

Chapter 1

Computer Abstractions and Technology

Adapted by Prof. Gheith Abandah

Content

- 1.2 Eight Great Ideas in Computer Architecture (*Review*)
- 1.5 Technologies for Building Processors and Memory
- 1.6 Performance (*Review*)
- 1.7 The Power Wall
- 1.8 The Sea Change: The Switch from Uniprocessors to Multiprocessors
- 1.9 Real Stuff: Benchmarking the Intel Core i7
- 1.10 Fallacies and Pitfalls
- 1.11 Concluding Remarks

Eight Great Ideas

- Design for *Moore's Law*
- Use *abstraction* to simplify design
- Make the *common case fast*
- Performance *via parallelism*
- Performance *via pipelining*
- Performance *via prediction*
- *Hierarchy* of memories
- *Dependability* *via* redundancy

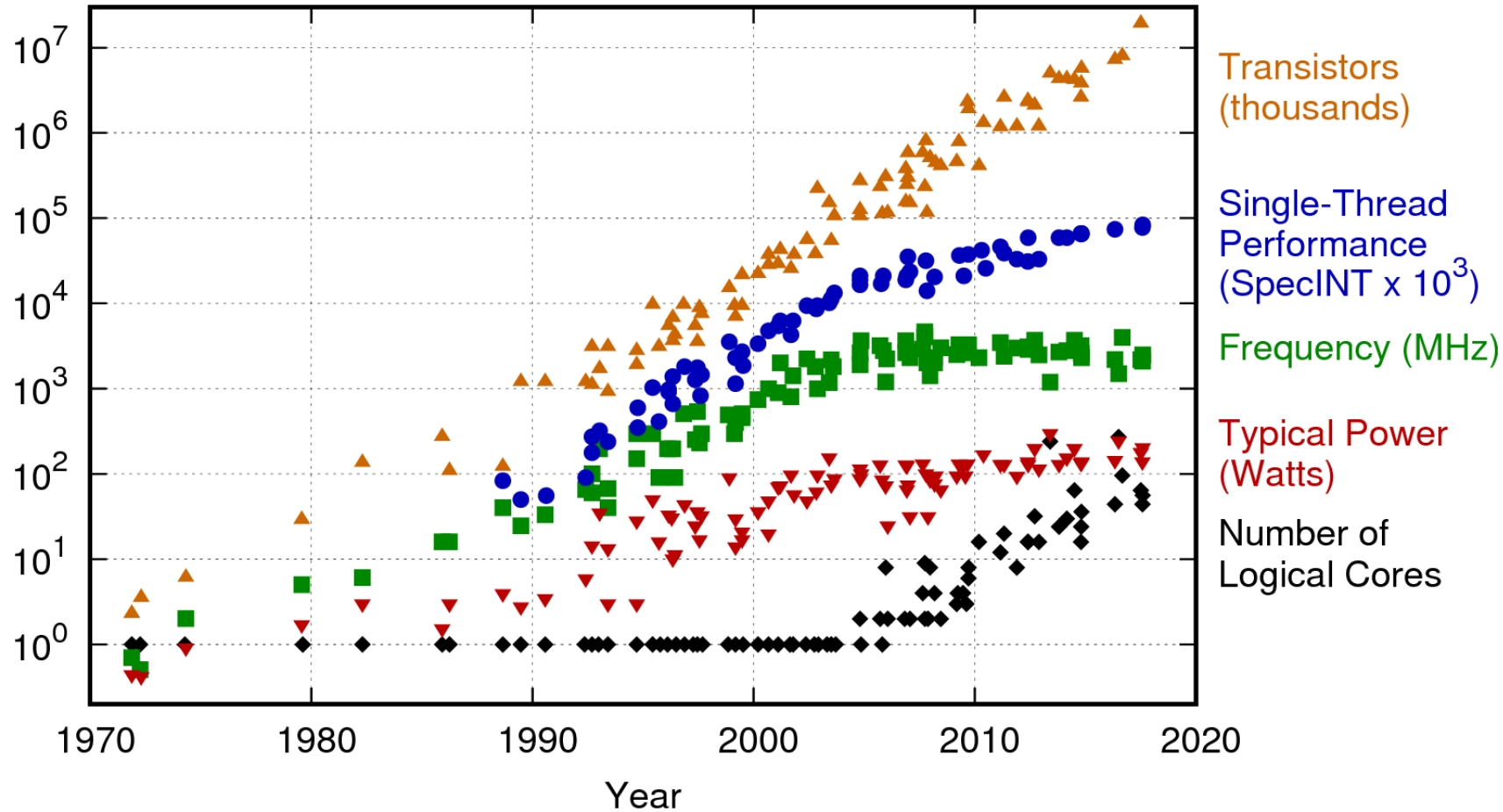


Content

- 1.2 Eight Great Ideas in Computer Architecture (*Review*)
- 1.5 Technologies for Building Processors and Memory
- 1.6 Performance (*Review*)
- 1.7 The Power Wall
- 1.8 The Sea Change: The Switch from Uniprocessors to Multiprocessors
- 1.9 Real Stuff: Benchmarking the Intel Core i7
- 1.10 Fallacies and Pitfalls
- 1.11 Concluding Remarks

