility**

$$\mu_{eff} = \frac{\mu_n}{1 + V_{ds}\big/\left(E_{sat}L_{eff}\right)}$$

$$Q_b \approx -\Upsilon C_{ox}\sqrt{2\phi_F + V_{sb} + V}$$

$$A_b = 1 + \frac{\Upsilon}{2\sqrt{2\phi_F + V_{sb}}}$$

❑ **Saturation current**

$$I_{dsat} = -Wv_{sat}Q_{sat} = Wv_{sat}C_{ox}\left(V_{gs} - V_t - A_b V_{dsat}\right) \quad (2)$$
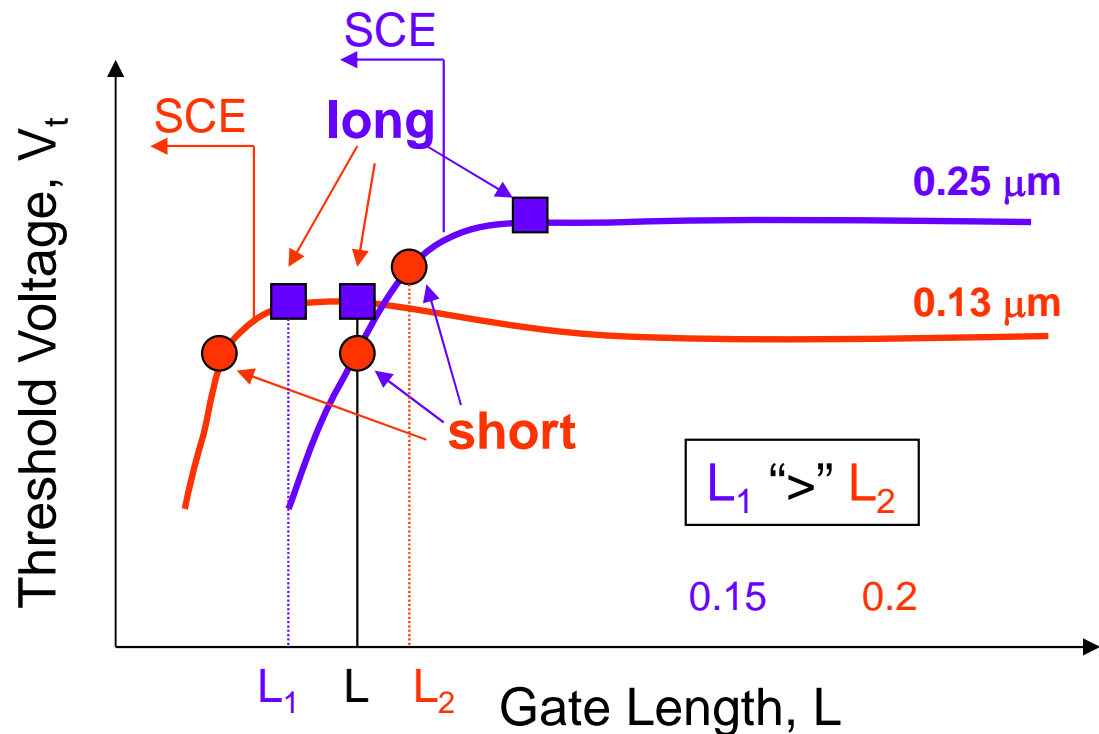
$(1)(V_{dsat}) = (2): \quad V_{dsat} = \dfrac{E_{sat}L_{eff}\left(V_{gs} - V_t\right)}{V_{gs} - V_t + A_b E_{sat} L_{eff}}$

$$I_{ds} = \mu_{eff}C_{ox}\frac{W}{L}\left(V_{gs} - V_t - \frac{1}{2}A_b V_{ds}\right)V_{ds} \quad (1)$$

$$I_{dsat} = Wv_{sat}C_{ox}\frac{\left(V_{gs} - V_t\right)^2}{V_{gs} - V_t + A_b E_{sat} L_{eff}} \xrightarrow{L \to 0} \propto \left(V_{gs} - V_t\right)$$

**Linear!**

# Long-Channel or Short-Channel?

❑ **Short-channel effect (SCE)** — technology dependent (depends on where the device "sits" on the $V_t$ – L curve, not the actual dimension)

