

CHAPTER 1: MEASURING & UNDERSTANDING PERFORMANCE

JOSÉ M^a GONZÁLEZ LINARES

DEPT. DE ARQUITECTURA DE COMPUTADORES



UNIVERSIDAD
DE MÁLAGA



INDEX

Introduction

Below the covers

Understanding Performance

The Power Wall

Concluding Remarks

INTRODUCTION

THE COMPUTER REVOLUTION

WHAT IS A COMPUTER?

A computer is a general purpose device that can be programmed to carry out a finite set of arithmetic or logical operations.

- We use programs, i.e., sequences of instructions, to solve problems.
- The Operating System help us to code and execute our programs.

CLASSES OF COMPUTERS

Desktop computers

- General purpose, variety of software, graphics display, keyboard and mouse
- Subject to cost/performance tradeoff

Server computers

- Network based
- High capacity, dependability, security

Supercomputers

- Hundreds to thousands processors, terabytes of memory, petabytes of storage
- High performance, cost

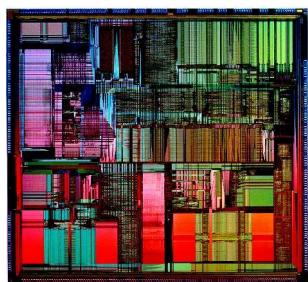
Embedded computers

- Hidden as components of systems
- Stringent power/performance/cost constraints

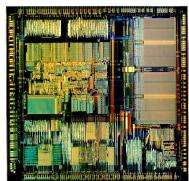
EVOLUTION OF PROCESSORS FOR DESKTOPS AND LAPTOPS



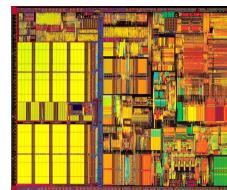
Intel 80286 (1982)
 134×10^3 transistors
12 Mhz, 68,7 mm²



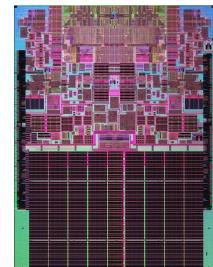
Intel Pentium (1993)
 3.1×10^6 transistors
66 Mhz, 264 mm²



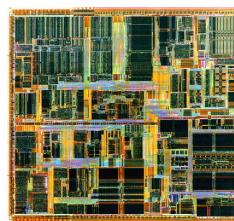
Intel 80386 (1985)
 275×10^3 transistors
33 Mhz, 104 mm



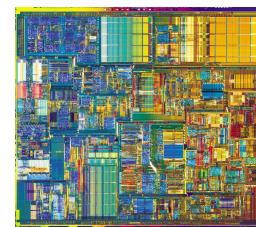
Intel Pentium III (1999)
 28×10^6 transistors
733 Mhz, 140 mm²



Intel Core (2006)
 291×10^6 transistors
1 Ghz, 143 mm²



Intel Pentium II (1997)
 7.5×10^6 transistors
300 Mhz, 209 mm²



Intel Pentium 4 (2000)
 42×10^6 transistors
1.5 Ghz, 224 mm²

CURRENT PROCESSORS FOR DESKTOPS AND LAPTOPS

Intel 12th Gen
Alder Lake



Process of 10nm and less, high-performance and power-efficient cores, Out-of-Order execution, TDP between 4 and 100 W, three levels of cache, integrated GPU, ...

SUPERCOMPUTERS

The top500 list in June 2022

Rank	System	Cores	Rmax (PFlop/s)	Rpeak (PFlop/s)	Power (kW)
1	Frontier - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE DOE/SC/Oak Ridge National Laboratory United States	8,730,112	1,102.00	1,685.65	21,100
2	Supercomputer Fugaku - Supercomputer Fugaku, A64FX 48C 2.2GHz, Tofu interconnect D, Fujitsu RIKEN Center for Computational Science Japan	7,630,848	442.01	537.21	29,899
3	LUMI - HPE Cray EX235a, AMD Optimized 3rd Generation EPYC 64C 2GHz, AMD Instinct MI250X, Slingshot-11, HPE EuroHPC/CSC Finland	1,110,144	151.90	214.35	2,942
4	Summit - IBM Power System AC922, IBM POWER9 22C 3.07GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, IBM DOE/SC/Oak Ridge National Laboratory United States	2,414,592	148.60	200.79	10,096
5	Sierra - IBM Power System AC922, IBM POWER9 22C 3.1GHz, NVIDIA Volta GV100, Dual-rail Mellanox EDR Infiniband, IBM / NVIDIA / Mellanox DOE/NNSA/LLNL United States	1,572,480	94.64	125.71	7,438

WHERE ELSE ARE EMBEDDED PROCESSORS FOUND?



THE COMPUTER REVOLUTION

Progress in computer technology

- Underpinned by Moore's Law

Makes novel applications feasible

- Computers in automobiles
- Cell phones
- Human genome project
- World Wide Web
- Search Engines

Computers are pervasive



MOORE'S LAW

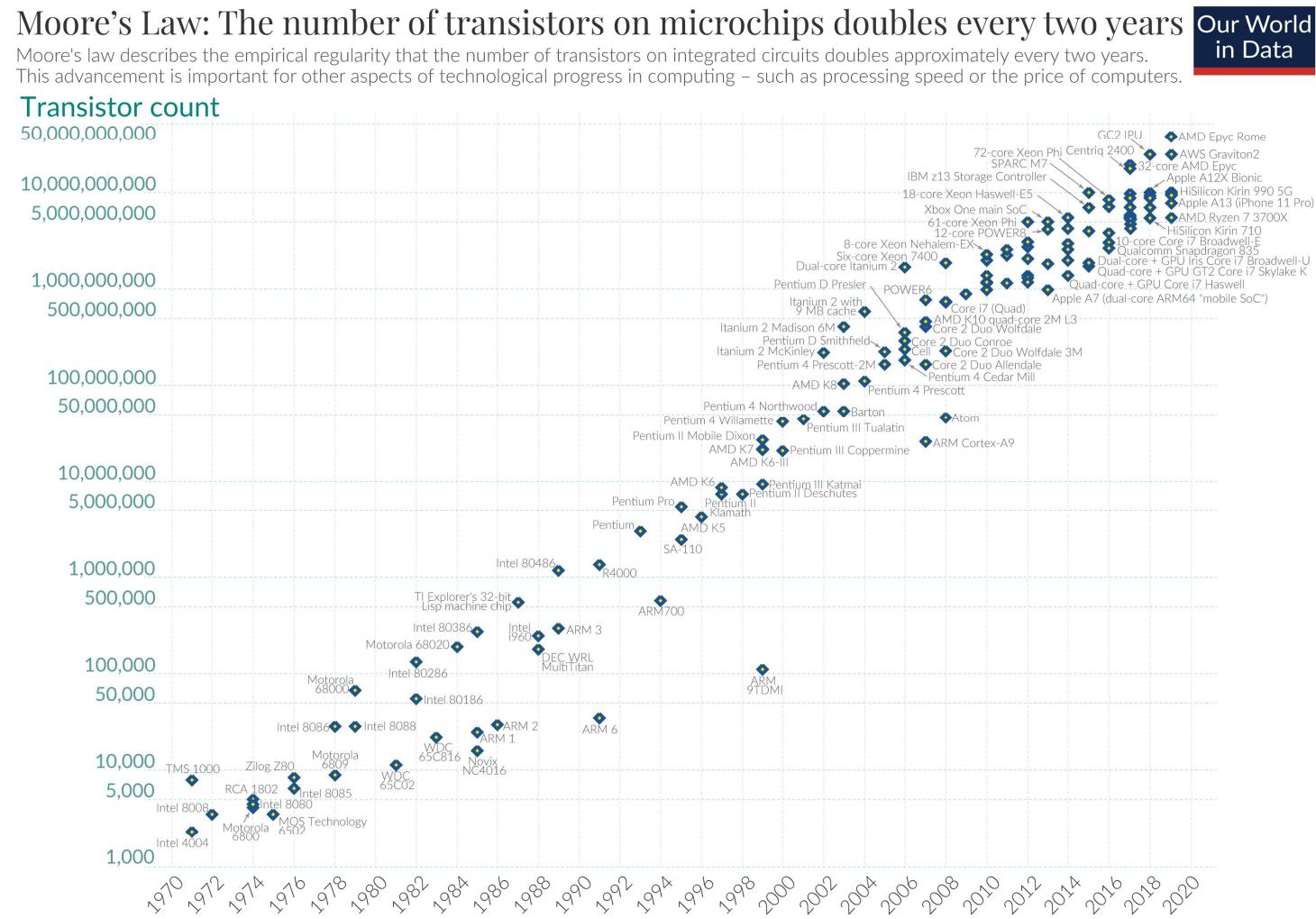
2017: Quad Core Kaby Lake with 3.2B transistors

