

CSE 560
Computer Systems Architecture

Technology

Technology Unit Overview

- **Technology basis**

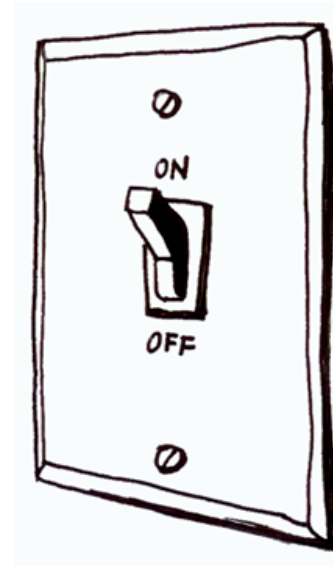
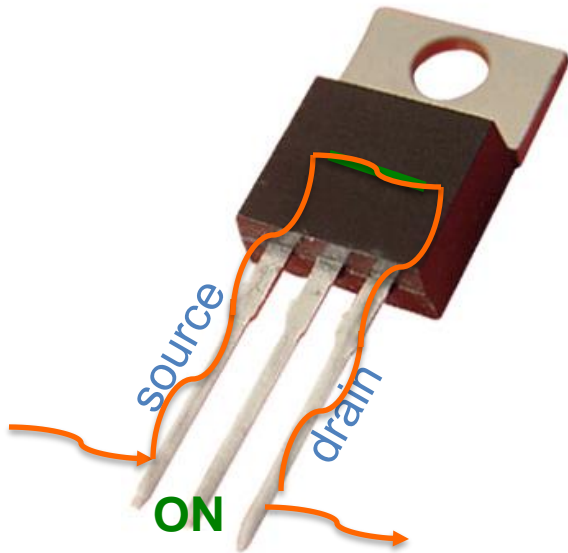
- Transistors
- Transistor scaling (Moore's Law)

- **The metrics**

- Cost
- Transistor speed
- Power
- Reliability

How do the metrics change with transistor scaling?

How do these changes affect the job of a computer architect?



The Transistor

Technology Generations

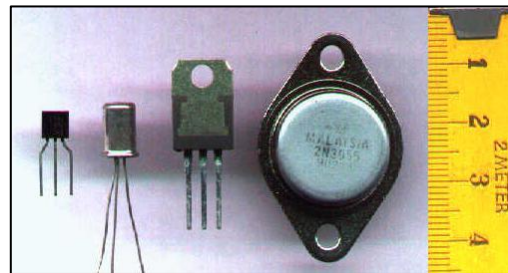
1950-1959 Vacuum Tubes

1960-1968 Transistors

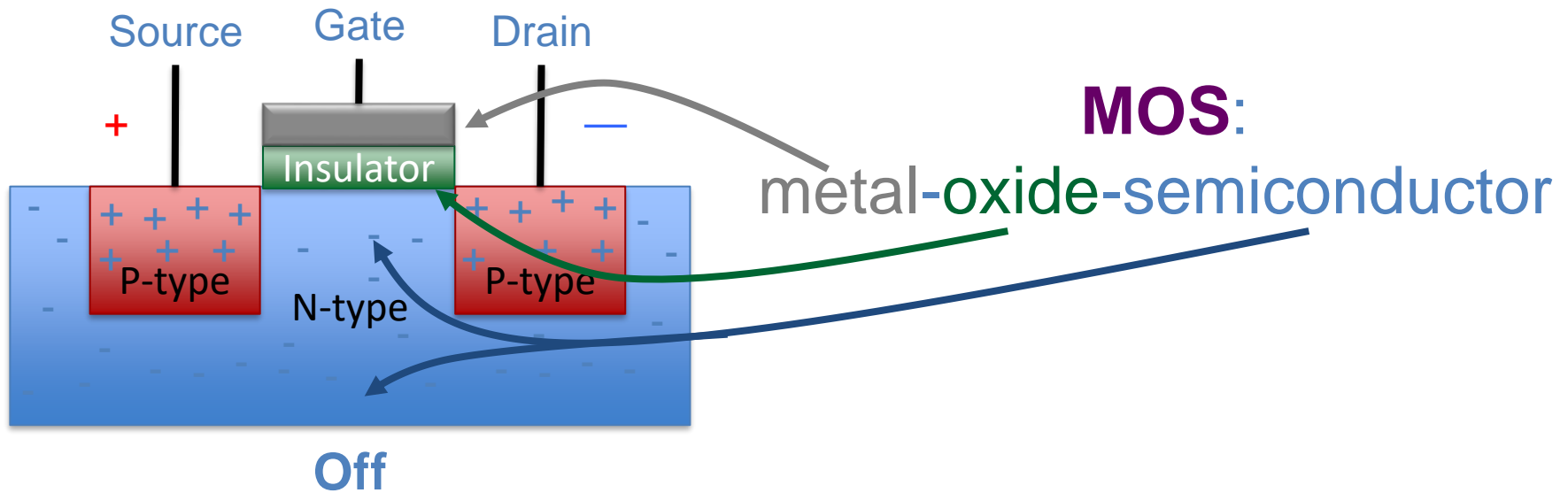
1969-1977 Integrated Circuit (multiple transistors on chip)

1978-1999 LSI & VLSI (10Ks & 100Ks transistors on chip)

2000-20xx VLSI (millions, now billions transistors on chip)



