# Computer Organization and Assembly Language
# CS / EE 320
# Spring 2024

Lecture 1

Shahid Masud

# Difference between Computer **Organization** and Computer **Architecture**

Architecture: View from outside, high level specs (Programs looking towards CPU Instructions)

Organization: View from Inside (CPU looking outwards towards buses, memory and peripherals)

Performance Graphs from the report of National Academy, USA

'The Future of Computer Performance

Game Over or Next Lavel
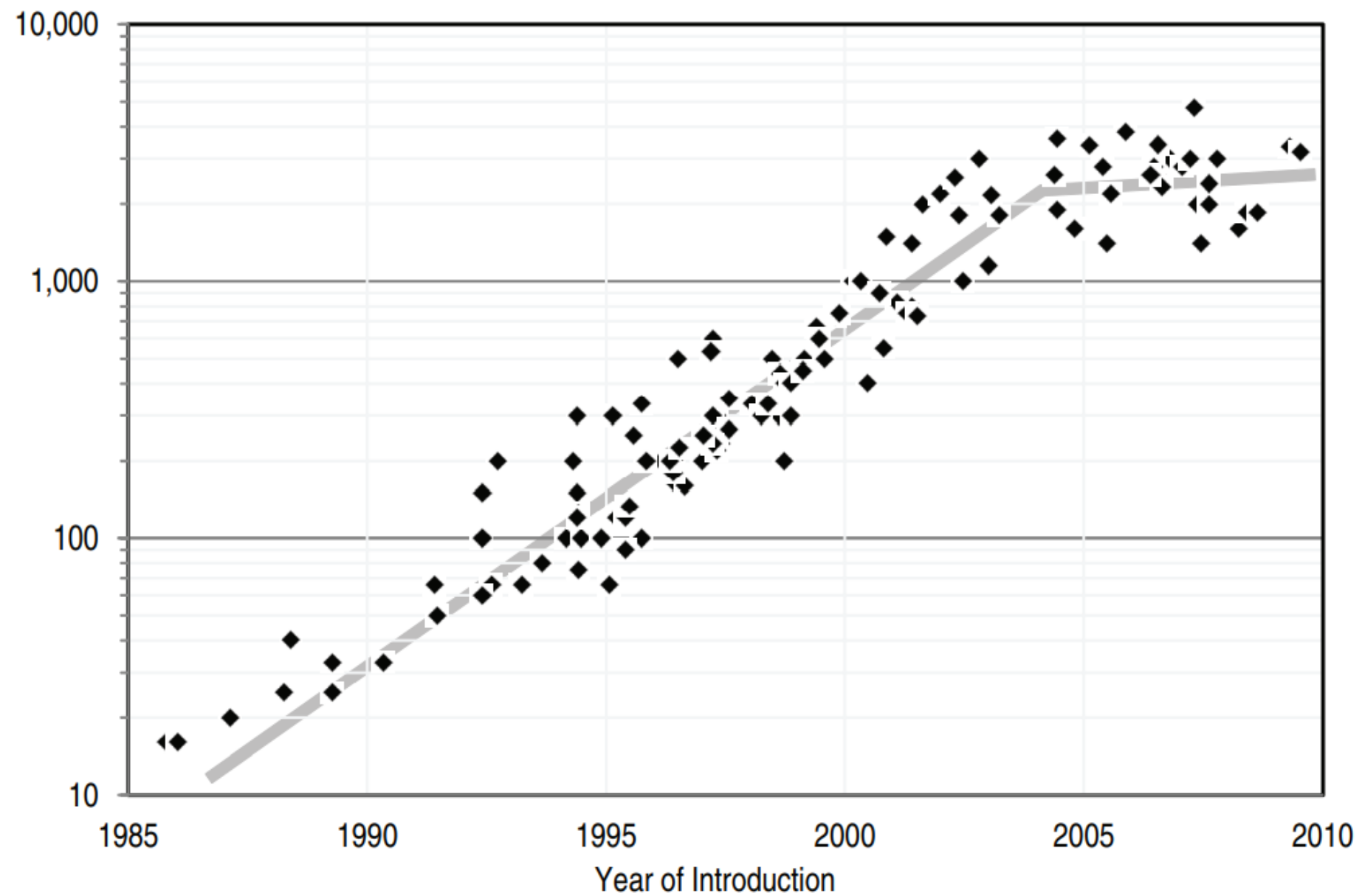
Ed. Samuel H Fuller

