



ADVANCED COMPUTER SYSTEM & ARCHITECTURE

Fundamentals of Quantitative Design and
Analysis

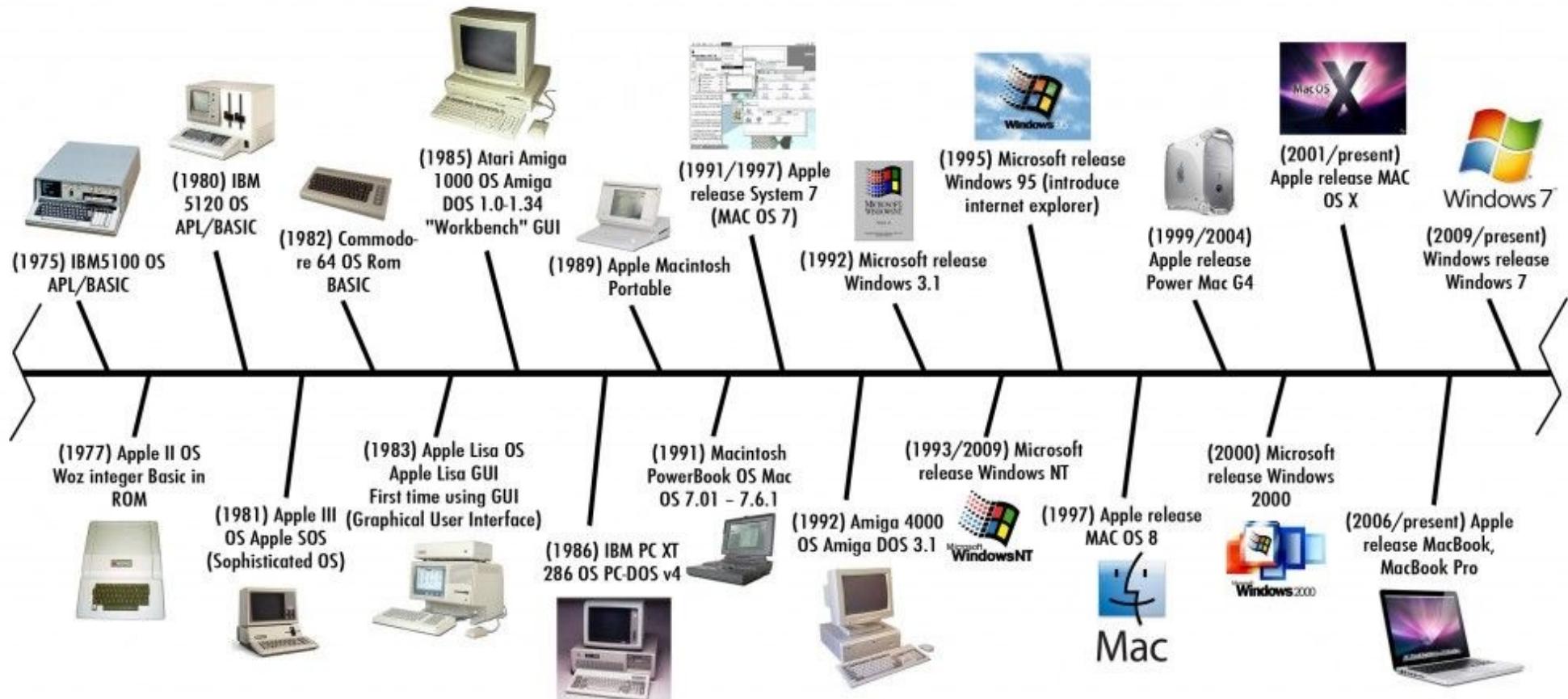
Assoc Prof Mohd Shahizan Othman

Innovating Solutions

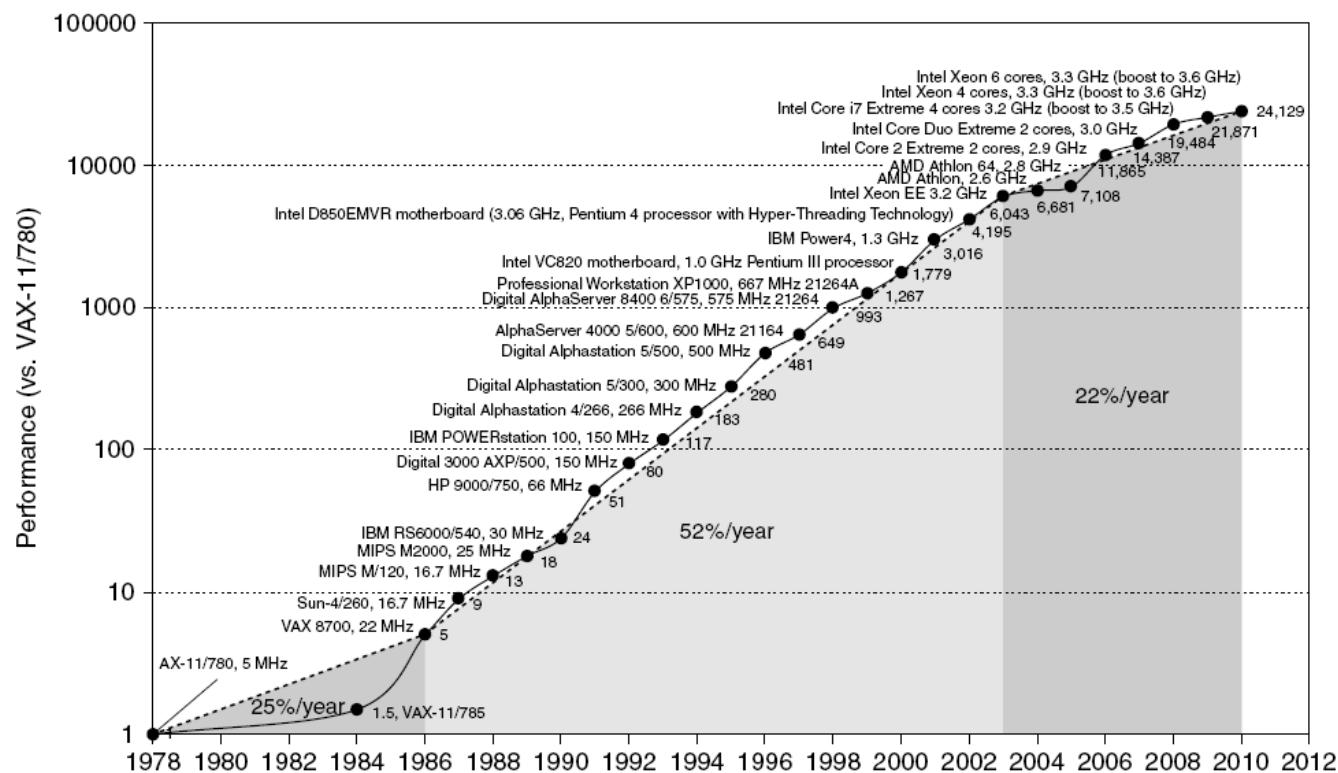
Computer Technology

- 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

History of Computer Technology



