

Chapter 1

Fundamentals of Quantitative Design and Analysis

Adapted by Prof. Gheith Abandah

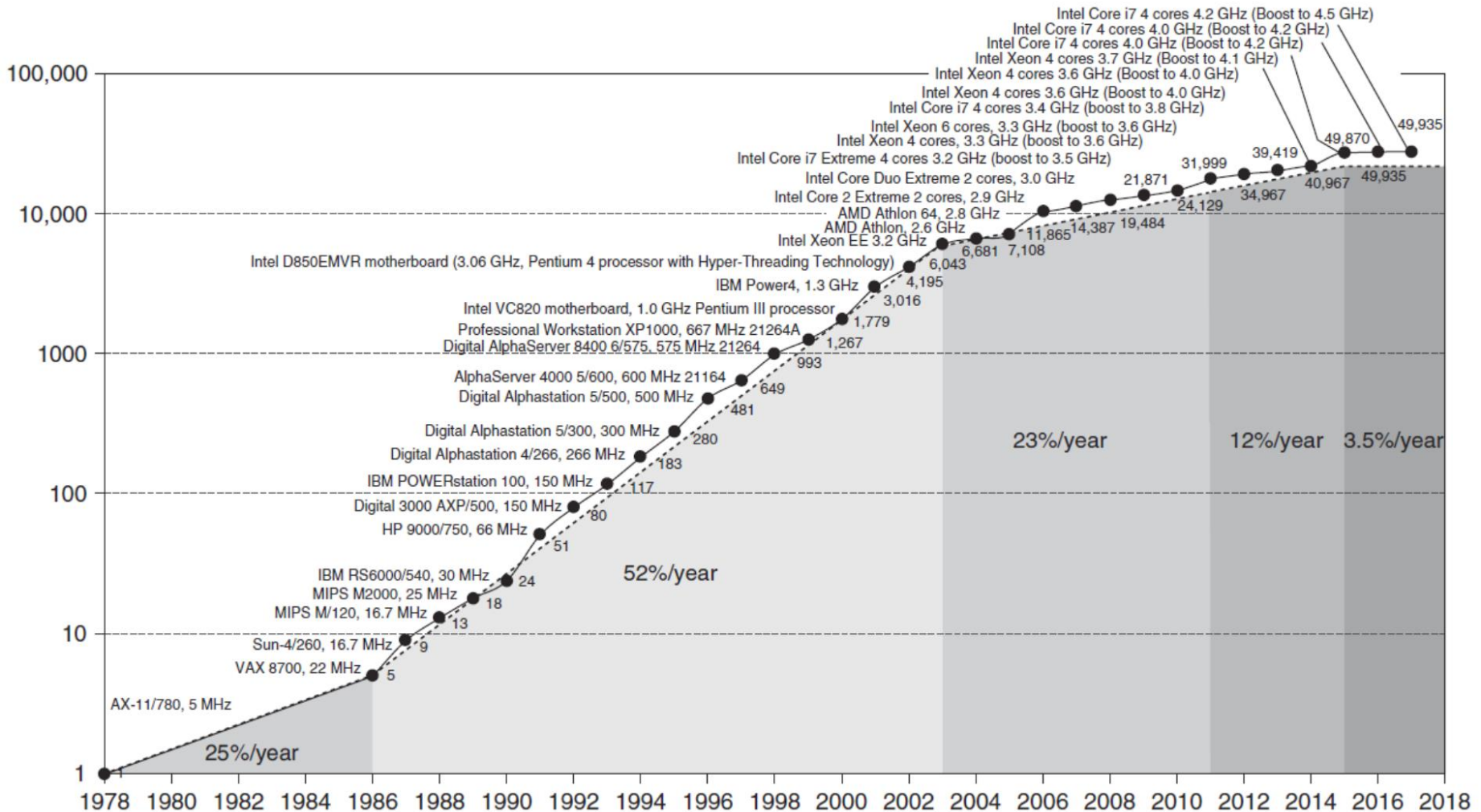
Contents

- Introduction
- Classes of Computers
- Defining Computer Architecture
- Trends in Technology
- Trends in Power and Energy in Integrated Circuits
- Trends in Cost
- Dependability
- Measuring, Reporting, and Summarizing Performance
- Quantitative Principles of Computer Design
- Fallacies and Pitfalls

