**ROM (Read Only Memory) – nMOS ROM**

PLA

decoder — $V_{DD}$ — $V_{DD}$ — $I_1$ (WL) — ROM (NOR) — $I_2$ — $I_3$ — $I_4$ (WL$_4$) — selector

$A_2$   $A_1$   $A_0$   $Z_1$ (BL$_1$)   (BL$_2$)  $Z_2$   $Z_3$ (BL$_3$)

---

$$I_1 = \overline{A_2 + \bar{A}_1}$$
$$I_2 = \overline{A_2 + \bar{A}_1}$$
$$I_3 = \overline{\bar{A}_2 + A_1}$$
$$I_4 = \overline{A_2 + A_1}$$

} NOR array

$$Z_1 = A_0 \,\overline{I_1 + I_2 + I_3} + \bar{A}_0\, I_2$$
$$Z_2 = A_0\, I_3 + \bar{A}_0\, I_4$$
$$Z_3 = A_0\, I_1 + \bar{A}_0\, I_2$$

} NOR array gated by $A_0, \bar{A}_0$

---

**NOR decoder – NOR ROM Array**

NOR row decoder  →  NOR ROM array

1, 2, 3, ... $2^N$ WLs

1 2 3 ··· N Address bits          1 2 3 ·· $2^M$ columns (BLs)

---

**4 × 4 NOR-based ROM**

$V_{DD}$

$R_1$
$R_2$
$R_3$
$R_4$

$C_1$   $C_2$   $C_3$   $C_4$

| $R_1$ | $R_2$ | $R_3$ | $R_4$ | $C_1$ | $C_2$ | $C_3$ | $C_4$ |
|---|---|---|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 |

---

**NAND-row decoder – NAND ROM Array**

NAND row decoder  →  NAND ROM array

1, 2, 3, ... $2^N$ WLs

1 2 3 ··· N Address bits          1 2 3 ··· $2^M$ columns (BLs)

---

**An Example of Column Decoder Circuit** (Binary Tree Decoder)

$C_1$ $C_2$ $C_3$ $C_4$   $C_5$ $C_6$ $C_7$ $C_8$

$B_1$
$\bar{B}_1$
$I_1$   $I_2$

$B_2$
$\bar{B}_2$
$J_1$        $J_2$

$B_3$
$\bar{B}_3$

Strength: compact, less area, less power

Weakness: many NMOS are in series slow!

$C_1$ or $C_2$ chosen       $I_1$ or $I_2$ chosen       Data Output      $J_1$ or $J_2$ chosen

## 4×4 NAND-based ROM

depletion NMOST

$C_1$ $C_2$ $C_3$ $C_4$

$R_1$ $R_2$ $R_3$ $R_4$

| $R_1$ | $R_2$ | $R_3$ | $R_4$ | $C_1$ | $C_2$ | $C_3$ | $C_4$ |
|---|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 |
| 1 | 0 | 1 | 1 | 0 | 0 | 1 | 1 |
| 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 |
| 1 | 1 | 1 | 0 | 0 | 1 | 1 | 0 |
| 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 | 1 | 1 | 1 |
| : | | | | | | | |

## Fabrication

n- diffusion

poly

threshold voltage implant to below 0 such that transistor is on (i.e. shorted)

First column   2nd column

## A Structure of CMOS Dynamic PLA (Layout)

$V_{DD}$

Decoder clock

Precharge

$A_0$

$A_1$

$A_2$

ground switch

$V_{SS}$ (ground)

$\alpha$   $\beta$   $\gamma$   $\delta$   $\epsilon$   $z_1$

$z_2$

$z_3$

ROM clock   Precharge

---

Precharge   $V_{DD}$   ROM clock   Word Line   Botline

$A_0$

$V_{SS}$

$V_{SS}$

domino CMOS (ground switched) NOR array   (domino circuit)

domino CMOS (ground switched NOR array)

Ref: R.H. Krambeck, C.M. Lee and H.F. Law, "High-Speed Compact Circuits with CMOS", IEEE J. of Solid-State Circuits, SC-17 (1982)

---

$$\tau_{HL} = \frac{b\,(C_{wire} + N\,C_{drain})\,V_{DD}}{I_{drive}} < \tau_{spec}$$

For $\tau_{spec} = 0.25\,ns$ (WL/BL delay)

$C_{drain} = 0.01\,fF$

$V_{dd} = 1.2\,V$

$b = 2$ (emperical constant)

