

An aerial photograph of a city grid, overlaid with various colored rectangles in shades of green, purple, blue, and red, representing different land use or zoning categories.

## CHAPTER 4

# Circuit Characterization and Performance Estimation III

# Outline

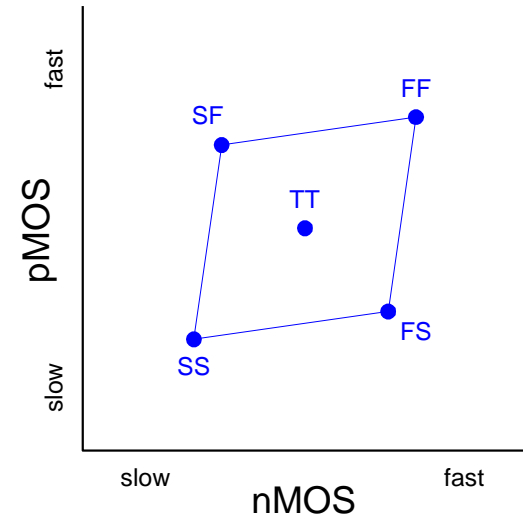
1. Delay Estimation
2. Logical Effort and Transistor Sizing
3. Power Dissipation
4. Interconnect
5. Wire Engineering
- 6. Design Margin**
7. Reliability
8. Scaling

# Design Margin

- Sources of variations
  - Supply voltage
  - Temperature
  - Process variation

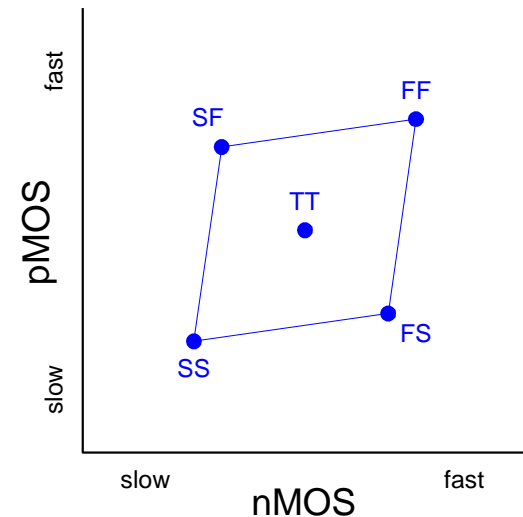
# Parameter Variation

- Transistors have uncertainty in parameters
  - Process:  $L_{\text{eff}}$ ,  $V_t$ ,  $t_{\text{ox}}$  of nMOS and pMOS
  - Vary around typical (T) values
- Fast (F)
  - $L_{\text{eff}}$ : \_\_\_\_\_
  - $V_t$ : \_\_\_\_\_
  - $t_{\text{ox}}$ : \_\_\_\_\_
- Slow (S): opposite
- Not all parameters are independent for nMOS and pMOS



# Parameter Variation

- Transistors have uncertainty in parameters
  - Process:  $L_{\text{eff}}$ ,  $V_t$ ,  $t_{\text{ox}}$  of nMOS and pMOS
  - Vary around typical (T) values
- Fast (F)
  - $L_{\text{eff}}$ : short
  - $V_t$ : low
  - $t_{\text{ox}}$ : thin
- Slow (S): opposite
- Not all parameters are independent for nMOS and pMOS



# Environmental Variation

- $V_{DD}$  and T also vary in time and space
- Fast:
  - $V_{DD}$ : \_\_\_\_\_
  - T: \_\_\_\_\_

Corner	Voltage	Temperature
F		
T	1.8	70 C
S		

# Environmental Variation

- $V_{DD}$  and T also vary in time and space
- Fast:
  - $V_{DD}$ : high
  - T: low

Corner	Voltage	Temperature
F	1.98	0 C
T	1.8	70 C
S	1.62	125 C

# Process Corners

- Process corners describe worst case variations
  - If a design works in all corners, it will probably work for any variation.
- Describe corner with four letters (T, F, S)
  - nMOS speed
  - pMOS speed
  - Voltage
  - Temperature



