# NBAY 6120

## Chip Making

## &

## Microprocessor Technology

NBAY 1620

March 2, 2017

Donald P. Greenberg

Lecture 2

- ***Required Reading:***

  - Craig R. Barrett. From Sand to Silicon: Manufacturing an Integrated Circuit, Scientific American, Special Issue, The Solid-State Century, January 1998, pp. 55-61. (Search: e-Journals/ Scientific American Archive Online/article (full text) http://www.library.cornell.edu/johnson/library/general/emba.html

  - Peter J. Denning and Ted G. Lewis. "Exponential Laws of Computing Growth." Communications of the ACM. January 2017. ACM.org.
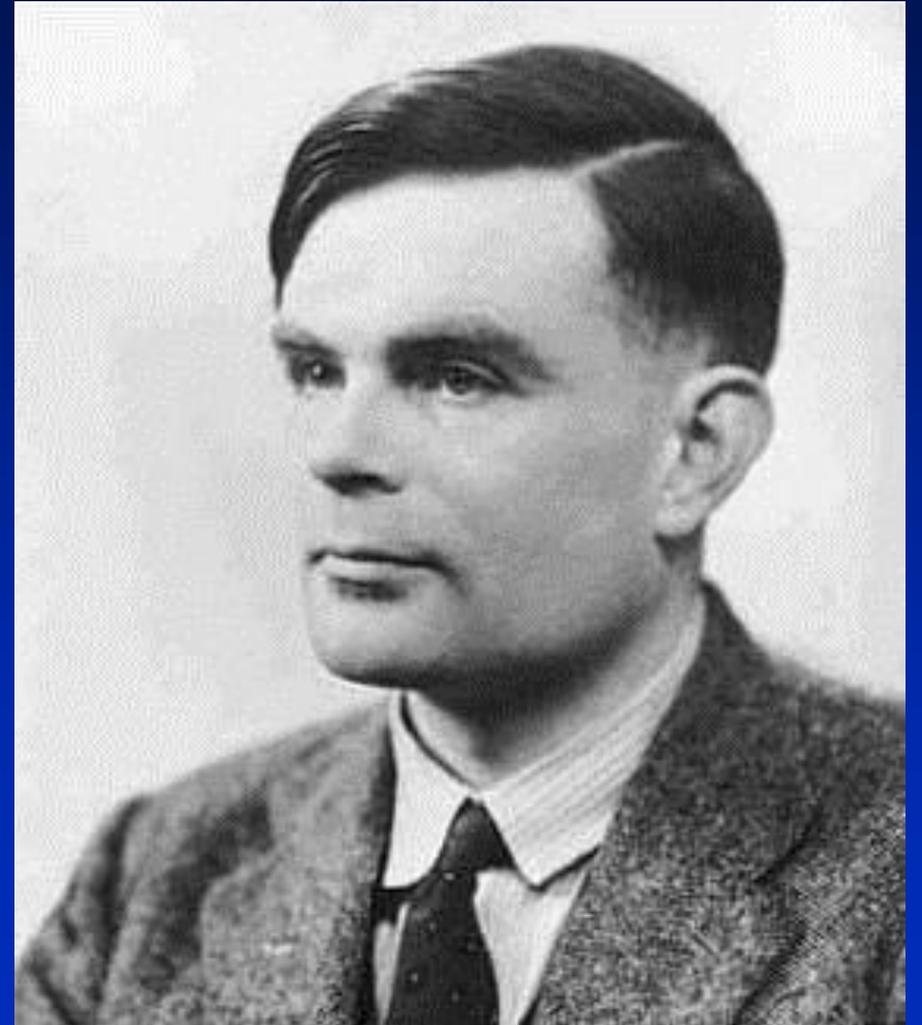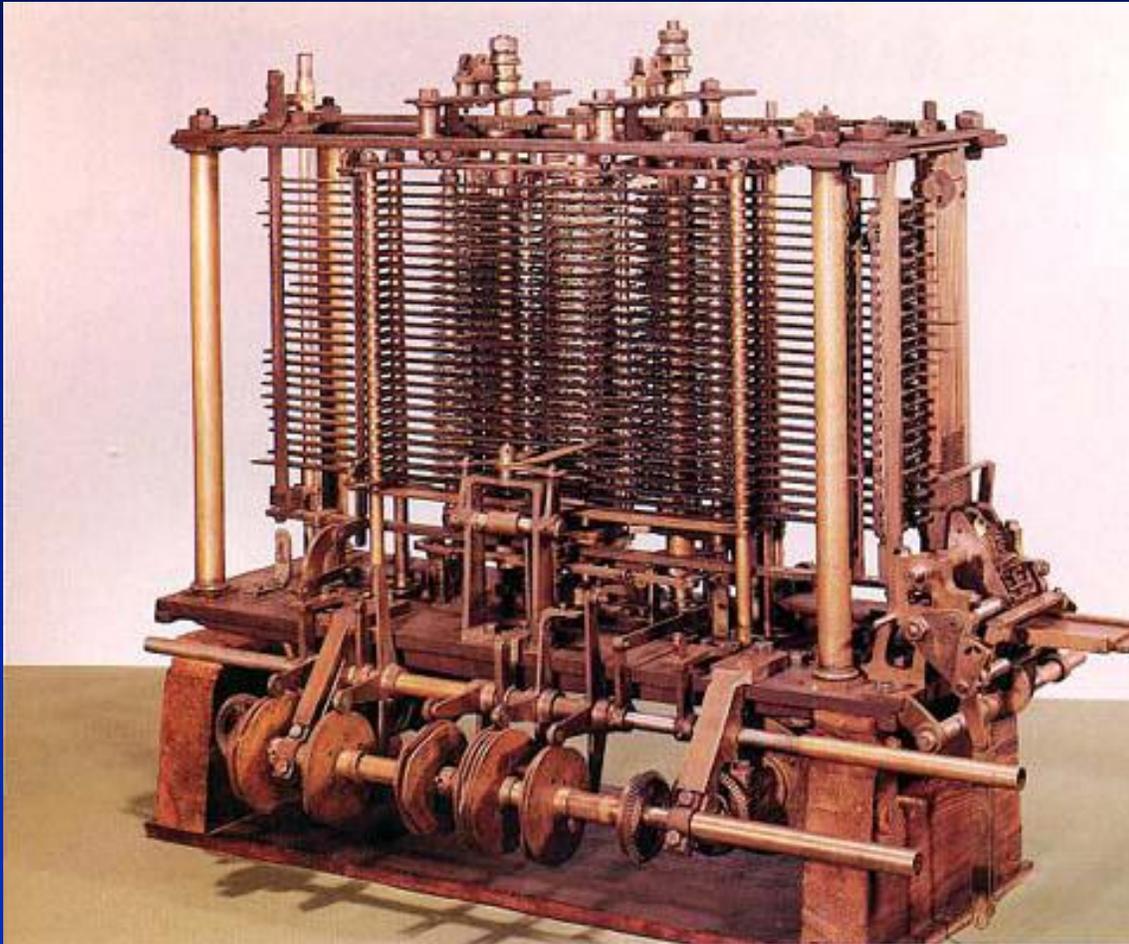
    Optional:

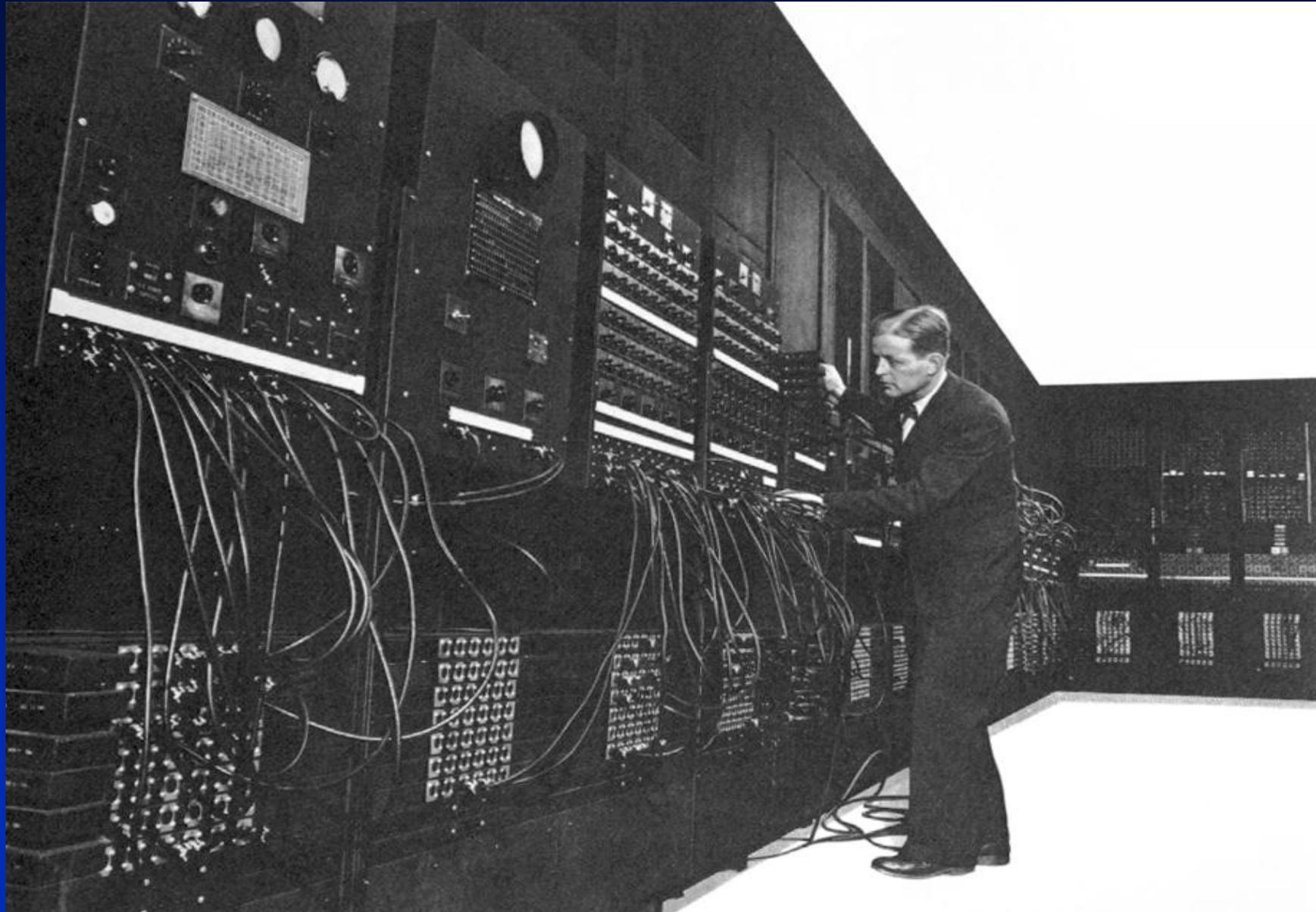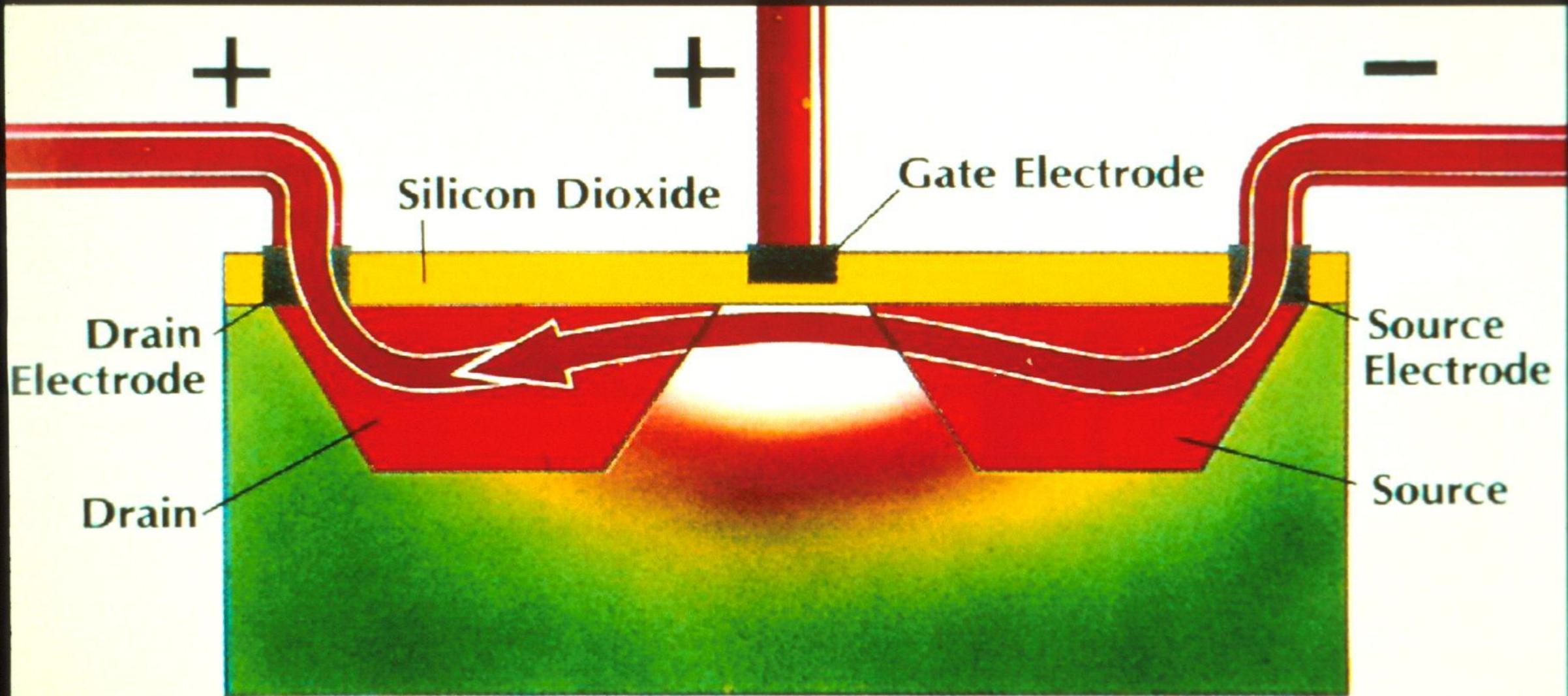    Mack, Chris. "The Multiple Lives of Moore's Law." *IEEE Spectrum* Apr. 2015: 30-37. *Cornell University Library*. Web. http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=7065415