Computer Technology

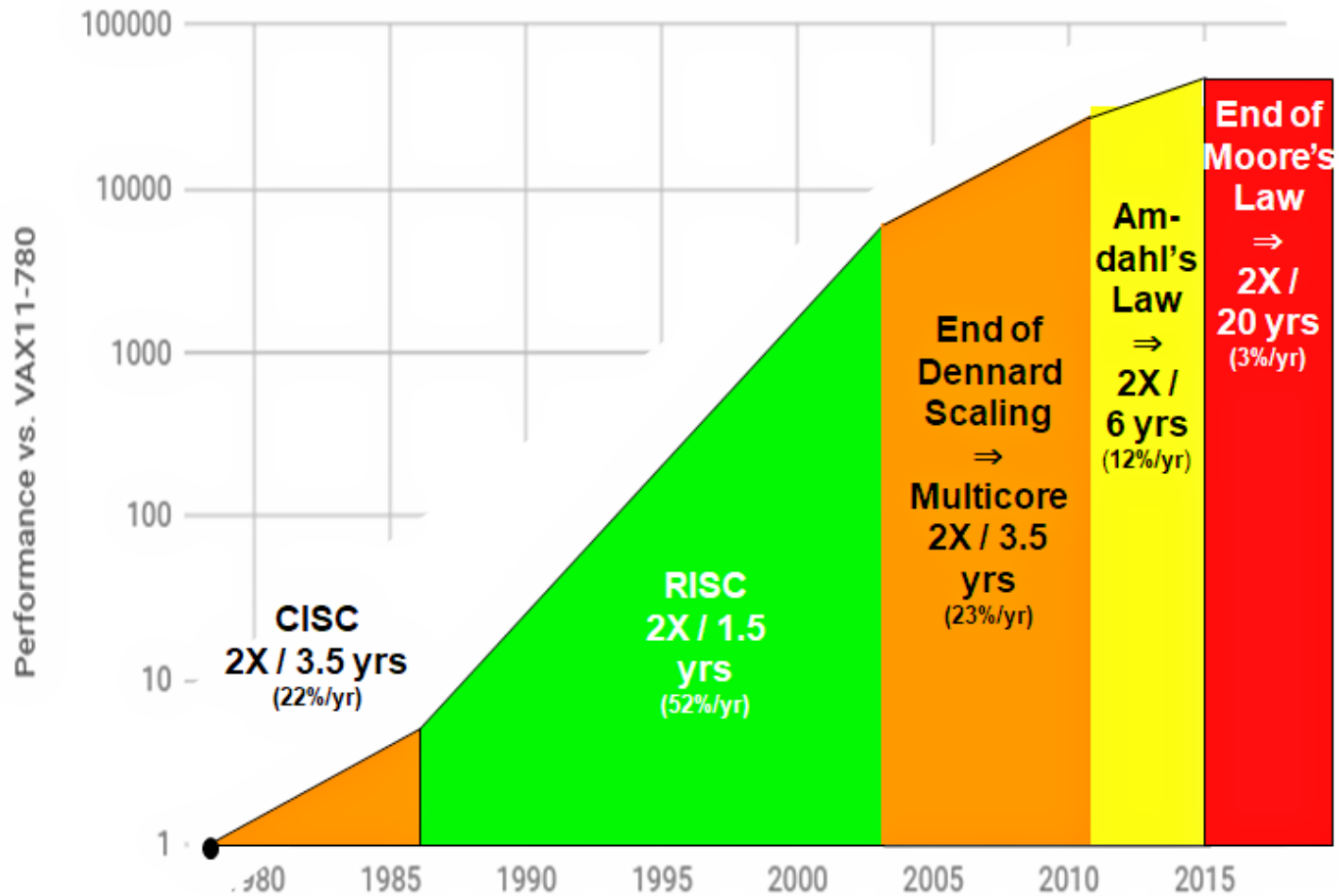
- Performance improvements:
 - Improvements in semiconductor technology
 - Feature size, clock speed
 - Improvements in computer architectures
 - Enabled by HLL compilers, UNIX
 - Lead to RISC architectures
- Together have enabled:
 - Lightweight computers
 - Productivity-based managed/interpreted programming languages

Single Processor Performance



Single Processor Performance

40 years of Processor Performance



Current Trends in Architecture

- Cannot continue to leverage instruction-level parallelism (ILP)
 - Single processor performance improvement ended in 2003
- New models for performance:
 - Data-level parallelism (DLP)
 - Thread-level parallelism (TLP)
 - Request-level parallelism (RLP)
- These require explicit restructuring of the application

Contents

- Introduction
- Classes of Computers
- Defining Computer Architecture
- Trends in Technology
- Trends in Power and Energy in Integrated Circuits
- Trends in Cost
- Dependability
- Measuring, Reporting, and Summarizing Performance
- Quantitative Principles of Computer Design
- Fallacies and Pitfalls

Classes of Computers

- Personal Mobile Device (PMD)
 - E.g. smart phones, tablet computers
 - Emphasis on energy efficiency and real-time
- Desktop Computing
 - Emphasis on price-performance
- Servers
 - Emphasis on availability, scalability, throughput
- Clusters / Warehouse Scale Computers
 - Used for “Software as a Service (SaaS)”
 - Emphasis on availability and price-performance
 - Sub-class: Supercomputers, emphasis: floating-point performance and fast internal networks
- Internet of Things/Embedded Computers
 - Emphasis: price