2019: Epyc Rome with 39.5B transistors

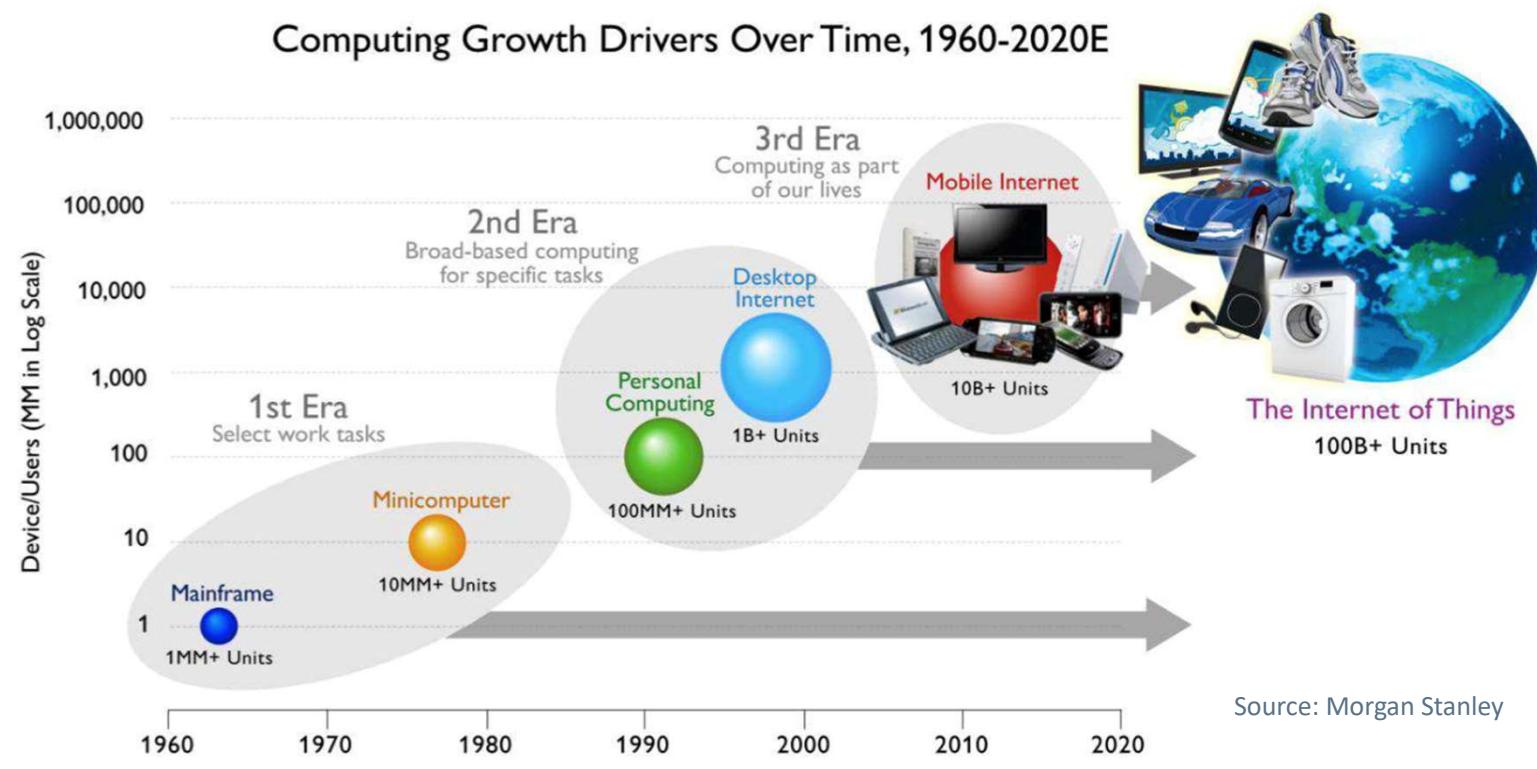
2022: Apple M1 Ultra with 114B transistors

2020: GA100 Ampere Nvidia GPU with 54B transistors

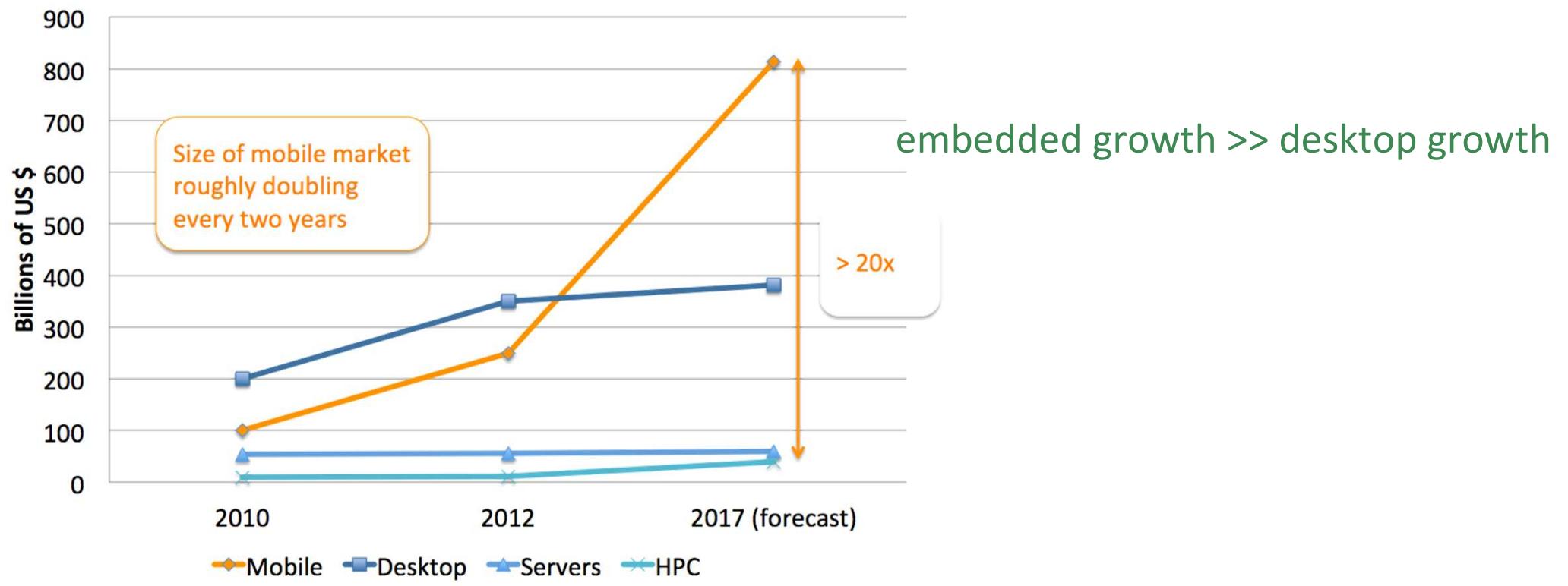
2022: GH100 Hopper Nvidia with 80B transistors