# Turing Machine

# Alan Turing

# Eniac 1946

# Cloud Computing - 2010

# Microprocessor Transistor Counts 1971-2011 & Moore's Law



Transistor count (y-axis): 2,600,000,000 · 1,000,000,000 · 100,000,000 · 10,000,000 · 1,000,000 · 100,000 · 10,000 · 2,300

Years (x-axis): 1971 · 1980 · 1990 · 2000 · 2011

**VAX-11/780 (1980)**

**IBM 360 Model 75 (1965)**

**Cray T3E (1995)**

**Cloud Computing (2010)**

curve shows transistor count doubling every two years

Curve shows 'Moore's Law': transistor count doubling every two years

Processor labels: 16-Core SPARC T3, Six-Core Core i7, Six-Core Xeon 7400, 10-Core Xeon Westmere-EX, Dual-Core Itanium 2, 8-core POWER7, Quad-core z196, Quad-Core Itanium Tukwila, 8-Core Xeon Nehalem-EX, AMD K10, POWER6, Itanium 2 with 9MB cache, Six-Core Opteron 2400, Core i7 (Quad), AMD K10, Itanium 2, Core 2 Duo, Cell, AMD K8, Barton, Atom, Pentium 4, AMD K7, AMD K6-III, AMD K6, Pentium III, Pentium II, AMD K5, Pentium, 80486, 80386, 80286, 68000, 80186, 8086, 8088, 8085, 6800, 6809, 8080, Z80, 8008, MOS 6502, 4004, RCA 1802
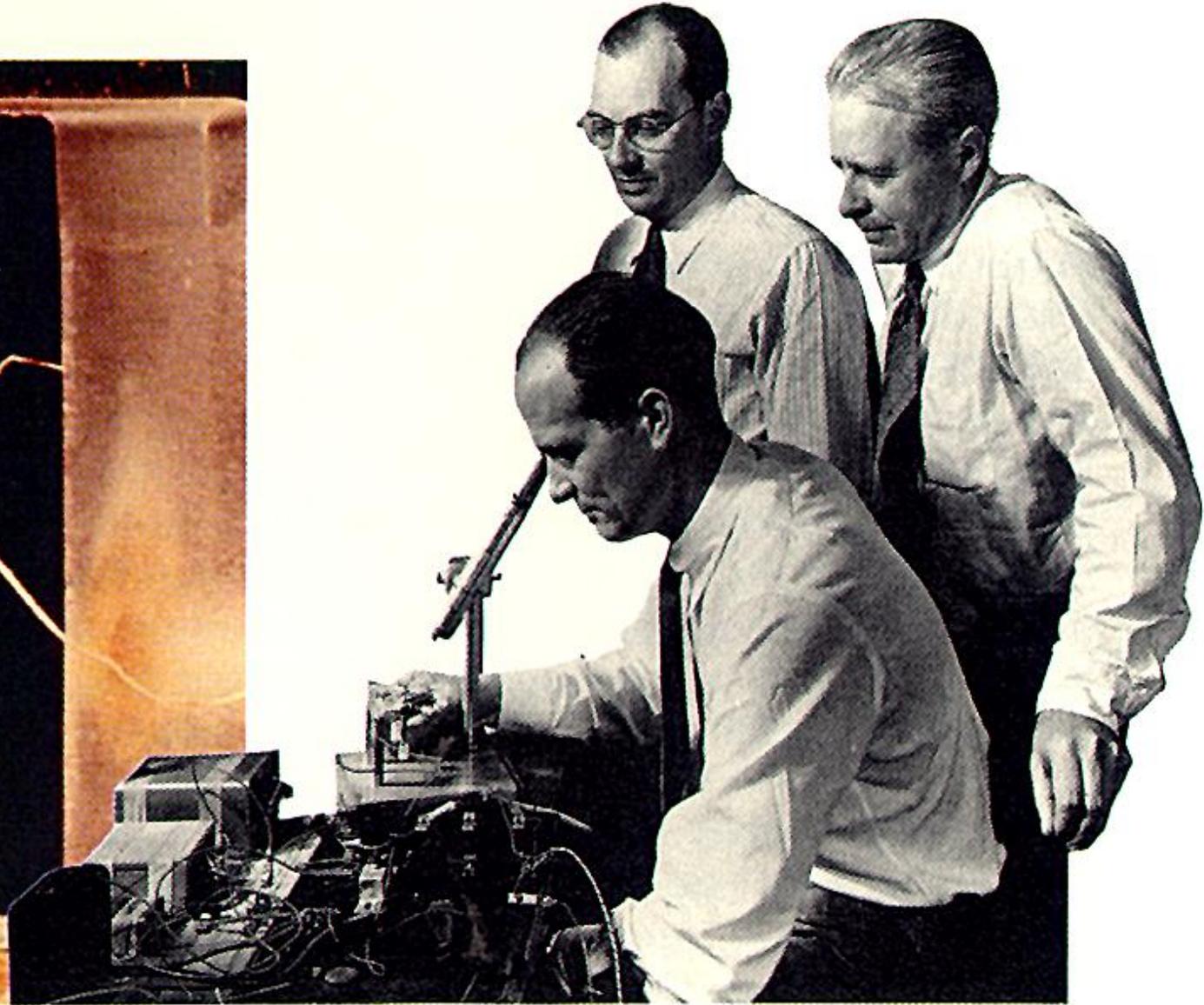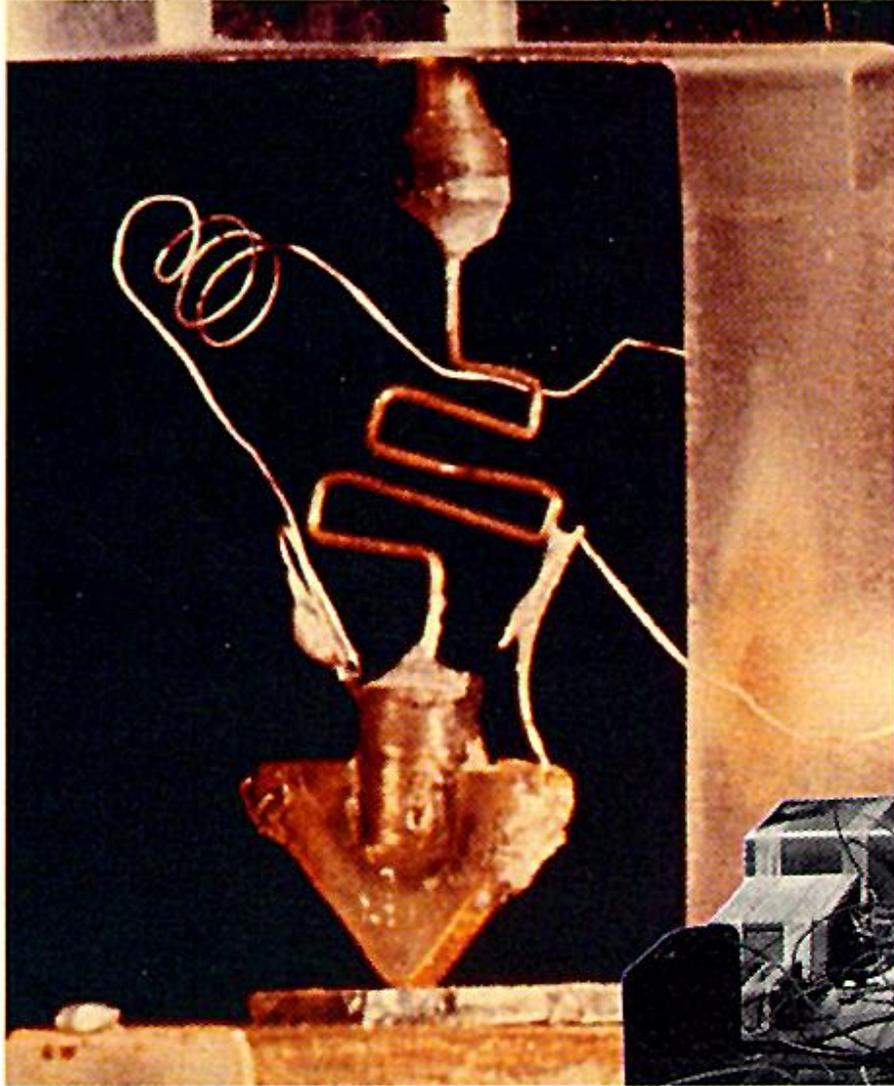
# Transistor

# Shockley, Bardeen & Brattain            1947

# Shockley & the Traitorous Eight

- William Shockley -  Receives the Nobel Prize in Physics with Bardeen and Brattain (1956)
leaves Bell Laboratory and forms Fairchild Semiconductor