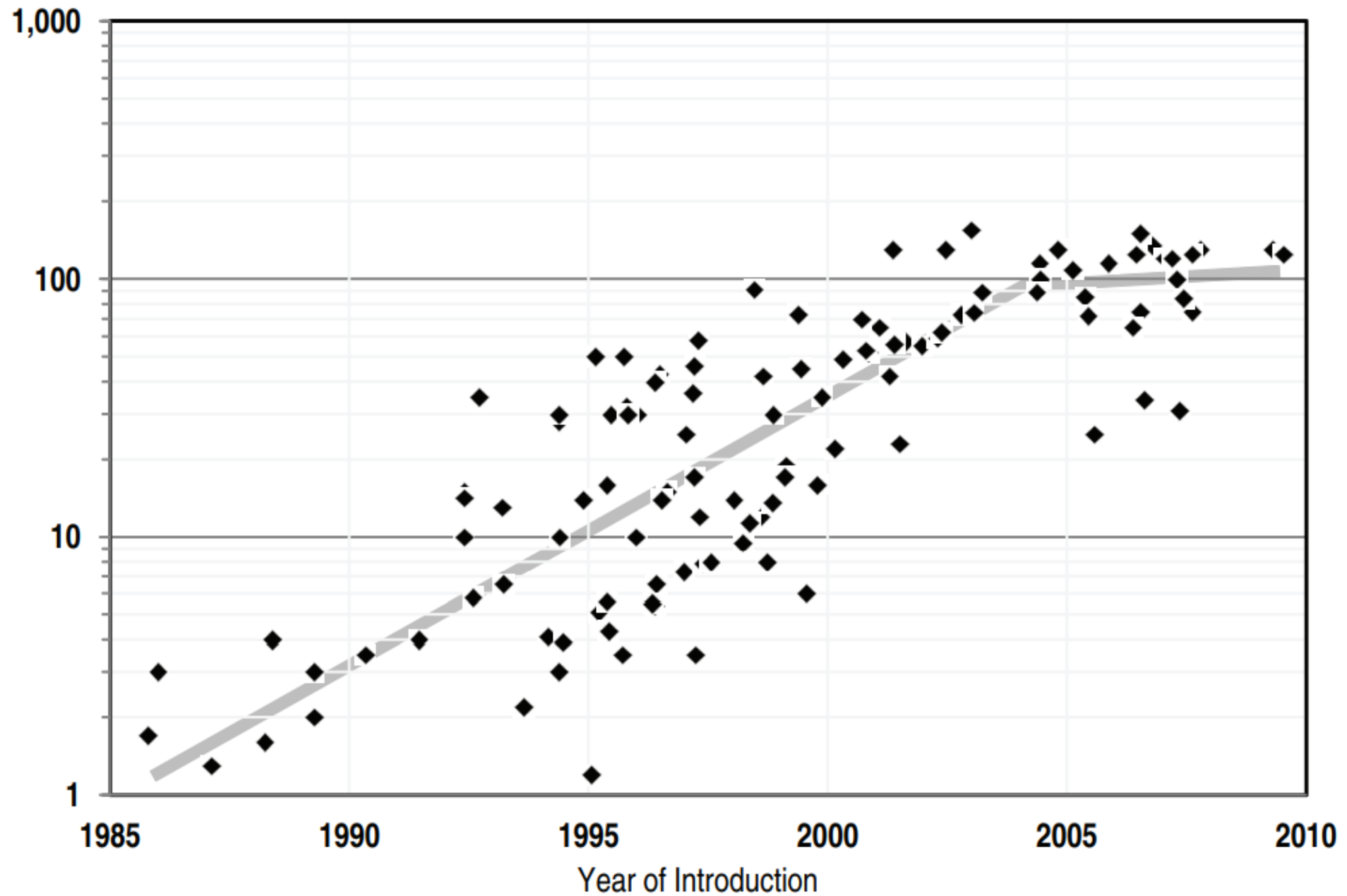FIGURE 3.3 Microprocessor-clock frequency (MHz) over time (1985-2010).