# Important Corners

- Some critical simulation corners include

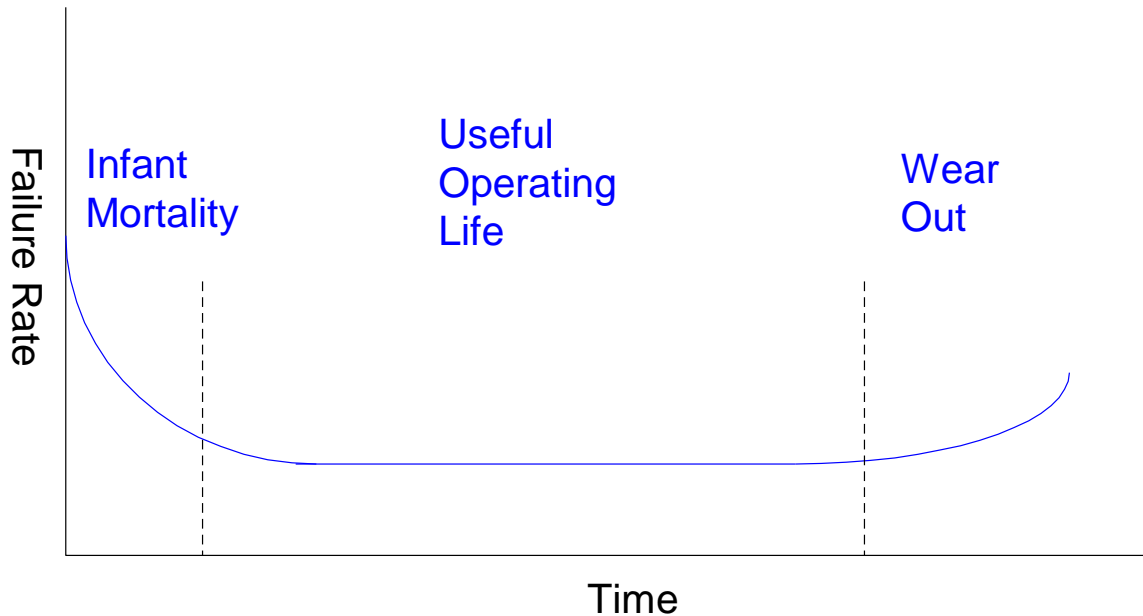
Purpose	nMOS	pMOS	V <sub>DD</sub>	Temp
Cycle time	S	S	S	S
Power	F	F	F	F
Subthreshold leakage	F	F	F	S
Pseudo-nMOS	S	F	?	?

# Outline

1. Delay Estimation
2. Logical Effort and Transistor Sizing
3. Power Dissipation
4. Interconnect
5. Wire Engineering
6. Design Margin
- 7. Reliability**
8. Scaling

# Reliability

- Hard Errors
- Soft Errors



# Electromigration

- “Electron wind” causes movement of metal atoms along wires
- Excessive electromigration leads to open circuits
- Most significant for unidirectional (DC) current
  - Depends on current density  $J_{dc}$  (current / area)
  - Exponential dependence on temperature

- Black’s Equation: 
$$MTTF \propto \frac{e^{\frac{E_a}{kT}}}{J_{dc}^n}$$

- Typical limits:  $J_{dc} < 1 - 2 \text{ mA} / \mu\text{m}^2$

- See videos

# Self-Heating

- Current through wire resistance generates heat
  - Oxide surrounding wires is a thermal insulator
  - Heat tends to build up in wires
  - Hotter wires are more resistive, slower
- Self-heating limits AC current densities for reliability