- Julius Blank -  founded Xicor

- Jean Hoerni -  invented the planar process
founded Amelco → Teledyne

- Jay Last -  founded Amelco → Teledyne
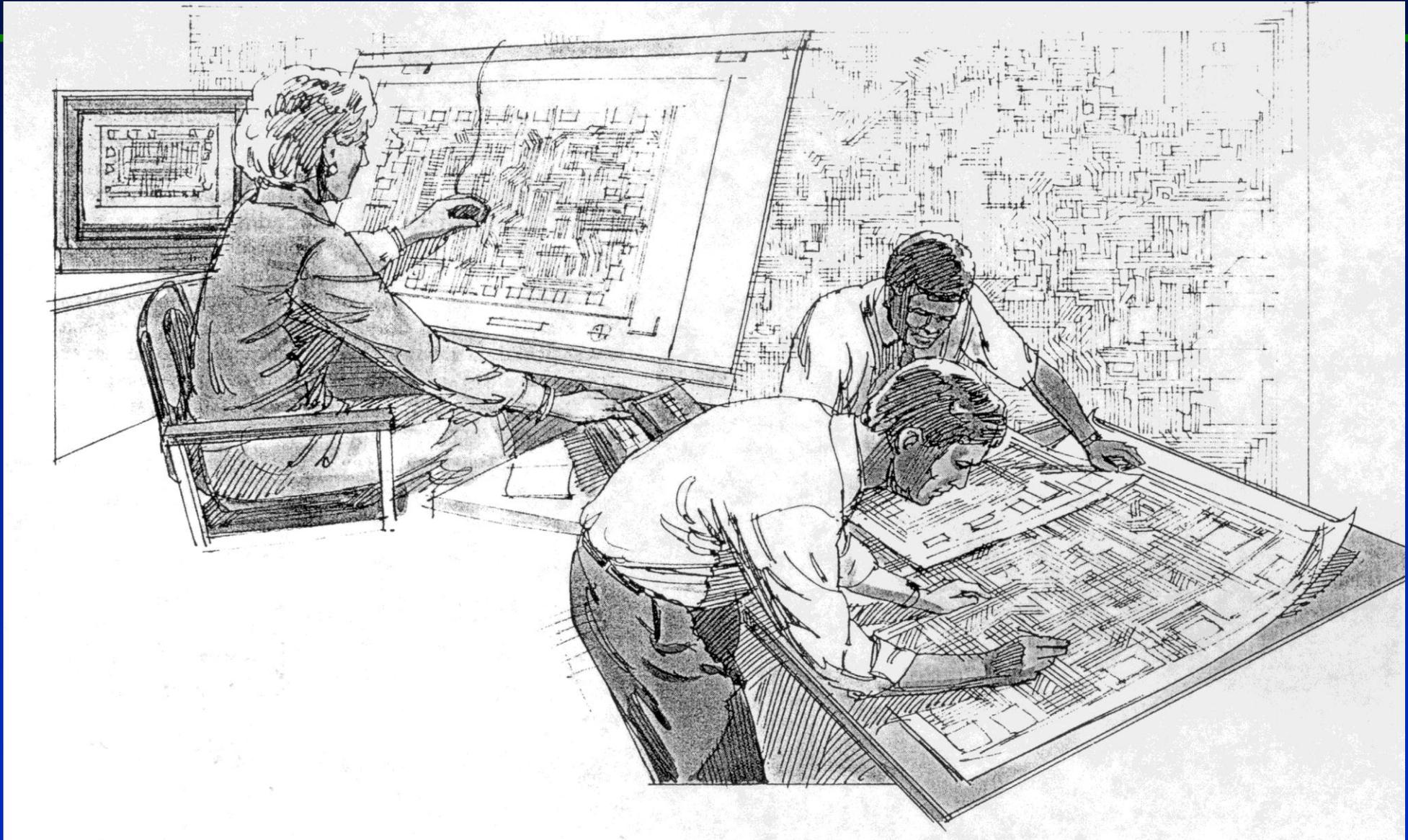
- Sheldon Roberts -  founded Amelco → Teledyne