Technology Trends

42 Years of Microprocessor Trend Data



Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2017 by K. Rupp

Technology Trends

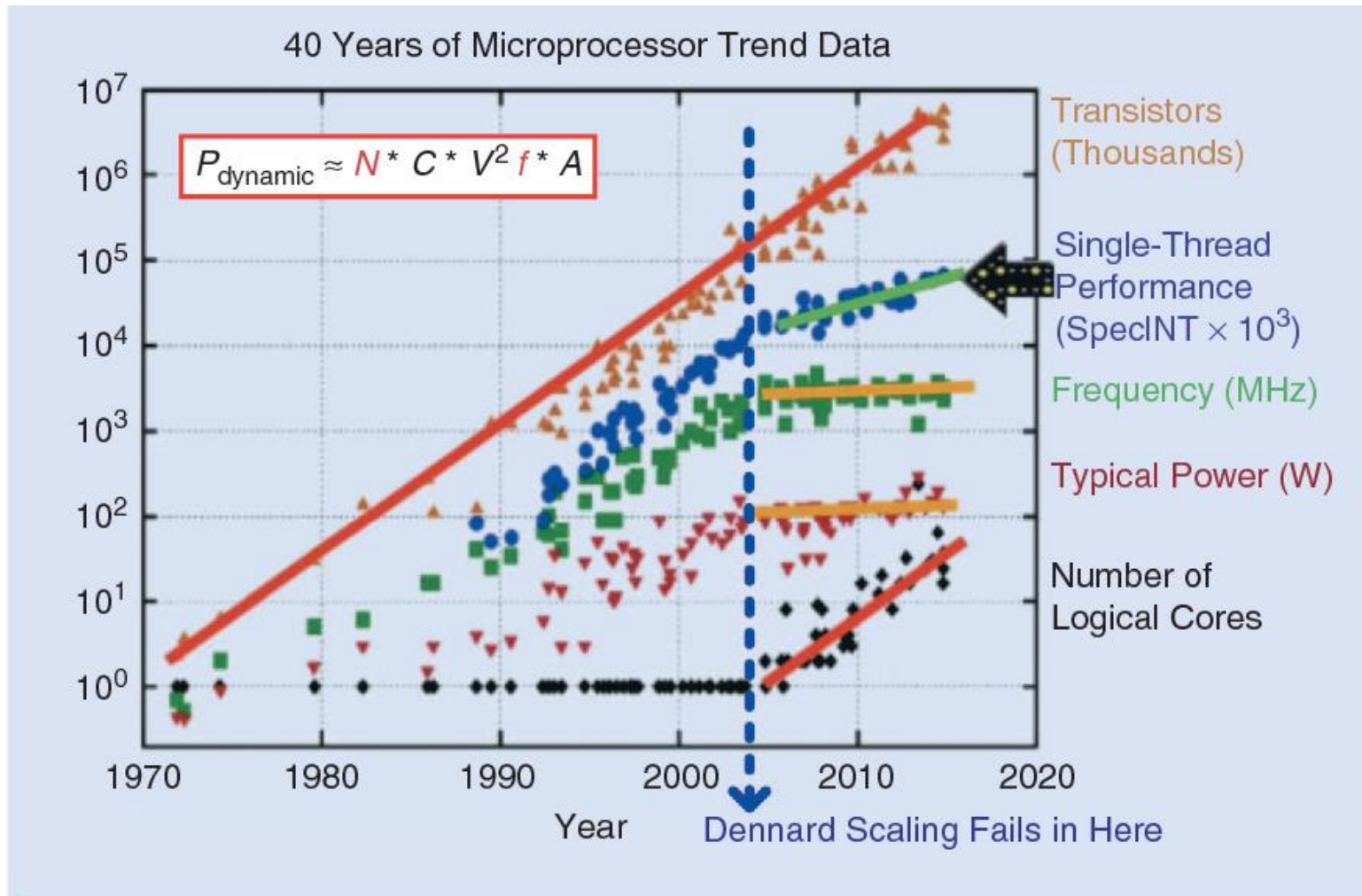
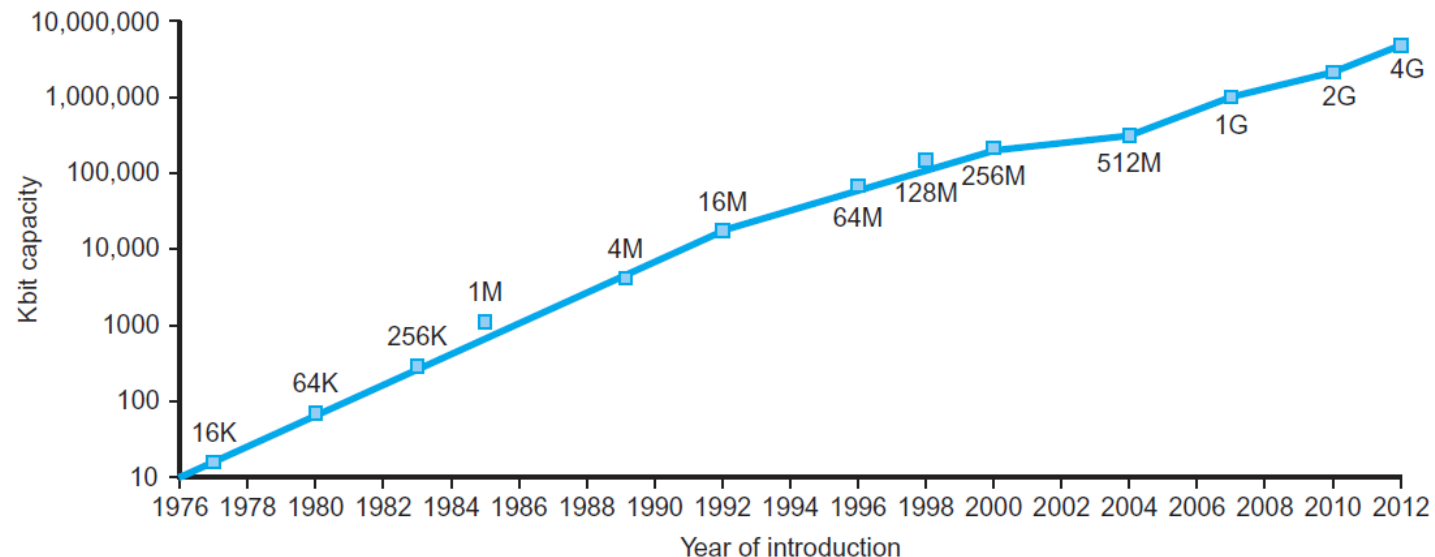