$$I_{rms} = \sqrt{\frac{\int_0^T I(t)^2 dt}{T}}$$

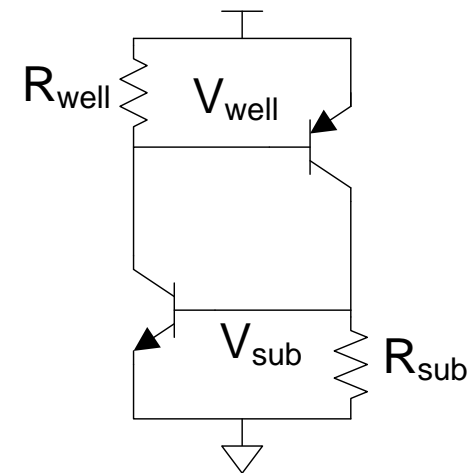
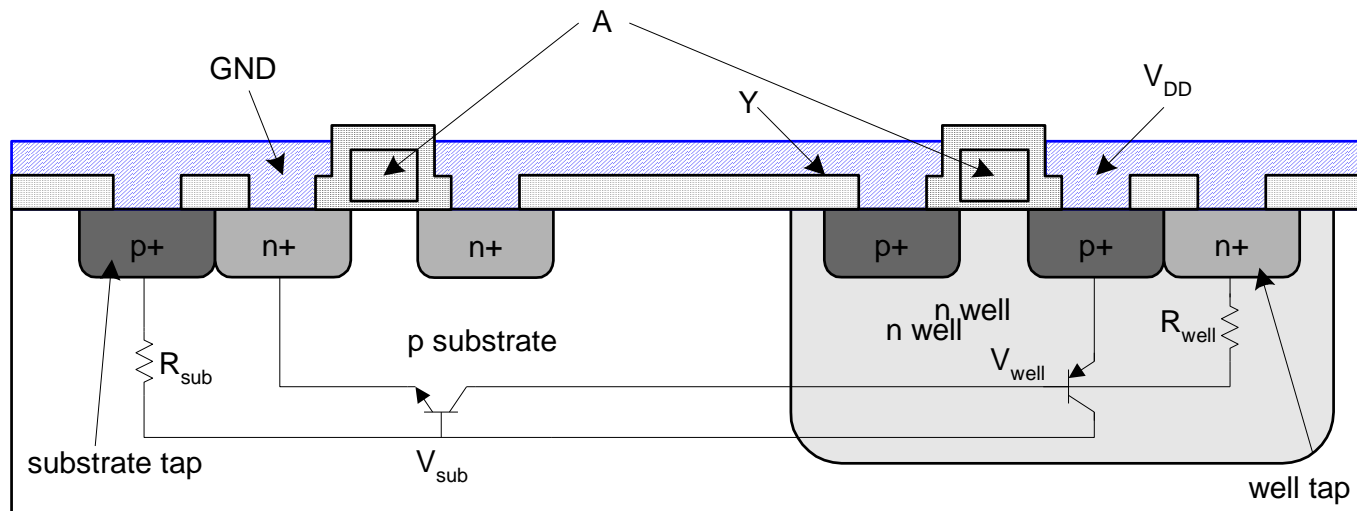
- Typical limits:  $J_{rms} < 15 \text{ mA} / \mu\text{m}^2$

# Hot Carriers

- Electric fields across channel impart high energies to some carriers
  - These “hot” carriers may be blasted into the gate oxide where they become trapped
  - Accumulation of charge in oxide causes shift in  $V_t$  over time
  - Eventually  $V_t$  shifts too far for devices to operate correctly
- Choose  $V_{DD}$  to achieve reasonable product lifetime
  - Worst problems for inverters and NORs with slow input risetime and long propagation delays