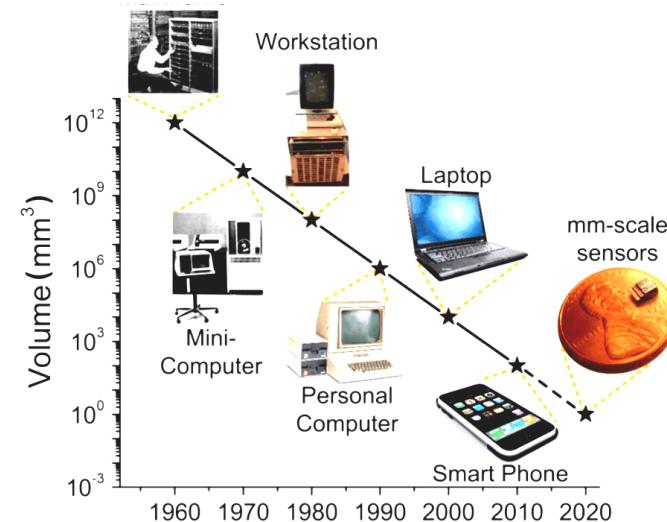
Historical Microprocessor Performance



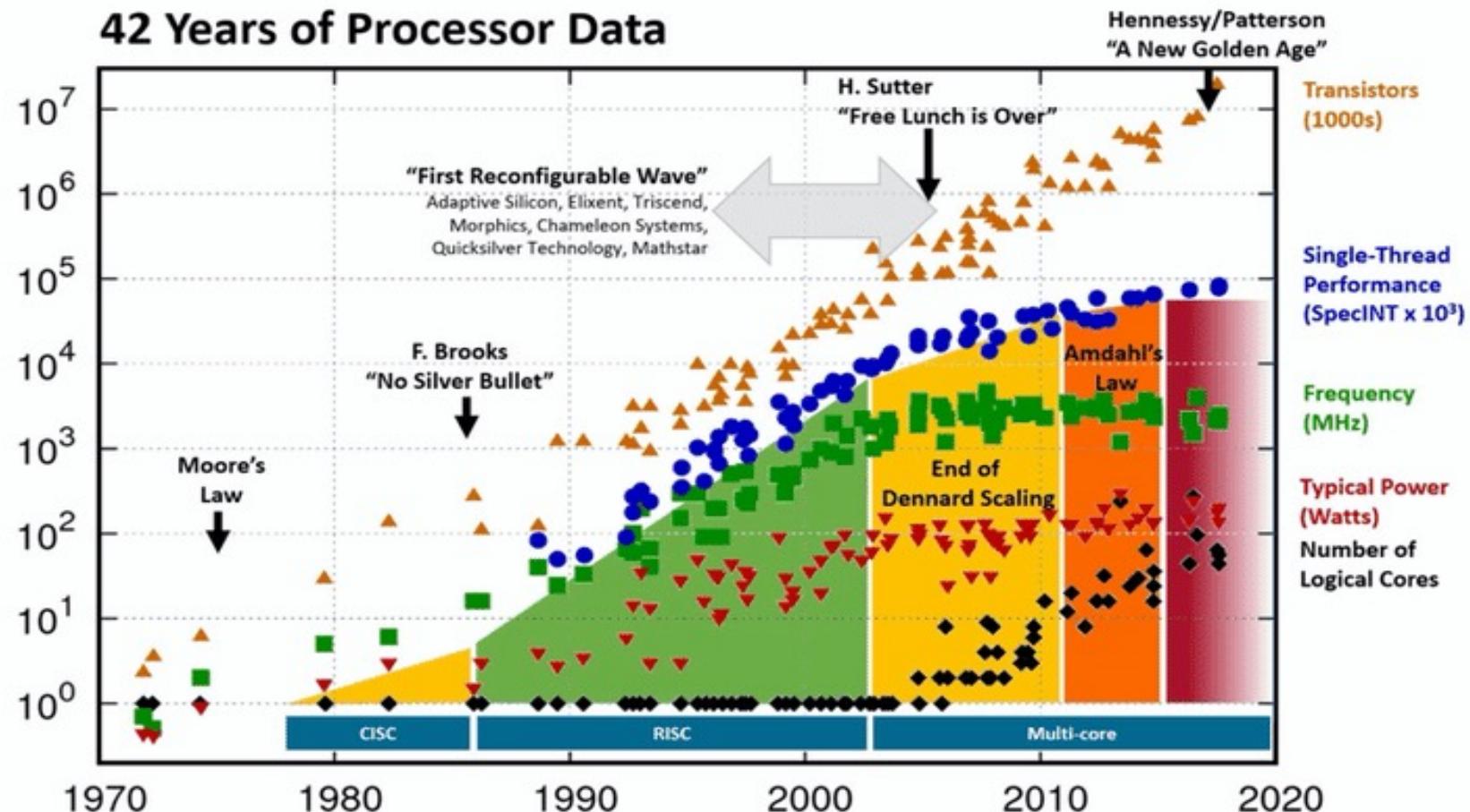
Source: H&P textbook

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



42 Years of Processor Data



Hennessy and Patterson, Turing Lecture 2018, overlaid over "42 Years of Processors Data"