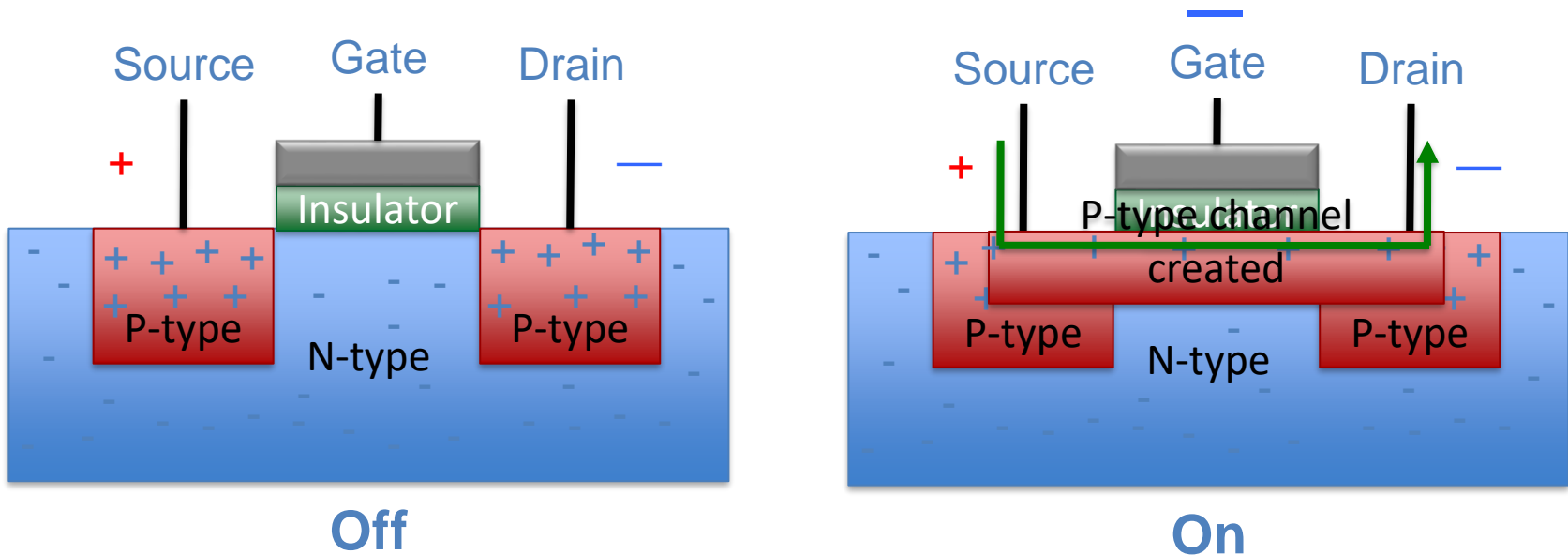
The **Silicon** in Silicon Valley



N-Type Silicon: negative free-carriers (free electrons)

P-Type Silicon: positive free-carriers (holes)

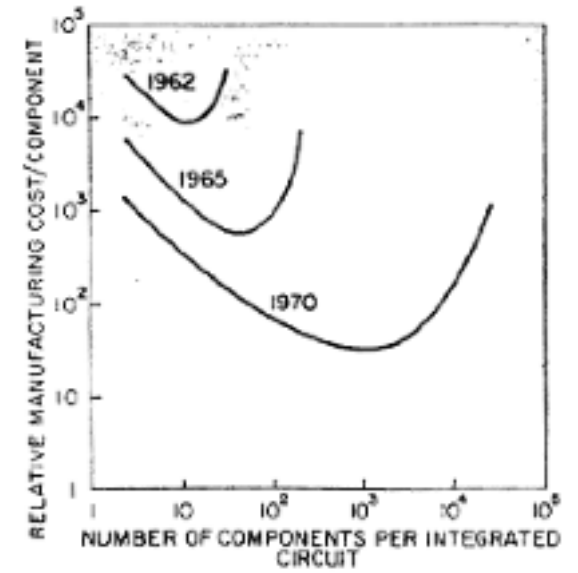
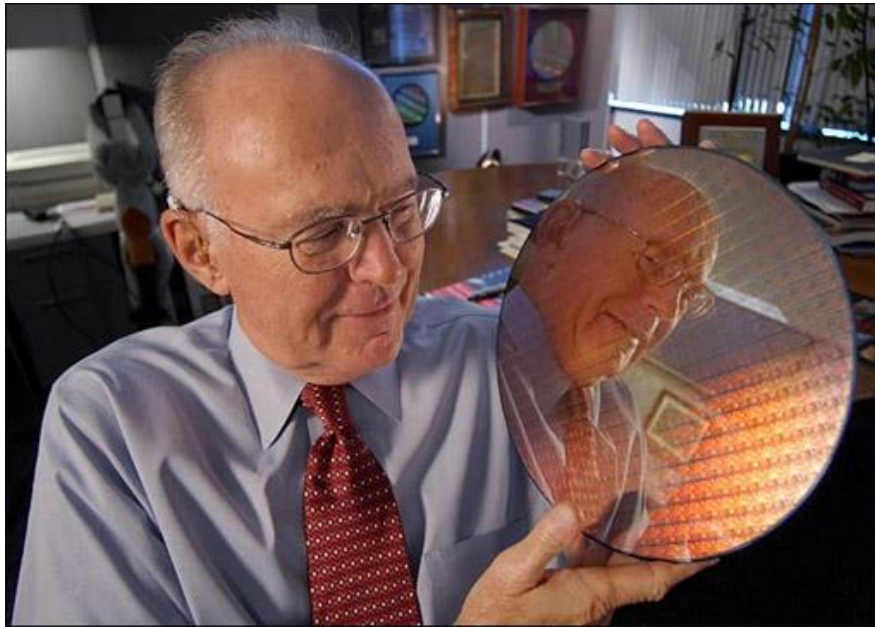
CMOS: Semiconductor Technology