❑ **Technology scaling** — optimization for each technology node

# Charge-Sharing Model: $V_t$ "Roll-Off"

❑ **Charge-sharing model**

➤ **Without charge-sharing**

$$V_t = V_{FB} - Q_{bm}/C_{ox} + 2\phi_F$$
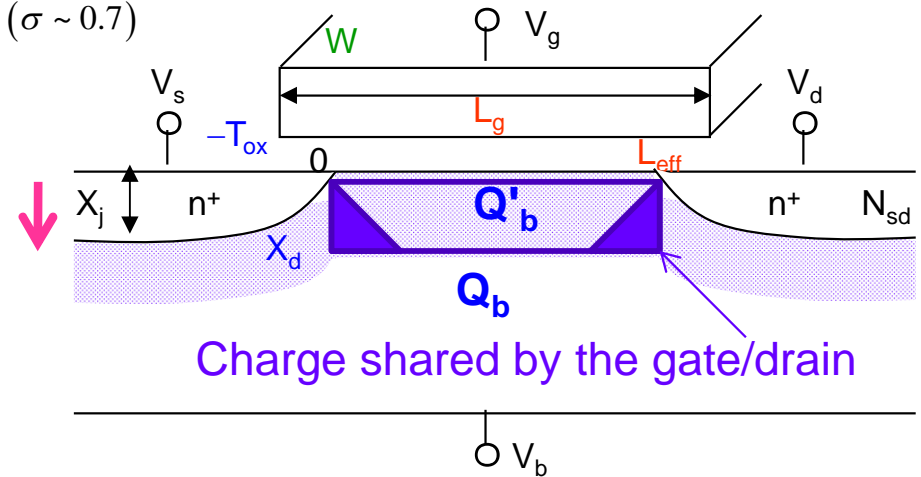
➤ **With charge-charing**

$$V_t' = V_{FB} - Q_{bm}'/C_{ox} + 2\phi_F$$

$$\boxed{L_{eff} = L_g - 2\sigma X_j} \quad (\sigma \sim 0.7)$$

$$C_{ox} = \varepsilon_{ox}/T_{ox} \qquad Q_{bm} = -qN_A X_{dm}$$

$$X_{dm} = \sqrt{2\varepsilon_{Si}(2\phi_F + V_{sb})/qN_A}$$

$$\Delta V_t \equiv V_t - V_t' = -\frac{Q_{bm}}{C_{ox}}\left(1 - \frac{Q_{bm}'}{Q_{bm}}\right) = -\frac{Q_{bm}}{C_{ox}}\frac{X_{dm}}{L_{eff}}$$

$$= \frac{qN_A X_{dm}}{\varepsilon_{ox}/T_{ox}}\frac{X_{dm}}{L_{eff}} = \frac{4\varepsilon_{Si}\phi_F}{\varepsilon_{ox}}\frac{T_{ox}}{L_g - 2\sigma X_j}$$

**Short-channel effect (SCE): $V_t$ "roll-off"**
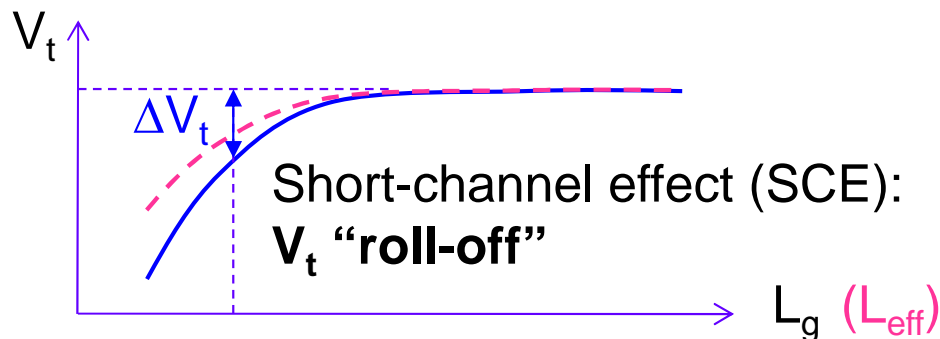
$V_t$ / $\Delta V_t$ / $L_g$ ($L_{eff}$)