$I_d = 0.6\,\mu A$ , $C_{wire} < N\,C_{drain}$

$$N < \frac{0.25\times 10^{-9} \times 0.6\times 10^{-6}\,A}{2\,(0.001\times 10^{-15})\times 1.2V} \quad \frac{0.15\times 10^{-15}}{2.4\times 10^{-15}\times 10^{-3}}$$

$\cong 62$ (number of nMOSTs per WL/BL)

## CMOS reference

F.M. Wanlass and C.T. Sah, "Nanowatt Logic using Field-Effect Metal-Oxide Semiconductor Triodes," IEEE Solid-State Circuits Conf. Philadelphia, PA (1963) (First CMOS paper)

## Nonvolatile Semiconductor Memory

Metal-Insulator-Semiconductor (MIS) Structure

metal
insulator
n+   n+   Semiconductor
P-substrate

metal
insulator
metal
n+   n+   insulator/
2-substrate

By D. Kahng and S.M. Sze
Bell System Technical Journal (1967)

**Slide 1:**

<u>Operation Principle of</u>
metal - insulator - metal - insulator - semiconductor
nonvolatile memory



semiconductor    metal 1
                        metal 2
insulator 1    insulator 2

<u>Electrons are injected into the floating gate (metal 1)</u>
by "tunneling" through insulator 1, and are stored
semi-permanently.

**Slide 2:**



$SiO_2$    floating gate
n+    ↑ thick    n+
p-type substrate

electrons are injected into the "floating gate"
by crossing the oxide ($SiO_2$) potential barrier
and are stored in the floating gate.
Thus, named FAMOS ( Floating-gate
Avalanche injection MOS)

**Slide 3:**

MNOS Structure ( Metal Nitride Oxide Semicon)



Metal
$Si_3N_4$ (~500Å)    FOX
FOX
$SiO_2$ (~20Å)
substrate

Hysteresis Loop

Erasing voltage
electrons are trapped
can be as high as + 60 V
threshold voltage
Writing voltage applied to the <u>metal gate</u>. ( $V_4$ )
electrons are released
changes A → B    B → A

**Slide 4:**

Energy Band Diagram

M (metal)    N (nitride)    O (oxide)    S (semiconductor)



Forward Bias
$V_C$
$E_{FM}$
$J_N$    $J_O$    $E_C$    $E_F$
Fowler-Nordheim tunneling

Current discontinuity model by
Frohman-Bentchkowsky

$E_{FM}$
$V_C$    $J_N$    $J_O$
$V_C$    $E_F$

**Slide 5:**

$$J_N = C_N E_N^2 \exp(- E_1/E_N)$$

$$J_{ox} = C_{ox} E_{ox}^2 \frac{\pi c kT/E_{ox}}{\sin(\pi kT/E_{ox})} \exp(- E_2/E_{ox})$$

$E_N$ = electric field in the nitride layer
$E_{ox}$ = "         "     "      oxide   layer
$k$ = Boltzman's constant
$T$ = temp in [°K]
$J_N$ = nitride layer current density
$J_0 (J_{ox})$ = oxide layer      "      $\boxed{V_4 = E_{ox} t_{ox} + E_N t_N}$

e.g. $C_{ox} = 10^{-5}$ A/V²    $C_N = 3.5 \times 10^{-10}$ A/V²
$E_2 = 2.54 \times 10^9$ V/cm    $E_1 = 1.2 \times 10^8$ V/cm
$t_{ox} = 50$ Å    $t_N = 1000$ Å

**Slide 6:**

Stacked Gate Tetrode Proposed by
H. G. Dill & T. N. Toombs (1969)



offset gate ($\phi_2$)
$Si_3N_4$
Source    Drain
$SiO_2$
Depletion Region
control gate ($\phi_1$)    substrate

## Flash Memory

Memory cell is a transistor with a **floating gate** whose threshold voltage can be programmed (changed) repeatedly by applying an electric field (through $V_G$ voltage) to its gate.