<https://www.karlrupp.net/2018/02/42-years-of-microprocessor-trend-data/>; "First Wave" added by Les Wilson, Frank Schirrmeister

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

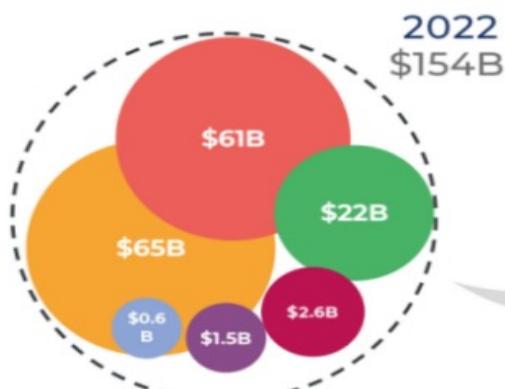
Processor Technology Trends

- Transistor density increases by 35% per year and die size increases by 10-20% per year... more functionality
- Transistor speed improves linearly with size (complex equation involving voltages, resistances, capacitances)
- Wire delays do not scale down at the same rate as logic delays
- The power wall: it is not possible to consistently run at higher frequencies without hitting power/thermal limits; fancy cooling required beyond ~150W

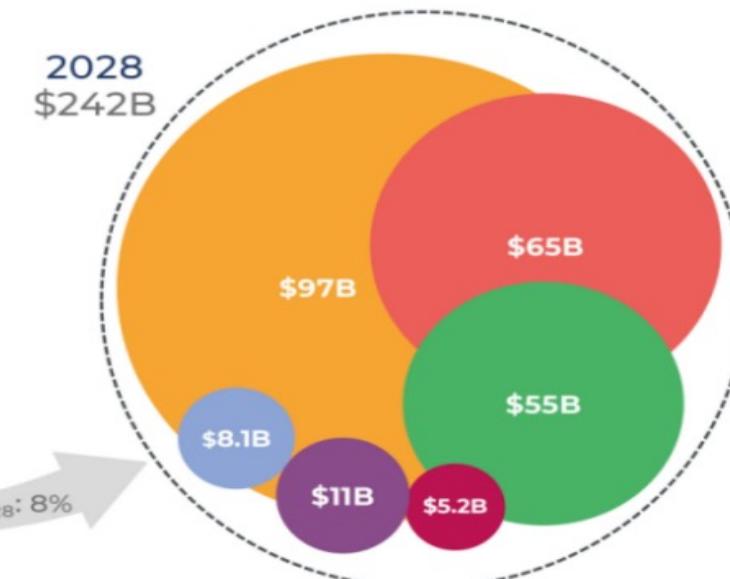
2022-2028 PROCESSOR REVENUE FORECAST BY TYPE OF PROCESSOR

Source: Status of the Processor Industry 2023 report, Yole Intelligence, 2023

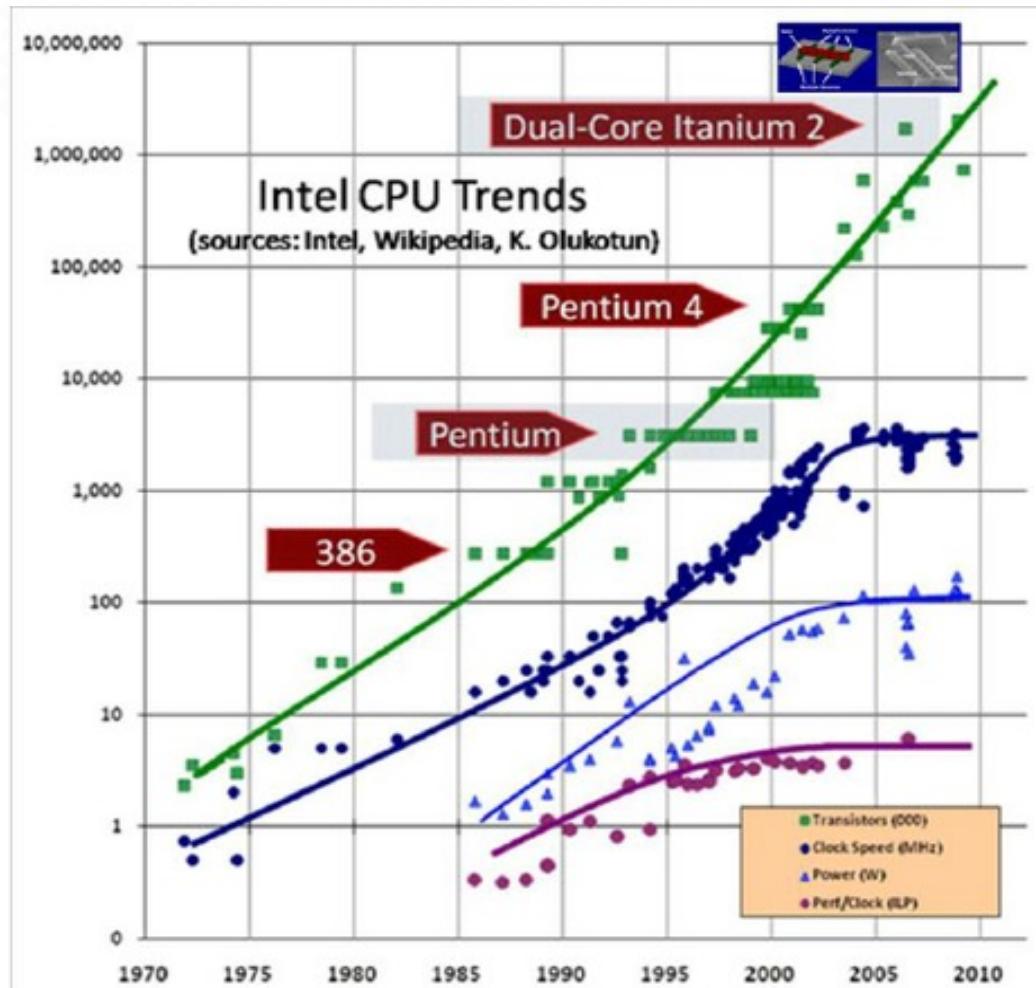
- Application Processing Unit (APU)
- Central Processing Unit (CPU)
- Graphics Processing Unit (GPU)
- Data Processing Unit (DPU)
- AI ASIC*
- SoC FPGA*