CONNECTIVITY DRIVING COMPUTING



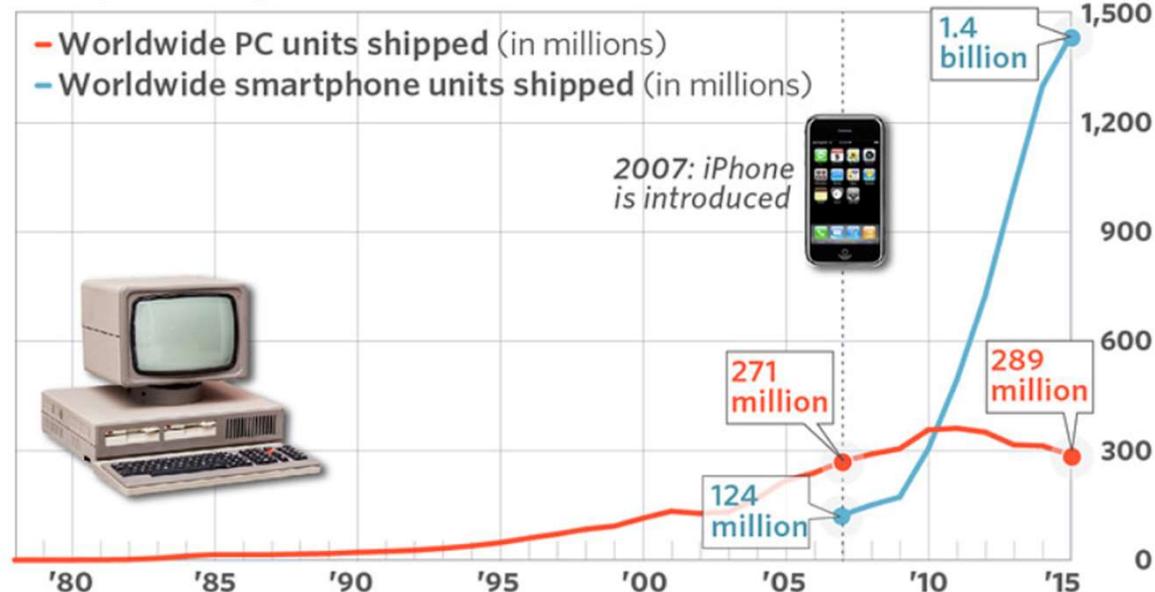
THE PROCESSOR MARKET



SMARTPHONE INDUSTRY EVOLUTION

How smartphones killed the PC

Smartphones outgrew PCs in 2011

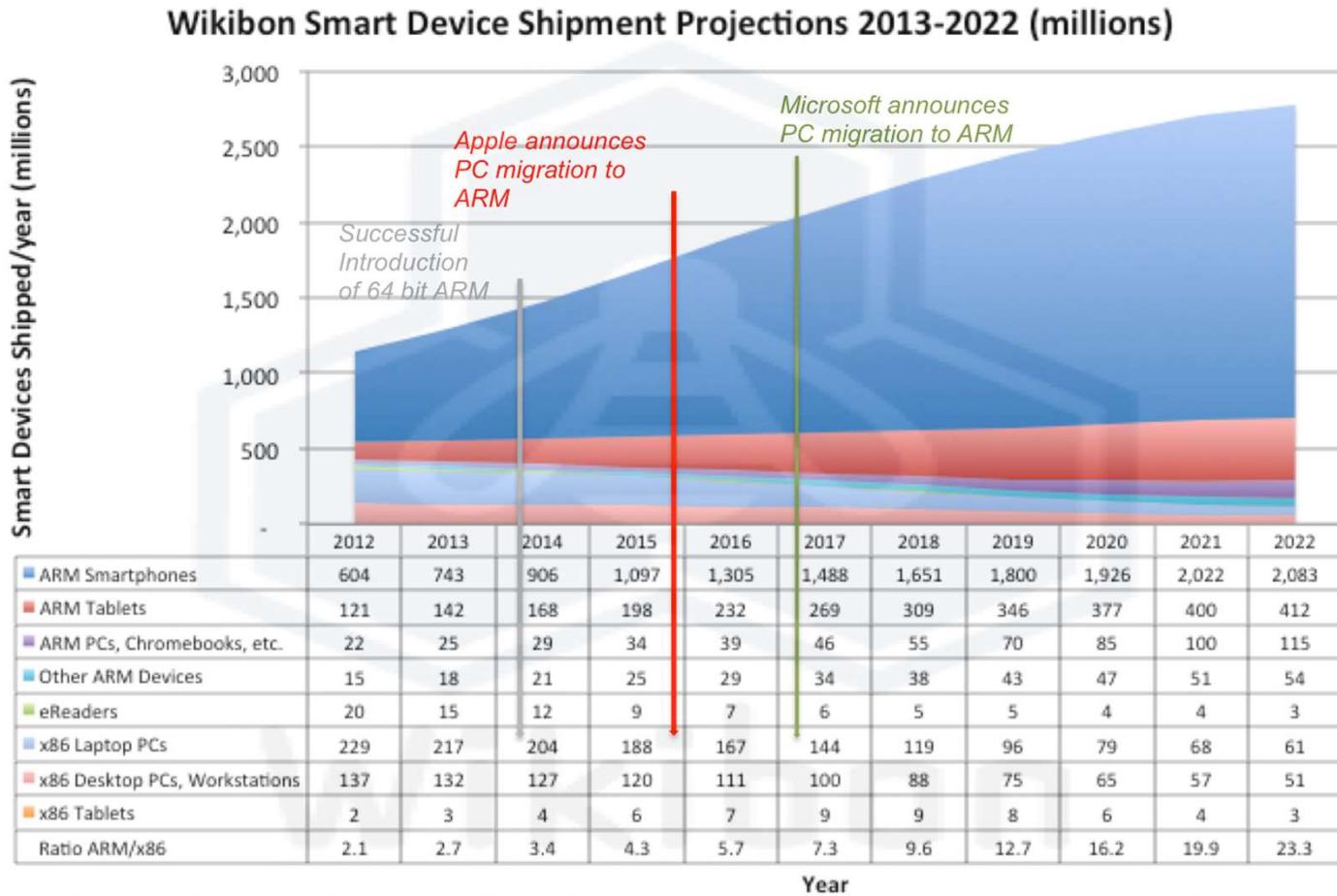


4bn people buying
every 2 years
instead of 1.6bn
every 5 years

Source: Gartner, IDC, Apple

TECNOLOGY SHIFT

ARM has cannibalized the processor market



WHAT WILL YOU LEARN?

How programs are translated into the machine language

- And how the hardware executes them

The hardware/software interface

What determines program performance

- And how it can be improved

How hardware designers improve performance

How software designers improve performance

- And what is parallel processing

IS IT USEFUL?

Marcos Fajardo wrote from scratch a rendering system, Arnold, that runs natively on x86 CPUs, where it tries to take advantage of all available threads and SIMD lanes for optimal parallelism.

It is used in

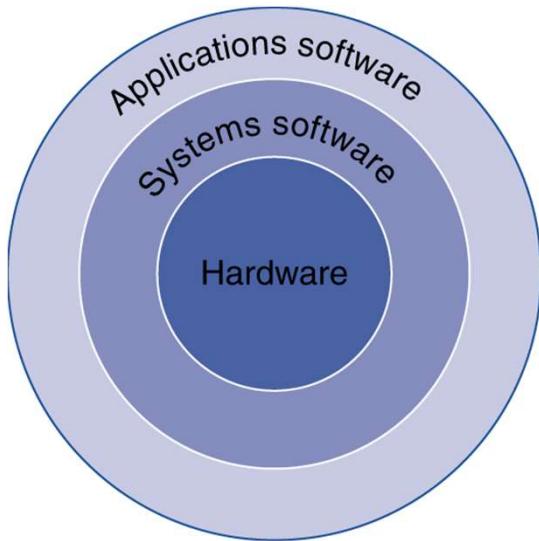
- Thor
- Captain America
- X-Men: First Class
- The Avengers
- Gravity
- Guardians of the Galaxy
- Star Wars: The Force Awakens
- Game of Thrones
- ...