Parallelism

- Classes of parallelism in applications:
 - Data-Level Parallelism (DLP)
 - Task-Level Parallelism (TLP)
- Classes of architectural parallelism:
 - Instruction-Level Parallelism (ILP)
 - Vector architectures/Graphic Processor Units (GPUs)
 - Thread-Level Parallelism
 - Request-Level Parallelism

Flynn's Taxonomy

- Single instruction stream, single data stream (SISD)
- Single instruction stream, multiple data streams (SIMD)
 - Vector architectures
 - Multimedia extensions
 - Graphics processor units
- Multiple instruction streams, single data stream (MISD)
 - No commercial implementation
- Multiple instruction streams, multiple data streams (MIMD)
 - Tightly-coupled MIMD
 - Loosely-coupled MIMD

Contents

- Introduction
- Classes of Computers
- Defining Computer Architecture
- Trends in Technology
- Trends in Power and Energy in Integrated Circuits
- Trends in Cost
- Dependability
- Measuring, Reporting, and Summarizing Performance
- Quantitative Principles of Computer Design
- Fallacies and Pitfalls

Defining Computer Architecture

- “Old” view of computer architecture:
 - Instruction Set Architecture (ISA) design
 - i.e. decisions regarding:
 - registers, memory addressing, addressing modes, instruction operands, available operations, control flow instructions, instruction encoding
- “Real” computer architecture:
 - Specific requirements of the target machine
 - Design to maximize performance within constraints: cost, power, and availability
 - Includes ISA, microarchitecture, hardware

Instruction Set Architecture

- Class of ISA
 - General-purpose registers
 - Register-memory vs load-store
- RISC-V registers
 - 32 g.p., 32 f.p.

Register	Name	Use	Saver
x0	zero	constant 0	n/a
x1	ra	return addr	caller
x2	sp	stack ptr	callee
x3	gp	gbl ptr	
x4	tp	thread ptr	
x5-x7	t0-t2	temporaries	caller
x8	s0/fp	saved/ frame ptr	callee
x9	s1	saved	callee
x10-x17	a0-a7	arguments	caller
x18-x27	s2-s11	saved	callee
x28-x31	t3-t6	temporaries	caller
f0-f7	ft0-ft7	FP temps	caller
f8-f9	fs0-fs1	FP saved	callee
f10-f17	fa0-fa7	FP arguments	callee
f18-f27	fs2-fs21	FP saved	callee
f28-f31	ft8-ft11	FP temps	caller

RISC-V ISA

- Memory addressing
 - byte addressed, aligned accesses faster
- Addressing modes
 - Register, immediate, displacement (base+offset)
- Types and size of operands
 - 8-bit, 32-bit, 64-bit

RISC-V ISA

- Operations

- Data transfer `ld x5, 0(x6)`
- Arithmetic `add x5, x6, x7`
- Logical `add x5, x6, x7`
- Floating point `fadd.d f5, f6, f7`

- Control flow instructions

- Use content of registers `beq x5, x6, Label`
- Return address in register `jal x1, F1`
`jalr x0, 0(x1)`

- Encoding

- Fixed, 32 bits

Contents

- Introduction
- Classes of Computers
- Defining Computer Architecture
- Trends in Technology
- Trends in Power and Energy in Integrated Circuits
- Trends in Cost
- Dependability
- Measuring, Reporting, and Summarizing Performance
- Quantitative Principles of Computer Design
- Fallacies and Pitfalls