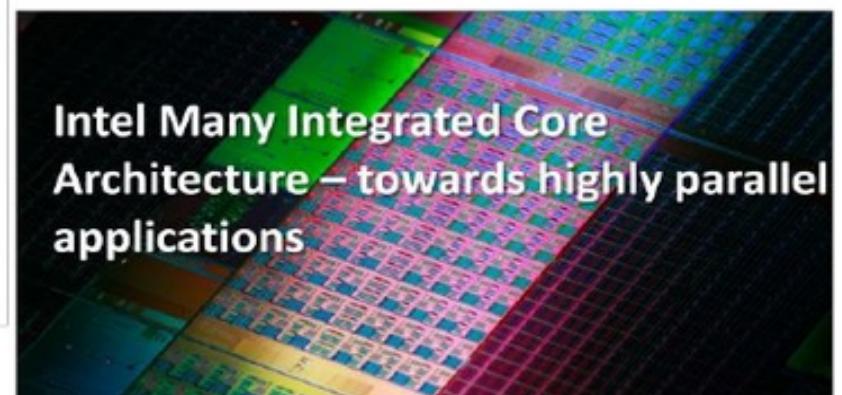
Note: Microcontroller Units are not included



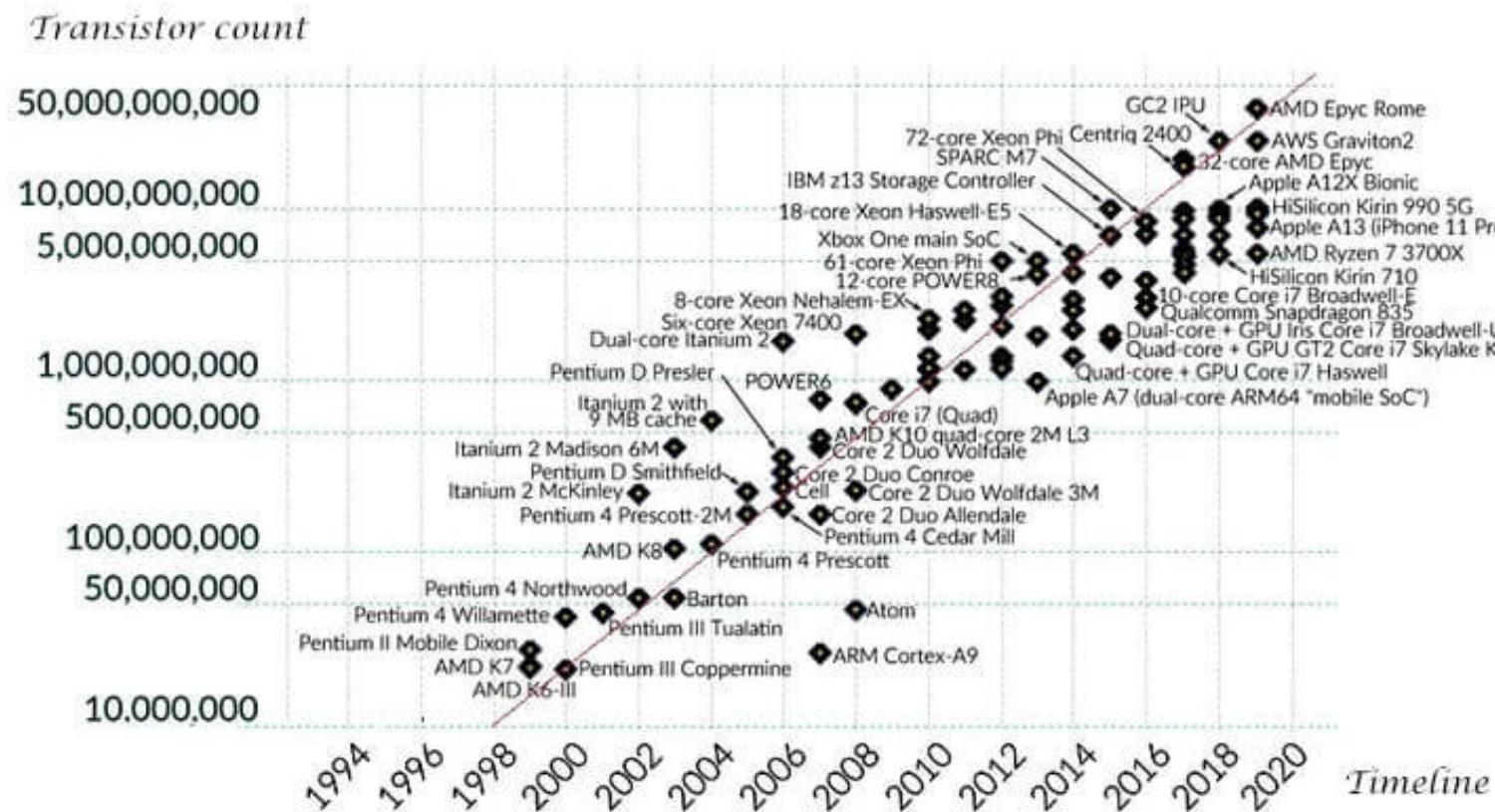
*AI ASIC : Application-Specific Integrated Circuit for Artificial Intelligence
SoC FPGA : System-on-Chip Field-Programmable Gate Array



Tesla k10 GPU (NVIDIA) – state of the art in GPU technology



Trends in Transistor Count within CPU



What Helps Performance?