FIGURE 1: The Dennard scaling failed around the middle of the 2000s [24].

Technology Trends

DRAM capacity

- Electronics technology continues to evolve
 - Increased capacity and performance
 - Reduced cost



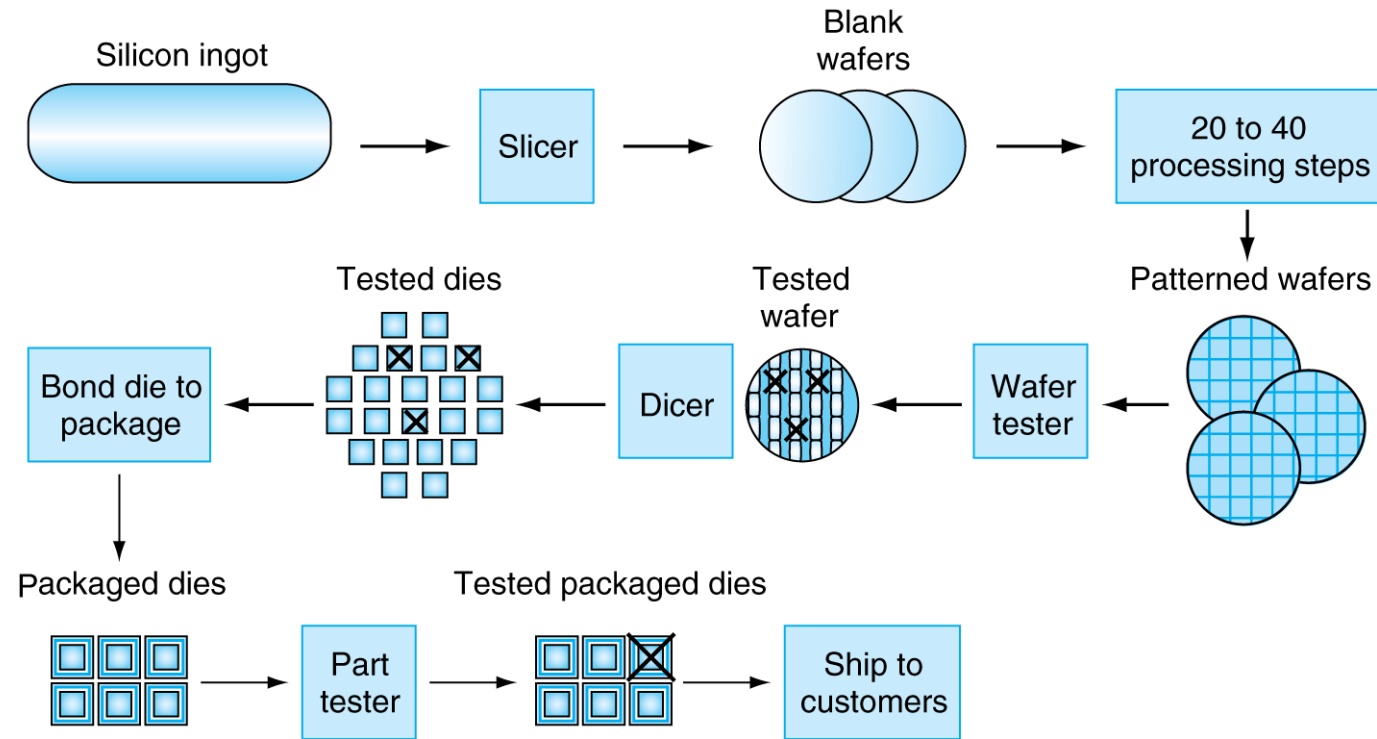
Year	Technology	Relative performance/cost
1951	Vacuum tube	1
1965	Transistor	35
1975	Integrated circuit (IC)	900
1995	Very large scale IC (VLSI)	2,400,000
2013	Ultra large scale IC	250,000,000,000

Semiconductor Technology

- Silicon: semiconductor
- Add materials to transform properties:
 - Conductors
 - Insulators
 - Switch
- YouTube Video: VLSI Fabrication Process