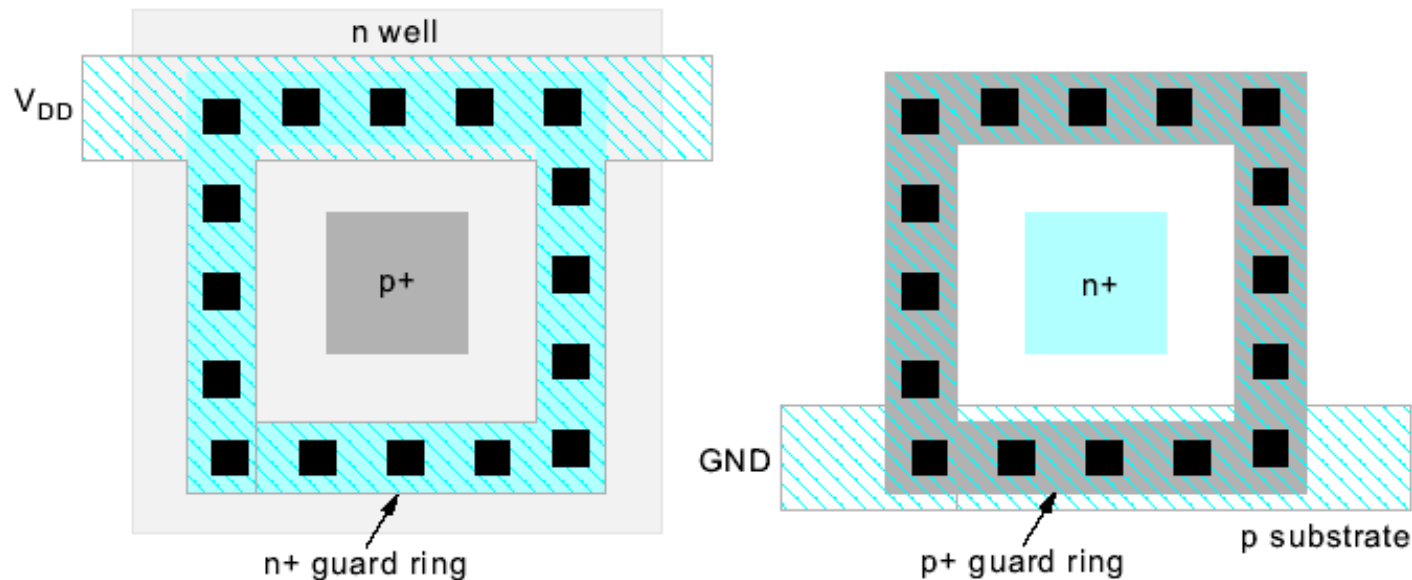
# Latchup

- Latchup: positive feedback leading to  $V_{DD}$  – GND short
  - Major problem for 1970's CMOS processes before it was well understood
- Avoid by minimizing resistance of body to GND /  $V_{DD}$ 
  - Use plenty of substrate and well taps



# Guard Rings

- Latchup risk greatest when diffusion-to-substrate diodes could become forward-biased
- Surround sensitive region with guard ring to collect injected charge





# Overvoltage

- High voltages can damage transistors
  - Electrostatic discharge
  - Oxide arcing
  - Punchthrough
  - Time-dependent dielectric breakdown (TDDB)
    - Accumulated wear from tunneling currents
- Requires low  $V_{DD}$  for thin oxides and short channels
- Use ESD protection structures where chip meets real world

# Summary

- Static CMOS gates are very robust
  - Will settle to correct value if you wait long enough
- Other circuits suffer from a variety of pitfalls
  - Tradeoff between performance & robustness
- Very important to check circuits for pitfalls
  - For large chips, you need an automatic checker.
  - Design rules aren't worth the paper they are printed on unless you back them up with a tool.

# Soft Errors

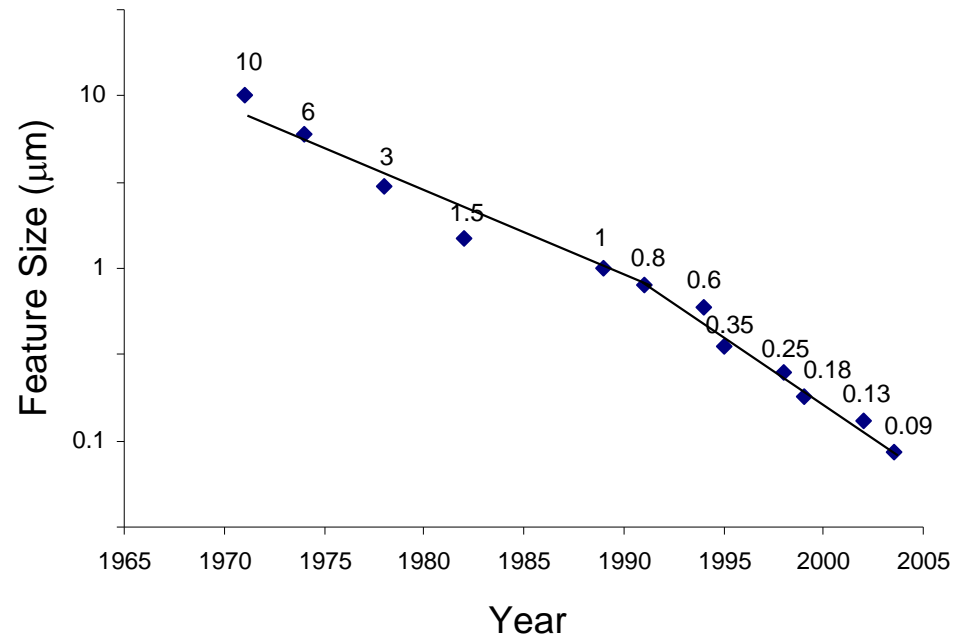
- In 1970's, DRAMs were observed to occasionally flip bits for no apparent reason
  - Ultimately linked to alpha particles and cosmic rays
- Collisions with particles create electron-hole pairs in substrate
  - These carriers are collected on dynamic nodes, disturbing the voltage
- Minimize soft errors by having plenty of charge on dynamic nodes
- Tolerate errors through ECC, redundancy

# Outline

1. Delay Estimation
2. Logical Effort and Transistor Sizing
3. Power Dissipation
4. Interconnect
5. Wire Engineering
6. Design Margin
7. Reliability
- 8. Scaling**

# Scaling

- The only constant in VLSI is constant change
- Feature size shrinks by 30% every 2-3 years
  - Transistors become cheaper
  - Transistors become faster
  - Wires do not improve (and may get worse)
- Scale factor  $S$ 
  - Typically  $S = \sqrt{2}$
  - Technology nodes