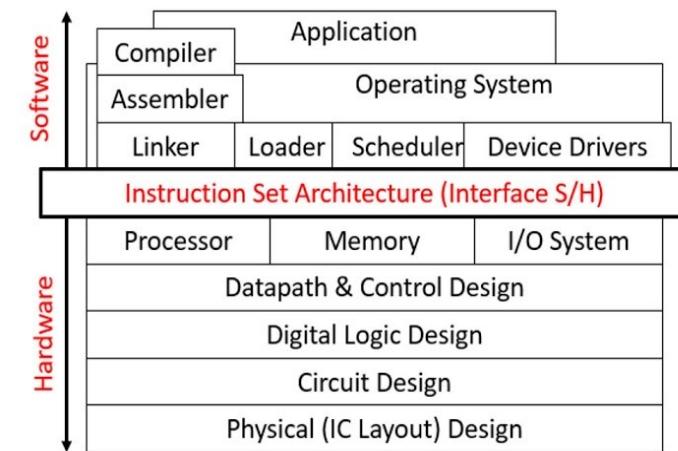
- In a clock cycle, can do more work -- since transistors are faster, transistors are more energy-efficient, and there's more of them
- Better architectures: finding more parallelism in one thread, better branch prediction, better cache policies, better memory organizations, more thread-level parallelism, moving computations to memory, accelerating some kernels, ...

Points to Note

- The 52% growth per year is because of faster clock speeds and architectural innovations (led to 25x higher speed)
- Clock speed increases have dropped to 1% per year in recent years
- The 22% growth includes the parallelization from multiple cores
- Moore's Law: transistors on a chip double every 18-24 months

Defining Computer Architecture

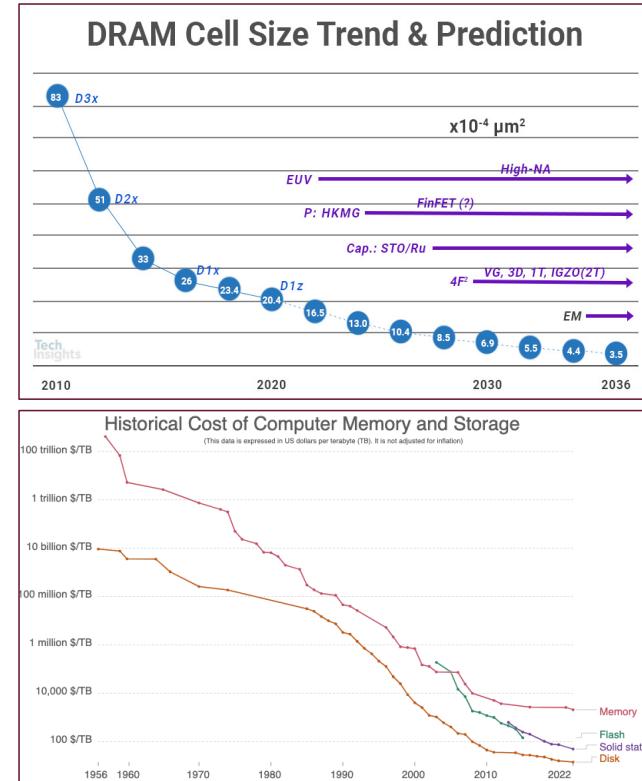
- “Old” view of computer architecture:
 - Instruction Set Architecture (ISA) design
 - i.e. decisions regarding:
 - registers, memory addressing, addressing modes, instruction operands, available operations, control flow instructions, instruction encoding
- “Real” computer architecture:
 - Specific requirements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



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

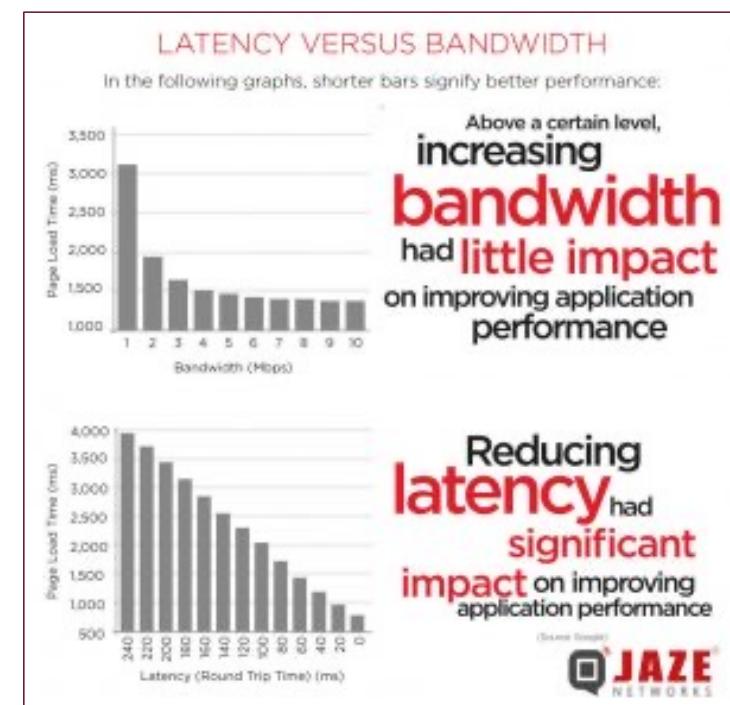
Copyright © 2019, Elsevier Inc. All rights reserved.