# Shockley & the Traitorous Eight

- Gordon Moore -        founded Intel in 1968

- Robert Noyce -        founded Intel in 1968

- Eugene Kleiner -      founded Kleiner-Perkins

- Victor Grinich -      only a poor professor at UC Berkeley & Stanford

# From Sand to Silicon – Manufacturing an Integrated Circuit

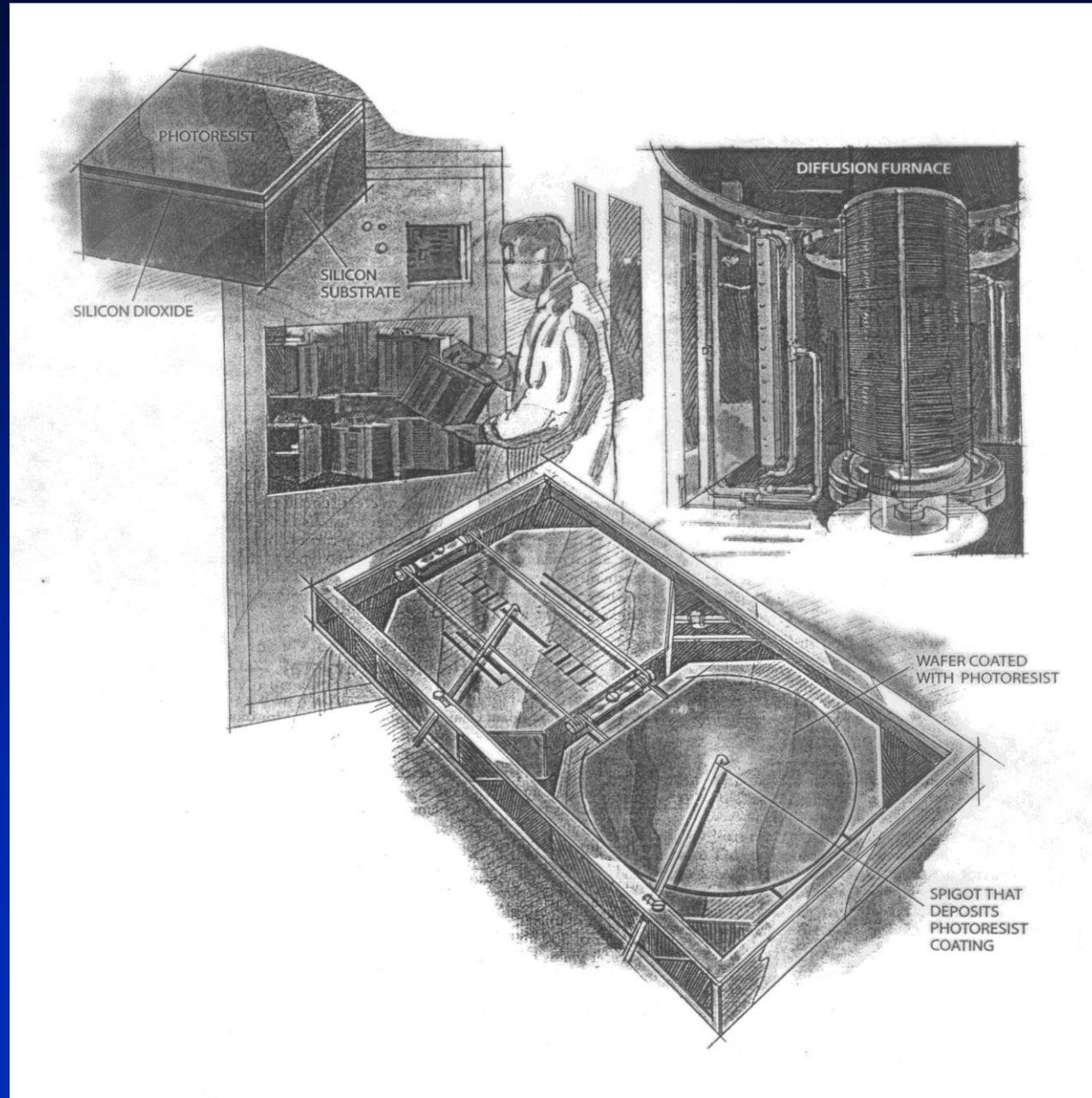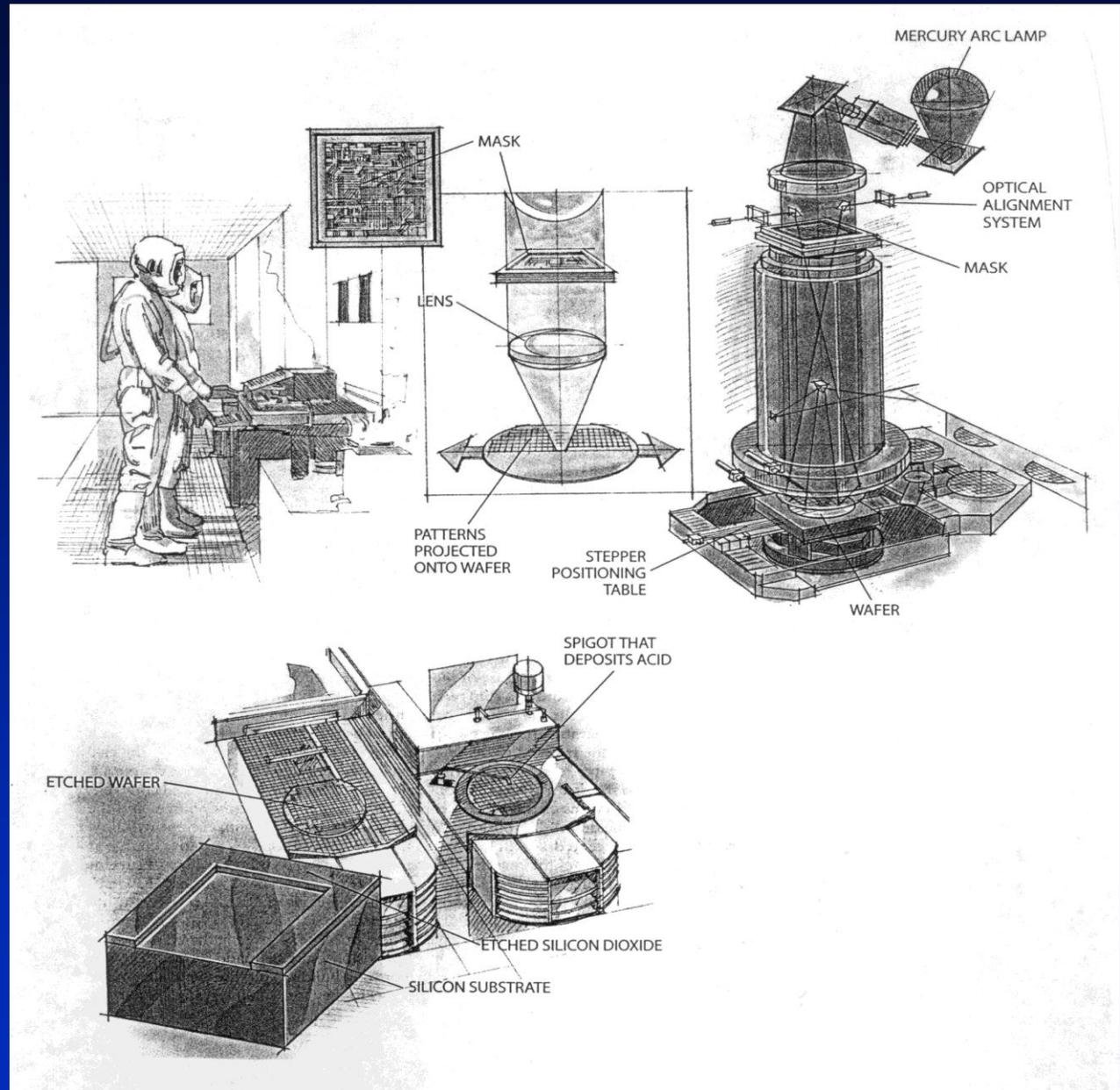## Scientific American: The Solid-State Century, Special Issue 1998
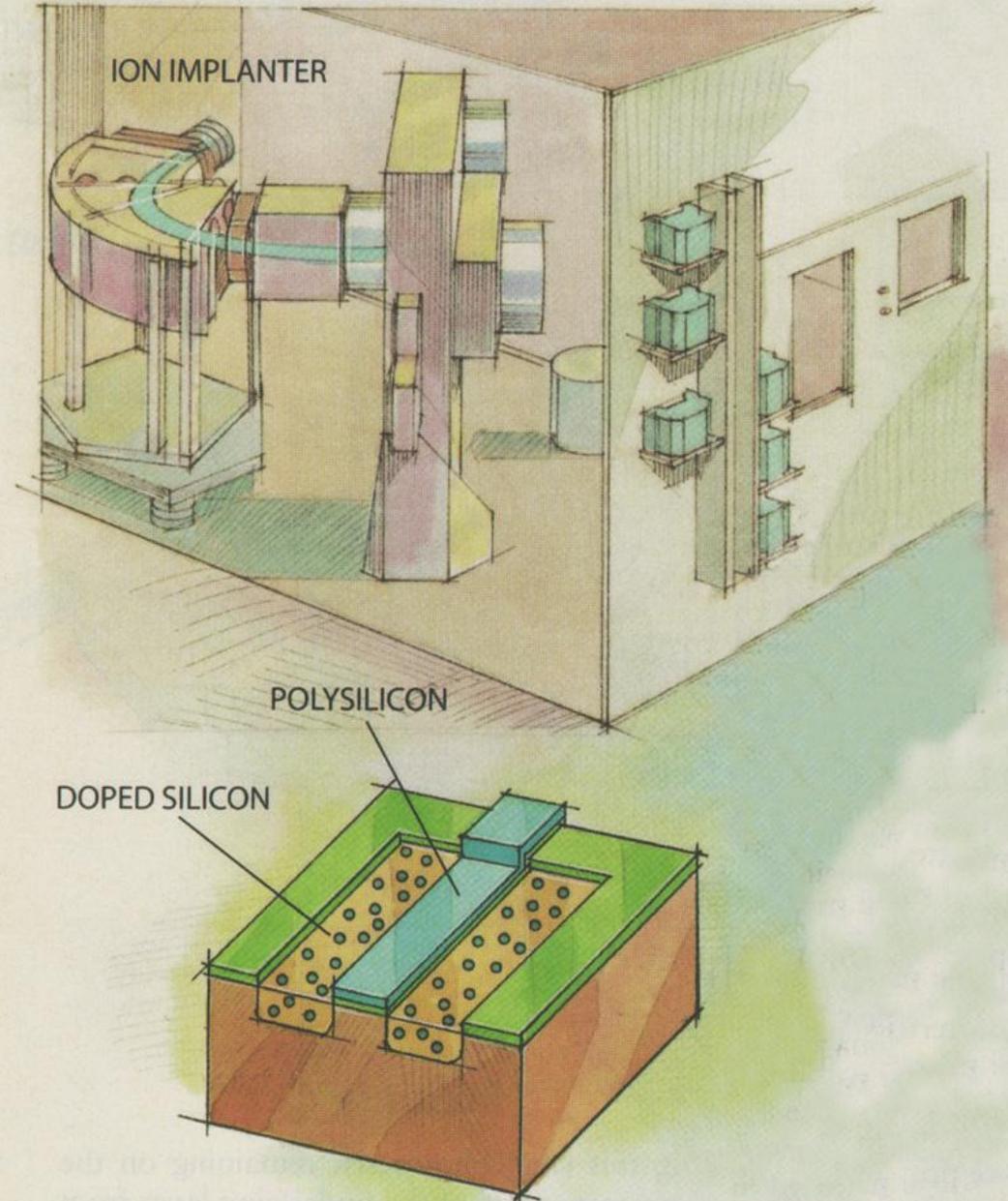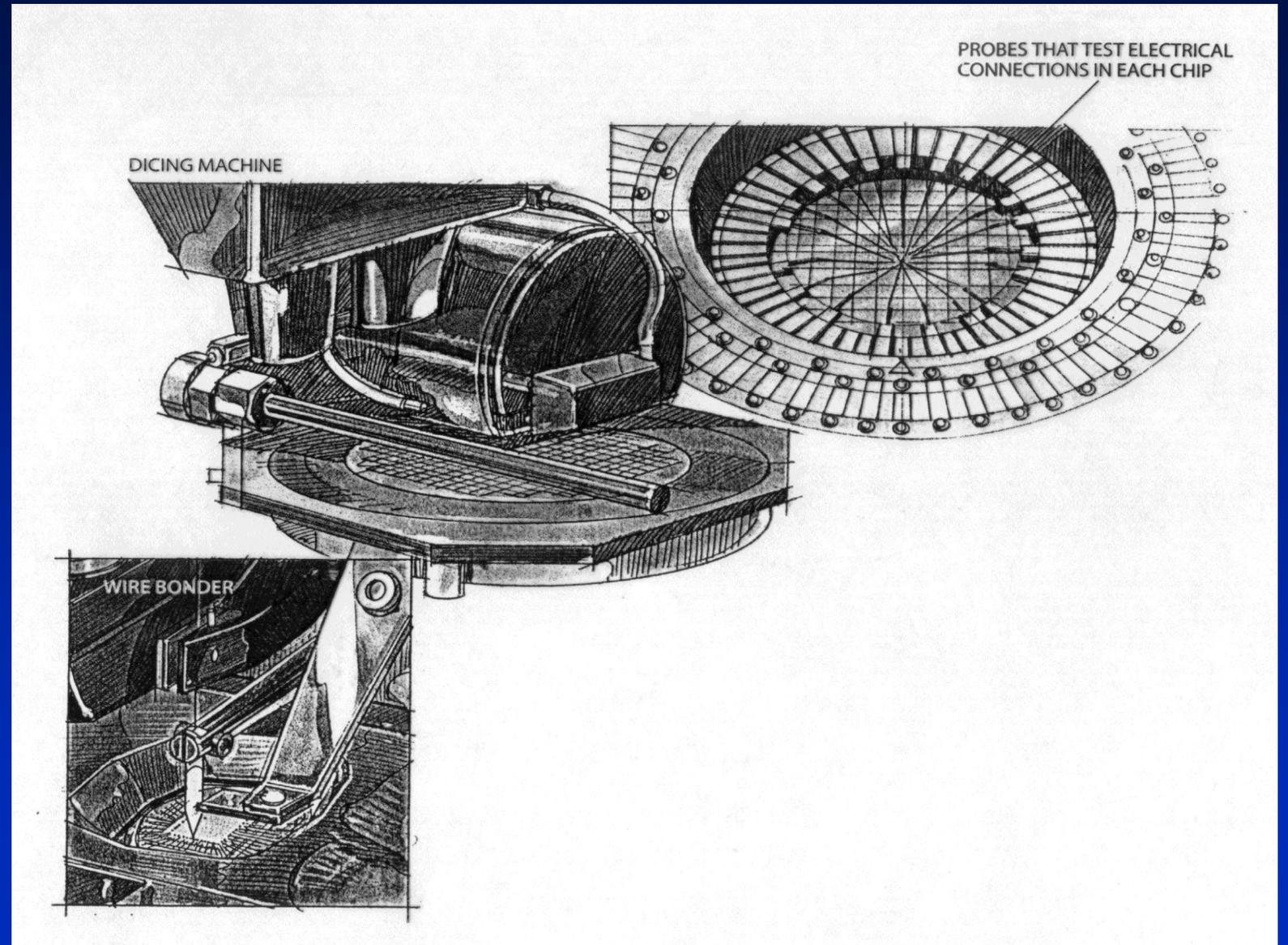
# Chip Design
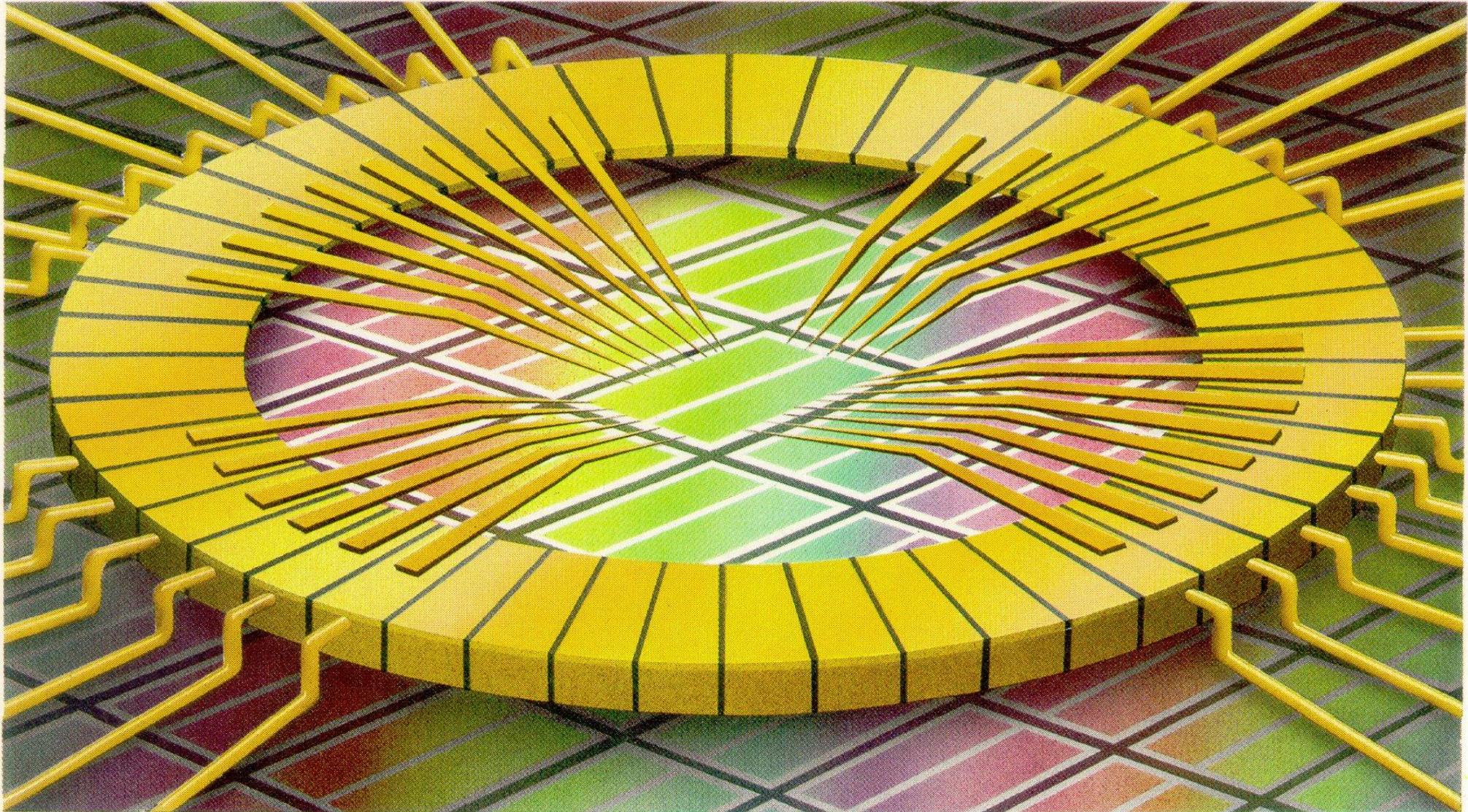
# Silicon Crystal
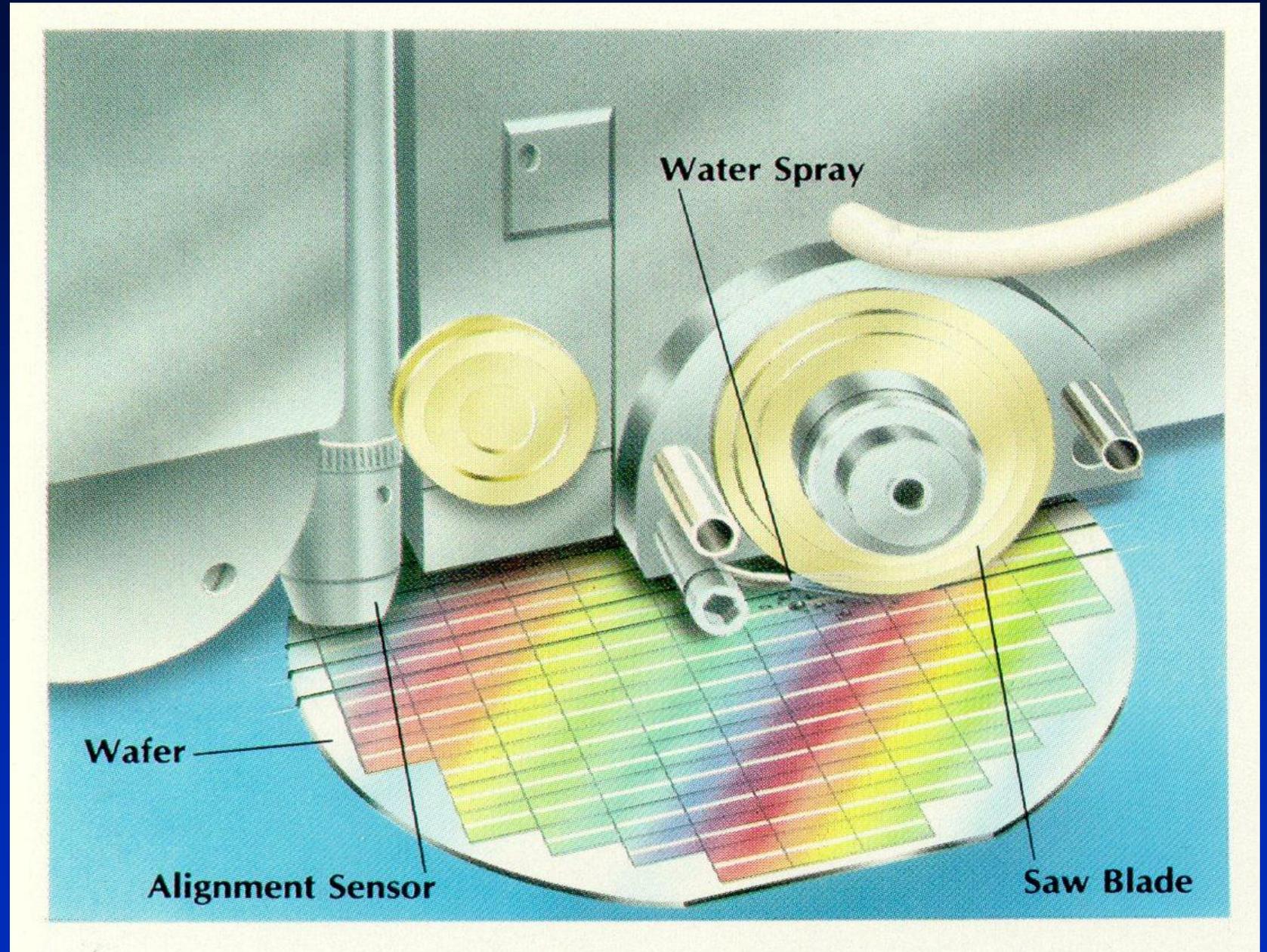
# Layering