BELLOW YOUR PROGRAM

UNDER THE COVERS

BELOW YOUR PROGRAM



Application software

- Written in high-level language (HLL)

System software

- Compiler: translates HLL code to machine code
- Operating System: service code
 - Handling input/output
 - Managing memory and storage
 - Scheduling tasks & sharing resources

Hardware

- Processor, memory, Input/Output (I/O) controllers

LEVELS OF PROGRAM CODE

High-level language

- Level of abstraction closer to problem domain
- Provides for productivity and portability

Assembly language

- Textual representation of instructions

Hardware representation

- Binary digits (bits)
- Encoded instructions and data

High-level
language
program
(in C)

```
swap(int v[], int k)
{int temp;
 temp = v[k];
 v[k] = v[k+1];
 v[k+1] = temp;
}
```

Assembly
language
program
(for MIPS)

```
swap:
    muli $2, $5,4
    add $2, $4,$2
    lw $15, 0($2)
    lw $16, 4($2)
    sw $16, 0($2)
    sw $15, 4($2)
    jr $31
```

Binary machine
language
program
(for MIPS)

```
0000000010100010000000000011000
000000000011000001100000100001
1000110001100010000000000000000
10001100011100100000000000000000
10101100111100100000000000000000
10101100011000100000000000000000
00000011110000000000000000000000
```

Compiler

Assembler

ABSTRACTION

Abstraction helps us deal with complexity

- Hide lower-level detail

Instruction set architecture (ISA)

- The hardware/software interface

Application binary interface

- The ISA plus system software interface

Implementation

- The details underlying and interface

UNDER THE COVERS

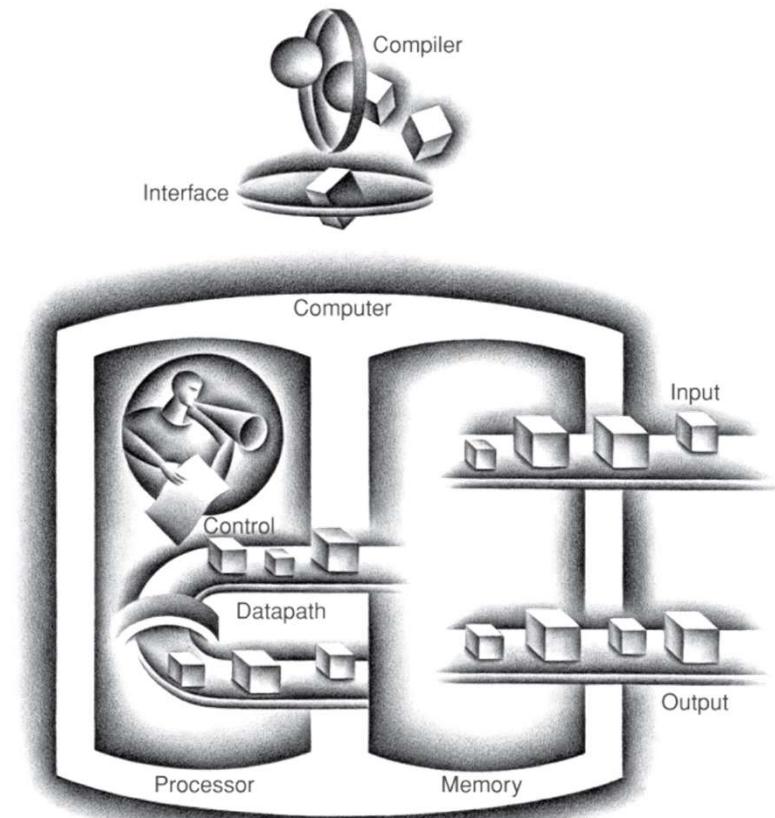
Same components for all kinds of computer

- Five classic components of a computer – input, output, memory, datapath, and control

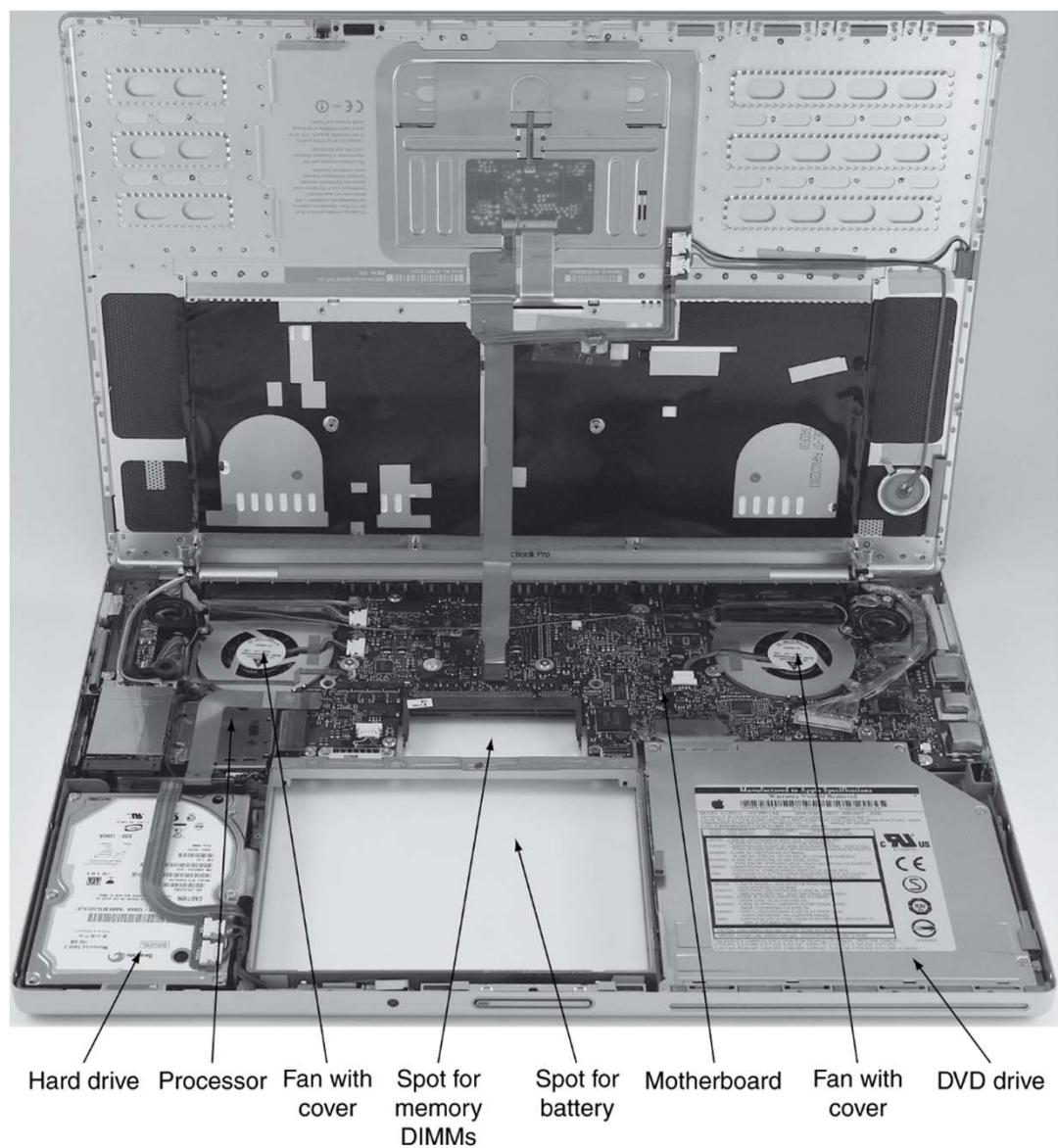
datapath + control = processor (CPU)

Input/output includes

- User-interface devices
 - Display, keyboard, mouse
- Storage devices
 - Hard disk, CD/DVD, flash
- Network adapters
 - For communicating with other computers



OPENING THE Box



A SAFE PLACE FOR DATA

Volatile main memory

- Loses instructions and data when power off

Non-volatile secondary memory

- Magnetic disk
- Optical disk (CDROM, DVD)
- SSD
- Flash memory

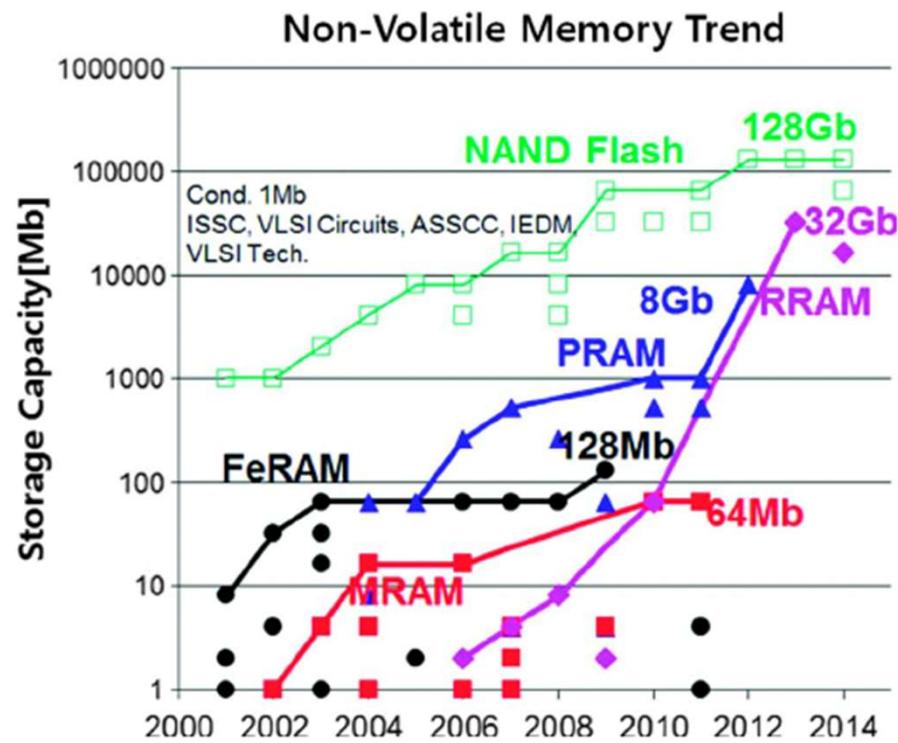


TECHNOLOGY TRENDS

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
2005	Ultra large scale IC	6,200,000,000