# Scaling Assumptions

- What changes between technology nodes?
- Constant Field Scaling
  - All dimensions ( $x, y, z \Rightarrow W, L, t_{ox}$ )
  - Voltage ( $V_{DD}$ )
  - Doping levels
- Lateral Scaling
  - Only gate length  $L$
  - Often done as a quick gate shrink ( $S = 1.05$ )

# Device Scaling

**Table 4.15** Influence of scaling on MOS device characteristics

Parameter	Sensitivity	Constant Field	Lateral
<b>Scaling Parameters</b>			
Length: $L$		$1/S$	$1/S$
Width: $W$		$1/S$	1
Gate oxide thickness: $t_{ox}$		$1/S$	1
Supply voltage: $V_{DD}$		$1/S$	1
Threshold voltage: $V_{tn}, V_{tp}$		$1/S$	1
Substrate doping: $N_A$		$S$	1
<b>Device Characteristics</b>			
$\beta$	$\frac{W}{L} \frac{1}{t_{ox}}$	$S$	$S$
Current: $I_{ds}$	$\beta(V_{DD} - V_t)^2$	$1/S$	$S$
Resistance: $R$	$\frac{V_{DD}}{I_{ds}}$	1	$1/S$
Gate capacitance: $C$	$\frac{WL}{t_{ox}}$	$1/S$	$1/S$
Gate delay: $\tau$	$RC$	$1/S$	$1/S^2$
Clock frequency: $f$	$1/\tau$	$S$	$S^2$
Dynamic power dissipation (per gate): $P$	$CV^2f$	$1/S^2$	$S$
Chip area: $A$		$1/S^2$	1
Power density	$P/A$	1	$S$
Current density	$I_{ds}/A$	$S$	$S$

# Observations

- Gate capacitance per micron is nearly independent of process
- But ON resistance \* micron improves with process
- Gates get faster with scaling (good)
- Dynamic power goes down with scaling (good)
- Current density goes up with scaling (bad)
- Velocity saturation makes lateral scaling unsustainable



# Example

- Gate capacitance is typically about  $2 \text{ fF}/\mu\text{m}$
- The FO4 inverter delay in the TT corner for a process of feature size  $f$  (in nm) is about  $0.5f \text{ ps}$
- Estimate the ON resistance of a unit ( $4 \lambda$  wide) nMOS transistor.
- $\text{FO4} = 5 \tau = 15 \text{ RC}$
- $\text{RC} = (0.5f) / 15 = (f/30) \text{ ps/nm}$
- If  $W = 2f$ ,  $C = 4f * (1e-3) \text{ fF}$ ,  $R = 8.33 \text{ k}\Omega$ 
  - Unit resistance is roughly independent of  $f$

# Scaling Assumptions

- Wire thickness
  - Hold constant vs. reduce in thickness
- Wire length
  - Local / scaled interconnect
  - Global interconnect
    - Die size scaled by  $D_c \approx 1.1$

# Interconnect Scaling

**Table 4.16** Influence of scaling on interconnect characteristics

Parameter	Sensitivity	Reduced Thickness	Constant Thickness
<b>Scaling Parameters</b>			
Width: $w$			$1/S$
Spacing: $s$			$1/S$
Thickness: $t$		$1/S$	1
Interlayer oxide height: $h$			$1/S$
<b>Characteristics Per Unit Length</b>			
Wire resistance per unit length: $R_w$	$\frac{1}{wt}$	$S^2$	$S$
Fringing capacitance per unit length: $C_{wf}$	$\frac{t}{s}$	1	$S$
Parallel plate capacitance per unit length: $C_{wp}$	$\frac{w}{h}$	1	1
Total wire capacitance per unit length: $C_w$	$C_{wf} + C_{wp}$	1	between 1, $S$
Unrepeated RC constant per unit length: $t_{wu}$	$R_w C_w$	$S^2$	between $S$ , $S^2$
Repeated wire RC delay per unit length: $t_{wr}$ (assuming constant field scaling of gates in Table 4.15)	$\sqrt{RCR_w C_w}$	$\sqrt{S}$	between 1, $\sqrt{S}$
Crosstalk noise	$\frac{t}{s}$	1	$S$