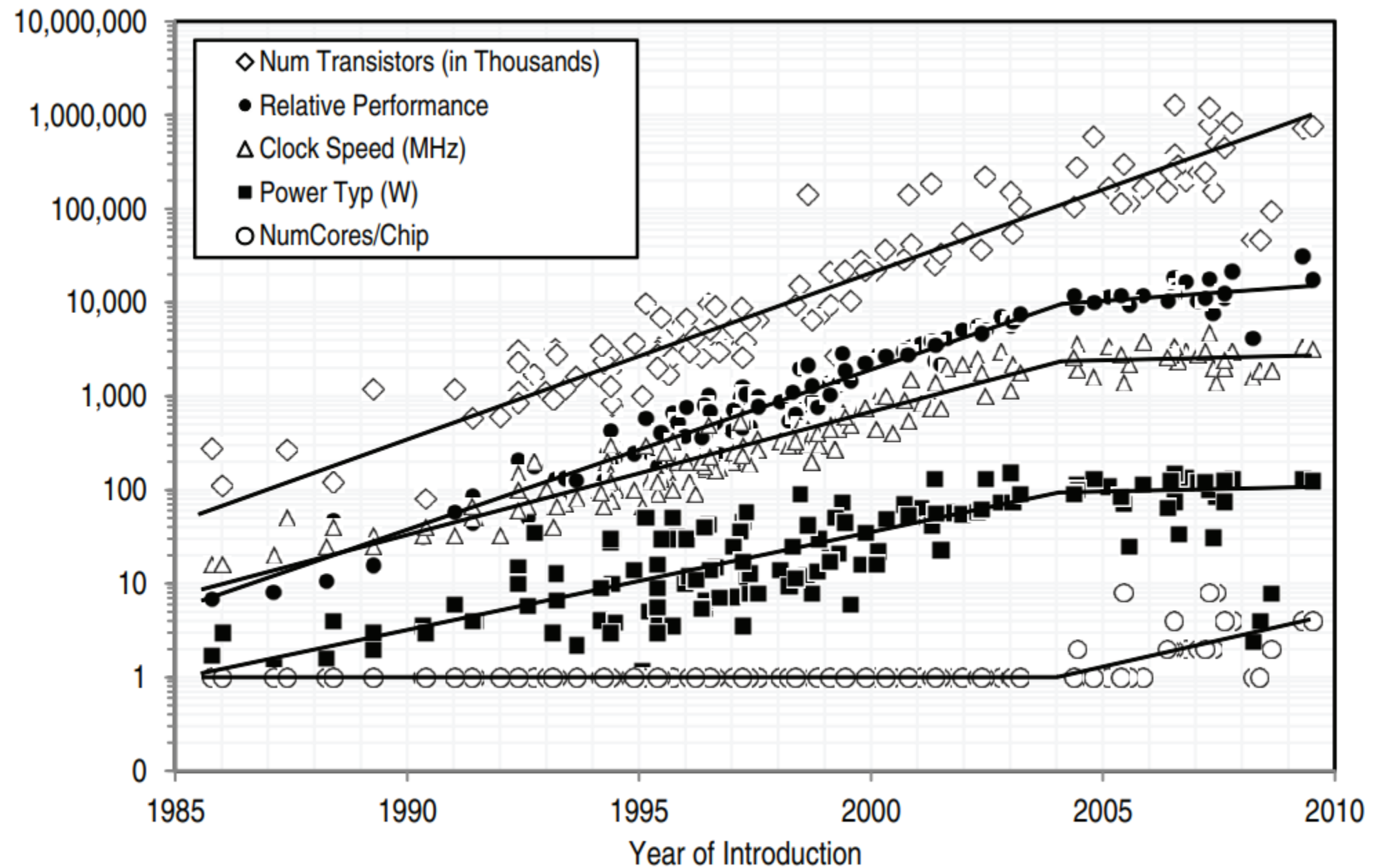FIGURE 3.1 Microprocessor power dissipation (watts) over time (1985-2010).