# Masking & Etching
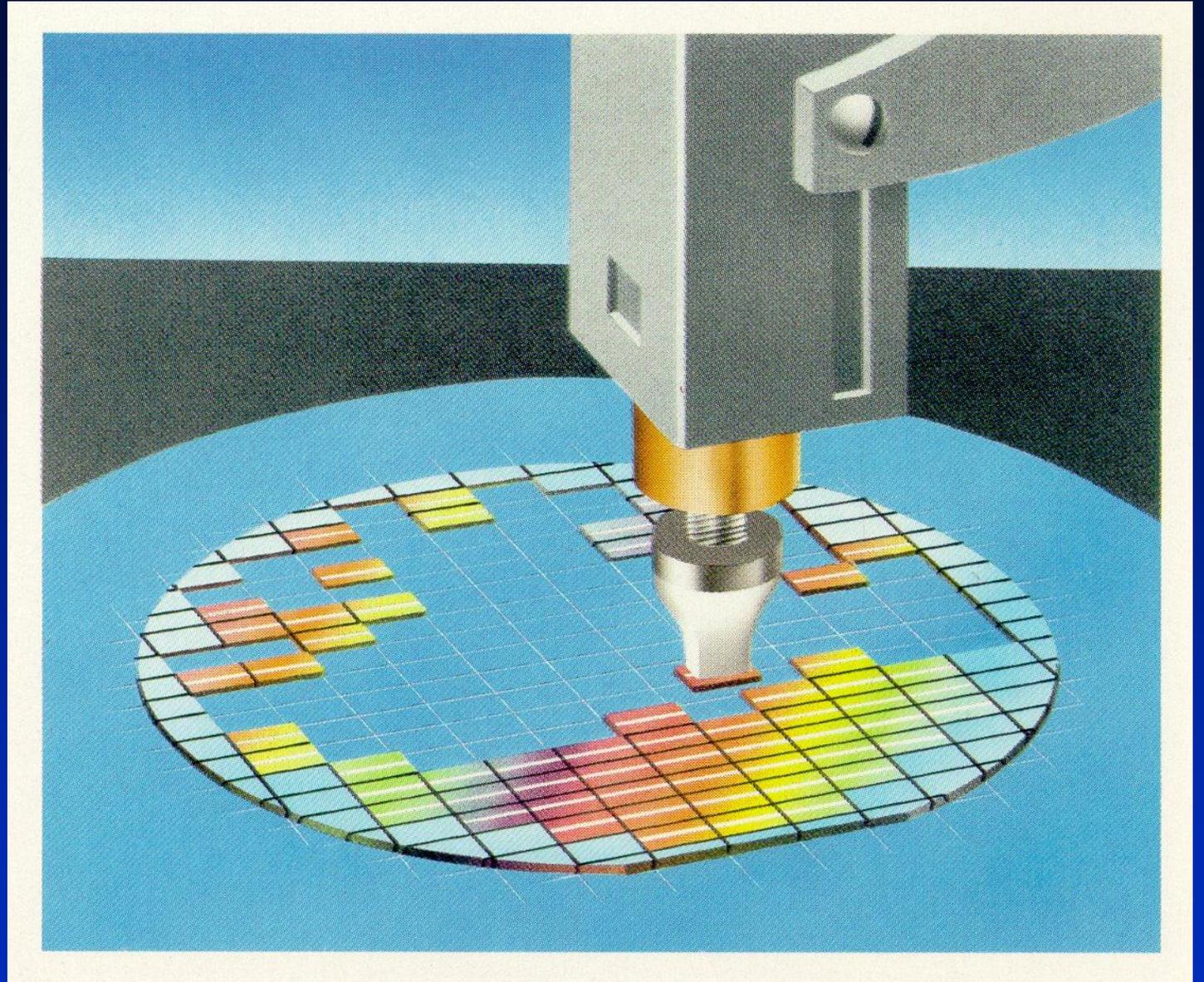
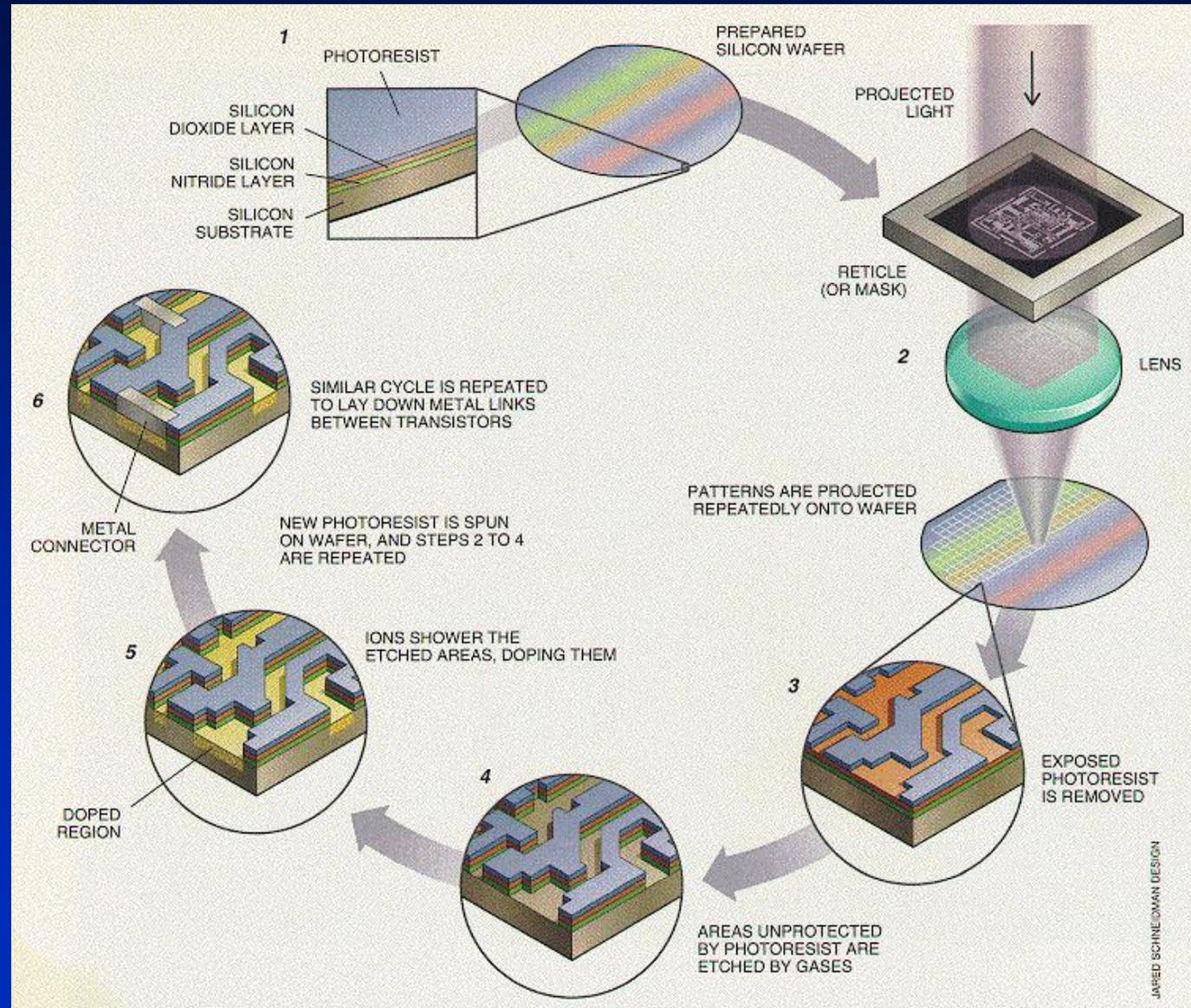# Doping

# Interconnections & Dicing
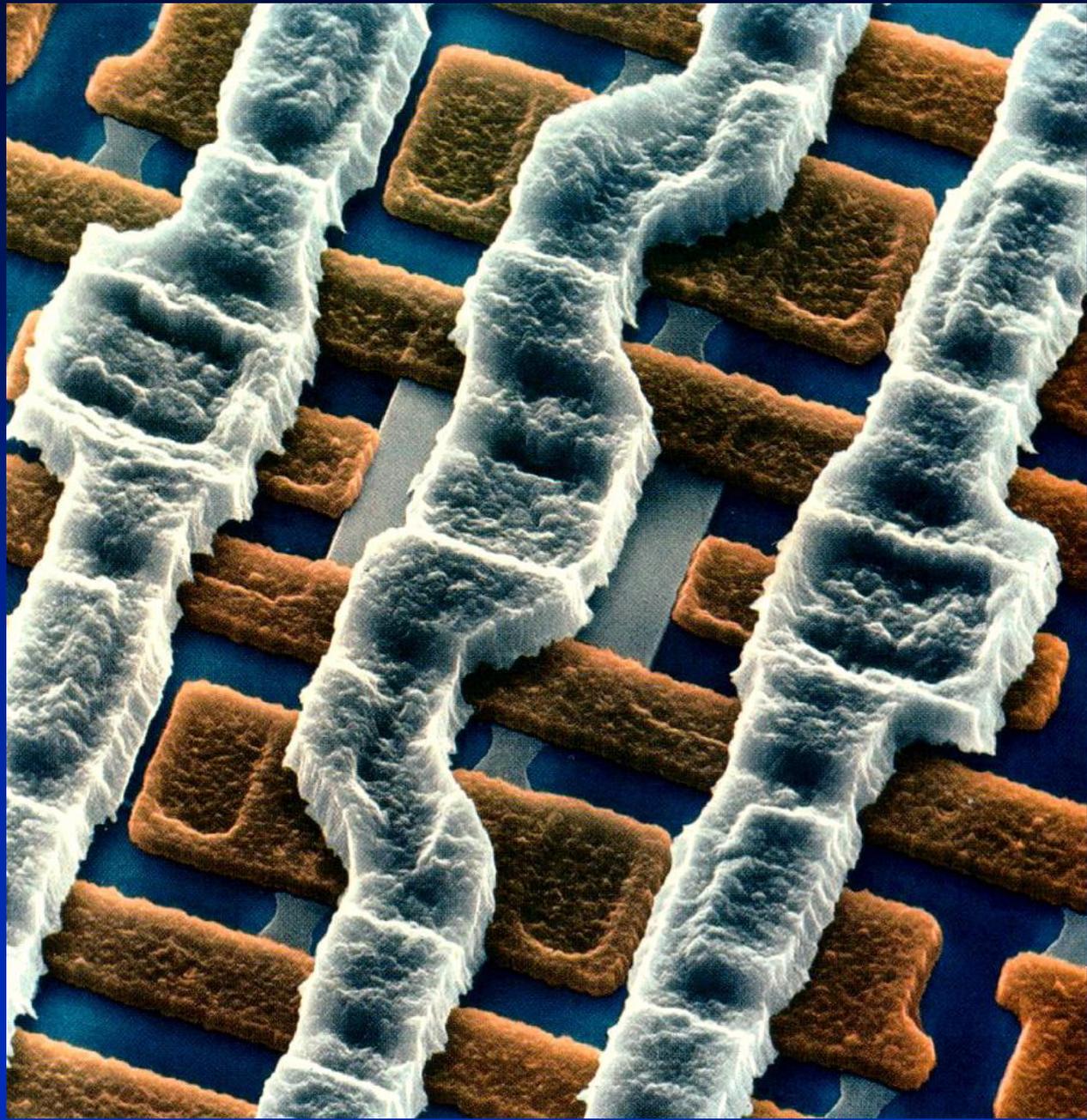
# Probing Electrical Connections
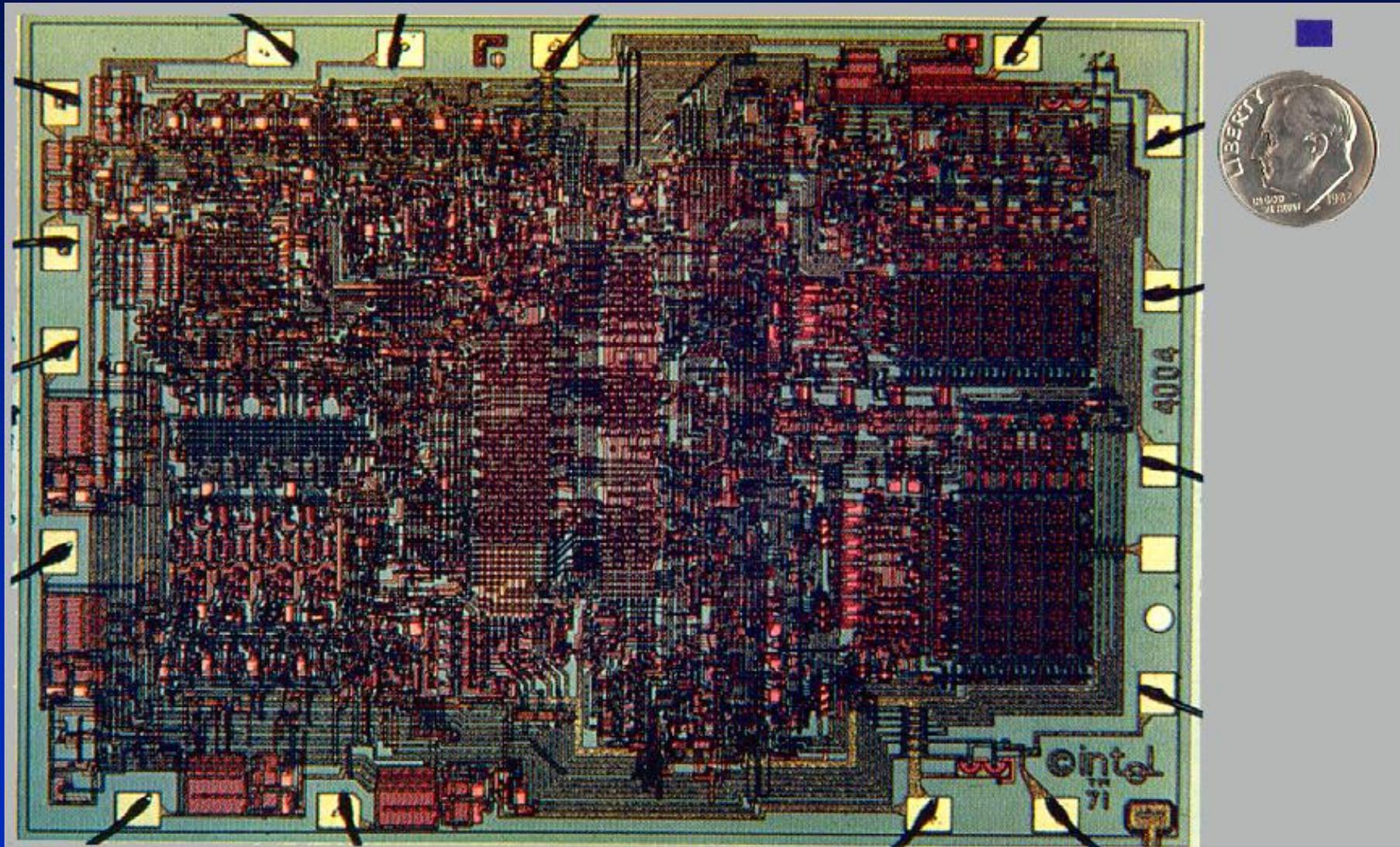
# Dicing

# Chip Selection

# Chip Fabrication

# International Technology Roadmap for Semiconductors

|  | 2001 | 2004 | 2007 | 2010 | 2013 | 2016 |
|---|---|---|---|---|---|---|
| **Technology (nanometers)** | **130nm** | **90nm** | **65nm** | **45nm** | **32nm** | **22nm** |
| **Functions per Chip (millions)** | 97 | 193 | 386 | 1546 | 3092 | 6184 |
| **Clock Speed (Ghz)** | 2.5Ghz | 4.1Ghz | 9.3Ghz | 15Ghz | 23Ghz | 40Ghz |
| **Wafer Size (millimeters)** | 200mm | 300mm | 300mm | 300mm | 450mm | 450mm |
| **Chip Size (mm$^2$)** | 140 mm$^2$ | 140 mm$^2$ | 140 mm$^2$ | 140 mm$^2$ | 140 mm$^2$ | 140 mm$^2$ |

Roughly 0.5 shrink every 3 years.  Intel released 22 nm chips in 2013

# Intel 4004

# November 1971

# Moore's Law – CPU Transistor Counts

| Processor | Transistor count | Date of introduction | Manufacturer | Process | Area |
|---|---|---|---|---|---|
| Intel 4004 | 2,300 | 1971 | Intel | 10 µm | |
| Intel 8008 | 3,500 | 1972 | Intel | 10 µm | |
| Intel 8080 | 4,500 | 1974 | Intel | 6 µm | |
| Intel 8088 | 29,000 | 1979 | Intel | 3 µm | |
| Intel 80286 | 134,000 | 1982 | Intel | 1.5 µm | |
| Intel 80386 | 275,000 | 1985 | Intel | 1.5 µm | |
| Intel 80486 | 1,180,000 | 1989 | Intel | 1 µm | |
| Pentium | 3,100,000 | 1993 | Intel | 0.8 µm | |
| AMD K5 | 4,300,000 | 1996 | AMD | 0.5 µm | |
| Pentium II | 7,500,000 | 1997 | Intel | 0.35 µm | |
| AMD K6 | 8,800,000 | 1997 | AMD | 0.35 µm | |
| Pentium III | 9,500,000 | 1999 | Intel | 0.25 µm | |
| AMD K6-III | 21,300,000 | 1999 | AMD | 0.25 µm | |
| AMD K7 | 22,000,000 | 1999 | AMD | 0.25 µm | |
| Pentium 4 | 42,000,000 | 2000 | Intel | 180 nm | |
| Atom | 47,000,000 | 2008 | Intel | 45 nm | |
| Barton | 54,300,000 | 2003 | AMD | 130 nm | |
| AMD K8 | 105,900,000 | 2003 | AMD | 130 nm | |
| Itanium 2 | 220,000,000 | 2003 | Intel | 130 nm | |

# Moore's Law – CPU Transistor Counts

| Processor | Transistor count | Date of introduction | Manufacturer | Process | Area |
|---|---|---|---|---|---|
| Core 2 Duo | 291,000,000 | 2006 | Intel | 65 nm | |
| AMD K10 | 463,000,000 | 2007 | AMD | 65 nm | |
| AMD K10 | 758,000,000 | 2008 | AMD | 45 nm | |
| Itanium 2 with 9MB cache | 592,000,000 | 2004 | Intel | 130 nm | |
| Core i7 (Quad) | 731,000,000 | 2008 | Intel | 45 nm | 263 mm² |
| POWER6 | 789,000,000 | 2007 | IBM | 65 nm | 341 mm² |
| Six-Core Opteron 2400 | 904,000,000 | 2009 | AMD | 45 nm | |
| Six-Core Core i7 | 1,170,000,000 | 2010 | Intel | 32 nm | |
| Dual-Core Itanium 2 | 1,700,000,000 | 2006 | Intel | 90 nm | 596 mm² |
| Six-Core Xeon 7400 | 1,900,000,000 | 2008 | Intel | 45 nm | |
| Quad-Core Itanium Tukwila | 2,000,000,000 | 2010 | Intel | 65 nm | |
| Six-Core Core i7 (Sandy Bridge-E) | 2,270,000,000 | 2011 | Intel | 32 nm | 434 mm² |
| 8-Core Xeon Nehalem-EX | 2,300,000,000 | 2010 | Intel | 45 nm | 684 mm² |
| 10-Core Xeon Westmere-EX | 2,600,000,000 | 2011 | Intel | 32 nm | 512 mm² |
| Six-core zEC12 | 2,750,000,000 | 2012 | IBM | 32 nm | 597 mm² |
| 8-Core Itanium Poulson | 3,100,000,000 | 2012 | Intel | 32 nm | 544 mm² |
| 15-Core Xeon Ivy Bridge-EX | 4,310,000,000 | 2014 | Intel | 22nm | 541 mm² |
| 62-Core Xeon Phi | 5,000,000,000 | 2012 | Intel | 22 nm | 350 mm² |
| Xbox One Main SoC | 5,000,000,000 | 2013 | Microsoft | 28 nm | 363 mm² |
| 18-core Xeon Haswell-E5 | 5,560,000,000 | 2014 | Intel | 22 nm | 661mm² |
| IBM z13 Storage Controller | 7,100,000,000 | 2015 | IBM | 22 nm | 678 mm² |

# Moore's Law – GPU Transistor Counts

| Processor | Transistor count | Date of introduction | Manufacturer | Process | Area |
|---|---|---|---|---|---|
| R520 | 321,000,000 | 2005 | AMD | 90 nm | 288 mm² |
| R580 | 384,000,000 | 2006 | AMD | 90 nm | 352 mm² |
| G80 | 681,000,000 | 2006 | NVIDIA | 90 nm | 480 mm² |
| R600 Pele | 700,000,000 | 2007 | AMD | 80 nm | 420 mm² |
| G92 | 754,000,000 | 2007 | NVIDIA | 65 nm | 324 mm² |
| RV790XT Spartan | 959,000,000 | 2008 | AMD | 55 nm | 282 mm² |
| GT200 Tesla | 1,400,000,000 | 2008 | NVIDIA | 65 nm | 576 mm² |
| Cypress RV870 | 2,154,000,000 | 2009 | AMD | 40 nm | 334 mm² |
| Cayman RV970 | 2,640,000,000 | 2010 | AMD | 40 nm | 389 mm² |
| GF100 Fermi | 3,200,000,000 | Mar 2010 | NVIDIA | 40 nm | 526 mm² |
| GF110 Fermi | 3,000,000,000 | Nov 2010 | NVIDIA | 40 nm | 520 mm² |
| GK104 Kepler | 3,540,000,000 | 2012 | NVIDIA | 28 nm | 294 mm² |
| Tahiti RV1070 | 4,312,711,873 | 2011 | AMD | 28 nm | 365 mm² |
| GK110 Kepler | 7,080,000,000 | 2012 | NVIDIA | 28 nm | 561 mm² |
| RV1090 Hawaii | 6,300,000,000 | 2013 | AMD | 28 nm | 438 mm² |
| GM204 Maxwell | 5,200,000,000 | 2014 | NVIDIA | 28 nm | 398 mm² |
| GM200 Maxwell | 8,100,000,000 | 2015 | NVIDIA | 28 nm | 601 mm² |
| Fiji | 8,900,000,000 | 2015 | AMD | 28 nm | 596 mm² |

2007

Paul S. Otellini
Intel Corporation's fifth CEO



Forget angels on pins: Paul Otellini gets 2,500 Silverthorne chips on a silicon wafer.