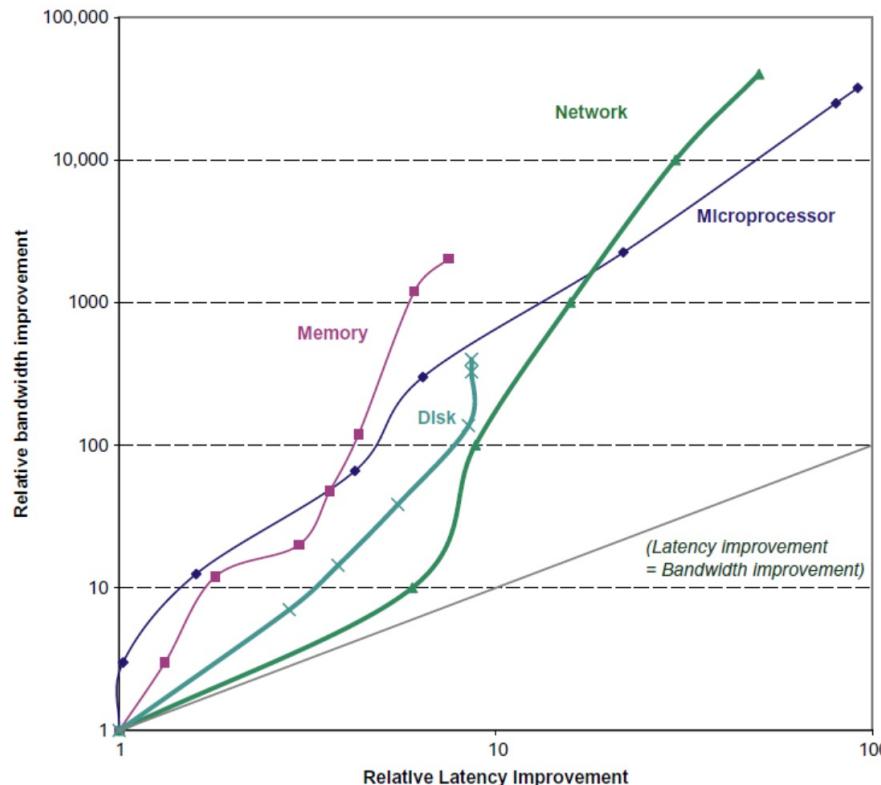
Bandwidth and Latency

- 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

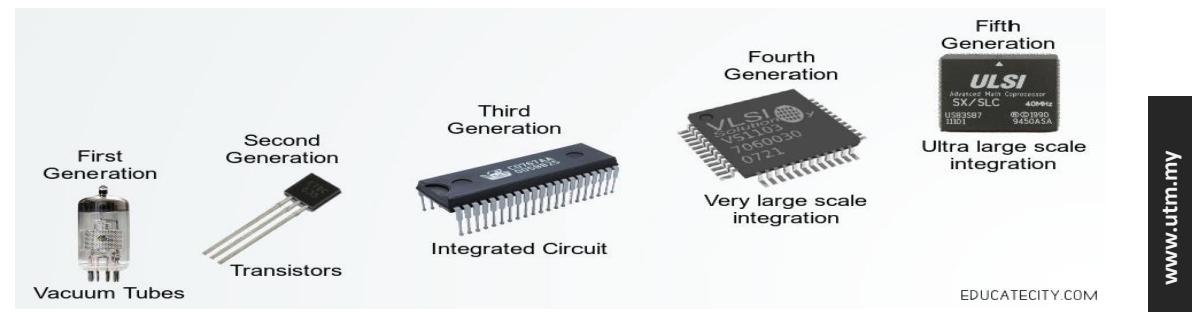


Log-log plot of bandwidth and latency milestones

Copyright © 2019, Elsevier Inc. All rights reserved.