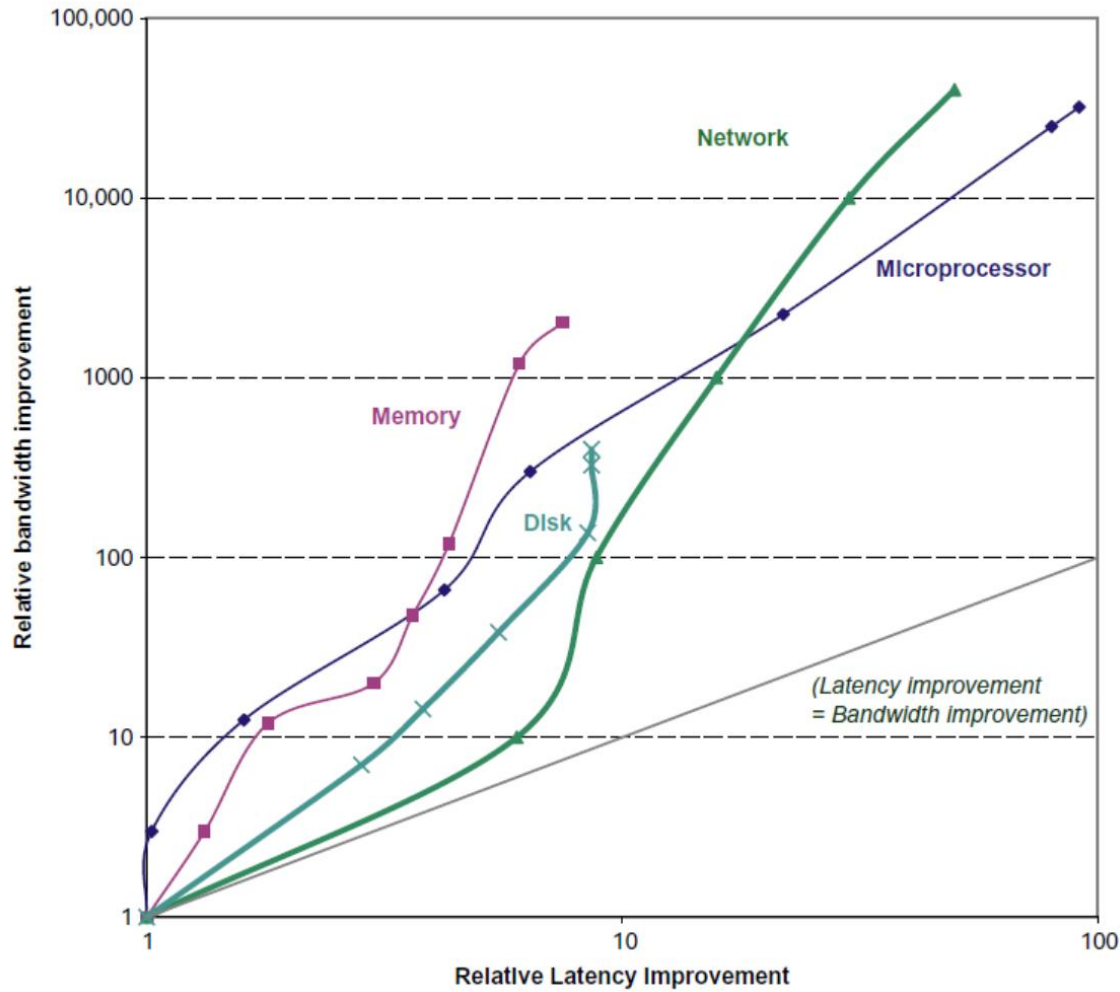
Trends in Technology

- Integrated circuit technology (Moore's Law)
 - Transistor density: 35%/year
 - Die size: 10-20%/year
 - Integration overall: 40-55%/year
- DRAM capacity: 25-40%/year (slowing)
 - 8 Gb (2014), 16 Gb (2019), possibly no 32 Gb
- Flash capacity: 50-60%/year
 - 8-10X cheaper/bit than DRAM
- Magnetic disk capacity: recently slowed to 5%/year
 - Density increases may no longer be possible, maybe increase from 7 to 9 platters
 - 8-10X cheaper/bit than Flash
 - 200-300X cheaper/bit than DRAM

Bandwidth and Latency

- For milestones in 25-40 years
- Bandwidth or throughput
 - Total work done in a given time
 - 32,000-40,000X improvement for processors
 - 300-1200X improvement for memory and disks
- Latency or response time
 - Time between start and completion of an event
 - 50-90X improvement for processors
 - 6-8X improvement for memory and disks

Bandwidth and Latency



Log-log plot of bandwidth and latency milestones

Transistors and Wires

- Feature size
 - Minimum size of transistor or wire in x or y dimension
 - 10 microns in 1971 to .011 microns (11 nm) in 2017
 - Transistor performance scales linearly
 - Wire delay does not improve with feature size!
 - Integration density scales quadratically

Contents

- Introduction
- Classes of Computers
- Defining Computer Architecture
- Trends in Technology
- Trends in Power and Energy in Integrated Circuits
- Trends in Cost
- Dependability
- Measuring, Reporting, and Summarizing Performance
- Quantitative Principles of Computer Design
- Fallacies and Pitfalls