$V_D < V_G$ — control gate, thick poly, floating gate, thin ($SiO_2$) tunneling oxide say 10 nm, $V_S$, $V_D$, $e^-$, $n+$ source, $n+$ drain, P-substrate

(high voltage) — control gate, $V_{pp}$, floating gate, open, Fowler-Nordheim tunneling, $n+$, $n+$ drain, P-substrate

---

control gate, $V_G = V_{CG}$, $C_{FC}$, floating gate, $V_{FG}$, $C_{FS}$, $C_{FB}$, $C_{FD}$, $V_D$, $n+$ source, $n+$ drain, P-substrate

$$C_{total} = C_{FC} + C_{FS} + C_{FB} + C_{FD}$$

$$V_{FG} = \frac{Q_{FG}}{C_{FG}} + \frac{C_{FC}}{C_{total}} V_{CG} + \frac{C_{FD}}{C_{total}} V_D$$

$Q_{FG}$ = charge stored in the floating gate

$$V_T(CG) = \frac{C_{total}}{C_{FC}} V_T(FG) - \frac{Q_{FG}}{C_{FC}} - \frac{C_{FD}}{C_{FC}} V_D$$

---

$$\Delta V_T(CG) = - \frac{\Delta Q_{FG}}{C_{FC}}$$
between '0' & '1'

$I_D$ vs $V_{CG}$ — Low $V_T$ data 0, High $V_T$ data 1, $\Delta V_T$, Control Gate voltage, $V_R$, control gate voltage $V_R$

### NOR Flash Memory

$WL_1$, $WL_2$, $BL_1$, $BL_2$, Source line

Operation

| Signal | Erase | Program | Read |
|---|---|---|---|
| $BL_1$ | open | 6V | 1V |
| | open | 0V | 0V |
| Source line | 12V | 0V | 0V |
| $WL_1$ | 0 | | |
| $WL_2$ | 0 | 12V | 5V |

---

NOR array faster but many contacts (area ↑)
*NAND array* slower but **compact** (area ↓)

$BL$, select1, $WL_1$, $WL_2$, $WL_3$ (selected), $WL_4$, $WL_5$, $WL_6$, select2, Source Line — selected → write/read

Operation

| Signal | Erase | Program | Read |
|---|---|---|---|
| $BL_1$ | open | 0V | 1V |
| $WL_1$ | 0 | 10V | 5V |
| 2 | 0 | 10V | 5V |
| 3 | 0 | 20V | 0V |
| 4 | 0 | 10V | 5V |
| 5 | 0 | 10V | 5V |
| 6 | 0 | 10V (Vpp) | 5V |
| Source line | 20V | 0V | 0V |
| select line 2 | open | 0V | 5V |
| P-well | 20V | 0V | 0V |
| n-tub | 20V | 0V | 0V |

---

### Charge Pumping for High Voltage Vpp generation

$V_{in}$, $M_1$, $M_2$, $M_3$, $M_4$, $M_{N-1}$, $M_N$, $V_{out}$, charge transfer

clock, clock, clock, clock

Voltage: $V_2$, $V_1$, $V_{in}$, $V_{in} - V_T(M_1)$

$$V_{out} = V_{in} + (\gamma V_{DD} - V_T(M_1)) + \cdots + (\gamma V_{DD} - V_T(M_N))$$

$\gamma$ = boosting efficiency factor

---

### Multi-level-cell Threshold Voltage Distribution in Flash (4 values) — 2 bits/cell

Population — state 0 (11), state 1 (01), state 2 (10), state 3 (00)

$R_1$, $R_2$, $R_3$, cell threshold voltage

## Slide 1 (top-left)

*Subthreshold operation*

### Opportunities for Ultra-Low Voltage

- Number of applications emerging that do not need high performance, only extremely low power dissipation
- Examples:
  - Standby operation for mobile components
  - Implanted electronics and artificial senses
  - Smart objects, fabrics, and e-textiles
- Need power levels below 1 mW (even µW in certain cases)

**Slide 11.4**
Although keeping the power density constant is one motivation for the continued search to lower the EOP, another, maybe even more important, reason is the exciting applications that only become feasible at very low energy/power levels. Consider, for instance, the digital wrist-watch. The concept, though straightforward, only became attractive once the power dissipation

$[R_1]$

## Slide 2 (top-right)

### Minimum Operational Voltage of Inverter

- Swanson, Meindl (April 1972)
- Further extended in Meindl (Oct 2000)