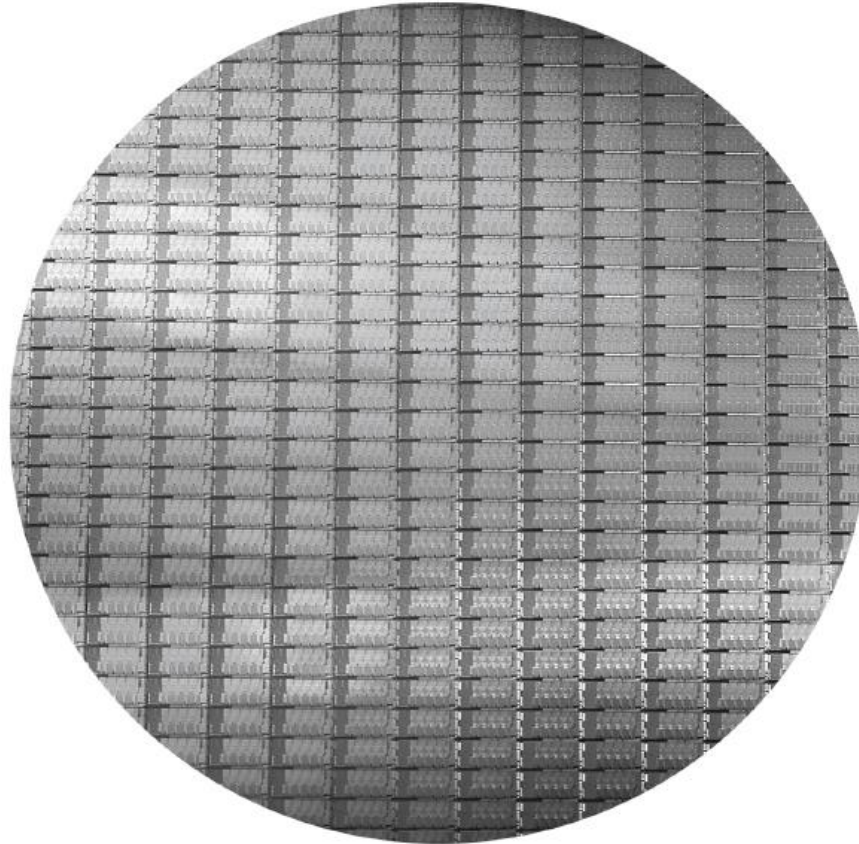
<https://youtu.be/fwNkg1fsqBY>

Manufacturing ICs



- Yield: proportion of working dies per wafer

Intel Core i7 Wafer



- 300mm wafer, 280 chips, 32nm technology
- Each chip is 20.7 x 10.5 mm

Integrated Circuit Cost

$$\text{Cost per die} = \frac{\text{Cost per wafer}}{\text{Dies per wafer} \times \text{Yield}}$$

$$\text{Dies per wafer} \approx \text{Wafer area} / \text{Die area}$$

$$\text{Yield} = \frac{1}{(1 + (\text{Defects per area} \times \text{Die area} / 2))^2}$$

- Nonlinear relation to area and defect rate
 - Wafer cost and area are fixed
 - Defect rate determined by manufacturing process
 - Die area determined by architecture and circuit design

Content

- 1.2 Eight Great Ideas in Computer Architecture (*Review*)
- 1.5 Technologies for Building Processors and Memory
- 1.6 Performance (*Review*)
- 1.7 The Power Wall
- 1.8 The Sea Change: The Switch from Uniprocessors to Multiprocessors
- 1.9 Real Stuff: Benchmarking the Intel Core i7
- 1.10 Fallacies and Pitfalls
- 1.11 Concluding Remarks

Response Time and Throughput

- Response time
 - How long it takes to do a task
- Throughput
 - Total work done per unit time
 - e.g., tasks/transactions/... per hour
- How are response time and throughput affected by
 - Replacing the processor with a faster version?
 - Adding more processors?
- We'll focus on response time for now...

Relative Performance

- Define Performance = $1/\text{Execution Time}$
- “X is n time faster than Y”

$$\begin{aligned} & \text{Performance}_X / \text{Performance}_Y \\ &= \text{Execution time}_Y / \text{Execution time}_X = n \end{aligned}$$

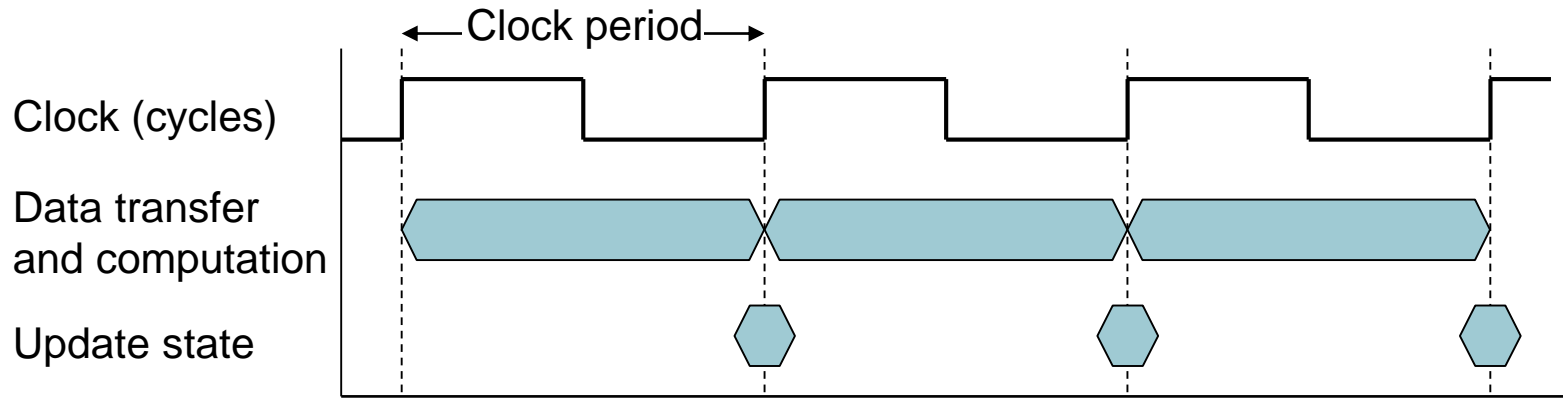
- Example: time taken to run a program
 - 10s on A, 15s on B
 - $\text{Execution Time}_B / \text{Execution Time}_A$
 $= 15\text{s} / 10\text{s} = 1.5$
 - So A is 1.5 times faster than B