Power and Energy

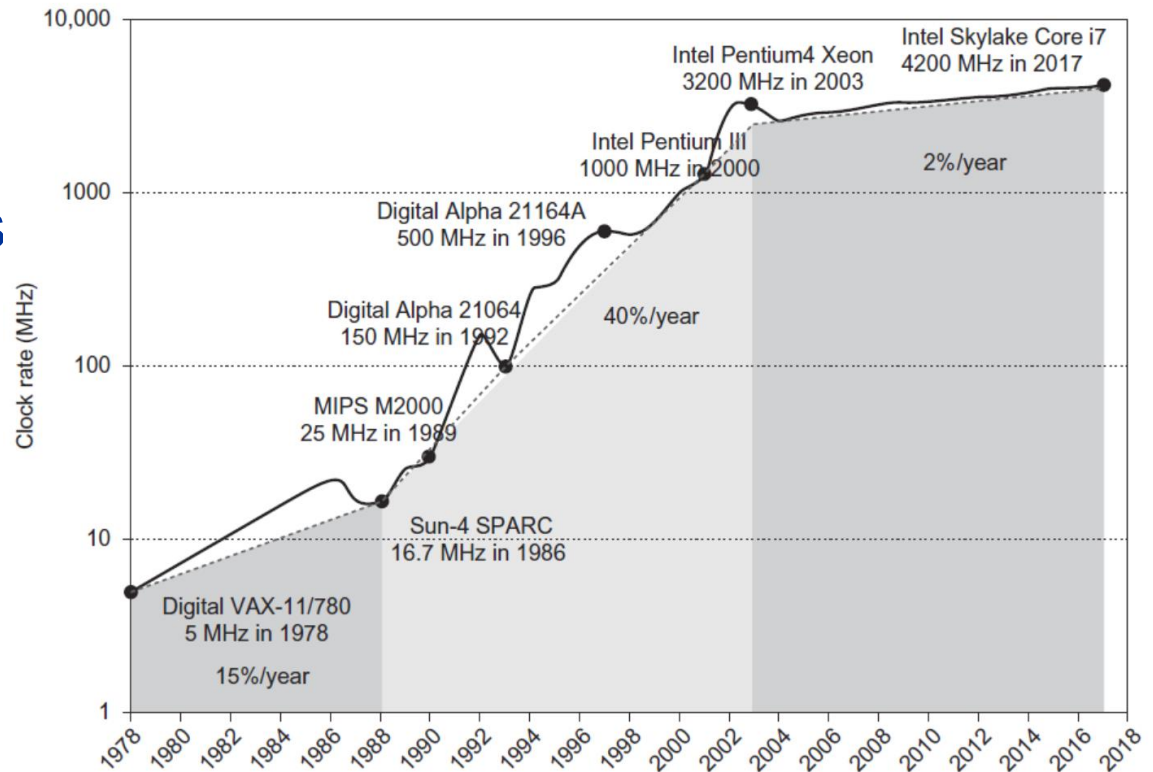
- Problem: Get power in, get power out
- Thermal Design Power (TDP)
 - Characterizes sustained power consumption
 - Used as target for power supply and cooling system
 - Lower than peak power (1.5X higher), higher than average power consumption
- Clock rate can be reduced dynamically to limit power consumption
- Energy per task is often a better measurement

Dynamic Energy and Power

- Dynamic energy
 - Transistor switch from 0 -> 1 or 1 -> 0
 - $\frac{1}{2} \times \text{Capacitive load} \times \text{Voltage}^2$
- Dynamic power
 - $\frac{1}{2} \times \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency switched}$
- Reducing clock rate reduces power, not energy

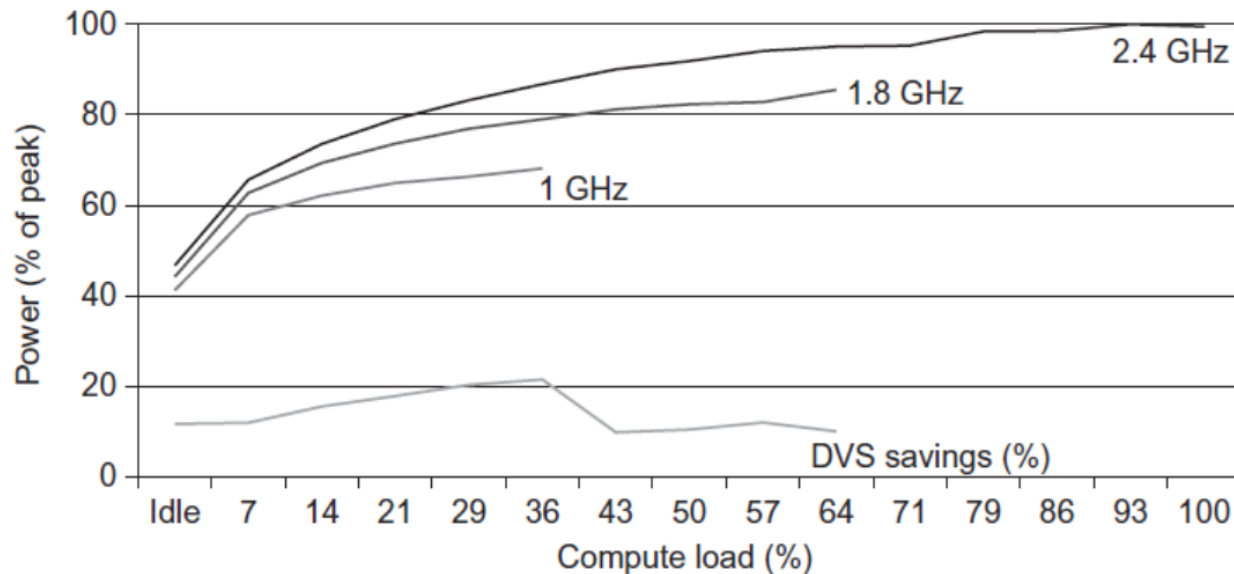
Power

- Intel 80386 consumed ~ 2 W
- 3.3 GHz Intel Core i7 consumes 130 W
- Heat must be dissipated from 1.5 x 1.5 cm chip
- This is the limit of what can be cooled by air