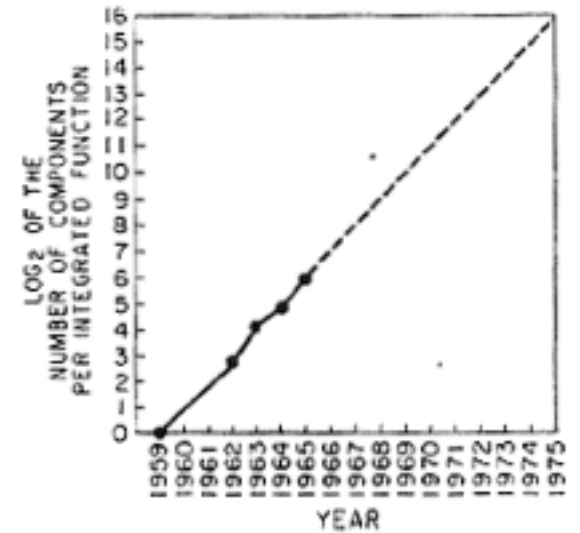
P-Transistor: negative charge on gate closes channel, connecting source & drain

(**N-Transistor** works the opposite way)

Complementary MOS (CMOS) Technology: uses p & n transistors



Transistor Scaling



Enter Gordon Moore

The experts look ahead

Cramming more components onto integrated circuits

With unit cost falling as the number of components per circuit rises, by 1975 economics may dictate squeezing as many as 65,000 components on a single silicon chip

By Gordon E. Moore

Director, Research and Development Laboratories, Fairchild Semiconductor division of Fairchild Camera and Instrument Corp.

The future of integrated electronics is the future of electronics itself. The advantages of integration will bring about a proliferation of electronics, pushing this science into many new areas.

Integrated circuits will lead to such wonders as home computers—or at least terminals connected to a central computer—automatic controls for automobiles, and personal portable communications equipment. The electronic wrist-watch needs only a display to be feasible today.

But the biggest potential lies in the production of large systems. In telephone communications, integrated circuits in digital filters will separate channels on multiplex equipment. Integrated circuits will also switch telephone circuits and perform data processing.

Computers will be more powerful, and will be organized in completely different ways. For example, memories built of integrated electronics may be distributed throughout the

machine instead of being concentrated in a central unit. In addition, the improved reliability made possible by integrated circuits will allow the construction of larger processing units. Machines similar to those in existence today will be built at lower costs and with faster turn-around.

Present and future

By integrated electronics, I mean all the various technologies which are referred to as microelectronics today as well as any additional ones that result in electronics functions supplied to the user as irreducible units. These technologies were first investigated in the late 1950's. The object was to miniaturize electronics equipment to include increasingly complex electronic functions in limited space with minimum weight. Several approaches evolved, including microassembly techniques for individual components, thin-film structures and semiconductor integrated circuits.

Each approach evolved rapidly and converged so that each borrowed techniques from another. Many researchers believe the way of the future to be a combination of the various approaches.

The advocates of semiconductor integrated circuitry are already using the improved characteristics of thin-film resistors by applying such films directly to an active semiconductor substrate. Those advocating a technology based upon films are developing sophisticated techniques for the attachment of active semiconductor devices to the passive film arrays.

Both approaches have worked well and are being used in equipment today.

The complexity for minimum component costs has increased at a rate of roughly a factor of two per year.... Certainly over the short term this rate can be expected to continue, if not to increase. Over the longer term, the rate of increase is a bit more uncertain, although there is no reason to believe it will not remain nearly constant for at least 10 years. That means by 1975, the number of components per integrated circuit for minimum cost will be 65,000. I believe that such a large circuit can be built on a single wafer.

(From the original 1965 Moore's Law paper)

"The number of transistors will double every year", 1965

("...or every two years", 1975)

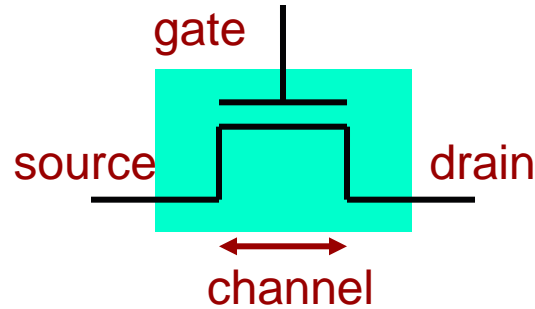
The author



Dr. Gordon E. Moore is one of the new breed of electronic engineers, schooled in the physical sciences rather than in electronics. He earned a B.S. degree in chemistry from the University of California and a Ph.D. degree in physical chemistry from the California Institute of Technology. He was one of the founders of Fairchild Semiconductor and has been director of the research and development laboratories since 1959.