Measuring Execution Time

- Elapsed time
 - Total response time, including all aspects
 - Processing, I/O, OS overhead, idle time
 - Determines system performance
- CPU time
 - Time spent processing a given job
 - Discounts I/O time, other jobs' shares
 - Comprises user CPU time and system CPU time
 - Different programs are affected differently by CPU and system performance

CPU Clocking

- Operation of digital hardware governed by a constant-rate clock



- Clock period: duration of a clock cycle
 - e.g., $250\text{ps} = 0.25\text{ns} = 250 \times 10^{-12}\text{s}$
- Clock frequency (rate): cycles per second
 - e.g., $4.0\text{GHz} = 4000\text{MHz} = 4.0 \times 10^9\text{Hz}$

CPU Time

$$\begin{aligned}\text{CPU Time} &= \text{CPU Clock Cycles} \times \text{Clock Cycle Time} \\ &= \frac{\text{CPU Clock Cycles}}{\text{Clock Rate}}\end{aligned}$$

- Performance improved by
 - Reducing number of clock cycles
 - Increasing clock rate
 - Hardware designer must often trade off clock rate against cycle count

Instruction Count and CPI

Clock Cycles = Instruction Count \times Cycles per Instruction

CPU Time = Instruction Count \times CPI \times Clock Cycle Time

$$= \frac{\text{Instruction Count} \times \text{CPI}}{\text{Clock Rate}}$$

- Instruction Count for a program
 - Determined by program, ISA and compiler
- Average cycles per instruction
 - Determined by CPU hardware
 - If different instructions have different CPI
 - Average CPI affected by instruction mix

CPI in More Detail

- If different instruction classes take different numbers of cycles

$$\text{Clock Cycles} = \sum_{i=1}^n (\text{CPI}_i \times \text{Instruction Count}_i)$$

- Weighted average CPI

$$\text{CPI} = \frac{\text{Clock Cycles}}{\text{Instruction Count}} = \sum_{i=1}^n \left(\text{CPI}_i \times \frac{\text{Instruction Count}_i}{\text{Instruction Count}} \right)$$

Relative frequency

Performance Summary

The BIG Picture

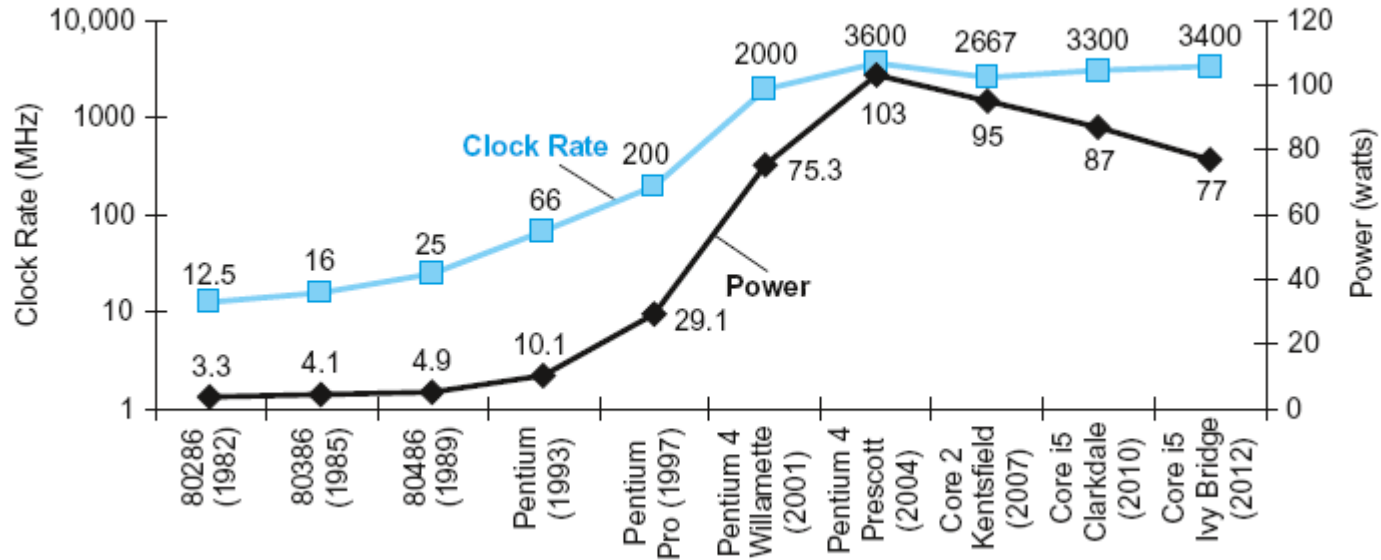
$$\text{CPU Time} = \frac{\text{Instructions}}{\text{Program}} \times \frac{\text{Clock cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{Clock cycle}}$$