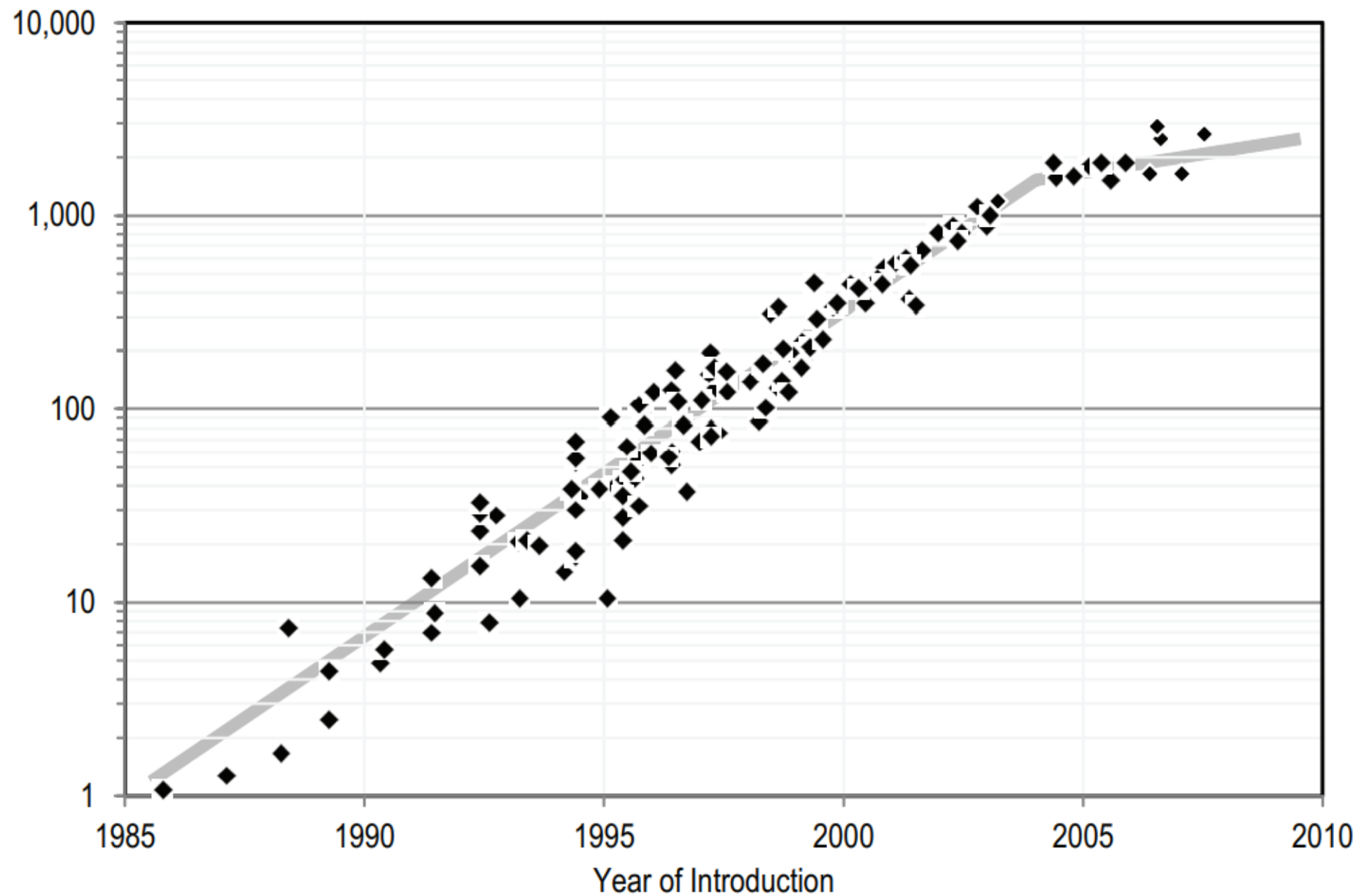FIGURE 2.1 Transistors, frequency, power, performance, and cores over time