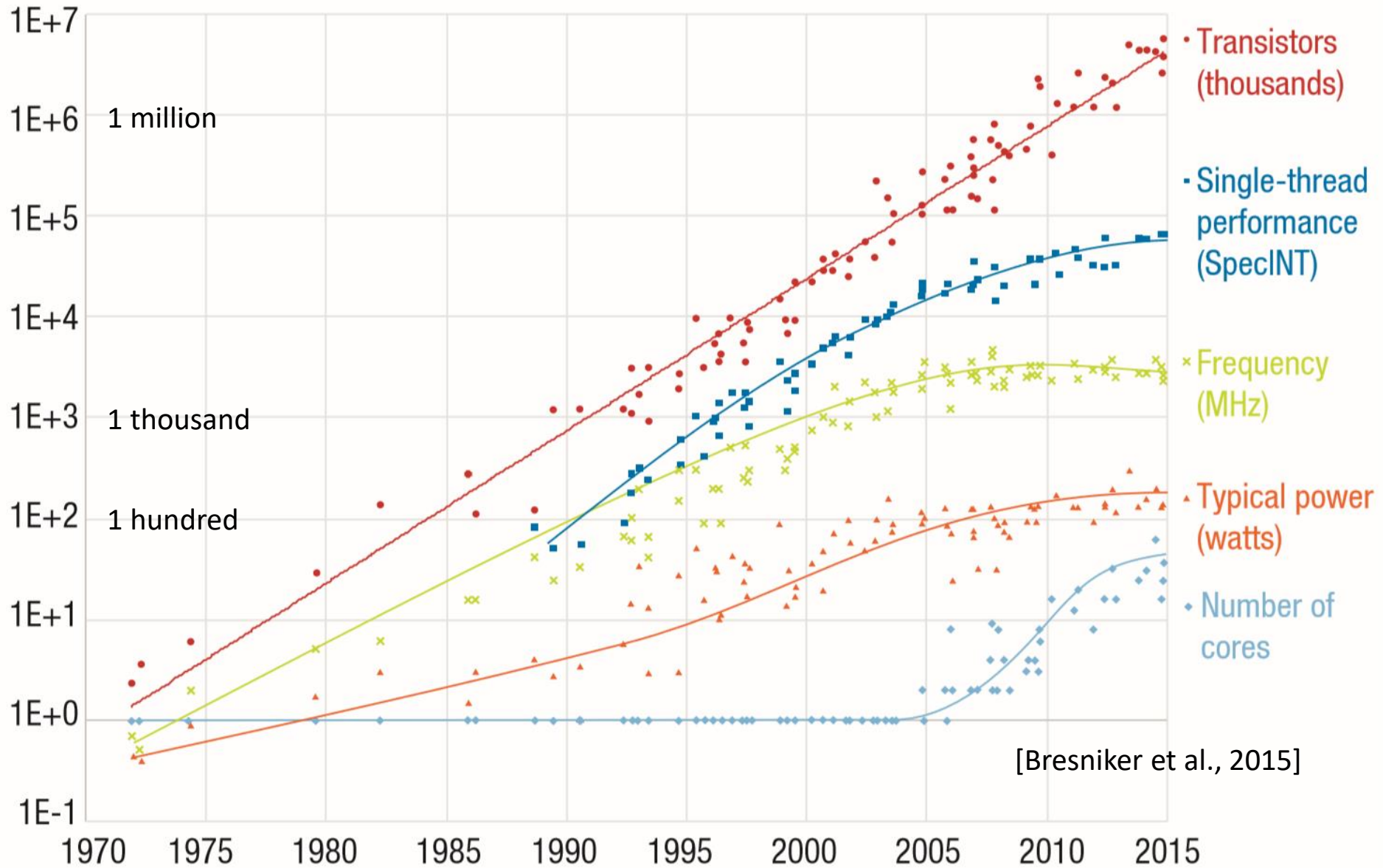
Electronics, Volume 38, Number 8, April 19, 1965

Moore's Law: Technology Scaling



- **Channel length:** characteristic parameter (short → fast)
 - Aka “feature size” or “technology”
 - Currently: 0.003 micron (μm), 3 nanometers (nm)
- **Moore's Law:** aka “technology scaling”
 - Continued miniaturization (\approx channel length)
 - + **Improves:** switching **speed**, **power**/transistor, area(**cost**)/transistor
 - **Reduces:** transistor **reliability**