# Interconnect Delay

<b>Table 4.16</b> Influence of scaling on interconnect characteristics			
Parameter	Sensitivity	Reduced Thickness	Constant Thickness
<b>Scaling Parameters</b>			
Width: $w$			$1/S$
Spacing: $s$			$1/S$
Thickness: $t$		$1/S$	1
Interlayer oxide height: $h$			$1/S$
<b>Local/Scaled Interconnect Characteristics</b>			
Length: $l$			$1/S$
Unrepeated wire RC delay	$l^2 t_{wu}$	1	between $1/S, 1$
Repeated wire delay	$lt_{wr}$	$\sqrt{1/S}$	between $1/S, \sqrt{1/S}$
<b>Global Interconnect Characteristics</b>			
Length: $l$			$D_c$
Unrepeated wire RC delay	$l^2 t_{wu}$	$S^2 D_c^2$	between $SD_c^2, S^2 D_c^2$
Repeated wire delay	$lt_{wr}$	$D_c \sqrt{S}$	between $D_c, D_c \sqrt{S}$

# Observations

- Capacitance per micron is remaining constant
  - About 0.2 fF/ $\mu\text{m}$
  - Roughly 1/10 of gate capacitance
- Local wires are getting faster
  - Not quite tracking transistor improvement
  - But not a major problem
- Global wires are getting slower
  - No longer possible to cross chip in one cycle

# ITRS

- Semiconductor Industry Association forecast
  - Intl. Technology Roadmap for Semiconductors

**Table 4.17** Predictions from the 2002 ITRS

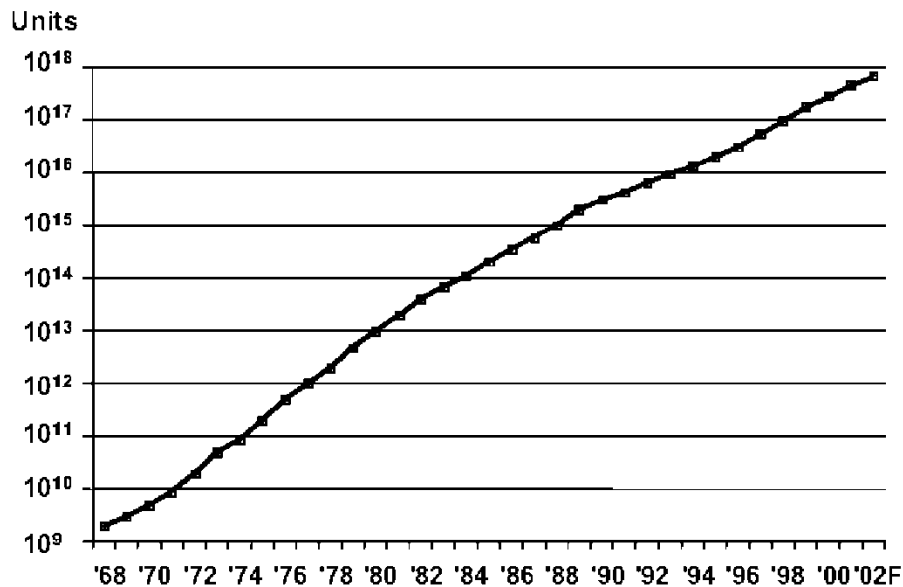
Year	2001	2004	2007	2010	2013	2016
Feature size (nm)	130	90	65	45	32	22
$V_{DD}$ (V)	1.1–1.2	1–1.2	0.7–1.1	0.6–1.0	0.5–0.9	0.4–0.9
Millions of transistors/die	193	385	773	1564	3092	6184
Wiring levels	8–10	9–13	10–14	10–14	11–15	11–15
Intermediate wire pitch (nm)	450	275	195	135	95	65
Interconnect dielectric constant	3–3.6	2.6–3.1	2.3–2.7	2.1	1.9	1.8
I/O signals	1024	1024	1024	1280	1408	1472
Clock rate (MHz)	1684	3990	6739	11511	19348	28751
FO4 delays/cycle	13.7	8.4	6.8	5.8	4.8	4.7
Maximum power (W)	130	160	190	218	251	288
DRAM capacity (Gbits)	0.5	1	4	8	32	64

# Scaling Implications

- Improved Performance
- Improved Cost
- Interconnect Woes
- Power Woes
- Productivity Challenges
- Physical Limits

# Cost Improvement

- In 2003, \$0.01 bought you 100,000 transistors
  - Moore's Law is still going strong