NETWORKS

Communication and resource sharing

Local area network (LAN): Ethernet

- Within a building

Wide area network (WAN): the Internet

Wireless network: WiFi, Bluetooth



NetFPGA 10 Gigabit Ethernet card

TECHNOLOGY SCALING ROAD MAP (ITRS)

Year	2006	2008	2010	2012	2014	2017
Feature size (nm)	65	45	32	22	14	10
Intg. Capacity (BT)	0.3	0.7	1.1	1.4	1.9	>3.2

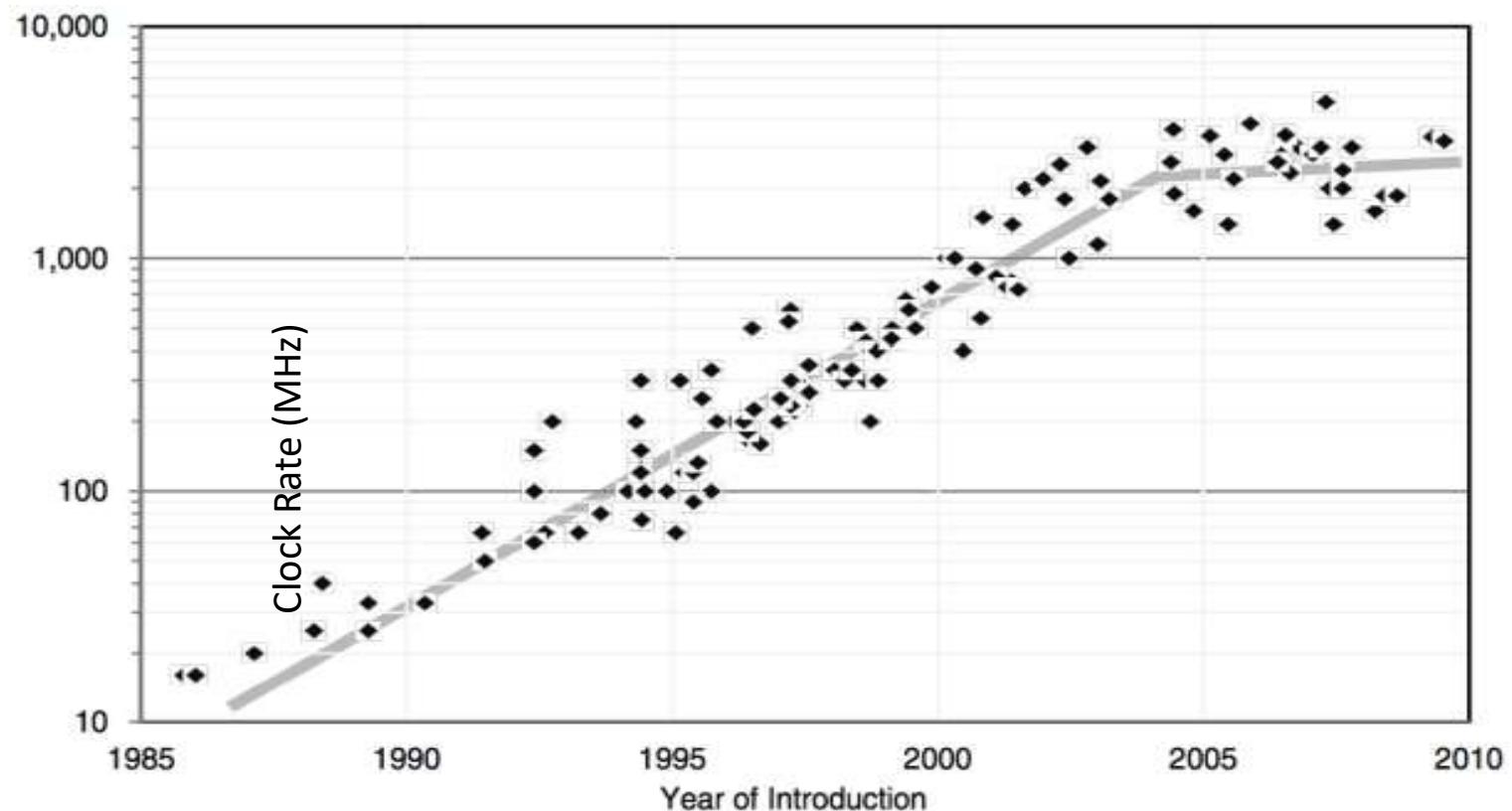
Fun facts about 45nm transistors

- 30 million can fit on the head of a pin
- You could fit more than 2,000 across the width of a human hair

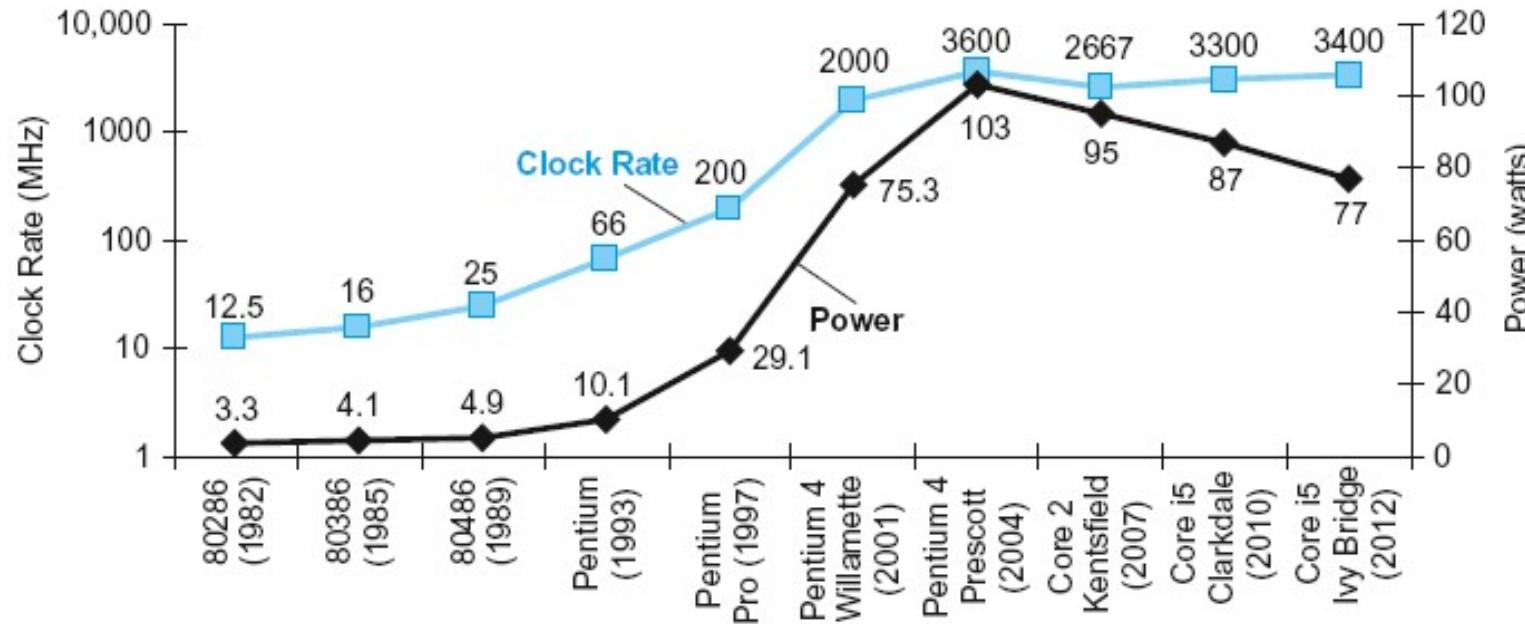
If car prices had fallen at the same rate as the price of a single transistor has since 1968, a new car today would cost about 1 cent

PROCESSOR FREQUENCY

Why processor clock rate have almost not increased in the last ten years?