Technology Trends



Cost



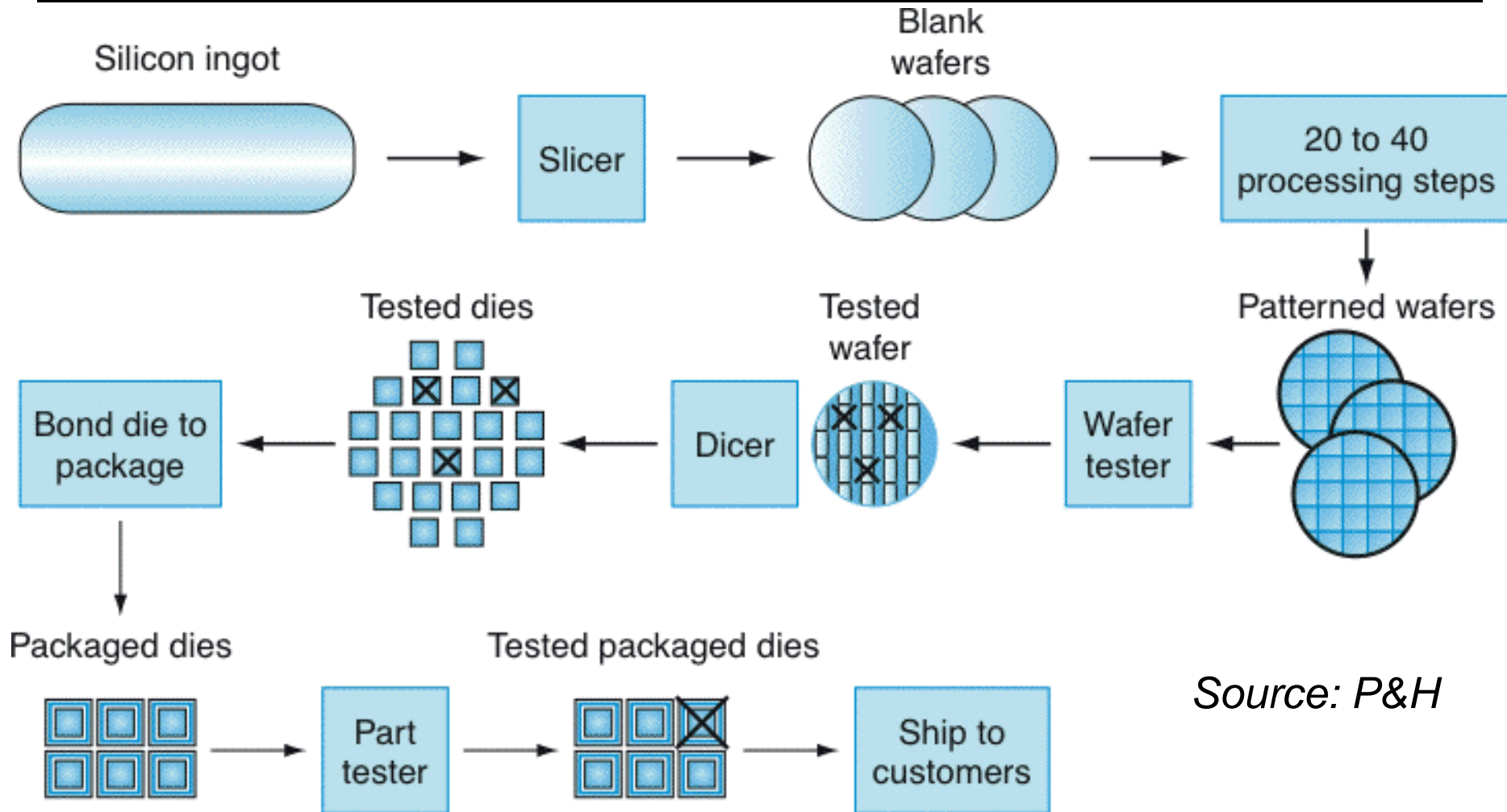
Cost

- Metric: \$
- CPU = fraction of cost, so is profit (Intel's, Dell's)

	Desktop	Laptop	Netbook	Phone
\$	\$100–\$300	\$150–\$350	\$50–\$100	\$10–\$20
% of total	10–30%	10–20%	20–30%	20-30%
Other costs	Memory, display, power supply/battery, storage, software			

- We are concerned about *chip cost*
 - **Unit cost:** costs to manufacture individual chips
 - **Startup cost:** cost to design chip, build the manufacturing facility

Unit Costs in Manufacturing Process

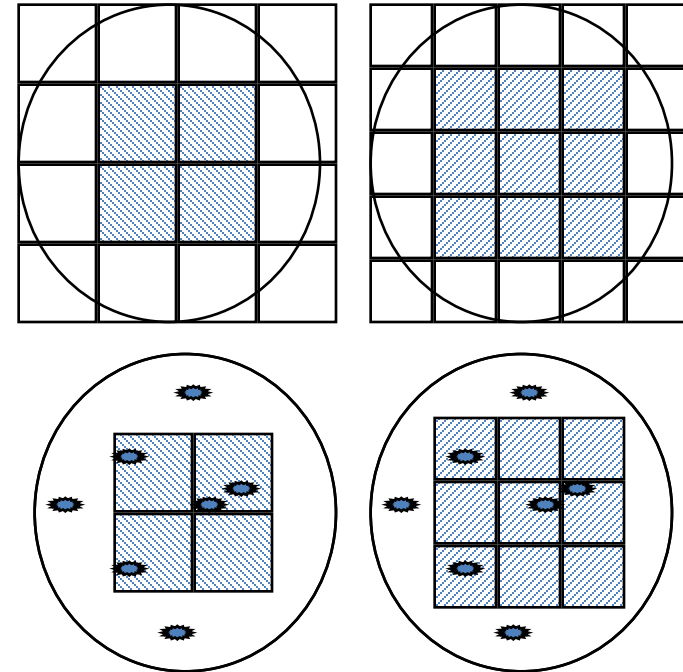


Source: P&H

Lots of unit testing, packaging, re-testing

Unit Cost: Integrated Circuit (IC)

- Cost / wafer is constant, $f(\text{wafer size, number of steps})$
- Chip (die) cost related to **area**
 - Larger chips \rightarrow fewer chips/wafer
 \rightarrow fewer *working* ones
 - Chip cost \sim chip area $^\alpha$
 - $\alpha = 2$ to 3
 - Why? random defects
- **Wafer yield:** % wafer that is chips
- **Die yield:** % chips that work
 - Yield is increasingly non-binary, fast vs. slow chips



Fixed Costs

- **For new chip design**
 - Design & verification: ~\$100M (500 person-years @ \$200K per)
 - Amortized over “proliferations”, *e.g.*, Xeon/Celeron cache variants
- **For new (smaller) technology generation**
 - ~\$3B for a new fab
 - Amortized over multiple designs
 - Amortized by “rent” from companies w/o their own fabs
- **Intel’s tick-tock** (smaller → better)



Transistor Speed

Moore's Speed Effect #1: Transistor Speed

Transistor length: "process generation"

45nm = transistor gate length

Shrink transistor length:

- + ↓ resistance of channel (shorter)
- + ↓ gate/source/drain capacitance

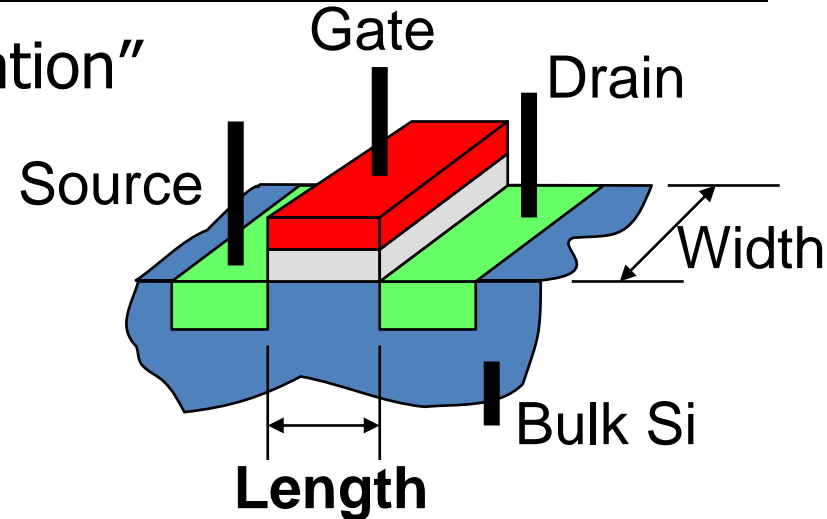
Result: switching speed ↑ linearly as gate length ↓

- much of past performance gains

But 2nd-order effects more complicated

– Process variation across chip increasing

- Some transistors slow, some fast
- Increasingly active research area: dealing with this



Moore's Speed Effect #2: More Transistors

Linear shrink in each of 2 dimensions

- 180 nm, 130 nm, 90 nm, 65 nm, 45 nm, 32 nm, 22 nm, 14 nm, 10 nm, 7 nm, 5 nm, 3 nm, ...
- Each generation is a 1.414 linear shrink
- Results in 2x more transistors (1.414×1.414)

More transistors → increased performance

- **Job of computer architect:** figure out what to do with the ever-increasing # of transistors
- *Examples:* caches, branch predictors, exploiting parallelism at all levels

Moore's Speed Effect #3: Psychological

Moore's Curve: common interpretation of Moore's Law

- "CPU performance doubles every 18 months"
- Self fulfilling prophecy: 2x in 18 months is $\sim 1\%$ per week
 - Q: Would you add a feature that improved performance 20% if it would delay the chip 8 months?
- Processors under Moore's Curve (arrive too late) fail spectacularly
 - *E.g.*, Intel's Itanium, Sun's Millennium

Power & Energy

Power/Energy Increasingly Important

- **Battery life** for mobile devices
 - Laptops, phones, cameras
- **Tolerable temperature** for devices without active cooling
 - Power means temperature, active cooling means **cost**
 - No fan in a cell phone, no market for a hot cell phone
- **Electric bill** for compute/data centers
 - Pay for power twice: once in, once out (to cool)
- **Environmental concerns**
 - “Computers” account for growing fraction of energy consumption

Energy & Power

Energy: total amount of energy stored/used

- Battery life, electric bill, environmental impact

Power: energy per unit time

- Related to “performance” (also a “per unit time” metric)
- Power impacts power supply, cooling needs (cost)
- Peak power vs. average power
 - E.g., camera power “spikes” when you take a picture

Two sources:

- **Dynamic power:** active switching of transistors
- **Static power:** transistors leak even when inactive

How to Reduce Dynamic Power

- Target each component: $P_{\text{dynamic}} \sim N \times C \times V^2 \times f \times A$
- **Reduce number of transistors (N)**
 - Use fewer transistors/gates
- **Reduce capacitance (C)**
 - Smaller transistors (Moore's law)
- **Reduce voltage (V)**
 - Quadratic reduction in energy consumption!
 - But also slows transistors (transistor speed is \sim to V)
- **Reduce frequency (f)**
 - Slow clock frequency – MacBook Air
- **Reduce activity (A)**
 - “Clock gating” disable clocks to unused parts of chip
 - Don't switch gates unnecessarily

How to Reduce Static Power

- Target each component: $P_{\text{static}} \sim N \times V \times e^{-Vt}$
- **Reduce number of transistors (N)**
 - Use fewer transistors/gates
- **Reduce voltage (V)**
 - Linear reduction in static energy consumption
 - But also slows transistors (transistor speed is \sim to V)
- **Disable transistors** (also targets N)
 - “Power gating” disable power to unused parts (long time to power up)
 - Power down units (or entire cores) not being used
- **Dual V_t** – use a mixture of high and low V_t transistors (slow for SRAM)
 - Requires extra fabrication steps (cost)
- **Low-leakage transistors**
 - High-K/Metal-Gates in Intel’s 45 nm process

Moore's Effect on Power

+ Reduces power/transistor

- Reduced sizes and surface areas reduce capacitance (C)

– Increases power density and total power

- By increasing transistors/area and total transistors
- Faster transistors → higher frequency → more power
- Hotter transistors leak more (thermal runaway)
- **What to do?** Reduce voltage [486 (5V) → Core2 (1.1V)]
 - + ↓ dynamic power quadratically, static power linearly
 - Keeping V_t the same and reducing frequency (F)
 - Lowering V_t and increasing leakage exponentially
 - **or** new techniques like high-K and dual- V_T