**Limitation: gain at midpoint** $\leftarrow$

$$V_{DD}(\text{Meindl}) = 2\frac{kT}{q}\ln(2 + \frac{C_{fix}}{C_{ox}})$$

or

$$V_{DD}(\text{Meindl}) = 2\frac{kT}{q}\ln(1 + n)$$

$C_{ox}$: gate capacitance
$C_r$: diffusion capacitance
$n$: slope factor

For ideal MOSFET (60 mV/decade slope):

$$V_{DD}(\text{Meindl}) = 2\ln(2)\frac{kT}{q} = 1.38\frac{kT}{q} = 0.036 \text{ V}$$

at a temperature of 300 K

[Ref: R. Swanson, JSSC'72, J. Meindl, JSSC'00]

$gain < -1$
$|gain| > 1$

**Slide 11.5**
The question of the minimum operational voltage of a CMOS inverter was addressed in a landmark paper [Swanson72] in the early 1970s – published even before CMOS integrated circuits came in vogue! For an inverter to be regenerative and to have two distinct steady-state operation points (a "1" and a "0"), it is essential that the absolute value of the gain of the gate in the transition region be larger than 1. Solving for those conditions leads to an expression for $V_{min}$ equal to $2(kT/q)\ln(1 + n)$, where n is the slope factor of the transistors. One important observation is that $V_{min}$ is proportional to the operational temperature $T$. Cooling down a CMOS circuit to temperatures close to absolute zero (e.g., liquid Helium), makes operation at mV levels possible. (Unfortunately, the energy going into the cooling more than often offsets the gains in operational energy.) Also, the closer the MOS transistor operating in sub-threshold mode gets to the ideal bipolar transistor behavior, the lower the minimum voltage. At room temperature, an ideal CMOS inverter (with a slope factor of 1) could marginally operate at as low as 36 mV.

## Slide 3 (middle-left)

$$I_{DS} = I_S e^{\frac{V_{GS}-V_{th}}{nkT/q}}\left(1 - e^{-\frac{V_{DS}}{kT/q}}\right) = I_0 e^{\frac{V_{GS}}{nkT/q}}\left(1 - e^{-\frac{V_{DS}}{kT/q}}\right)$$
where $I_0 = I_S e^{-\frac{V_{th}}{nkT/q}}$

### Sub-threshold Modeling of CMOS Inverter

- From Chapter 2:

$$I_{DS} = I_0 e^{\frac{V_{GS}}{nkT/q}}\left[1 - e^{-\frac{V_{DS}}{kT/q}}\right] = I_0 e^{\frac{V_{GS}}{nkT/q}}\left[1 - e^{-\frac{V_{DS}}{kT/q}}\right]$$

where

$$I_0 = I_S e^{-\frac{V_{th}}{nkT/q}}$$

(DIBL can be ignored at low voltages)

**Slide 11.6**
Given the importance of this expression, a quick derivation is worth undertaking. We assume that at these low operational voltages, the transistors operate only in the sub-threshold regime, which is often also called the *weak-inversion* mode. The current-voltage relationship for a MOS transistor in sub-threshold mode was presented in Chapter 2, and is repeated here for the sake of clarity. For low values of $V_{DS}$ the DIBL effect can be ignored.

## Slide 4 (middle-right)

$x_i = V_i/\phi_T$   $x_o = V_o/\phi_T$   $x_D = V_{DD}/\phi_T$
thermal voltage
$x_o = x_D + \ln\left(\frac{1 - G + \sqrt{(G-1)^2 + 4Ge^{x_o}}}{2}\right)$, $G = e^{\frac{x_i - x_o}{n}}$

### Sub-threshold DC model of CMOS Inverter

Assume NMOS and PMOS are fully symmetrical and all voltages normalized to the thermal voltage $\phi_T = kT/q$
$(x_i = V_i/\phi_T; x_o = V_o/\phi_T; x_D = V_{DD}/\phi_T)$
The VTC of the inverter for NMOS and PMOS in sub-threshold can be derived:

$$x_o = x_D + \ln\frac{1 - G + \sqrt{(G-1)^2 + 4Ge^{x_o}}}{2} \text{ where } G = e^{\frac{x_i - x_o}{n}}$$

so that

$$A_v = \frac{2(1 - e^{x_o - x_D}) - e^{-x_o} - e^{-x_D}}{n(2e^{x_o - x_D} - e^{(x_o - x_D)} - e^{-x_o})} \text{ and } A_{vmax} = -(e^{x_D/2} - 1)/n$$
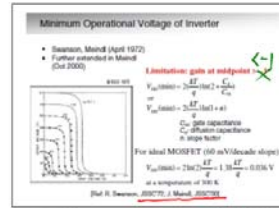
For $|A_{vmax}| = 1$: $x_D = 2n(n+1)$