POWER TRENDS



In CMOS IC technology

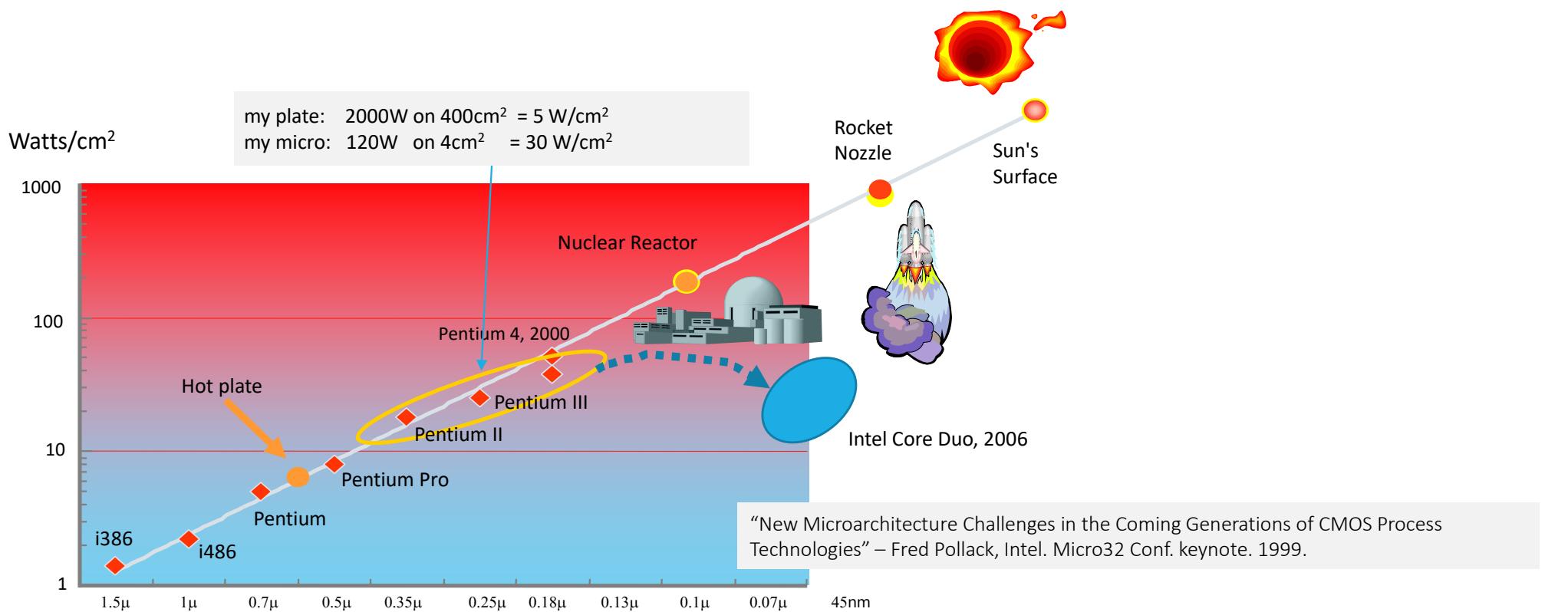
$$\text{Power} = \text{Capacitive load} \times \text{Voltage}^2 \times \text{Frequency}$$

×30

5V → 1V

×1000

WHAT HAPPENED TO CLOCK RATES AND WHY?



REDUCING POWER

Suppose a new CPU has

- 85% of capacitive load of old CPU
- 15% voltage and 15% frequency reduction

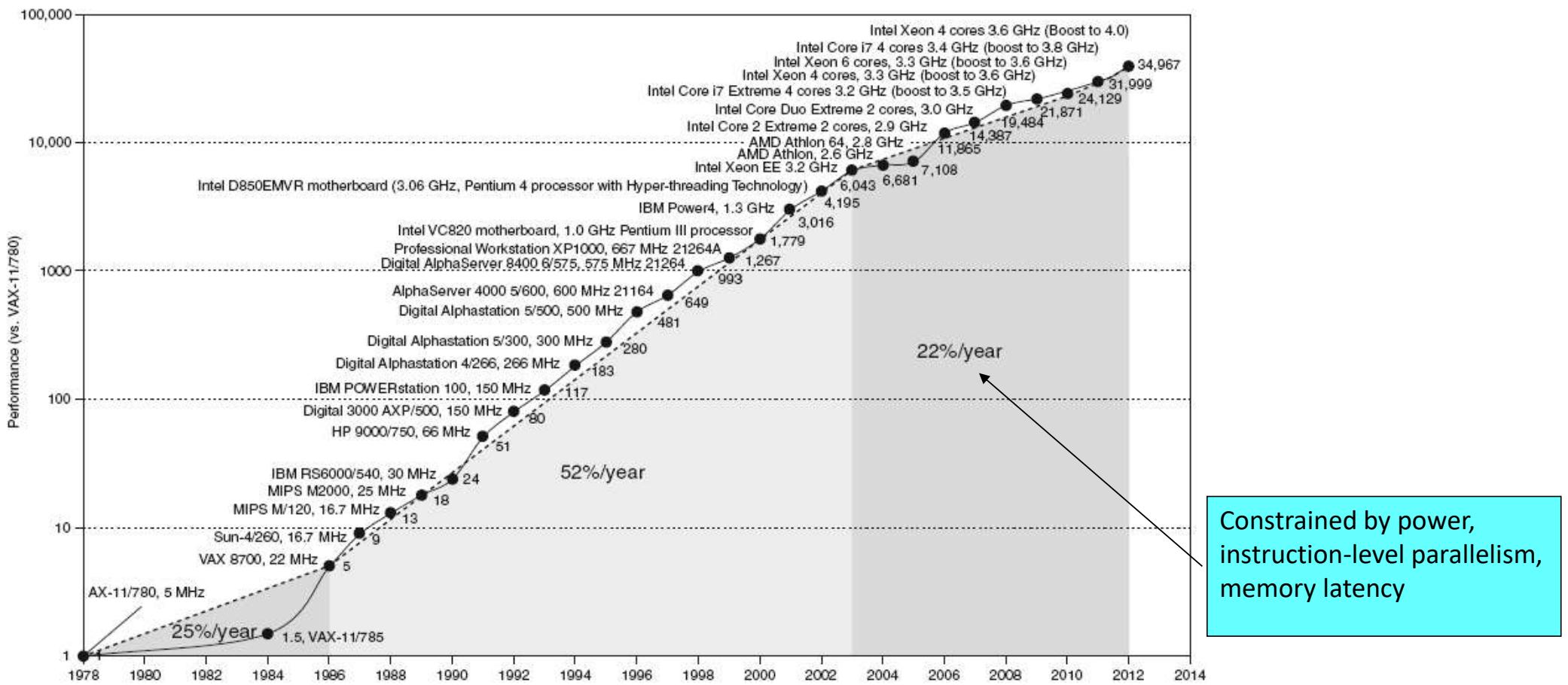
$$\frac{P_{\text{new}}}{P_{\text{old}}} = \frac{C_{\text{old}} \times 0.85 \times (V_{\text{old}} \times 0.85)^2 \times F_{\text{old}} \times 0.85}{C_{\text{old}} \times V_{\text{old}}^2 \times F_{\text{old}}} = 0.85^4 = 0.52$$