Continuation of Moore's Law

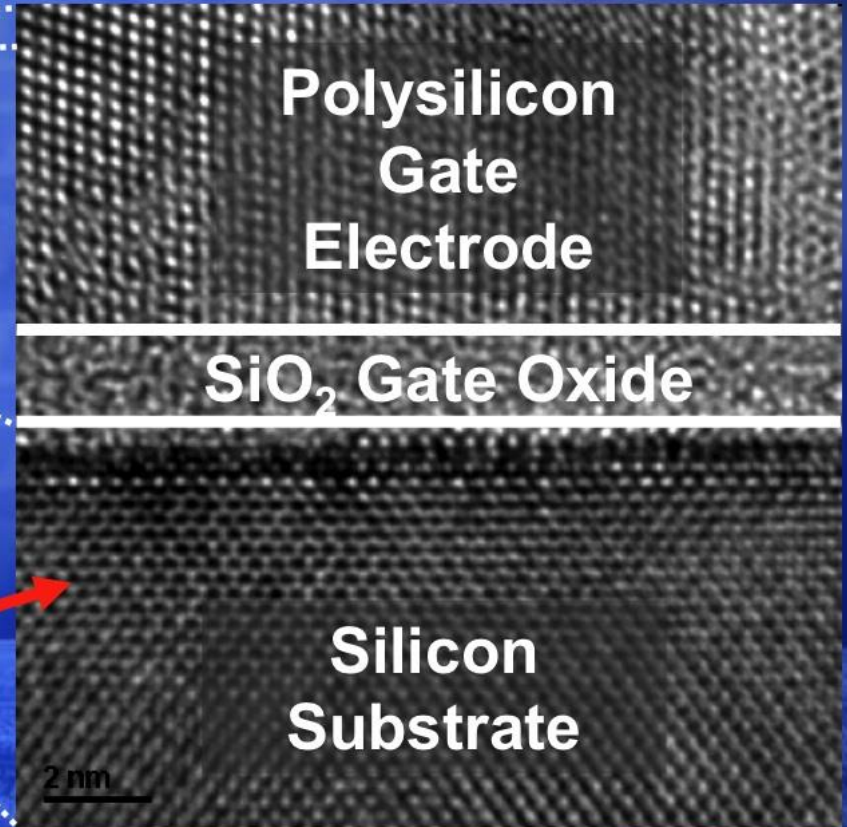
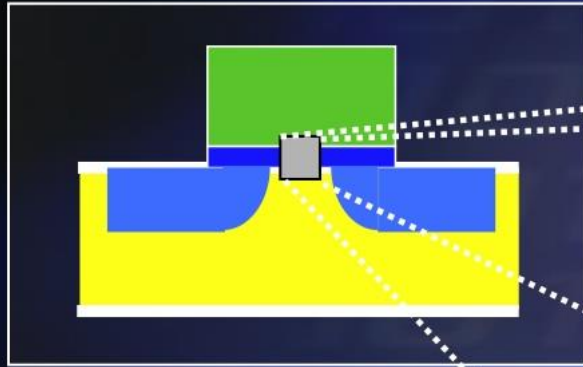
Process Name	P856	P858	Px60	P1262	P1264	P1266	P1268	P1270
1st Production	1997	1999	2001	2003	2005	2007	2009	2011
Process Generation	0.25 μ m	0.18 μ m	0.13 μ m	90 nm	65 nm	45 nm	32 nm	22 nm
Wafer Size (mm)	200	200	200/300	300	300	300	300	300
Inter-connect	Al	Al	Cu	Cu	Cu	Cu	Cu	?
Channel	Si	Si	Si	Strained Si	Strained Si	Strained Si	Strained Si	Strained Si
Gate dielectric	SiO ₂	SiO ₂	SiO ₂	SiO ₂	SiO ₂	High-k	High-k	High-k
Gate electrode	Poly-silicon	Poly-silicon	Poly-silicon	Poly-silicon	Poly-silicon	Metal	Metal	Metal

Introduction targeted at this time

Subject to change

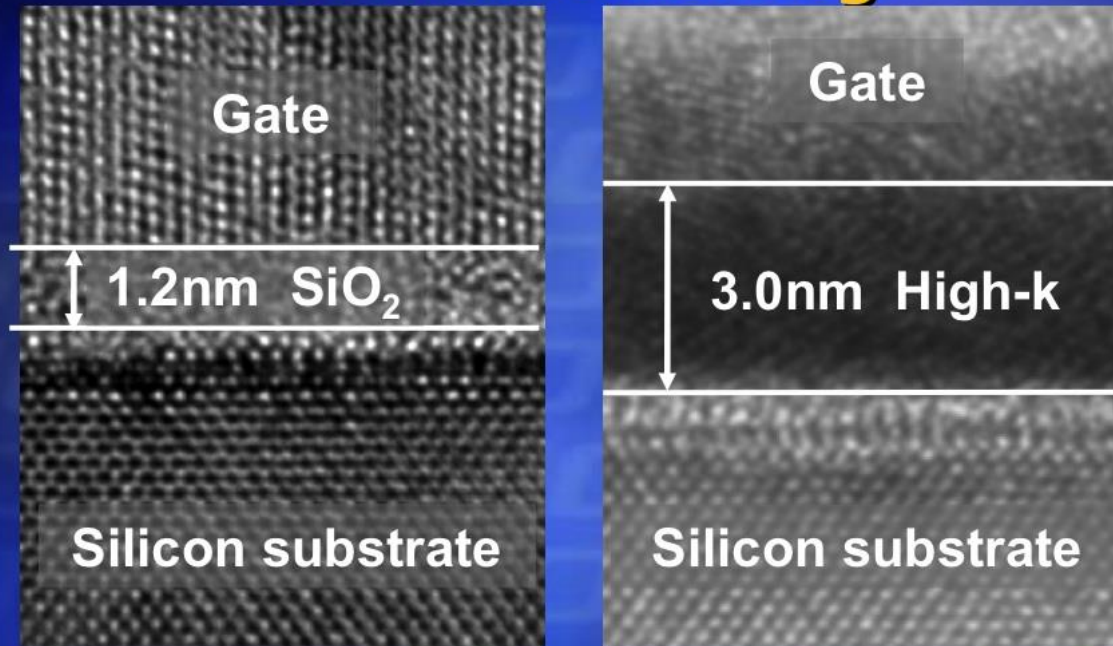
Intel found a solution for High-k and metal gate