Reducing Power

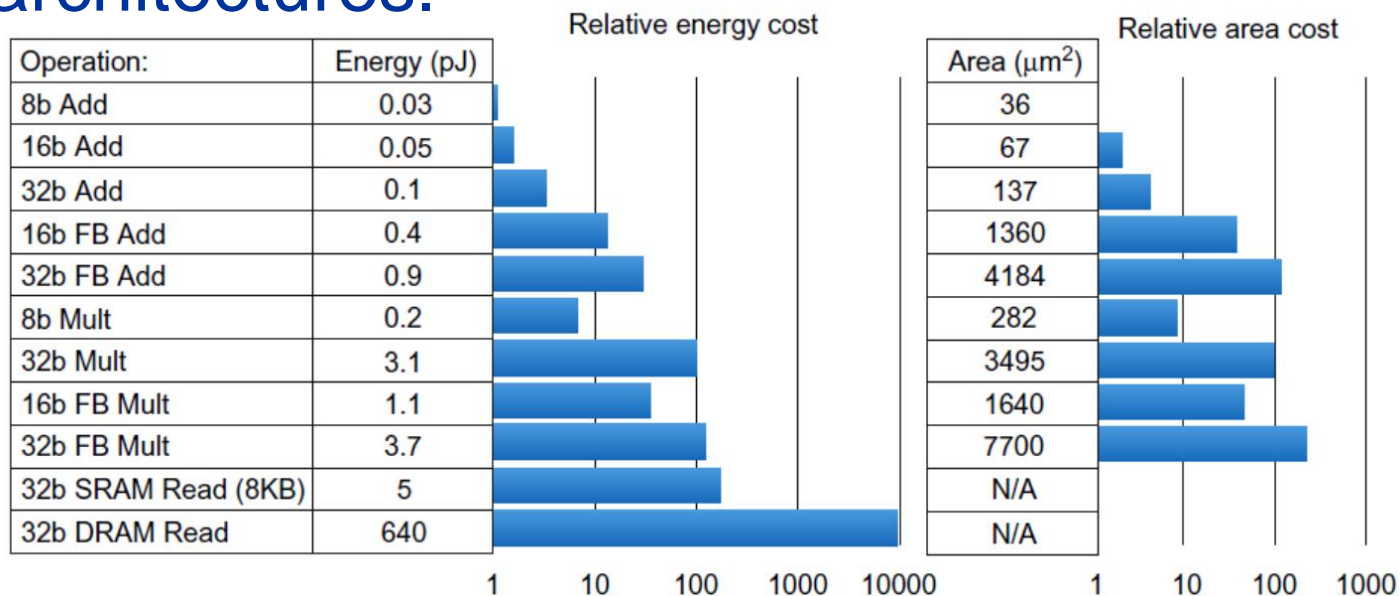
- Techniques for reducing power:
 - Do nothing well (turn clock off)
 - Dynamic Voltage-Frequency Scaling



- Low power state for DRAM, disks
- Overclocking, turning off cores

Static Power

- Static power consumption
 - 25-50% of total power
 - $\text{Current}_{\text{static}} \times \text{Voltage}$
 - Scales with number of transistors
 - To reduce: power gating
- Energy and die area motivate domain-specific architectures.