But a power wall has been reached

- We can't reduce voltage further
- We can't remove more heat

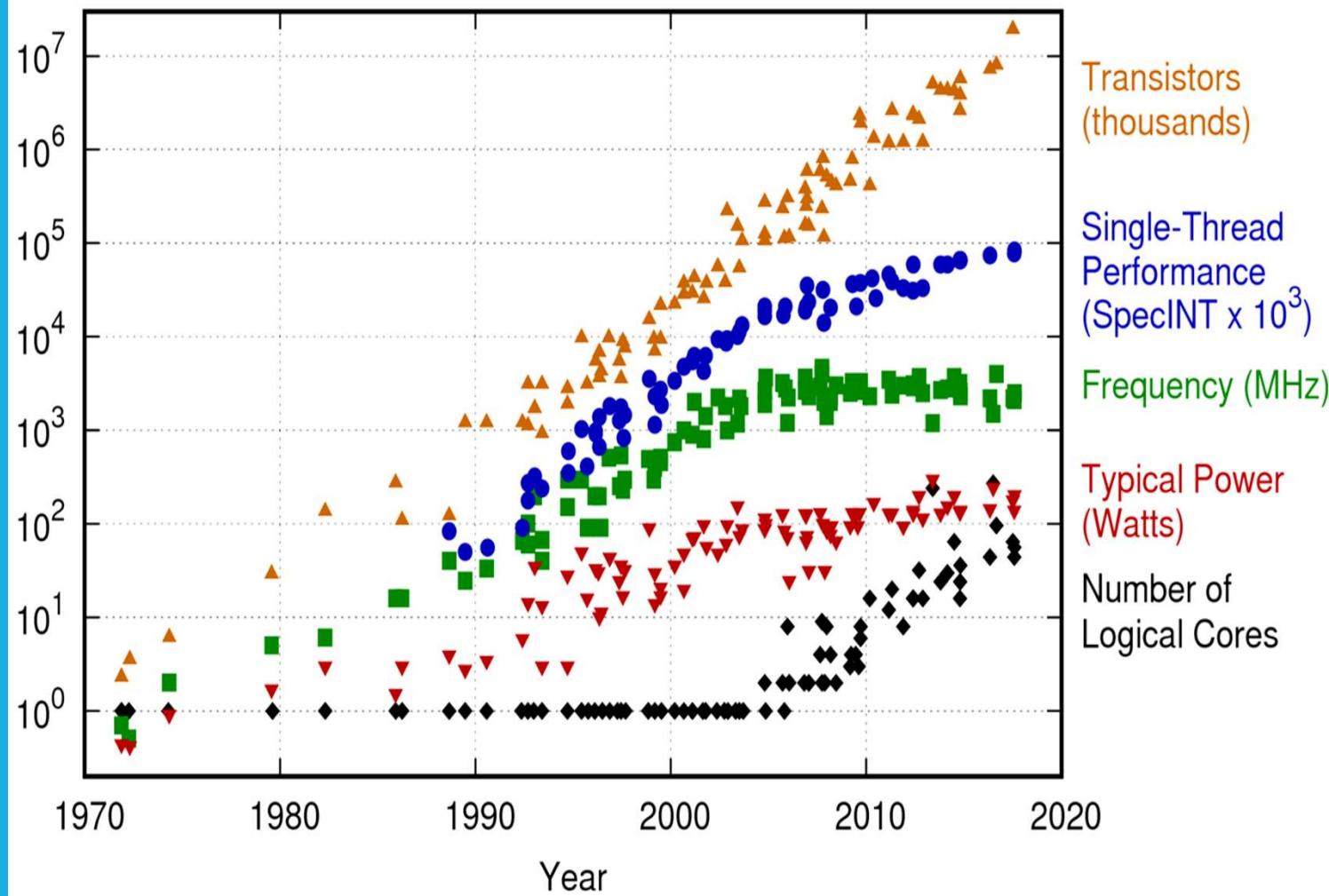
How else can we improve performance?

UNIPROCESSOR PERFORMANCE



THE SEA CHANGE: MULTIPROCESSORS

42 Years of Microprocessor Trend Data (Feb. 15, 2018)



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

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

THE SEA CHANGE: MULTIPROCESSORS

The power challenge has forced a change in the design of microprocessors

- Since 2002 the rate of improvement in the response time of programs on desktop computers has slowed from a factor of 1.5 per year to less than a factor of 1.2 per year

As of 2006 all desktop and server companies are shipping microprocessors with multiple processors – cores – per chip

Product	AMD Ryzen	Intel Kaby Lake	Samsung Exynos Octa	Qualcomm Snapdragon
Cores per chip	4-16 + GPU	4 + GPU	4 + 4	4 + GPU +DSP
Clock rate	2.0-3.8 GHz	1.8-3.0 GHz	1.3-1.8 GHz	2.3 GHz
Power	~12-95 W	~4.5-91W	< 10 W	< 10 W

Plan of record is to double the number of cores per chip per generation (about every two years)

UNDERSTANDING PERFORMANCE

POWER, THROUGHPUT AND RESPONSE TIME

UNDERSTANDING PERFORMANCE

Algorithm

- Determines number of operations executed

Programming language, compiler, architecture

- Determine number of machine instructions executed per operation

Processor and memory system

- Determine how fast instructions are executed

I/O system (including OS)

- Determines how fast I/O operations are executed

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?

RELATIVE PERFORMANCE

Define performance as $1/\text{ExecutionTime}$

X is n times faster than Y means

$$\frac{\text{Performance}_X}{\text{Performance}_Y} = \frac{\text{ExecutionTime}_Y}{\text{ExecutionTime}_X} = n$$

Example: time taken to run a program

- 10s on A, 15s on B
- $\frac{\text{ExecutionTime}_Y}{\text{ExecutionTime}_X} = \frac{15s}{10s} = 1.5$
- So A is 1.5 times faster than B

MEASURING EXECUTION TIME

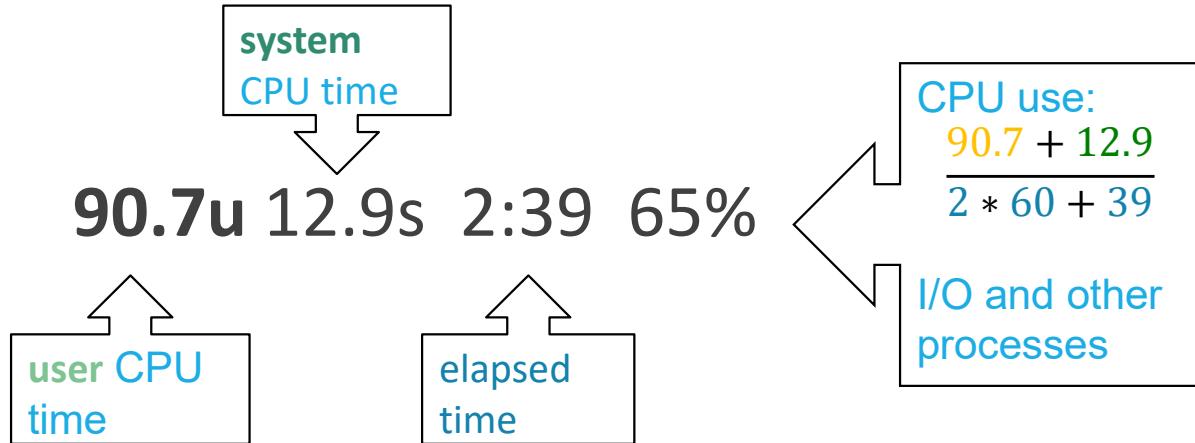
Elapsed time

- Total response time, including all aspects
 - Processing, I/O, OS overhead, idle time
- Determines system performance

CPU time

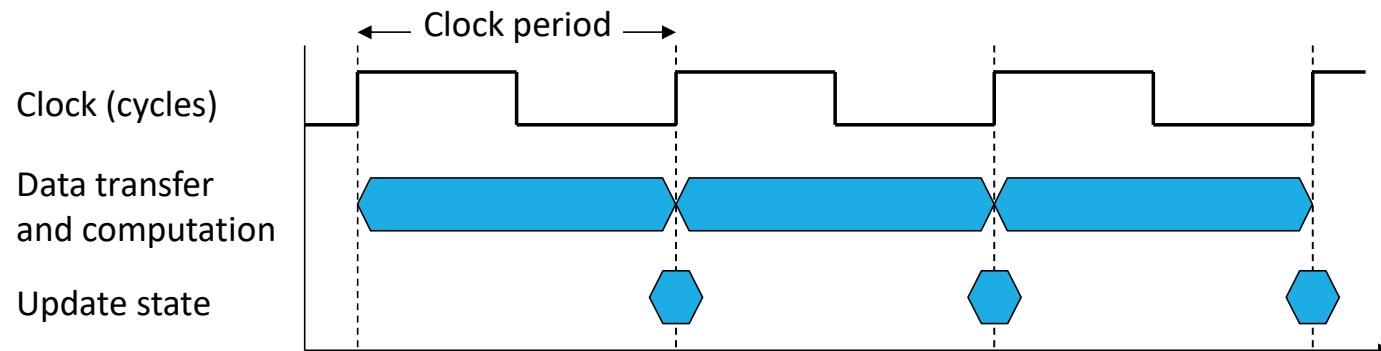
- Time spent processing a given job
 - Discounts I/O time, other jobs' shares
- Comprises user CPU time and system CPU time
- Different programs are affected differently by CPU and system performance