[Ref: E. Vittoz, CRC'05]

$K_p = K_n$
$K_x = \mu \frac{W}{L} \times C_{ox}$

$$A_N = -\frac{2(1 - e^{x_o - x_D}) - e^{-x_o} - e^{-x_D}}{n(2e^{x_o - x_D} - e^{(x_o - x_D)} - e^{-x_o})}$$

$|Avmax| = 1$, $x_D = 2\ln(n+1)$

**Slide 11.7**
The (static) voltage transfer characteristic (VTC) of the inverter is derived by equating the current through the NMOS and PMOS transistors. The derivation is substantially simplified if we assume that two devices have exactly the same strength when operating in sub-threshold. Also, normalizing all voltages with respect to the thermal voltage $\phi_T$ leads to more elegant expressions. Setting the gain to –1 yields the same expression for the minimum voltage as was derived by Swanson.

## Slide 5 (bottom-left)

### Results from Analytical Model

**Sub-threshold Inverter**

Normalized VTC for n=1.5 as a function of $V_{DD}$ ($x_D$)

Minimum supply voltage for a given maximum gain as a function of the slope factor n

$n = 1.5$

$x_D = 4$ sufficient for reliable operation

$x_{Dmin} = 2\ln(2.5) = 1.83$ for $n = 1.5$

[Ref: E. Vittoz, CRC'05]

**Slide 11.8**
thermal voltage leads to reasonable noise margins (assuming n = 1.5). This is approximately equal to 100 mV.

$x_D = V_{DD}/\phi_T$   $\frac{kT}{q}$ (thermal voltage)
$x_D = 4 \approx 100 \text{ mV}$   $= 26 \text{ mV at } T = 300 \text{ K (room temp.)}$

## Slide 6 (bottom-right)

### Confirmed by Simulation (at 90 nm)

Minimum operational supply voltage

For n = 1.5,
$V_{DDmin} = 1.83 \cdot \phi_T$
$= 48 \text{ mV}$

Observe: non-symmetry of VTC increases $V_{DDmin}$

**Slide 11.9**
Simulations (for a 90 nm technology) confirm these results. When plotting the minimum supply voltage as a function of the PMOS/NMOS ratio, a minimum can be observed when the inverter is completely symmetrical, that is when the PMOS and NMOS transistors have identical drive strengths. Any deviation from the symmetry causes $V_{min}$ to rise. This implies that transistor sizing will play a role in the design of minimum-voltage circuits.
Also worth noticing is that the simulated minimum voltage of 60 mV is slightly higher than the theoretical value of 48 mV. This is mostly owing to the definition of "operational" point. At 48 mV, the inverter is only marginally functional. In the simulation, we assume a small margin of approximately 25%.

## Also Holds for More Complex Gates

Minimum operational supply voltage (two-input NOR)

Degradation due to asymmetry

When only one input is switched (and the other fixed to "0") the built-in asymmetry leads to a higher minimum voltage than when both are switched simultaneously. This leads to a useful design rule-of-thumb when designing logic networks for minimum-voltage operation, *one should strive to make the gate topology symmetrical over all input conditions.*

---

## Minimum Energy per Operation

Predicted by von Neumann $kT\ln(2)$

J. von Neumann,
[Theory of Self-Reproducing Automata, 1966]

- Moving one electron over $V_{DDmin}$:
  - $Emin = CV_{DD}^2 = 2(n2)kT/q = k$ Fin(2)
  - Also called the Von Neumann–Landauer–Shannon bound
  - At room temperature (300 K): Emin = 0.29 $\times 10^{-19}$ J
- Minimum sized CMOS inverter at 90 nm operating at 1V
  $E = CV_{DD}^2 = 0.8 \times 10^{-15}$ J, or 5 orders of magnitude larger!
  How close can one get?

$0.29 \times 10^{-19}$ J
$0.8 \times 10^{-15}$ J

[Ref: J. Von Neumann, 1966]

**Slide 11.11**
Now that the issue of the minimum voltage is settled, the question of the minimum energy per operation (EOP), can be tackled. In a follow-up to the 1972 paper by Swanson, Meindl [Meindl, JSSC'00] Jagaud that moving a single electron over the minimum voltage requires an energy equal to $kT\ln 2/\lambda$. This result is remarkable in a number of ways.

