

# Scaling & Economics

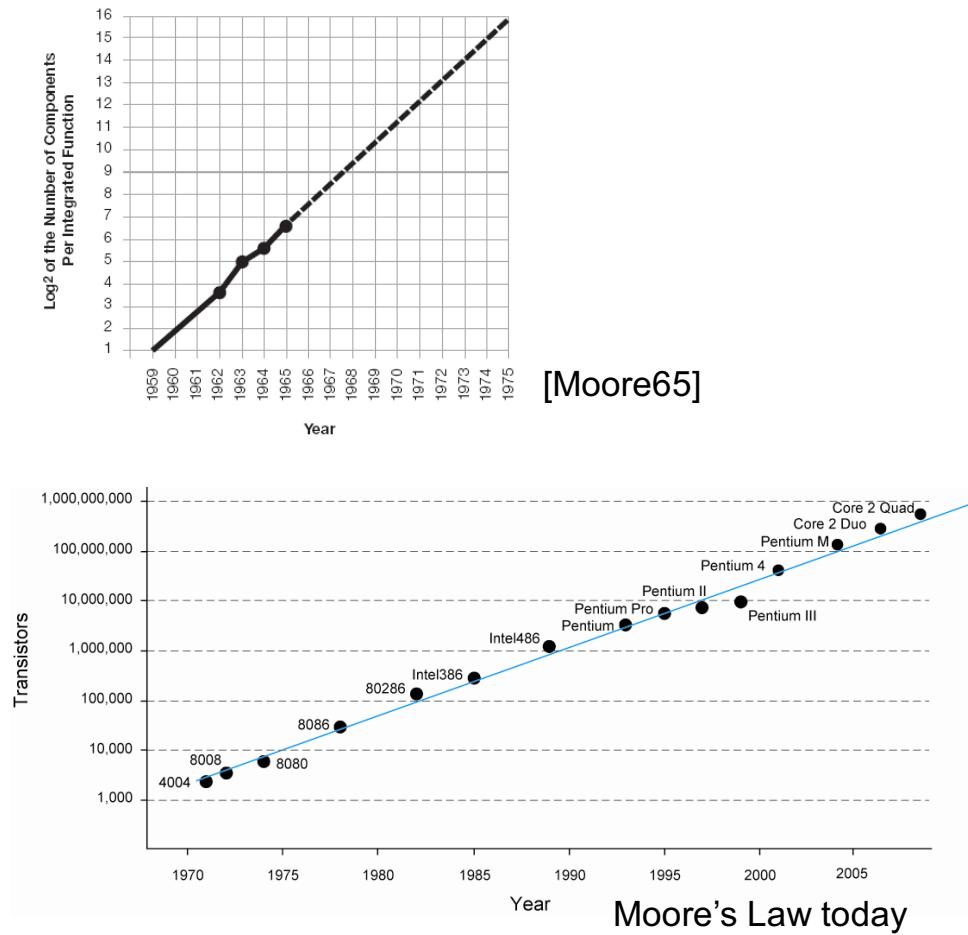
# Outline

---

- ❑ Scaling
    - Transistors
    - Interconnect
    - Future Challenges
  - ❑ Economics
-

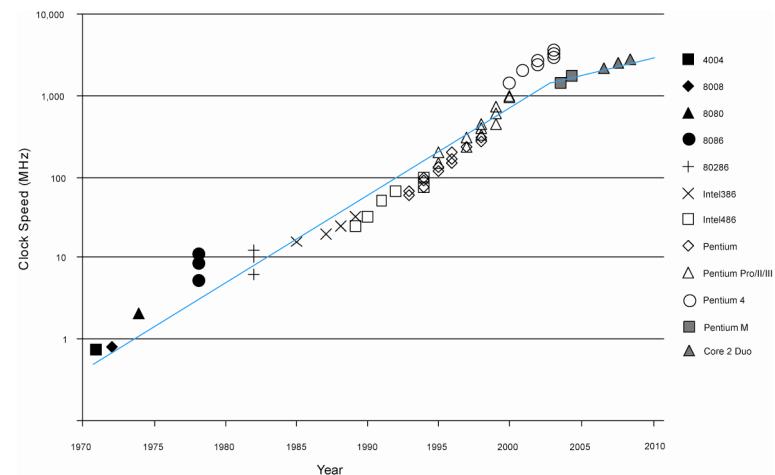
# Moore's Law

□ Recall that Moore's Law has been driving CMOS



[Moore65]