- Performance depends on
 - Algorithm: affects IC, possibly CPI
 - Programming language: affects IC, CPI
 - Compiler: affects IC, CPI
 - Instruction set architecture: affects IC, CPI, T_c

Content

- 1.2 Eight Great Ideas in Computer Architecture (*Review*)
- 1.5 Technologies for Building Processors and Memory
- 1.6 Performance (*Review*)
- 1.7 The Power Wall
- 1.8 The Sea Change: The Switch from Uniprocessors to Multiprocessors
- 1.9 Real Stuff: Benchmarking the Intel Core i7
- 1.10 Fallacies and Pitfalls
- 1.11 Concluding Remarks

Power Trends



- In CMOS IC technology

$$\text{Power} = \frac{1}{2} \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency}$$

×30

5V → 1V

×1000

Reducing Power

- Suppose a new CPU has
 - 85% of capacitive load of old CPU
 - 15% voltage and 15% frequency reduction

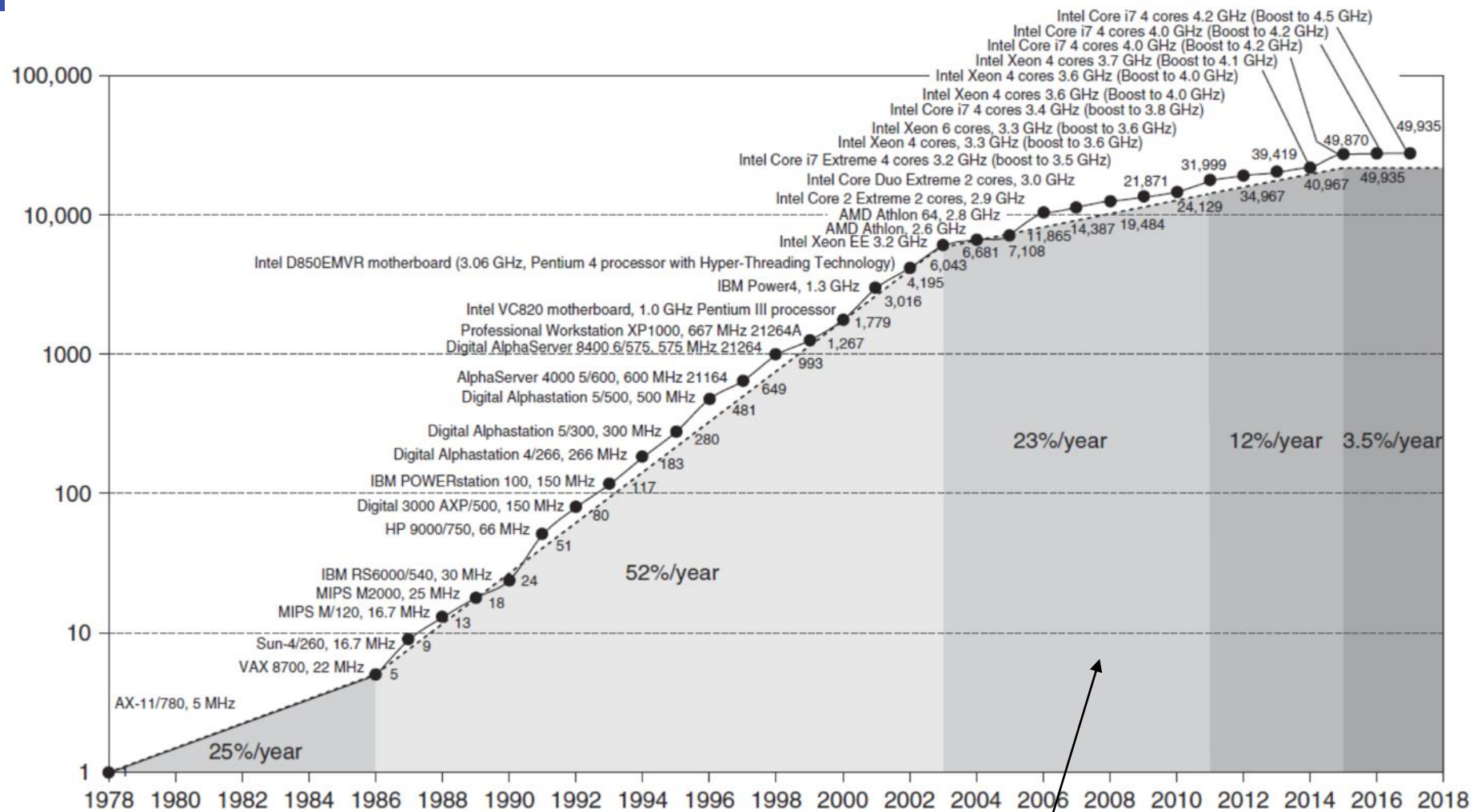
$$\frac{P_{\text{new}}}{P_{\text{old}}} = \frac{C_{\text{old}} \times 0.85 \times (V_{\text{old}} \times 0.85)^2 \times F_{\text{old}} \times 0.85}{C_{\text{old}} \times V_{\text{old}}^2 \times F_{\text{old}}} = 0.85^4 = 0.52$$

- The power wall
 - We can't reduce voltage further
 - We can't remove more heat
- How else can we improve performance?