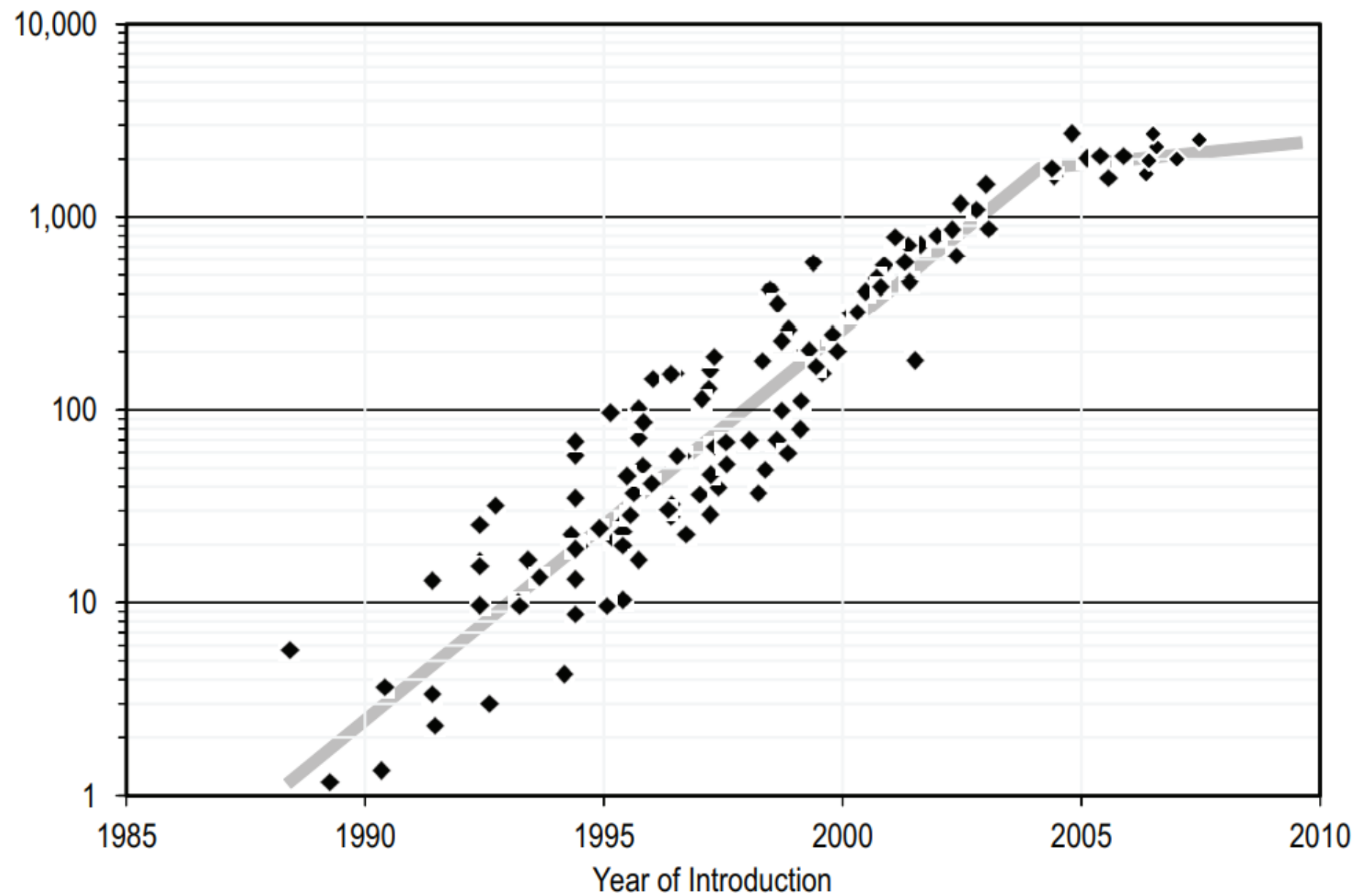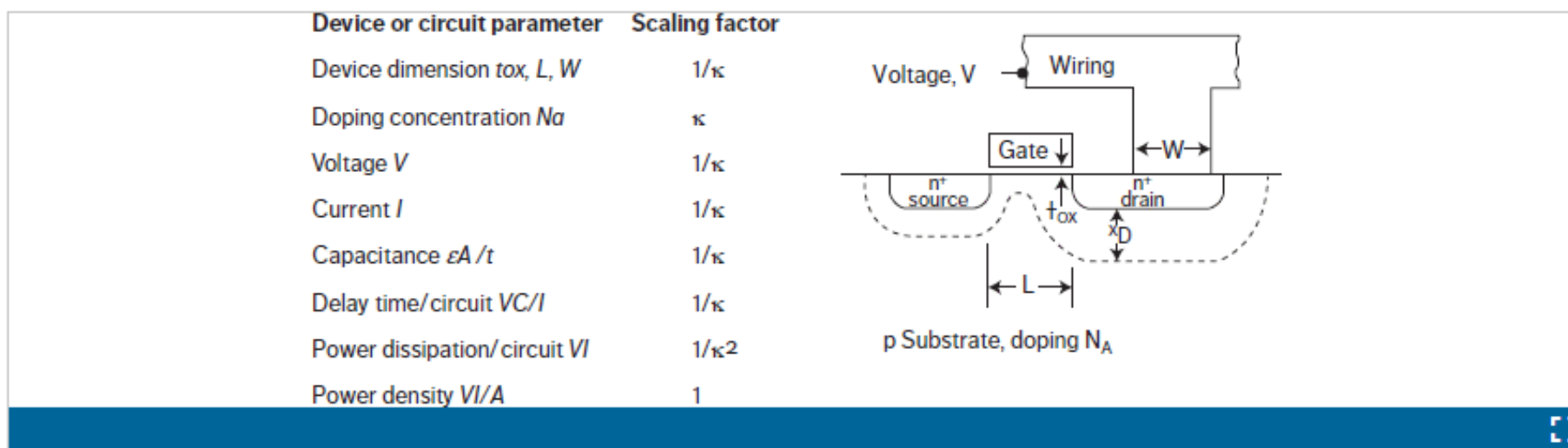FIGURE A.1 Integer application performance (SPECint2000) over time (1985-2010).