Moore's Law today



Corollary: clock speeds have improved

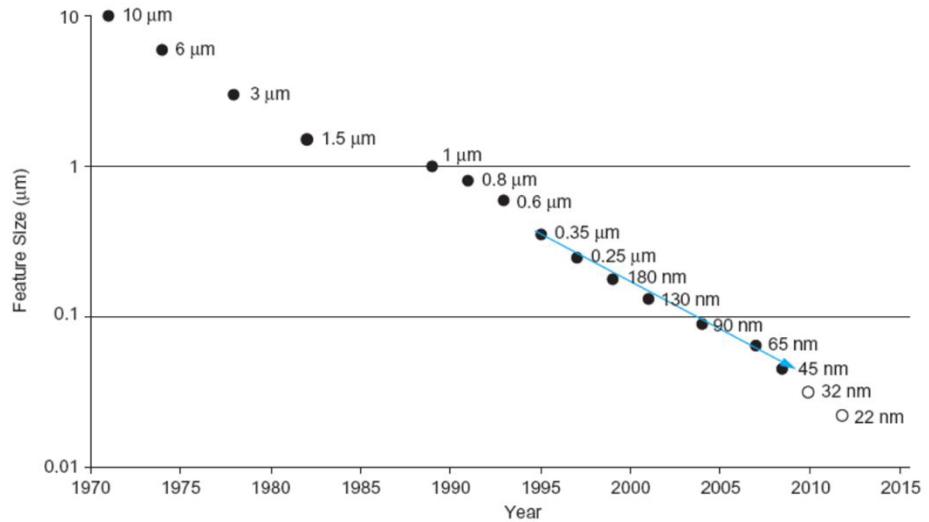
# Why?

---

- Why more transistors per IC?
    - Smaller transistors
    - Larger dies
  - Why faster computers?
    - Smaller, faster transistors
    - Better microarchitecture (more IPC)
    - Fewer gate delays per cycle
-

# Scaling

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



# Dennard Scaling

---

- Proposed by Dennard in 1974
- Also known as *constant field scaling*
  - Electric fields remain the same as features scale
- Scaling assumptions
  - All dimensions ( $x, y, z \Rightarrow W, L, t_{ox}$ )
  - Voltage ( $V_{DD}$ )
  - Doping levels

# Device Scaling

Parameter	Sensitivity	Dennard Scaling
L: Length		1/S
W: Width		1/S
$t_{ox}$ : gate oxide thickness		1/S
$V_{DD}$ : supply voltage		1/S
$V_t$ : threshold voltage		1/S
NA: substrate doping		S
$\beta$	$W/(Lt_{ox})$	S
$I_{on}$ : ON current	$\beta(V_{DD}-V_t)^2$	1/S
R: effective resistance	$V_{DD}/I_{on}$	1
C: gate capacitance	$WL/t_{ox}$	1/S
$\tau$ : gate delay	RC	1/S
f: clock frequency	$1/\tau$	S
E: switching energy / gate	$CV_{DD}^2$	1/S <sup>3</sup>
P: switching power / gate	Ef	1/S <sup>2</sup>
A: area per gate	WL	1/S <sup>2</sup>
Switching power density	P/A	1
Switching current density	$I_{on}/A$	S

# Observations

---

- Gate capacitance improves with process (good)
- Gates get faster with scaling (good)
- Dynamic power goes down with scaling (good)
- Current density goes up with scaling (bad)

# Real Scaling

---

- $t_{ox}$  scaling has slowed since 65 nm
    - Limited by gate tunneling current
    - Gates are only about 4 atomic layers thick!
    - High-k dielectrics have helped continue scaling of effective oxide thickness
  - $V_{DD}$  scaling has slowed since 65 nm
    - SRAM cell stability at low voltage is challenging
  - Dennard scaling predicts cost, speed, power all improve
    - Below 65 nm, some designers find they must choose just two of the three
-