# Why are we continuing to strive for smaller and smaller technology?

- More transistors/chip $\rightarrow$ increased functionality and performance

- Higher speeds $\rightarrow$ partially depends on how close together the components are placed

- Cheaper – more chips/wafer, greater yields

# Yield Ratio

$$yield = \frac{n_w}{n_t}$$

$$n_w = yield \bullet n_t$$

$n_w = number\ of\ working\ chips/wafer$

$n_t = total\ number\ of\ chips/wafer$

$Old\ fab\ lines,\ yield \rightarrow > 90\%$

$New\ fab\ lines,\ yield \rightarrow < 40\%$

# Yield per Wafer

# Yield Ratio

Number of defects/unit area depends on the process

$$\therefore Yield \approx \frac{1}{Chip\ Area}$$

Total chips ($n_t$) for a given wafer size is also inversely proportional to the chip area

# Why does the shrinking technology make the cost of manufacturing cheaper per component?

# Example:

For a 10% shrink in feature size :

$$n_{w_{new}} = n_{w_{old}} \left( \frac{1}{.9} \right)^2 \left( \frac{1}{.9} \right)^2$$

$\uparrow$ New yield $\quad$ $\uparrow$ New $n_t$

$$n_{w_{new}} = 1.52 n_{w_{old}}$$

# Can this Shrinking Technology Continue?

**Human Hair**

50μm
(50,000 nm)



**Integrated Chip**

32nm

400nm

700nm

**Visible Light**

# Keeping Up with Moore's Law

Remarkably, Moore's Law-the number of transistors that can fit on a microchip will double every 18-24 months-has held true for many years. Keeping up with this famous prediction by Intel founder Gordon Moore is getting harder.

# Getting Wafers Wet



By adding a thin layer of water between the projection lens and the wafer, the immersion system can create features 30 percent smaller.

# Photolithography

- Defining the smallest components requires short wavelengths of light.

- Currently, most fabrication processors use extreme ultra-violet light at 193nm.

- Can pass the light through water.  The water slows the light (less velocity) shrinking its wavelength.  It is estimated that this technique will meet demands for 7 more years.

- On February 20, 2006 IBM Almaden & JSR Micro demonstrated a system using an "unidentified" light slowing liquid yielding patterns 29.9nm wide.

*Science News, March 2, 2006*

# The HIGH-k SOLUTION

**By Mark T. Bohr, Robert S. Chau, Tahir Ghani & Kaizad Mistry –** October 2007 IEEE Spectrum

In 2007 new 45nm Microprocessors were the result of the first big material redesign in CMOS transistors since the late 1960s

# Intel 4th Generation i7 Chip (Haswell) - 2013



- 4 core/8 threads for desktop and mobile solutions

- 22nm Hi-K+ process

- 3D tri-gate transistors

# Technology Outlook

| High Volume Manufacturing | 2008 | 2010 | 2012 | 2014 | 2016 | 2018 | 2020 | 2022 |
|---|---|---|---|---|---|---|---|---|
| Technology Node (nm) | 45 | 32 | 22 | 16 | 11 | 8 | 6 | 4 |
| Integration Capacity (BT) | 8 | 16 | 32 | 64 | 128 | 256 | 512 | 1024 |
| Delay Scaling | >0.7 | | | ~1? | | | | |
| Energy Scaling | ~0.5 | | | >0.5 | | | | |
| Transistors | Planar | | | 3G, FinFET | | | | |
| Variability | High | | | Extreme | | | | |
| ILD | ~3 | | | towards 2 | | | | |
| RC Delay | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| Metal Layers | 8-9 | | | 0.5 to 1 Layer per generation | | | | |

# 10 Nanometer Technology

- Nov. 15, 2012, Samsung unveiled a 64 gigabyte (GB) multimedia card (eMMC) based on 10 nm technology.

- April 11, 2013, Samsung announced it was mass-producing High-Performance 128-gigabit NAND Flash Memory with 10 nm and 20 nm technology.

- April 2015, TSMC announced that 10 nm production would begin at the end of 2016.

- May 23rd 2015, Samsung Electronics showed off a wafer of 10nm FinFET chips.

# Moore's Law

# Factors Contributing to Advancing Microprocessor Performance

- Shrinking Component Size

- Increasing Speed

- Reducing Circuit Resistance

- New Materials

# Factors Contributing to Advancing Microprocessor Performance

- RISC vs. CISC

- VLIW

- Multi-level Cache

- Parallelism & Pipelining

# Factors Contributing to Advancing Microprocessor Performance

- RISC vs. CISC

- VLIW

- Multi-level Cache

- Parallelism & Pipelining

- Multi-core Technology

# Multicore Craze

- For years, the trend was to make chips faster

  Today $\rightarrow$ 3 Ghz

- But the power required (Watts) and the heat generated is proportional to the frequency squared.

- Therefore, put more computers on the chip but run at slower speeds.

**Human Hair**

50μm

**Integrated Chip**

32nm

**Visible Light**

400nm

700nm

95.84 pm

104.45°

H   O   H

- How long can Moore's Law continue?

- What are the limits to this integrated circuit technology?

*"There are two constraints:*

    –   *The finite velocity of light*

    –   *The atomic nature of materials"*

- Stephen Hawking

# Miniaturized Data Storage at Atomic Scale



IBM researchers have stored and retrieved digital data from an array of just 12 atoms

New York Times, 1/12/12

# IBM's Chip Stacking Technology

# Single-phase, miniaturized convective cooling



Distributed return architecture with cross section showing inlet jets with neighboring drainage holes.

# Single-phase, miniaturized convective cooling



SEM section of two-level jet plate. Water flow is indicated by blue arrows.

# Interior Structure of 3-D chips

| 3-D Volatile Memory [Matrix Semiconductor] | 2-D Random-Access Memory [IBM 256-Megabit] | 3-D Logic Circuit [Lab Prototype] | 2-D Microprocessor [Advanced Micro Devices Athlon] |



NOT TO SCALE

**Monosilicon substrate** · **Insulators** · **Aluminum wires** · **Polysilicon** · **Tungsten plugs** · **Ion-doped silicon** · **Isolation oxides** · **Silicide**

# Intel's 3D Transistor                    2011

# Intel's 22nm 3D tri-gate transistor
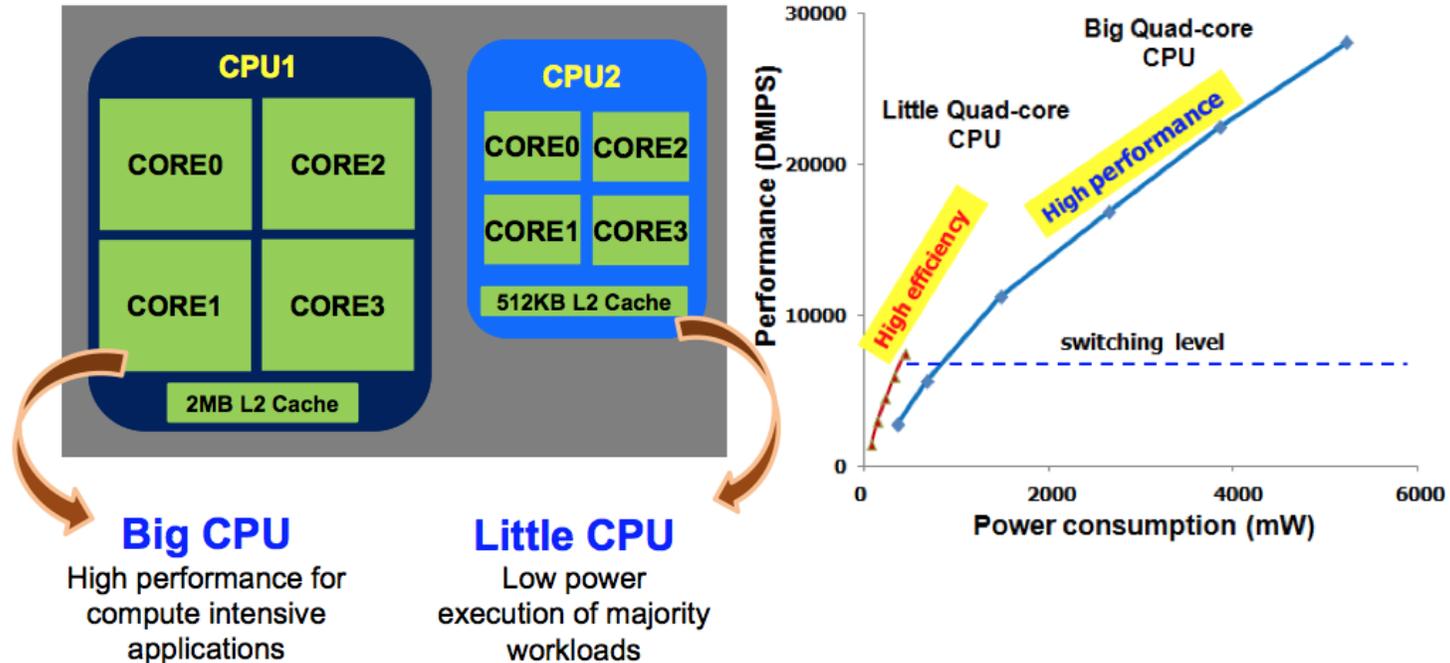
# Intel 5th Generation i7 Chip (Broadwell) - 2014



- 4 core/8 threads for desktop and mobile solutions

- 14nm Hi-K+ process

- 3D tri-gate transistors

- Predicted 30% improvement in power consumption

# Samsung Galaxy Exynos 5 Octa

# Samsung Galaxy S5 — April 2014



- Samsung Exynos 5 Octa 5422, (8 cores)

- Heterogeneous ARM architecture (2.1 GHz Cortex-A15 quad, 1.5 GHz Cortex-A7 quad core)

- 2 GB DDR3 Ram

- 64 GB storage

# Samsung Galaxy S6                    2016

- Screen Size: 5.1 inch screen

- Screen Resolution: 1440 x 2560 pixels (~577 ppi)