FIGURE A.2 Floating-point application performance (SPECfp2000) over time (1985-2010).
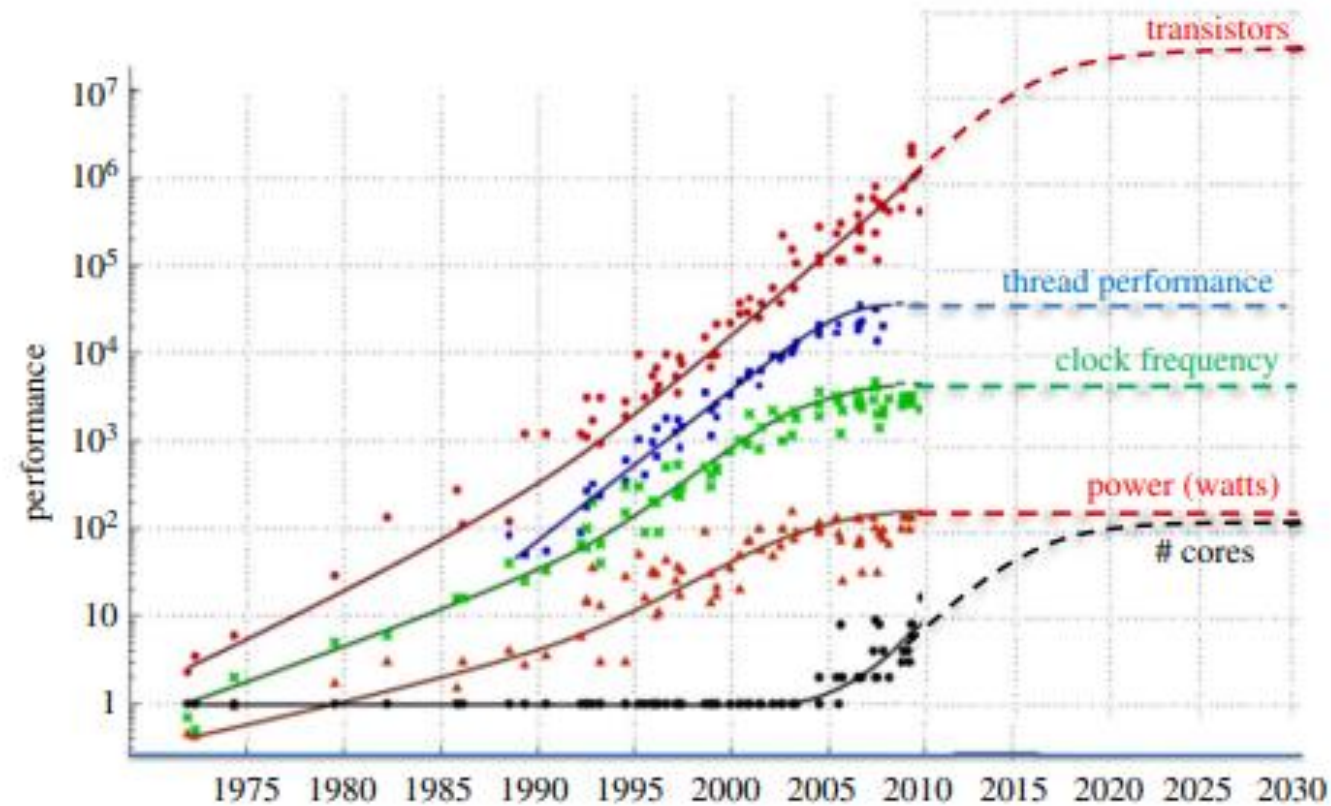
# Dennard Scaling

Robert Dennard and colleagues described in 1974 a scaling methodology for metal-oxide-semiconductor field-effect transistors (MOSFETs) that would deliver consistent improvements in transistor area, performance, and power reduction.[3] The methodology called for the scaling of transistor gate length, gate width, gate oxide thickness, and supply voltage all by the same scale factor, and increasing channel doping by the inverse of the same scale factor (see Figure 1). The result would be transistors with smaller area, higher drive current (higher performance), and lower parasitic capacitance (lower active power). This method for scaling MOSFET transistors is generally referred to as "classic" or "traditional" scaling and was very successfully used by the industry up until the 130-nm generation in the early 2000s.