Source: Dataquest/Intel



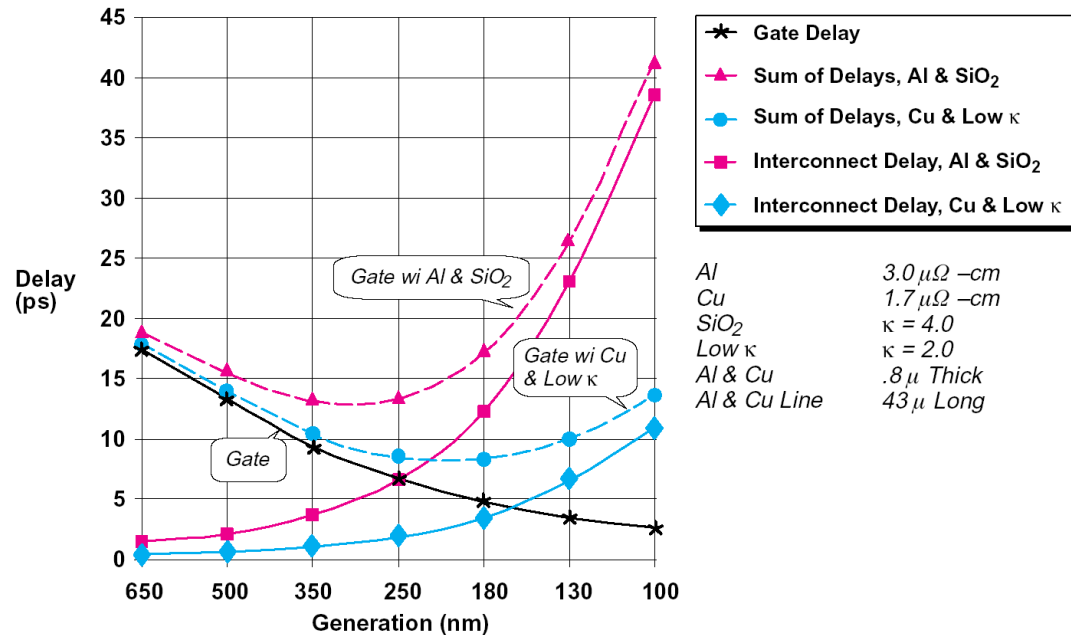
Source: Dataquest/Intel

[Moore03]



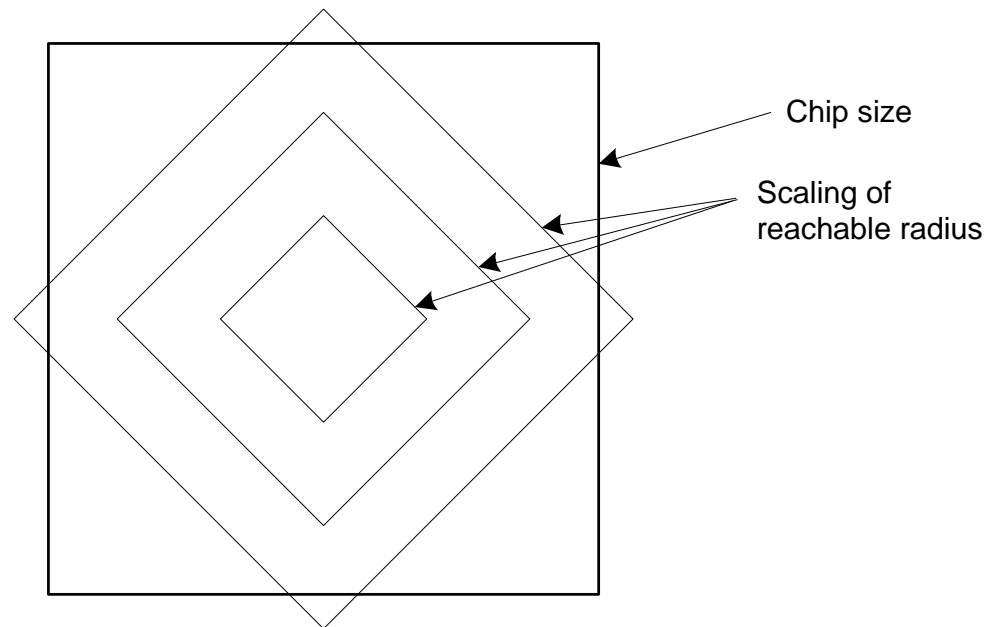
# Interconnect Woes

- SIA made a gloomy forecast in 1997
  - Delay would reach minimum at 250 – 180 nm, then get worse because of wires
- But...
  - Misleading scale
  - Global wires
- 100 kgate blocks ok