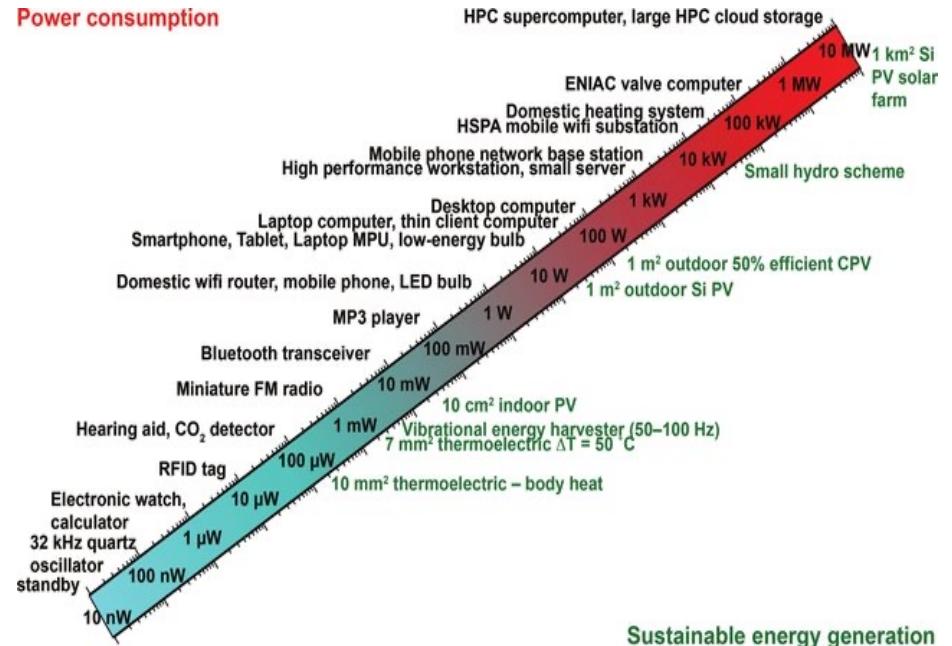
Transistors and Wires

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

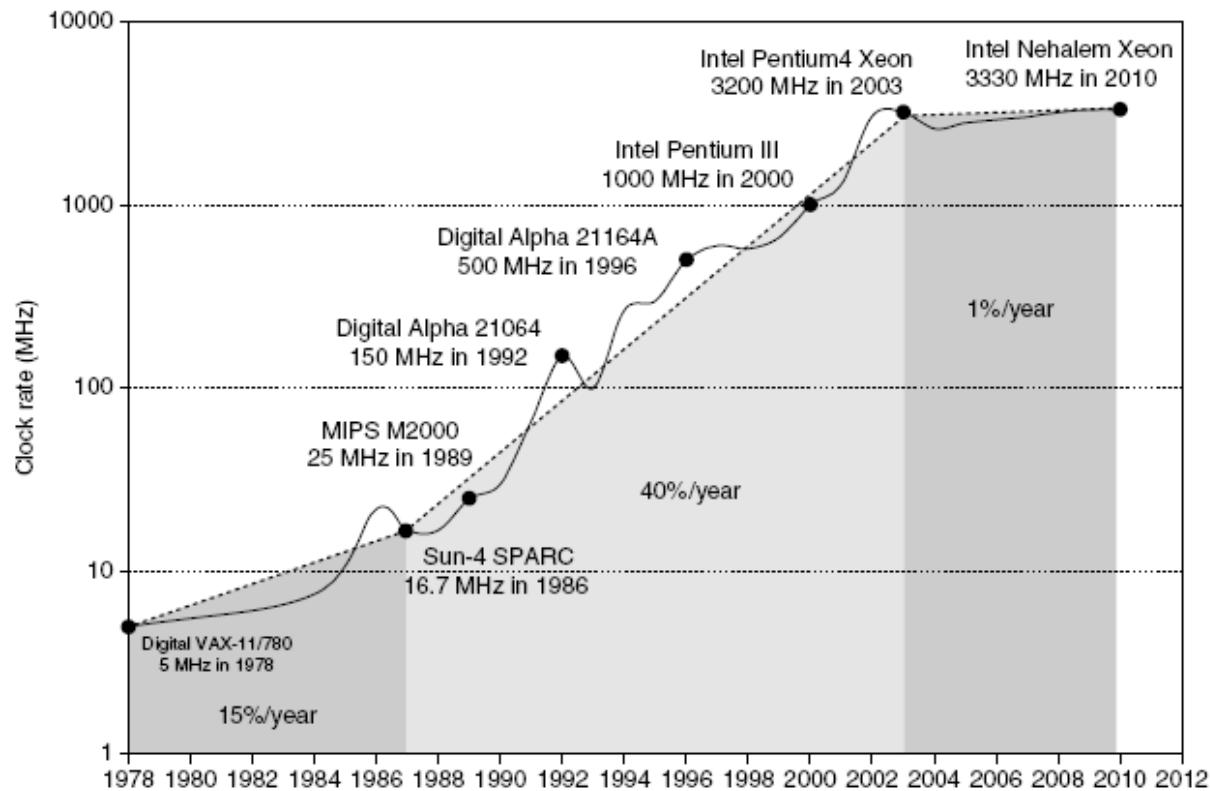


Power and Energy

- 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



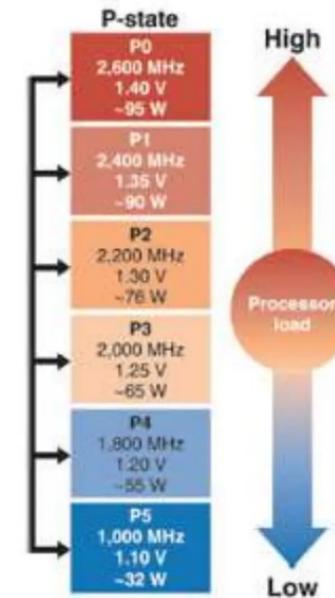
Clock Speed Increases



Source: H&P textbook

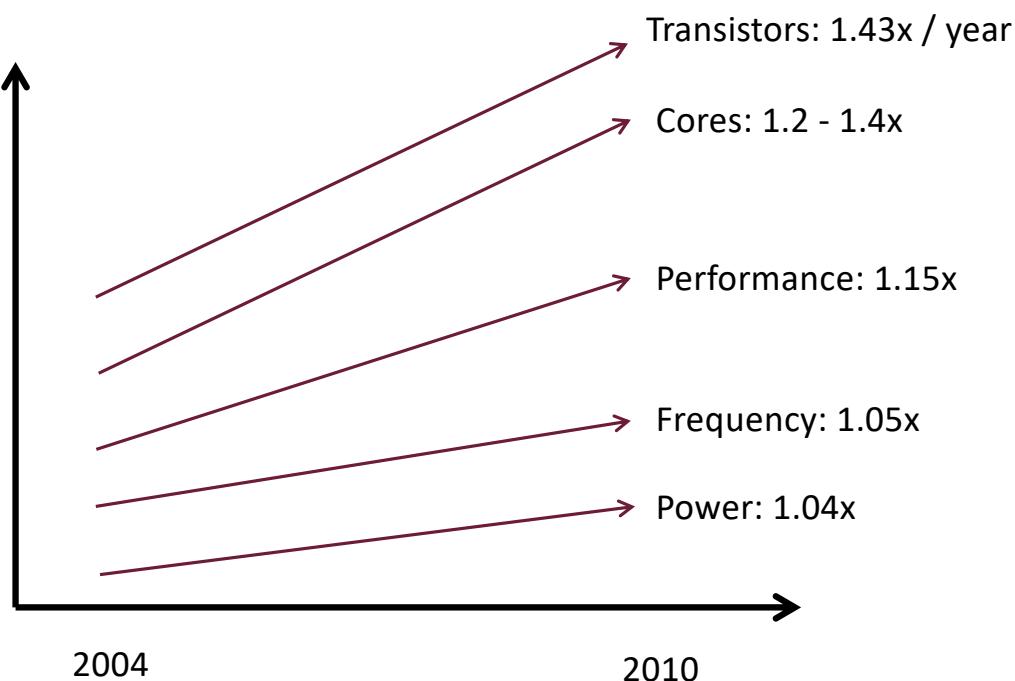
Reducing Power

- Techniques for reducing power:
 - Do nothing well
 - Dynamic Voltage-Frequency Scaling
 - e.g, AMD Opteron
 - Low power state for DRAM, disks
 - sleep mode
 - Overclocking, turning off cores
 - Intel i7, AMD Ryzen



Source: AMD

Recent Microprocessor Trends



Source: Micron University Symp.