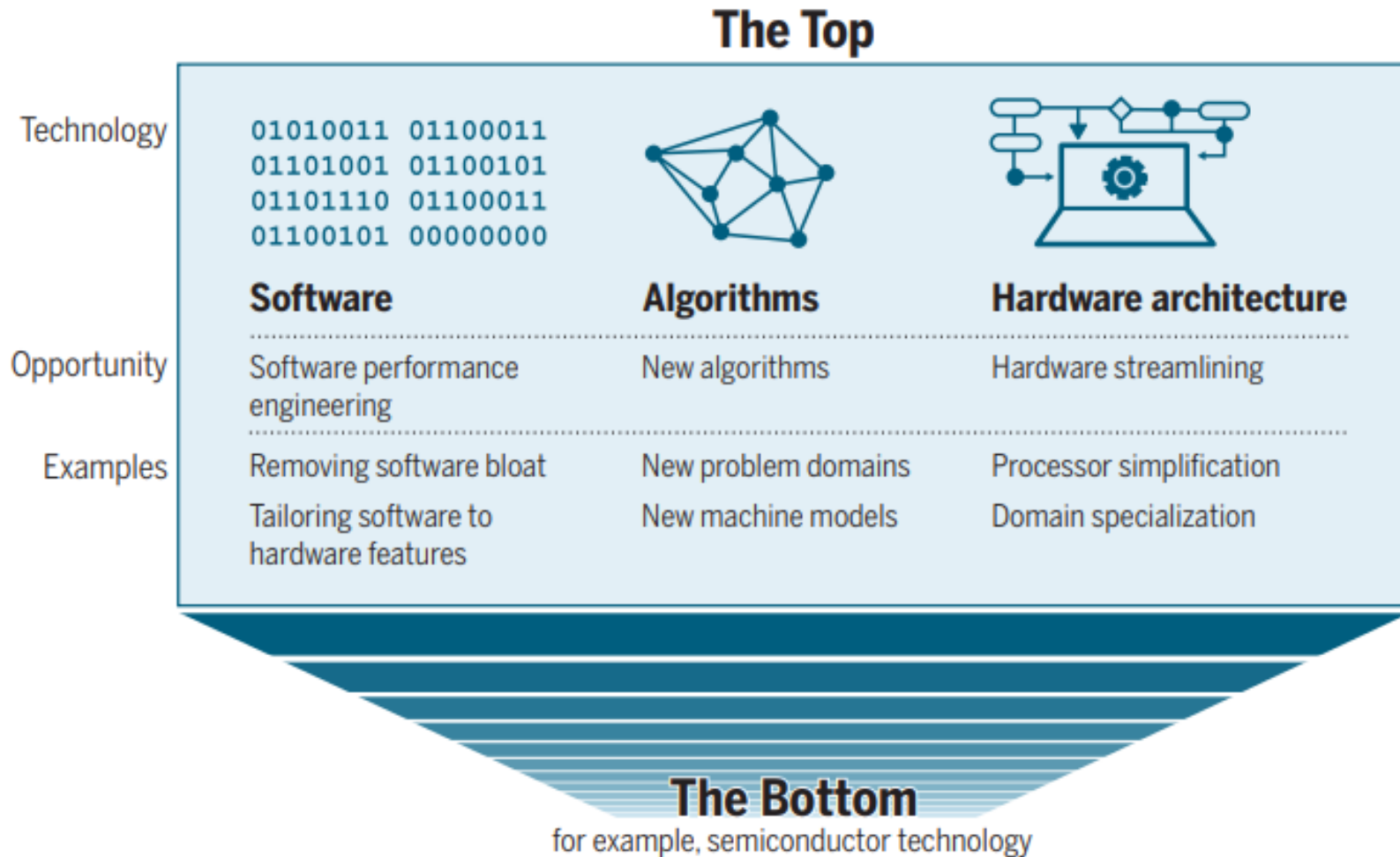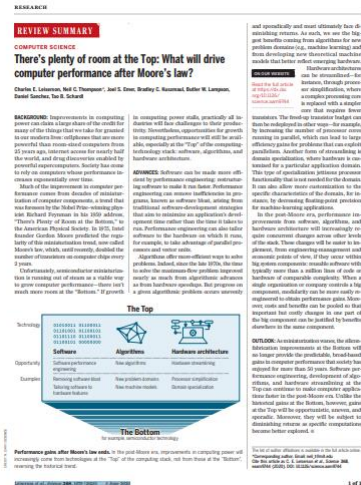| Device or circuit parameter | Scaling factor |
|---|---|
| Device dimension $t_{ox}$, $L$, $W$ | $1/\kappa$ |
| Doping concentration $N_a$ | $\kappa$ |
| Voltage $V$ | $1/\kappa$ |
| Current $I$ | $1/\kappa$ |
| Capacitance $\varepsilon A/t$ | $1/\kappa$ |
| Delay time/circuit $VC/I$ | $1/\kappa$ |
| Power dissipation/circuit $VI$ | $1/\kappa^2$ |
| Power density $VI/A$ | $1$ |