**Charge shared by the gate/drain**

## Simple "Triangle" Model

Total bulk charge:
$$Q_B' = -qN_A W\left(L_{eff} X_d - X_d^2\right)$$
$$Q_B = -qN_A W\left(L_{eff} X_d\right)$$

Bulk charge per unit area:
$$\therefore \frac{Q_b'}{Q_b} = \frac{Q_B'}{Q_B} = 1 - \frac{X_d}{L_{eff}}$$

# DIBL and Reverse SCE: $V_t$ "Roll-Up"

$V_t$

$\Delta V_{t0}$

$\Delta V_{DIBL}$

$V_{t0}$

$V_{ts}$

Drain-induced barrier lowering: **DIBL**

$L_g$

"**Halo**"

$W$  $V_g$  $V_d$

$V_s$

$-T_{ox}$  $0$  $L_g$  $L_{eff}$  $y$

$X_j$  $n^+$  $Q'_b$  $n^+$  $N_{sd}$

$X_d$  $Q_b$

Charge shared by the drain

$V_b$

Reverse SCE: "**Halo**"

$V_t$

$V_{t0}$

$L_g$

$N_A$

$y$

# Summary of Important Equations

❑ **Threshold voltage**

➢ *Long-channel (1D theoretical model)*

$$\phi_F = \frac{kT}{q} \ln\left(\frac{N_A}{n_i}\right) \qquad \Upsilon = \sqrt{2q\varepsilon_{Si}N_A}\big/C_{ox} \qquad C_{ox} = \varepsilon_{ox}\big/T_{ox}$$

$$\boxed{V_t \equiv V_{gs}\Big|_{\psi_s = 2\phi_F + V_{sb}} = V_{FB} + \Upsilon\sqrt{2\phi_F + V_{sb}} + 2\phi_F}$$

$$V_{FB} \equiv \phi_{MS} - Q_{ox}\big/C_{ox} = \Phi_M - \left(\chi + E_g\big/2 + \phi_F\right) - Q_{ox}\big/C_{ox}$$

➢ *Short-channel (triangle charge-sharing model)*

➢ **Short-channel DIBL**

$$\boxed{V_{t0}\left(L_g\right) \equiv V_{t0\_long} - \Delta V_{t0} = V_{t0\_long} - \frac{4\varepsilon_{Si}\phi_F}{\varepsilon_{ox}}\frac{T_{ox}}{L_g - 2\sigma X_j}}$$

$$\boxed{\Delta V_{DIBL}\left(L_g\right) \equiv V_{t0}\left(L_g\right) - V_{ts}\left(L_g\right)}$$

❑ **Drain current**

➢ *Linear*

➢ *Subthreshold*

$$n = 1 + C_d\big/C_{ox}$$

$$\boxed{I_{ds} = \mu_{eff}C_{ox}\frac{W}{L}\left(V_{gs} - V_t - \frac{1}{2}A_b V_{ds}\right)V_{ds}}$$

$$\boxed{I_{ds} = \mu_n C_d v_{th}^2 \frac{W}{L} e^{(V_{gs}-V_t)/(nv_{th})}\left(1 - e^{-V_{ds}/v_{th}}\right)}$$

$$C_d = \varepsilon_{Si}\big/X_{dm}$$

➢ *Saturation*

$$= \frac{\Upsilon C_{ox}}{2\sqrt{2\phi_F + V_{sb}}}$$

$$\boxed{I_{dsat} = W v_{sat} C_{ox}\frac{\left(V_{gs} - V_t\right)^2}{V_{gs} - V_t + A_b E_{sat} L_{eff}}} \Rightarrow \propto \begin{cases} \left(V_{gs} - V_t\right)^2 & \left(L_{eff} \to \infty; \text{long-channel: quadratic}\right) \\ \left(V_{gs} - V_t\right) & \left(L_{eff} \to 0; \text{short-channel: linear}\right) \end{cases}$$

# Gate-Controlled Drift ("ON") and Diffusion ("OFF")

# Threshold Voltage Definition ($I_{crit}$@$V_{t0}$)