EXAMPLE: 'TIME' IN UNIX



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

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

Performance improved by

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

CPU TIME EXAMPLE

Computer A: 2GHz clock, 10s CPU time

Designing Computer B

- Aim for 6s CPU time
- Can do faster clock, but causes $1.2 \times$ clock cycles

How fast must Computer B clock be?

$$\text{Clock Rate}_B = \frac{\text{Clock Cycles}_B}{\text{CPU Time}_B} = \frac{1.2 \times \text{Clock Cycles}_A}{6s}$$

$$\begin{aligned}\text{Clock Cycles}_A &= \text{CPU Time}_A \times \text{Clock Rate}_A \\ &= 10s \times 2\text{GHz} = 20 \times 10^9\end{aligned}$$

$$\text{Clock Rate}_B = \frac{1.2 \times 20 \times 10^9}{6s} = \frac{24 \times 10^9}{6s} = 4\text{GHz}$$

INSTRUCTION COUNT AND CPI

$$\text{ClockCycles} = \text{InstructionCount} \times \underbrace{\text{CyclesPerInstruction}}_{\text{CPI}}$$

$$\begin{aligned}\text{CPU Time} &= \text{InstructionCount} \times \text{CPI} \times \text{ClockCycleTime} \\ &= \frac{\text{InstructionCount} \times \text{CPI}}{\text{ClockRate}}\end{aligned}$$

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 EXAMPLE

Computer A: ClockCycleTime = 250ps, CPI = 2.0

Computer B: ClockCycle Time = 500ps, CPI = 1.2

Same ISA, which is faster, and by how much?

$$\begin{aligned}\text{CPU Time}_A &= \text{Instruction Count} \times \text{CPI}_A \times \text{ClockCycleTime}_A \\ &= I \times 2.0 \times 250\text{ps} = I \times 500\text{ps}\end{aligned}$$

A is faster...

$$\begin{aligned}\text{CPU Time}_B &= \text{Instruction Count} \times \text{CPI}_B \times \text{ClockCycleTime}_B \\ &= I \times 1.2 \times 500\text{ps} = I \times 600\text{ps}\end{aligned}$$

$$\frac{\text{CPU Time}_B}{\text{CPU Time}_A} = \frac{I \times 600\text{ps}}{I \times 500\text{ps}} = 1.2$$

...by this much

CPI IN MORE DETAIL

If different instruction classes take different numbers of cycles

$$\text{ClockCycles} = \sum_{i=1}^n (\text{CPI}_i \times \text{InstructionCount}_i)$$

The weighted average CPI is

$$\text{CPI} = \frac{\text{ClockCycles}}{\text{InstructionCount}} = \sum_{i=1}^n \left(\text{CPI}_i \times \frac{\text{InstructionCount}_i}{\text{InstructionCount}} \right)$$

Relative frequency

CPI EXAMPLE

Alternative compiled code sequences using instructions in classes A, B, C

Class	A	B	C
CPI for class	1	2	3
IC in sequence 1	2	1	2
IC in sequence 2	4	1	1

- Sequence 1: IC = 5
 - Clock Cycles
 $= 2 \times 1 + 1 \times 2 + 2 \times 3$
 $= 10$
 - Avg. CPI = $10/5 = 2.0$
- Sequence 2: IC = 6
 - Clock Cycles
 $= 4 \times 1 + 1 \times 2 + 1 \times 3$
 $= 9$
 - Avg. CPI = $9/6 = 1.5$

PERFORMANCE SUMMARY

$$\text{CPU Time} = \frac{\text{InstructionCount}}{\text{Program}} \times \frac{\text{Clock Cycles}}{\text{Instruction}} \times \frac{\text{Seconds}}{\text{ClockCycleTime}}$$

IC CPI CC ($= T_c$)

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

PITFALL: AMDAHL'S LAW

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

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

Example: if multiply accounts for 80s of 100s,

- How much improvement in multiply performance is needed to get 5x overall?
- $\frac{100s}{5} = 20s = \frac{80}{n} + 20 \rightarrow \text{can't be done!}$

TOTAL SPEED-UP

How much speed-up is obtained by improving a part?

$$S = \frac{T_{CPU\ original}}{T_{CPU\ improved}} = \frac{1}{\frac{F_e}{S_e} + (1 - F_e)}$$

where F_e is the fraction of enhancement and S_e is the factor of enhancement.

TOTAL SPEED-UP

How much speed-up is obtained by improving a part?

$$S = \frac{T_{CPU\ original}}{T_{CPU\ improved}} = \frac{1}{\frac{F_e}{S_e} + (1 - F_e)}$$

where F_e is the fraction of enhancement and S_e is the factor of enhancement.

The opportunity for improvement is affected by how much time the modified event consumes

- Corollary 1: make the common case fast

The performance enhancement possible with a given improvement is limited by the amount the improved feature is used

- Corollary 2: the bottleneck will limit the improvement

PITFALL: MIPS AS A PERFORMANCE METRIC

MIPS: Millions of Instructions Per Second

$$MIPS = \frac{InstructionCount}{ExecutionTime \times 10^6} = \frac{ClockRate}{CPI \times 10^6}$$

Doesn't account for

- Differences in ISAs between computers
- Differences in complexity between instructions

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)

$$\sqrt[n]{\prod_{i=1}^n \text{ExecutionTimeRatio}_i}$$

SPEC CPU 2017

Since July 2017 a new benchmark is available

- A new metric, SPECrate, is added
 - Former metric is called SPECspeed
- It includes 4 new suites:
 - Pairs of int/fp with speed/rate
 - Base and peak metrics can be chosen

Name	Description	IC×10 ⁹	CPI	Tc (ns)	Tcpu(s)	Ref time(s)	SPECratio
perl	Interpreted string processing	2118	0.75	0.40	637	9777	15.3
bzip2	Block-sorting compression	2389	0.85	0.40	817	9650	11.8
gcc	GNU C Compiler	1050	1.72	0.40	722	8050	11.1
mcf	Combinatorial optimization	336	10.00	0.40	1345	9120	6.8
go	Go game (AI)	1658	1.09	0.40	721	10490	14.6
hmmer	Search gene sequence	2783	0.80	0.40	890	9330	10.5
sjeng	Chess game (AI)	2176	0.96	0.40	835	12100	14.5
libquantum	Quantum computer simulation	1623	1.61	0.40	1047	20720	19.8
h264avc	Video compression	3102	0.80	0.40	993	22130	22.3
omnetpp	Discrete event simulation	587	2.94	0.40	690	6250	9.1
astar	Games/path finding	1082	1.79	0.40	773	7020	9.1
xalancbmk	XML parsing	1058	2.70	0.40	1143	6900	6.0
Geometric mean	High cache miss rates						11.7

CINT2006 FOR OPTERON X4 2356

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

CINT2006 FOR INTEL CORE i7 920

SPEC POWER BENCHMARK

Power consumption of server at different workload levels

- Performance: ssj_ops/sec
- Power: Watts (Joules/sec)

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

SPECPOWER_SSJ2008 FOR X4

Target Load %	Performance (ssj_ops/sec)	Average Power (Watts)
100%	231,867	295
90%	211,282	286
80%	185,803	275
70%	163,427	265
60%	140,160	256
50%	118,324	246
40%	920,35	233
30%	70,500	222
20%	47,126	206
10%	23,066	180
0%	0	141
Overall sum	1,283,590	2,605
$\Sigma ssj_ops / \Sigma power$		493

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
$\Sigma ssj_ops / \Sigma power =$		2,490

FALLACY: Low POWER AT IDLE

Look back at X4 power benchmark

- At 100% load: 295W
- At 50% load: 246W (83%)
- At 10% load: 180W (61%)

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

CONCLUDING REMARKS

Cost/performance is improving

- Due to underlying technology development

Hierarchical layers of abstraction

- In both hardware and software

Instruction set architecture

- The hardware/software interface

Execution time: the best performance measure

Power is a limiting factor

- Use parallelism to improve performance