**Figure 1.**
Traditional MOSFET scaling as described by robert dennard.

# Depiction of Dennard Scaling Effects



**Figure 2.** Sources of computing performance have been challenged by the end of Dennard scaling in 2004. All additional approaches to further performance improvements end in approximately 2025 due to the end of the roadmap for improvements to semiconductor lithography. Figure from Kunle Olukotun, Lance Hammond, Herb Sutter, Mark Horowitz and extended by John Shalf. (Online version in colour.)

**The Top**

Technology

| Software | Algorithms | Hardware architecture |
| --- | --- | --- |
| 01010011 01100011<br>01101001 01100101<br>01101110 01100011<br>01100101 00000000 | | |

| | Software | Algorithms | Hardware architecture |
| --- | --- | --- | --- |
| Opportunity | Software performance engineering | New algorithms | Hardware streamlining |
| Examples | Removing software bloat | New problem domains | Processor simplification |
| | Tailoring software to hardware features | New machine models | Domain specialization |

**The Bottom**
for example, semiconductor technology

**Performance gains after Moore's law ends.** In the post-Moore era, improvements in computing power will increasingly come from technologies at the "Top" of the computing stack, not from those at the "Bottom", reversing the historical trend.

Ref: Leiserson et al., Science 368, 1079 (2020) 5 June 2020

# Example of Performance Gains through Architecture

**Table 1. Speedups from performance engineering a program that multiplies two 4096-by-4096 matrices.** Each version represents a successive refinement of the original Python code. "Running time" is the running time of the version. "GFLOPS" is the billions of 64-bit floating-point operations per second that the version executes. "Absolute speedup" is time relative to Python, and "relative speedup," which we show with an additional digit of precision, is time relative to the preceding line. "Fraction of peak" is GFLOPS relative to the computer's peak 835 GFLOPS. See Methods for more details.

| Version | Implementation | Running time (s) | GFLOPS | Absolute speedup | Relative speedup | Fraction of peak (%) |
|---|---|---|---|---|---|---|
| 1 | Python | 25,552.48 | 0.005 | 1 | — | 0.00 |
| 2 | Java | 2,372.68 | 0.058 | 11 | 10.8 | 0.01 |
| 3 | C | 542.67 | 0.253 | 47 | 4.4 | 0.03 |
| 4 | Parallel loops | 69.80 | 1.969 | 366 | 7.8 | 0.24 |
| 5 | Parallel divide and conquer | 3.80 | 36.180 | 6,727 | 18.4 | 4.33 |
| 6 | plus vectorization | 1.10 | 124.914 | 23,224 | 3.5 | 14.96 |
| 7 | plus AVX intrinsics | 0.41 | 337.812 | 62,806 | 2.7 | 40.45 |

AVX = Intel Advanced Vector Extension

# Look at some recent research papers for inspiration

# Topics

1. Discussion on Moore's Law, through paper 'MORE THAN MOORE' by M. Mitchell Waldrop, published In Nature, February 2016

2. Computer Performance Graphs from National Academy, USA, report

3. Introduction to the specialization area of 'Computer Architecture' through paper by Hennessy and Patterson, 'A New Golden Age for Computer Architecture', published in Communications of the ACM, February 2019, pages 48 to 60.

    Important Directions for Computers of the future:

    a. Moore's Law failing to keep up

    b. Domain Specific Architectures

    c. Enhanced security features inside microprocessors

    d. Open Instruction Sets and Extension of Instruction Sets through Customized Accelerators

    e. Agile Hardware Design for Microprocessors

    f. Combination of Domain Specific Languages and Domain Specific Architectures

4. Intel Processor Timeline

5. Reviewed course outline including detailed topics and grading breakup of lectures and labs.

# Video of Turing Lecture by Patterson, 2019

David Patterson - A New Golden Age for Computer Architecture: History, Challenges and Opportunities – YouTube