- This expression for the minimum energy for a digital operation was already predicted much earlier by John von Neumann (as reported in [von Neumann, 1966]). Landauer later established that this is only the case for "logically irreversible" operations in a physical computer that dissipate energy by generating a corresponding amount of entropy for each bit of information that then gets irreversibly erased. This bound hence does not hold for reversible computers (if such could be built) [Landauer, 1961].

---

## Propagation Delay of Sub-threshold Inverter

$$t_p = \frac{CV_{DD}}{I_{on}} = \frac{CV_{DD}}{I_S e^{\frac{V_{DD}}{n\phi_T}}} \quad (\text{for } V_{DD} \gg \phi_T)$$

Normalizing $t_p$ to $\tau_0 = C\phi_T/I_S$

$$\frac{t_p}{\tau_0} = x_D e^{-x_D/n}$$

$\tau_0 = \frac{C\phi_T}{I_S}$

$\frac{\tau}{\tau_0} = x_D e^{-x_D/n}$

$x_D = \frac{V_{DD}}{\phi_T}$

$E = p \cdot \tau$

Comparison between curve-fitted model and simulations (FO4, 90 nm)

$\frac{\tau}{\tau_0} = \frac{C}{I_S} \frac{V_{DD}}{e^{V_{DD}/n\phi_T}} \Big/ \frac{C\phi_T}{I_S} = \frac{x_D}{e^{x_D/n}} = x_D e^{-x_D/n}$  normalize delay

**Slide 11.12**
The above analysis, though useful in setting absolute bounds, ignores some practical aspects, such as leakage. Hence, lowering the voltage as low as we can may not necessarily be the right answer to minimize energy.

Operating an inverter in the sub-threshold region, however, may be one way to get closer to the minimum-energy bound. Yet, as should be no surprise, this comes at a substantial cost in performance. Following the common practice of this book, we again map the design task as an optimization problem in the E–D space.

One interesting by-product of operating in the sub-threshold region is that the equations are quite simple and are exponentials (as used to be the case for bipolar transistors). Under the earlier assumptions of symmetry, an expression of the inverter delay is readily derived. Observe again that a reduction in supply voltage has an exponential effect on the delay!

---

*Table ORTC1     Summary Table of ITRS Technology Trend Targets*

| Year of Production | 2013 | 2015 | 2017 | 2019 | 2021 | 2023 | 2025 | 2028 |
|---|---|---|---|---|---|---|---|---|
| Logic Industry "Node Name" Label | "16/14" | "10" | "7" | "5" | "3.5" | "2.5" | "1.8" | |
| Logic ½ Pitch (nm) | 40 | 32 | 25 | 20 | 16 | 13 | 10 | 7 |
| Flash ½ Pitch (2D) (nm) | 18 | 15 | 13 | 11 | 9 | 8 | 8 | 8 |
| DRAM ½ Pitch (nm) | 28 | 24 | 20 | 17 | 14 | 12 | 10 | 7.7 |
| FinFET Fin Halfpitch (new) (nm) | 30 | 24 | 19 | 15 | 12 | 9.5 | 7.5 | 5.3 |
| FinFET Fin Width (new) (nm) | 7.6 | 7.2 | 6.8 | 6.4 | 6.1 | 5.7 | 5.4 | 5.8 |
| 6-t SRAM Cell Size(nm2) (W#02) | 0.096 | 0.061 | 0.038 | 0.024 | 0.015 | 0.010 | 0.0060 | 0.0030 |
| MPU/ASIC HighPed to NAND Gate Sizes(nm2) | 0.248 | 0.157 | 0.099 | 0.062 | 0.039 | 0.025 | 0.016 | 0.008 |
| 2-input NAND Gate Density (Equivalent) (W135/2) | 4.03E+03 | 6.37E+03 | 1.01E+04 | 1.61E+04 | 2.55E+04 | 4.05E+04 | 6.42E+04 | 1.28E+05 |
| Flash Generations Label (bits per chip) (SLC/MLC) | 64G /128G | 128G /256G | 256G / 512G | 512G / 1T | 512G / 1T | 1T / 2T | 2T / 4T | 4T / 8T |
| Flash 2D Number of Layer targets (at relaxed Poly half pitch) | 16-32 | 16-32 | 16-32 | 32-64 | 48-96 | 64-128 | 96-192 | 192-384 |
| Flash 3D Layer half-pitch targets (nm) | 64nm | 54nm | 45nm | 30nm | 28nm | 27nm | 25nm | 22nm |
| DRAM Generations Label (bits per chip) | 4G | 8G | 8G | 16G | 32G | 32G | 32G | 32G |
| ½Vmin Production High Volume Manufacturing Begins (1000copies) | | | | 2016 | | | | |
| Vdd (High Performance, high Vdd transistors)[**] | 0.86 | 0.83 | 0.80 | 0.77 | 0.74 | 0.71 | 0.68 | 0.64 |
| Ldk^Hf 111(3)(new)[**] | 1.13 | 1.53 | 1.75 | 1.97 | 2.19 | 2.29 | 2.52 | 3.17 |
| On-chip local clock MPU/ HP (at 4% CAGR) | 5.50 | 5.95 | 6.44 | 6.96 | 7.53 | 8.14 | 8.8 | 9.9 |
| Maximum number wiring levels (unchanged) | 13 | 13 | 14 | 14 | 15 | 15 | 16 | 17 |
| MPU High-Performance (HP) Printed Gate Length (GLpr) (nm) [**] | 28 | 22 | 18 | 14 | 11 | 9 | 7 | 5 |
| MPU High-Performance Physical Gate Length (GLph) (nm) [**] | 20 | 17 | 14 | 12 | 10 | 8 | 7 | 5 |
| ASIC/Low Standby-Power (LP) Physical Gate Length (nm) (GLph)[**] | 23 | 19 | 16 | 13 | 11 | 9 | 8 | 6 |