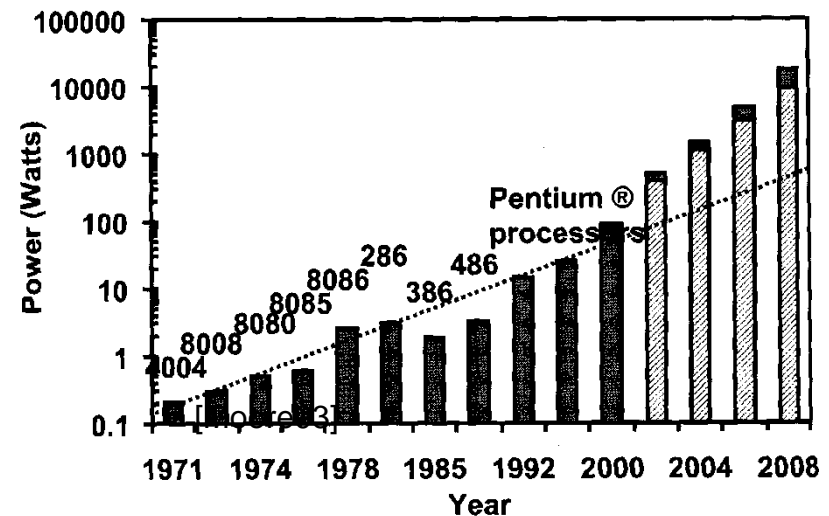
# Reachable Radius

- We can't send a signal across a large fast chip in one cycle anymore
- But the microarchitect can plan around this
  - Just as off-chip memory latencies were tolerated



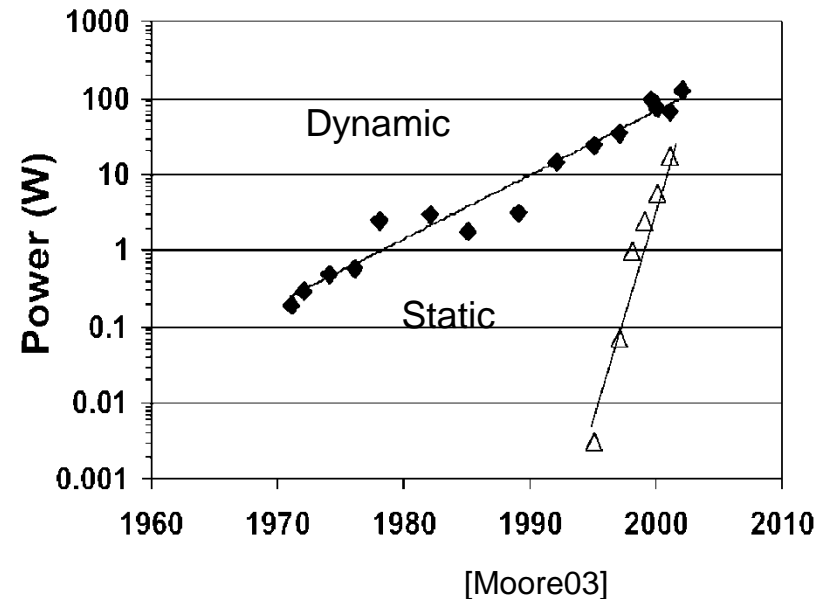
# Dynamic Power

- Intel VP Patrick Gelsinger (ISSCC 2001)
  - If scaling continues at present pace, by 2005, high speed processors would have power density of nuclear reactor, by 2010, a rocket nozzle, and by 2015, surface of sun.
  - “Business as usual will not work in the future.”
- Intel stock dropped 8% on the next day
- But attention to power is increasing



# Static Power

- $V_{DD}$  decreases
  - Save dynamic power
  - Protect thin gate oxides and short channels
  - No point in high value because of velocity sat.
- $V_t$  must decrease to keep device performance
- But this causes exponential increase in OFF leakage
- Major future challenge



# Productivity

- Transistor count is increasing faster than designer productivity (gates / week)
  - Bigger design teams
    - Up to 500 for a high-end microprocessor
  - More expensive design cost
  - Pressure to raise productivity
    - Rely on synthesis, IP blocks
  - Need for good engineering managers

# Physical Limits

- Will Moore's Law run out of steam?
  - Can't build transistors smaller than an atom...
- Many reasons have been predicted for end of scaling
  - Dynamic power
  - Subthreshold leakage, tunneling
  - Short channel effects
  - Fabrication costs
  - Electromigration
  - Interconnect delay
- Rumors of demise have been exaggerated