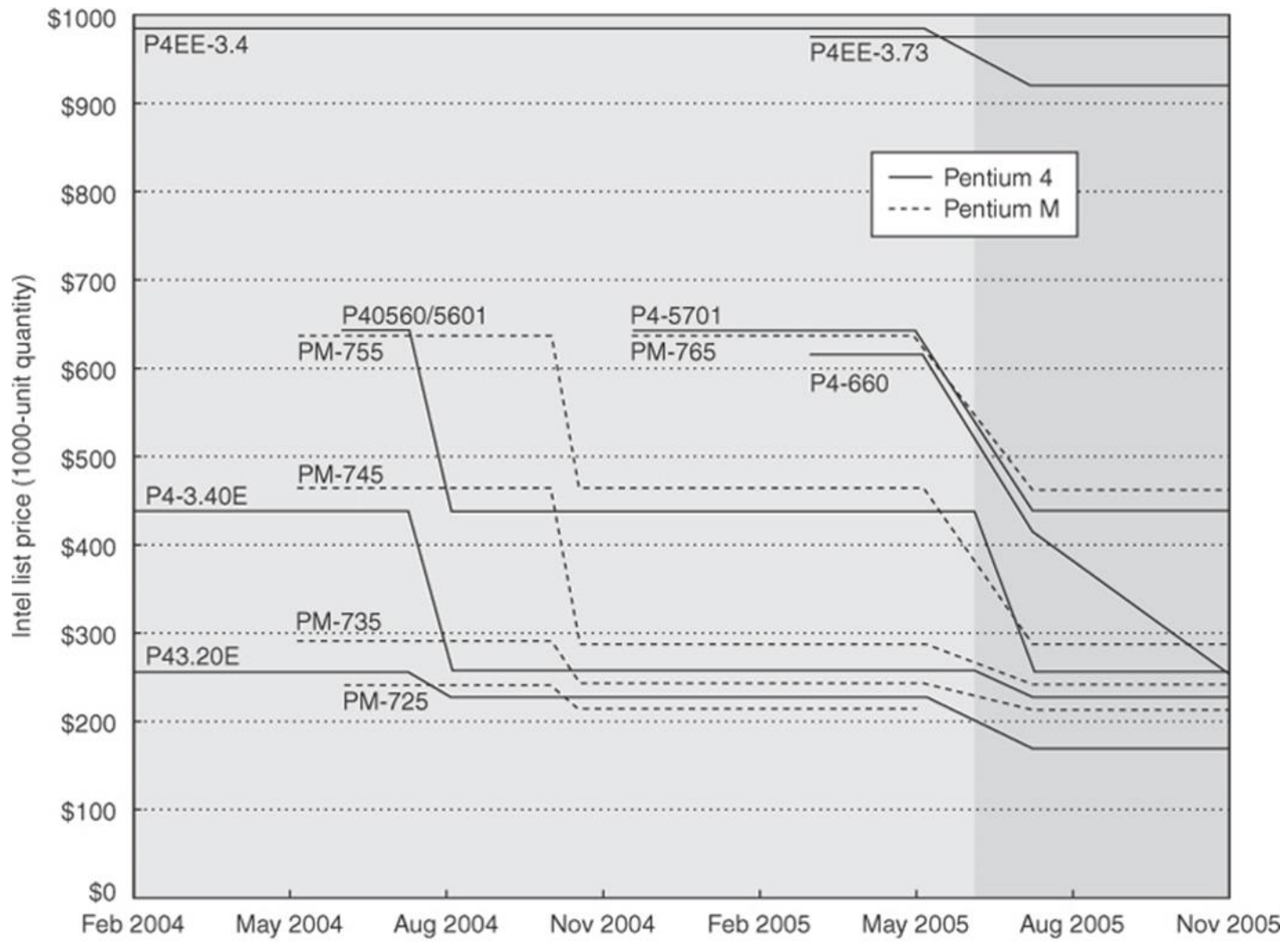
Contents

- Introduction
- Classes of Computers
- Defining Computer Architecture
- Trends in Technology
- Trends in Power and Energy in Integrated Circuits
- Trends in Cost
- Dependability
- Measuring, Reporting, and Summarizing Performance
- Quantitative Principles of Computer Design
- Fallacies and Pitfalls

Trends in Cost

- Cost driven down by learning curve
 - Yield
- DRAM: price closely tracks cost
- Microprocessors: price depends on volume
 - 10% less for each doubling of volume

Cost over Time



© 2007 Elsevier, Inc. All rights reserved.

Integrated Circuit Cost

$$\text{Cost of integrated circuit} = \frac{\text{Cost of die} + \text{Cost of testing die} + \text{Cost of packaging and final test}}{\text{Final test yield}}$$

$$\text{Cost of die} = \frac{\text{Cost of wafer}}{\text{Dies per wafer} \times \text{Die yield}}$$

$$\text{Dies per wafer} = \frac{\pi \times (\text{Wafer diameter} / 2)^2}{\text{Die area}} - \frac{\pi \times \text{Wafer diameter}}{\sqrt{2} \times \text{Die area}}$$

- Bose-Einstein formula:

$$\text{Die yield} = \text{Wafer yield} \times 1 / (1 + \text{Defects per unit area} \times \text{Die area})^N$$

- Defects per unit area = 0.016-0.057 defects per square cm (2010)
- N = process-complexity factor = 11.5-15.5 (40 nm, 2010)
- IC cost \propto (die area)^x: x ~ 2

Contents

- Introduction
- Classes of Computers
- Defining Computer Architecture
- Trends in Technology
- Trends in Power and Energy in Integrated Circuits
- Trends in Cost
- Dependability
- Measuring, Reporting, and Summarizing Performance
- Quantitative Principles of Computer Design
- Fallacies and Pitfalls

Dependability

- Module reliability
 - Mean time to failure (MTTF)
 - Mean time to repair (MTTR)
 - Mean time between failures (MTBF) = $MTTF + MTTR$
 - Availability = $MTTF / MTBF$
- FIT = Failure in Time = $10^9 / MTTF$
- Example 1: Find the system MTTF and availability for:

Component	QTY	MTTF (hrs)	MTTR (hrs)
C1	10	10^6	10
C2	1	10^5	10

Dependability

- Dependability can be improve through redundancy.
- MTTF for redundant pair =
$$(MTTF/2) / (MTTR/MTTF) = MTTF^2 / 2MTTR$$
- Example 2: Solve Example 1 assuming a redundant pair.