** Note:  from the PIDS working group data; however, the calibration of Vdd, GLph, and LCV is ongoing for improved targets in 2014 ITRS work.

---

## GRAND CHALLENGES IN THE NEAR-TERM (THROUGH 2020) AND LONG-TERM (2021 AND BEYOND)

### LOGIC DEVICE SCALING [PROCESS INTEGRATION, DEVICES, AND STRUCTURES, EMERGING RESEARCH DEVICES, FRONT END PROCESSES, MODELING AND SIMULATION, AND METROLOGY]

The conventional path of scaling planar CMOS will face significant challenges set by performance and power consumption requirements.

Reduction of the equivalent gate oxide thickness (EOT) will continue to be a difficult challenge in the near term despite the introduction of high-κ metal gate (HKMG). Integration of higher-κ materials while limiting the fundamental increase in gate tunneling currents due to band-gap narrowing are also challenges to be faced. The complete gate stack material systems need to be optimized together for best device characteristics (power and performance) and cost.

New device architecture such as multiple-gate MOSFETs (e.g., finFETs) and ultra-thin body FD-SOI are expected. A particularly challenging issue is the control of the thickness, including its variability, of these ultra-thin MOSFETs. The solutions for these issues should be pursued concurrently with circuit design and system architecture improvements.

High mobility channel materials such as Ge and III-V have been considered as an enhancement or replacement for Si channel for CMOS logic applications. High-κ metal gate dielectric with low interface trap density (DIT), low bulk traps and leakage, unpinned Fermi level and low ohmic contact resistances are major challenges.

### MEMORY DEVICE SCALING [PROCESS INTEGRATION, DEVICES, AND STRUCTURES, EMERGING RESEARCH DEVICES, FRONT-END PROCESSES, MODELING AND SIMULATION, AND METROLOGY]

The challenges for DRAM devices are adequate storage capacitance with reduced feature size, high-κ dielectrics implementation, low leakage access device design, and low sheet resistance materials for bit and word lines. The drive to $4F^2$ type cell to increase bit density and to lower production cost will require high aspect ratio and non-planar FET structures.

Flash memory has become a new FEOL technology driver for critical dimension scaling, materials and processing (lithography, etching, etc.) technology, ahead of DRAM and logic. Continued Flash density improvements in the near term rely on the thickness scaling of the tunnel oxide and the intergate dielectric. To guarantee the charge retention and endurance requirements, the introduction of high-κ materials will be necessary. Cost effective implementation of 3-D NAND flash beyond 256 Gb with MLC and acceptable reliability performance remains a difficult challenge. New challenges also include the inception into mainstream manufacturing of new memory types and storage concepts such as magnetic RAM (MRAM), phase-change memory (PCM), Resistive RAM ReRAM and ferroelectric RAM (FeRAM).

6