Content

- 1.2 Eight Great Ideas in Computer Architecture (*Review*)
- 1.5 Technologies for Building Processors and Memory
- 1.6 Performance (*Review*)
- 1.7 The Power Wall
- 1.8 The Sea Change: The Switch from Uniprocessors to Multiprocessors
- 1.9 Real Stuff: Benchmarking the Intel Core i7
- 1.10 Fallacies and Pitfalls
- 1.11 Concluding Remarks

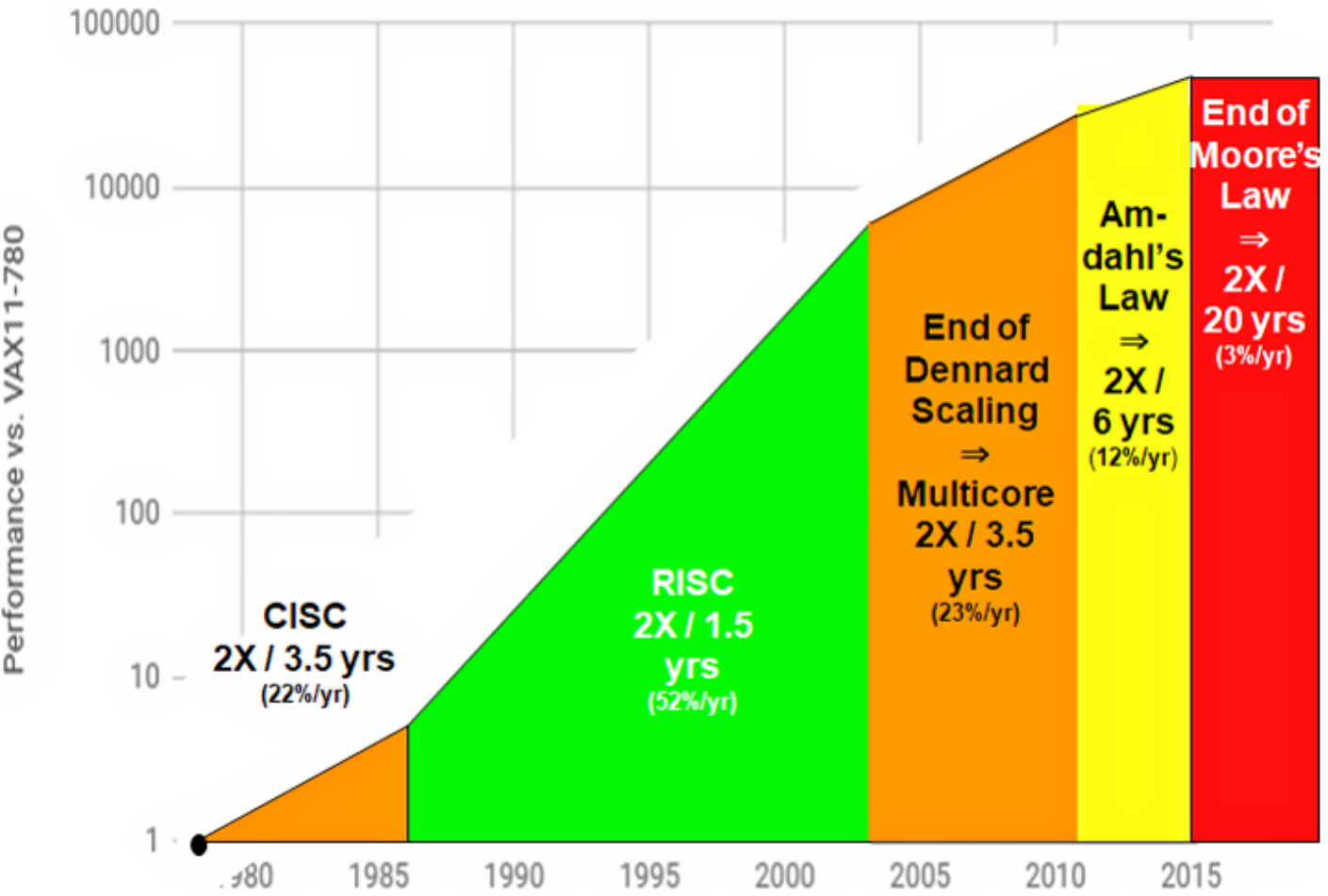
Uniprocessor Performance



Constrained by power, instruction-level parallelism, memory latency

Uniprocessor Performance

40 years of Processor Performance



Multiprocessors

- Multicore microprocessors
 - More than one processor per chip
- Requires explicitly parallel programming
 - Compare with instruction level parallelism
 - Hardware executes multiple instructions at once
 - Hidden from the programmer
 - Hard to do
 - Programming for performance
 - Load balancing
 - Optimizing communication and synchronization

Content

- 1.2 Eight Great Ideas in Computer Architecture (*Review*)
- 1.5 Technologies for Building Processors and Memory
- 1.6 Performance (*Review*)
- 1.7 The Power Wall
- 1.8 The Sea Change: The Switch from Uniprocessors to Multiprocessors
- 1.9 Real Stuff: Benchmarking the Intel Core i7**
- 1.10 Fallacies and Pitfalls**
- 1.11 Concluding Remarks**