Contents

- Introduction
- Classes of Computers
- Defining Computer Architecture
- Trends in Technology
- Trends in Power and Energy in Integrated Circuits
- Trends in Cost
- Dependability
- Measuring, Reporting, and Summarizing Performance
- Quantitative Principles of Computer Design
- Fallacies and Pitfalls

Measuring Performance

- Typical performance metrics:
 - Response time
 - Throughput
- Speedup of X relative to Y
 - $\text{Performance}_X / \text{Performance}_Y$
 - $\text{Execution time}_Y / \text{Execution time}_X$
- Execution time
 - Wall clock time: includes all system overheads
 - CPU time: only computation time
- Benchmarks
 - Kernels (e.g. matrix multiply)
 - Toy programs (e.g. sorting)
 - Synthetic benchmarks (e.g. Dhrystone)
 - Benchmark suites (e.g. SPEC CPU2017, TPC-C)

Contents

- Introduction
- Classes of Computers
- Defining Computer Architecture
- Trends in Technology
- Trends in Power and Energy in Integrated Circuits
- Trends in Cost
- Dependability
- Measuring, Reporting, and Summarizing Performance
- **Quantitative Principles of Computer Design**
- **Fallacies and Pitfalls**

Principles of Computer Design

- Take Advantage of Parallelism
 - e.g. multiple processors, disks, memory banks, pipelining, multiple functional units
- Principle of Locality
 - Reuse of data and instructions
- Focus on the Common Case
 - Amdahl's Law

$$Execution\ time_{new} = Execution\ time_{old} \times \left((1 - Fraction_{enhanced}) + \frac{Fraction_{enhanced}}{Speedup_{enhanced}} \right)$$

$$Speedup_{overall} = \frac{1}{(1 - Fraction_{enhanced}) + \frac{Fraction_{enhanced}}{Speedup_{enhanced}}}$$

Principles of Computer Design

- The Processor Performance Equation

CPU time = CPU clock cycles for a program × Clock cycle time

$$CPU\ time = \frac{CPU\ clock\ cycles\ for\ a\ program}{Clock\ rate}$$

$$CPI = \frac{CPU\ clock\ cycles\ for\ a\ program}{Instruction\ count}$$

CPU time = Instruction count × Cycles per instruction × Clock cycle time

$$\frac{Instructions}{Program} \times \frac{Clock\ cycles}{Instruction} \times \frac{Seconds}{Clock\ cycle} = \frac{Seconds}{Program} = CPU\ time$$

Principles of Computer Design

- Different instruction types having different CPIs

$$\text{CPU clock cycles} = \sum_{i=1}^n IC_i \times CPI_i$$

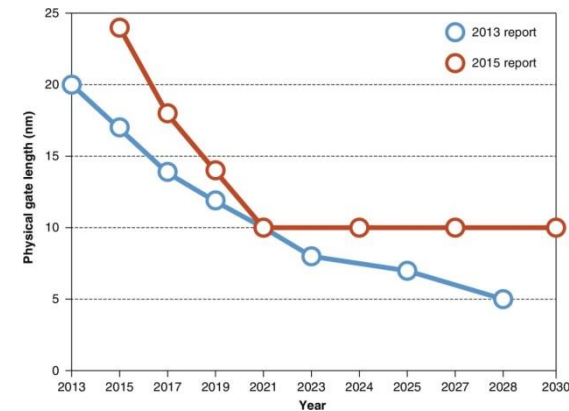
$$\text{CPU time} = \left(\sum_{i=1}^n IC_i \times CPI_i \right) \times \text{Clock cycle time}$$

Contents

- Introduction
- Classes of Computers
- Defining Computer Architecture
- Trends in Technology
- Trends in Power and Energy in Integrated Circuits
- Trends in Cost
- Dependability
- Measuring, Reporting, and Summarizing Performance
- Quantitative Principles of Computer Design
- **Fallacies and Pitfalls**

Fallacies and Pitfalls

- P: All exponential laws must come to an end
 - Dennard scaling (constant power density)
 - Stopped by threshold voltage
 - Disk capacity
 - 30-100% per year to 5% per year
 - Moore's Law
 - Most visible with DRAM capacity
 - ITRS disbanded
 - Only four foundries left producing state-of-the-art logic chips
 - 11 nm, 3 nm might be the limit



Fallacies and Pitfalls

- F: Microprocessors are a silver bullet
 - Performance is now a programmer's burden
- P: Falling prey to Amdahl's Law
- P: A single point of failure
 - The system is as strong as its weakest part.
- F: Hardware enhancements that increase performance also improve energy efficiency, or are at worst energy neutral
- F: Benchmarks remain valid indefinitely
 - Compiler optimizations target benchmarks

Fallacies and Pitfalls

- F: The rated mean time to failure of disks is 1,200,000 hours or almost 140 years, so disks practically never fail
 - MTTF value from manufacturers assume regular replacement every 5 years
- P: Fault detection can lower availability
 - Not all operations are needed for correct execution

Fallacies and Pitfalls

- F: Peak performance tracks observed performance