Gate dielectric today is only a few molecular layers thick



Individual
Atoms

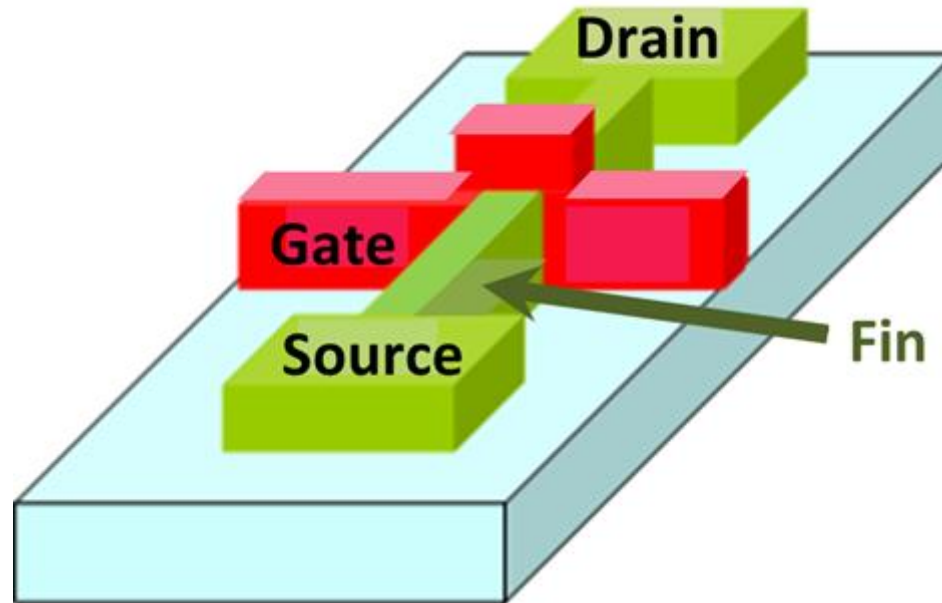
High-k Dielectric reduces leakage substantially



Benefits compared to current process technologies

	High-k vs. SiO ₂	Benefit
Capacitance	60% greater	<i>Much faster transistors</i>
Gate dielectric leakage	> 100x reduction	<i>Far cooler</i>

FinFET



By Irene Ringworm at the English language Wikipedia, CC BY-SA 3.0,
<https://commons.wikimedia.org/w/index.php?curid=3833512>

Reliability

Technology Basis for Reliability

- **Transient faults**
 - A bit “flips” randomly, **temporarily**
 - Cosmic rays etc. (more common at higher altitudes!)
- **Permanent (hard) faults**
 - A gate or memory cell wears out, **breaks and stays broken**
 - Temperature & electromigration slowly deform components
- Solution for both: **redundancy** to detect and tolerate

Moore's Bad Effect on Reliability

– **Transient faults:**

- Small (low charge) transistors are more easily flipped
- Even low-energy particles can flip a bit now

– **Permanent faults:**

- Small transistors and wires deform and break more quickly
- Higher temperatures accelerate the process

Wasn't a problem until ~10 years ago (except in satellites)

- Memory (DRAM): these dense, small devices hit first
- Then on-chip memory (SRAM)
- Logic is starting to have problems...

Moore's Good Effect on Reliability

- Scaling makes devices less reliable
- + Scaling increases device density to enable **redundancy**
- Examples
 - Error correcting code for memory (DRAM), β s (SRAM)
 - Core-level redundancy: paired-execution, hot-spare, etc.
 - Intel's Core i7 (Nehalem) uses 8 transistor SRAM cells
 - Versus the standard 6 transistor cells
- Big open questions
 - Can we protect logic efficiently? (w/o 2-3x overhead)
 - Can architectural techniques help hardware reliability?
 - Can software techniques help?

Summary

Moore's Law in the Future

- Won't last forever, approaching physical limits
 - But betting against it has proved foolish in the past
 - Likely to "slow" rather than stop abruptly
- Transistor count will likely continue to scale
 - "Die stacking" is on the cusp of becoming main stream
 - Uses the third dimension to increase transistor count
- But transistor performance scaling?
 - Running into physical limits
 - Example: gate oxide is less than 10 silicon atoms thick!
 - Can't decrease it much further
 - Power is becoming a limiting factor

Summary of Device Scaling

- + Reduces unit cost
 - But increases startup cost
- + Increases performance
 - Reduces transistor/wire delay
 - Gives us more transistors with which to increase performance
- + Reduces local power consumption
 - Quickly undone by increased integration, frequency
 - Aggravates power-density and temperature problems
- Aggravates reliability problem
 - + But gives us the transistors to solve it via redundancy