SPEC CPU Benchmark

- Programs used to measure performance
 - Supposedly typical of actual workload
- Standard Performance Evaluation Corp (SPEC)
 - Develops benchmarks for CPU, I/O, Web, ...
- SPEC CPU2006
 - Elapsed time to execute a selection of programs
 - Negligible I/O, so focuses on CPU performance
 - Normalize relative to reference machine
 - Summarize as geometric mean of performance ratios
 - CINT2006 (integer) and CFP2006 (floating-point)

$$\sqrt[n]{\prod_{i=1}^n \text{Execution time ratio}_i}$$

CINT2006 for Intel Core i7 920

Description	Name	Instruction Count x 10 ⁹	CPI	Clock cycle time (seconds x 10 ⁻⁹)	Execution Time (seconds)	Reference Time (seconds)	SPECratio
Interpreted string processing	perl	2252	0.60	0.376	508	9770	19.2
Block-sorting compression	bzip2	2390	0.70	0.376	629	9650	15.4
GNU C compiler	gcc	794	1.20	0.376	358	8050	22.5
Combinatorial optimization	mcf	221	2.66	0.376	221	9120	41.2
Go game (AI)	go	1274	1.10	0.376	527	10490	19.9
Search gene sequence	hmmer	2616	0.60	0.376	590	9330	15.8
Chess game (AI)	sjeng	1948	0.80	0.376	586	12100	20.7
Quantum computer simulation	libquantum	659	0.44	0.376	109	20720	190.0
Video compression	h264avc	3793	0.50	0.376	713	22130	31.0
Discrete event simulation library	omnetpp	367	2.10	0.376	290	6250	21.5
Games/path finding	astar	1250	1.00	0.376	470	7020	14.9
XML parsing	xalancbmk	1045	0.70	0.376	275	6900	25.1
Geometric mean	-	-	-	-	-	-	25.7

SPEC Power Benchmark

- Power consumption of server at different workload levels
 - Performance: ssj_ops/sec
 - Power: Watts (Joules/sec)

$$\text{Overallssj_ops per Watt} = \left(\sum_{i=0}^{10} \text{ssj_ops}_i \right) / \left(\sum_{i=0}^{10} \text{power}_i \right)$$

SPECpower_ssj2008 for Xeon X5650

Target Load %	Performance (ssj_ops)	Average Power (Watts)
100%	865,618	258
90%	786,688	242
80%	698,051	224
70%	607,826	204
60%	521,391	185
50%	436,757	170
40%	345,919	157
30%	262,071	146
20%	176,061	135
10%	86,784	121
0%	0	80
Overall Sum	4,787,166	1,922
Σ ssj_ops/ Σ power =		2,490

Content

- 1.2 Eight Great Ideas in Computer Architecture (*Review*)
- 1.5 Technologies for Building Processors and Memory
- 1.6 Performance (*Review*)
- 1.7 The Power Wall
- 1.8 The Sea Change: The Switch from Uniprocessors to Multiprocessors
- 1.9 Real Stuff: Benchmarking the Intel Core i7
- 1.10 Fallacies and Pitfalls**
- 1.11 Concluding Remarks**

Pitfall: Amdahl's Law

- Improving an aspect of a computer and expecting a proportional improvement in overall performance

$$T_{\text{improved}} = \frac{T_{\text{affected}}}{\text{improvement factor}} + T_{\text{unaffected}}$$

- Example: multiply accounts for 80s/100s
 - How much improvement in multiply performance to get 5× overall?

$$20 = \frac{80}{n} + 20 \quad \text{■ Can't be done!}$$

- Corollary: make the common case fast

Fallacy: Low Power at Idle

- Look back at i7 power benchmark
 - At 100% load: 258W
 - At 50% load: 170W (66%)
 - At 10% load: 121W (47%)
- Google data center
 - Mostly operates at 10% – 50% load
 - At 100% load less than 1% of the time
- Consider designing processors to make power proportional to load

Pitfall: MIPS as a Performance Metric

- MIPS: Millions of Instructions Per Second
 - Doesn't account for
 - Differences in ISAs between computers
 - Differences in complexity between instructions

$$\begin{aligned} \text{MIPS} &= \frac{\text{Instruction count}}{\text{Execution time} \times 10^6} \\ &= \frac{\text{Instruction count}}{\frac{\text{Instruction count} \times \text{CPI}}{\text{Clock rate}} \times 10^6} = \frac{\text{Clock rate}}{\text{CPI} \times 10^6} \end{aligned}$$

- CPI varies between programs on a given CPU

Content

- 1.2 Eight Great Ideas in Computer Architecture (*Review*)
- 1.5 Technologies for Building Processors and Memory
- 1.6 Performance (*Review*)
- 1.7 The Power Wall
- 1.8 The Sea Change: The Switch from Uniprocessors to Multiprocessors
- 1.9 Real Stuff: Benchmarking the Intel Core i7
- 1.10 Fallacies and Pitfalls
- 1.11 Concluding Remarks**

Concluding Remarks

- Cost/performance is improving
 - Due to underlying technology development
- Execution time: the best performance measure
- Power is a limiting factor
 - Use parallelism to improve performance