# Wire Scaling

---

- Wire cross-section
  - w, s, t all scale
- Wire length
  - Local interconnects: run within functional units => scale with feature size
  - Semiglobal (scaled) interconnects: run among functional units => scale with feature size
  - Global interconnects: run across entire die => may get longer because:
    - Die size scales by  $D_c \approx 1.1$

# Interconnect Scaling

Parameter	Sensitivity	Scale Factor
w: width		1/S
s: spacing		1/S
t: thickness		1/S
h: height		1/S
D <sub>c</sub> : die size		D <sub>c</sub>
R <sub>w</sub> : wire resistance/unit length	1/wt	S <sup>2</sup>
C <sub>wf</sub> : fringing capacitance / unit length	t/s	1
C <sub>wp</sub> : parallel plate capacitance / unit length	w/h	1
C <sub>w</sub> : total wire capacitance / unit length	C <sub>wf</sub> + C <sub>wp</sub>	1
t <sub>wu</sub> : unrepeated RC delay / unit length	R <sub>w</sub> C <sub>w</sub>	S <sup>2</sup>
t <sub>wr</sub> : repeated RC delay / unit length	sqrt(RCR <sub>w</sub> C <sub>w</sub> )	sqrt(S)
Crosstalk noise	w/h	1
E <sub>w</sub> : energy per bit / unit length	C <sub>w</sub> V <sub>DD</sub> <sup>2</sup>	1/S <sup>2</sup>

# Observations

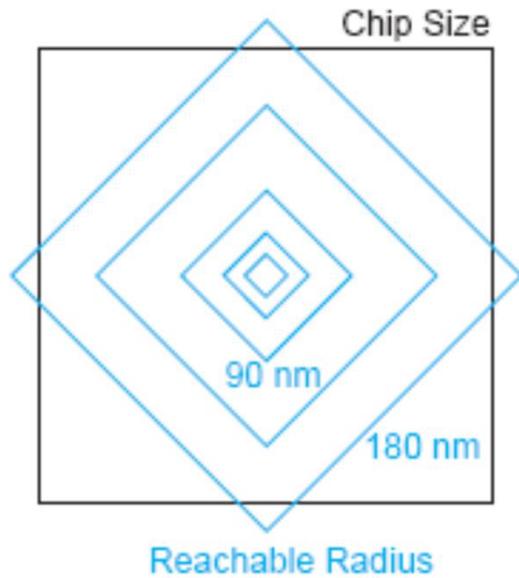
---

- Capacitance per unit length is remaining constant
  - About  $0.2 \text{ fF}/\mu\text{m}$
  - Roughly 1/5 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

# Reachable Radius

---

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



# 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

# VLSI Economics

---

- ❑ Selling price  $S_{\text{total}}$ 
  - $S_{\text{total}} = C_{\text{total}} / (1-m)$
- ❑  $m$  = profit margin
- ❑  $C_{\text{total}}$  = total cost
  - Nonrecurring engineering cost (NRE)
  - Recurring cost
  - Fixed cost

# NRE

---

- ❑ Engineering cost
  - Depends on size of design team
  - Include benefits, training, computers
  - CAD tools:
    - Digital front end: \$10K
    - Analog front end: \$100K
    - Digital back end: \$1M
- ❑ Prototype manufacturing
  - Mask costs: \$5M in 45 nm process
  - Test fixture and package tooling

# Fixed Costs

---

- Data sheets and application notes
- Marketing and advertising
- Yield analysis

# Example

---

- You want to start a company to build a wireless communications chip. How much venture capital must you raise?
  
  - Because you are smarter than everyone else, you can get away with a small team in just two years:
    - Seven digital designers
    - Three analog designers
    - Five support personnel
-

# Solution

- Digital designers:
  - \$70k salary
  - \$30k overhead
  - \$10k computer
  - \$10k CAD tools
  - Total:  $\$120k * 7 = \$840k$
- Analog designers
  - \$100k salary
  - \$30k overhead
  - \$10k computer
  - \$100k CAD tools
  - Total:  $\$240k * 3 = \$720k$
- Support staff
  - \$45k salary
  - \$20k overhead
  - \$5k computer
  - Total:  $\$70k * 5 = \$350k$
- Fabrication
  - Back-end tools: \$1M
  - Masks: \$5M
  - Total: \$6M / year
- Summary
  - 2 years @ \$7.91M / year
  - \$16M design & prototype