- Battery Life: 23 hour talk time

- Feature: 4G LTE

- Operating System: Android

- Camera Resolution: 16 MP

- Weight: 4.9 ounce

- Memory: 32 GB

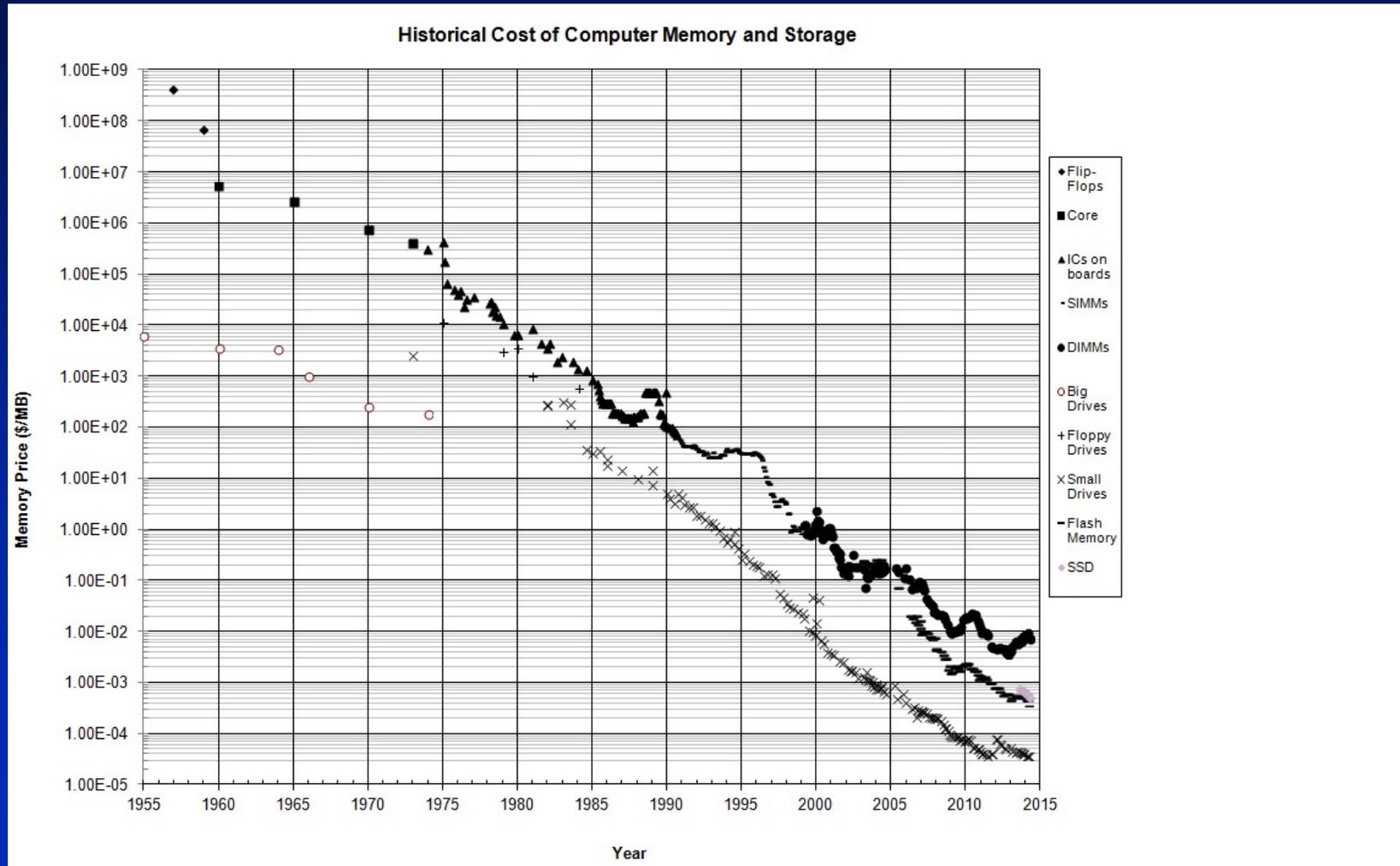- 64-bit Exynos 7420 Processor

- 14 nm FinFET technology

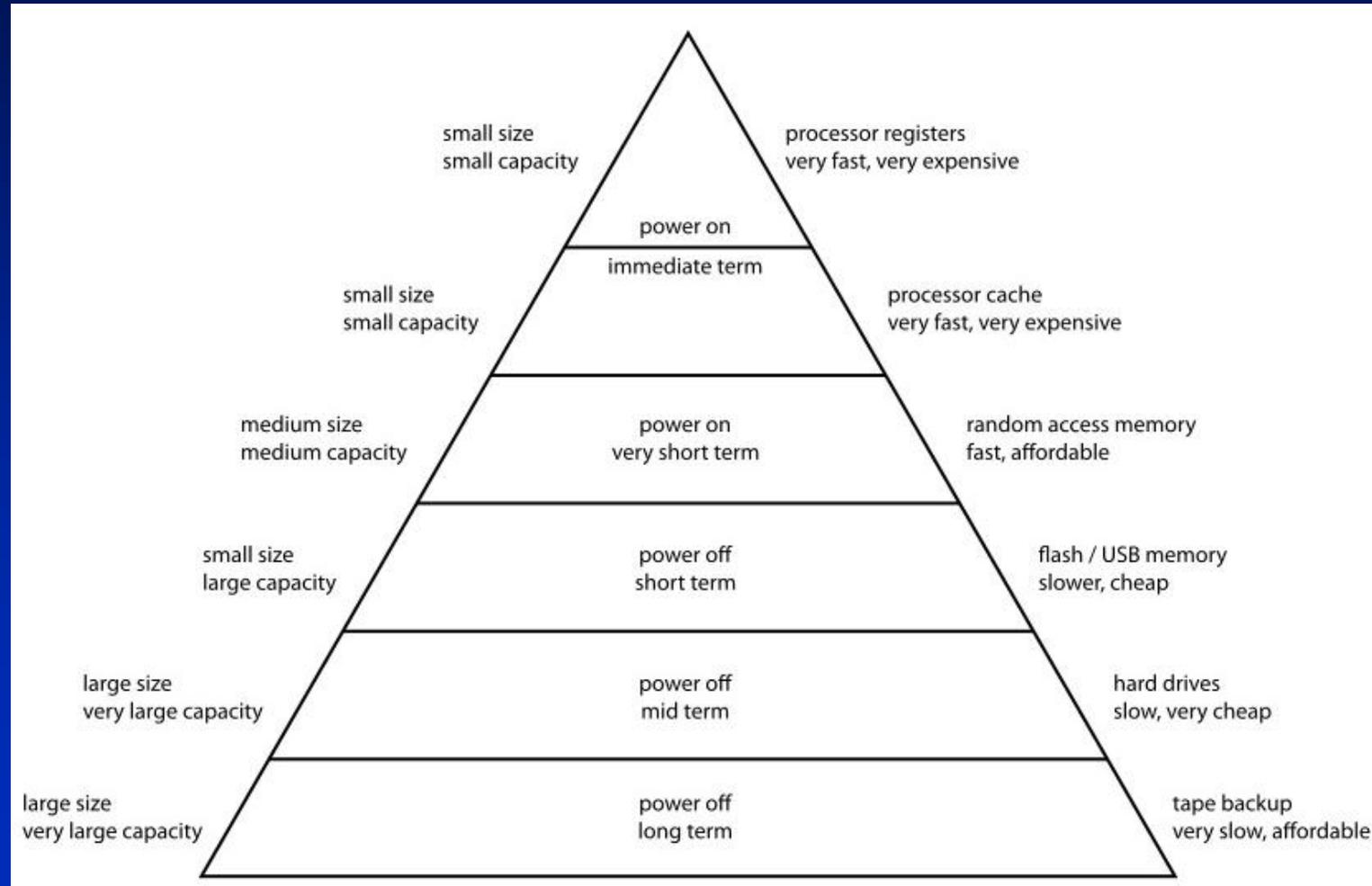# Intel's $7B Investment                    2017

# Cost of Computer Memory and Storage
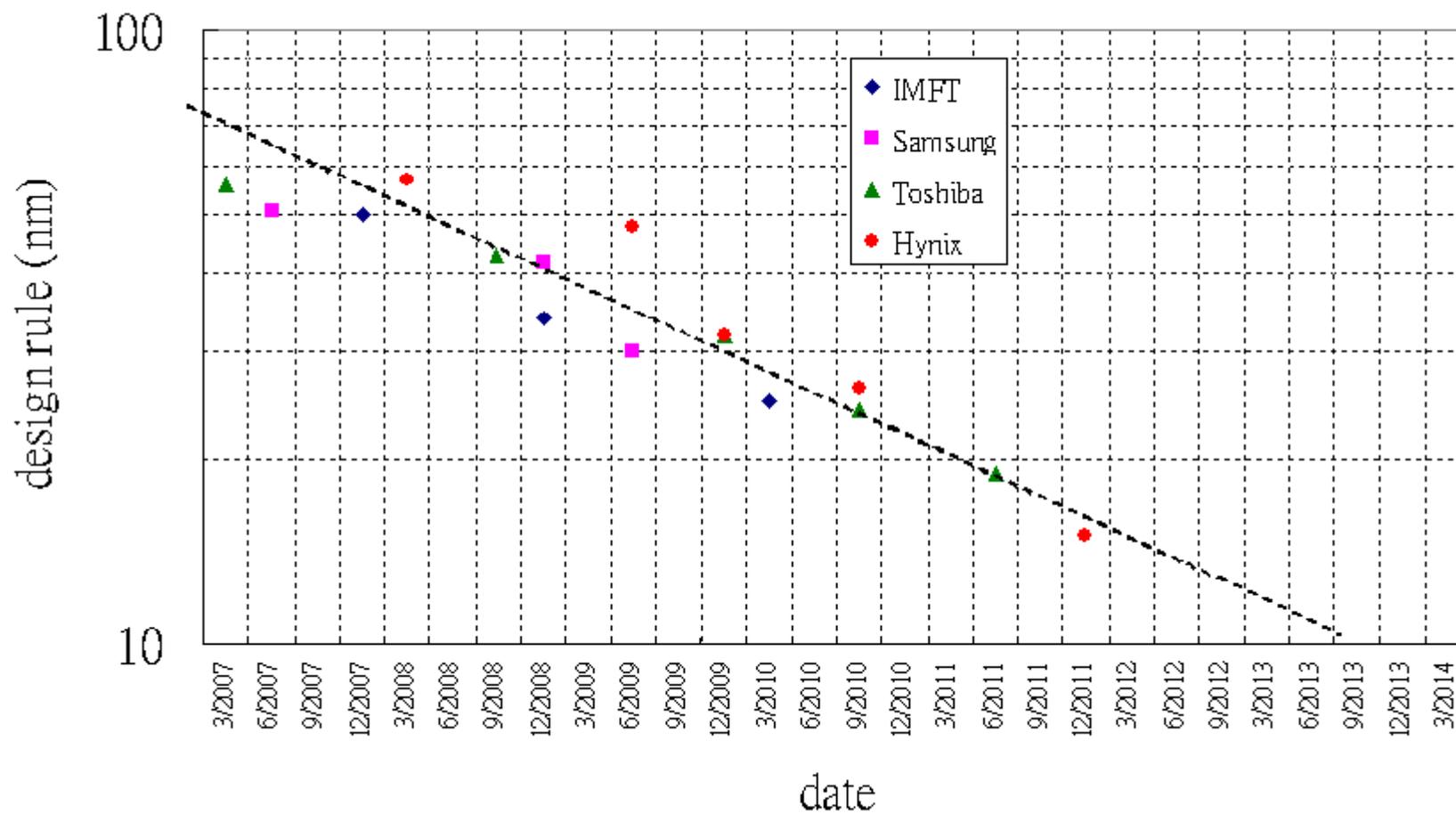


Historical Cost of Computer Memory and Storage

# Computer Memory Hierarchy
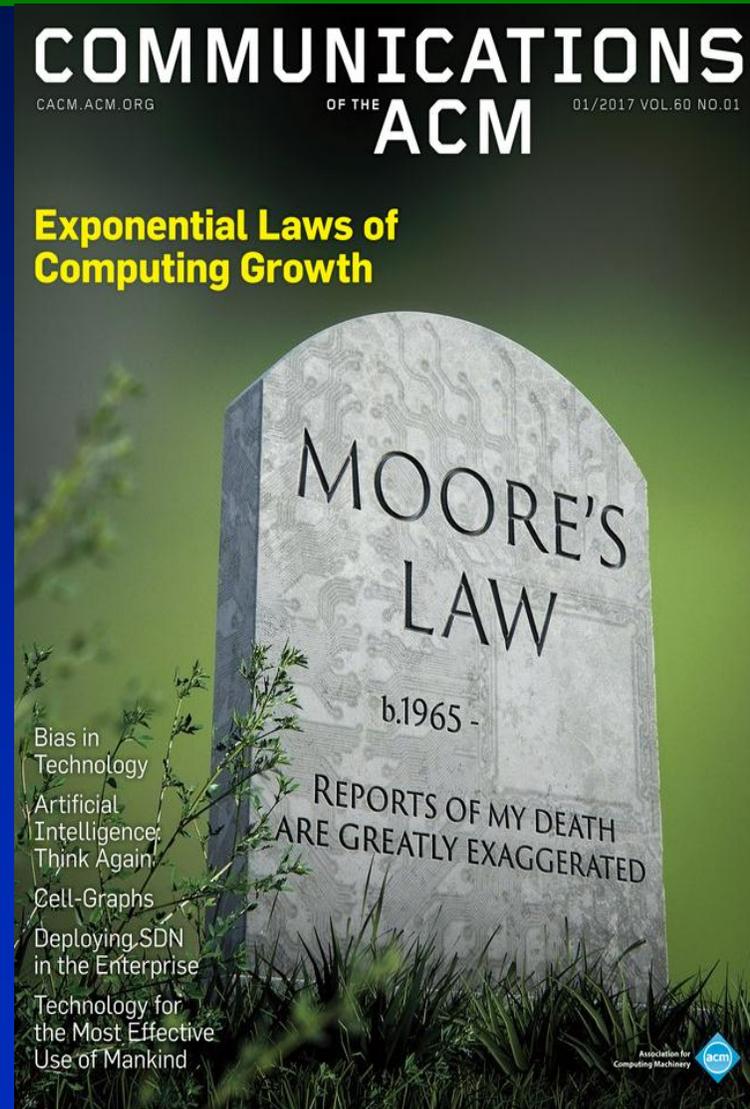
# Flash Scalability

"Every economic era is based on a key <u>abundance</u> and a <u>key</u> scarcity."

*George Gilder,*
*Forbes ASAP, 1992*

What are the key scarcities?

# Exponential Laws of Computing Growth



Dennings and Lewis, Jan. 2016

# End