Integrated Circuit Cost

- Integrated circuit

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

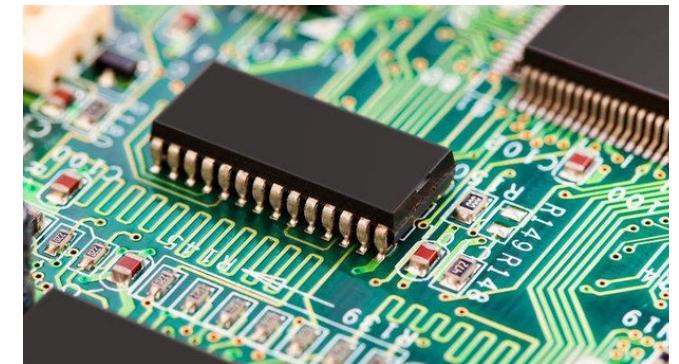
$$\text{Cost of die} = \frac{\text{Cost of wafer}}{\text{Dies per wafer} \times \text{Die yield}}$$

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

- 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)



More Diverse Platforms

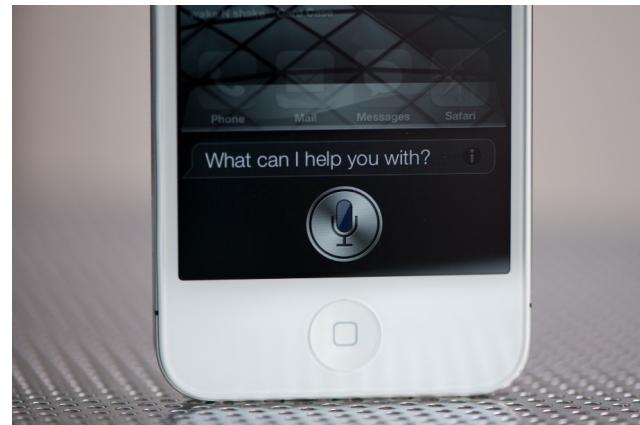
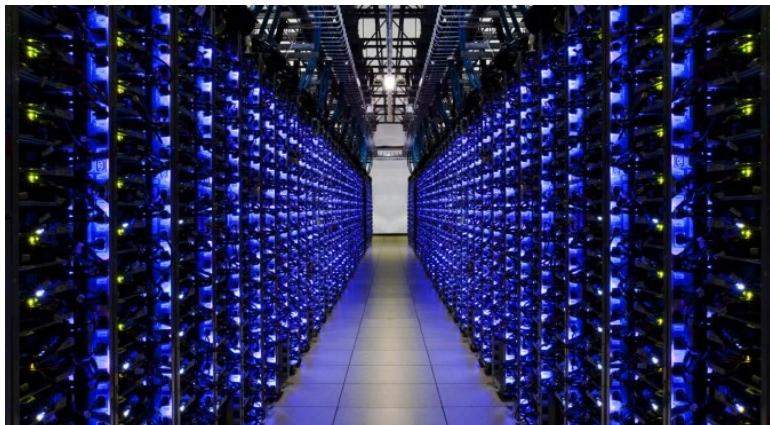


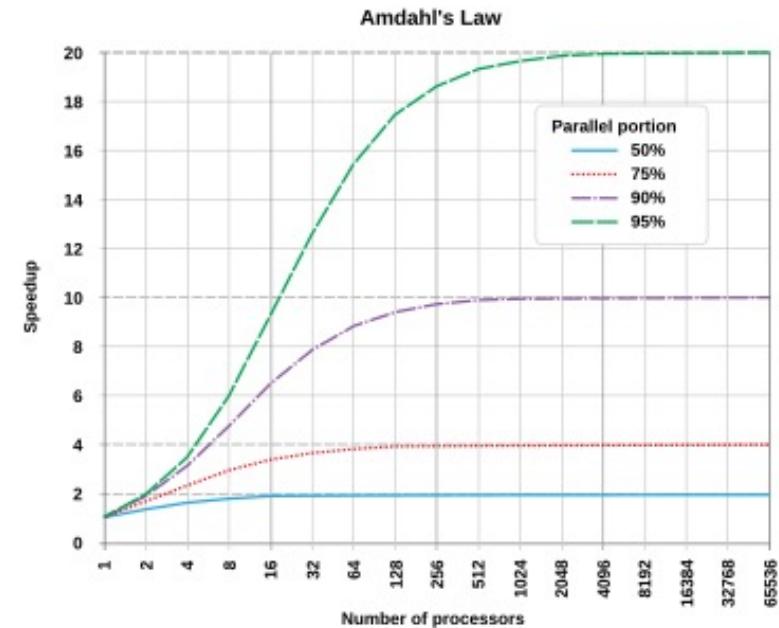
Image credits: uber, extremetech, anandtech

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

$$\text{Execution time}_{\text{new}} = \text{Execution time}_{\text{old}} \times \left((1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}} \right)$$

$$\text{Speedup}_{\text{overall}} = \frac{1}{(1 - \text{Fraction}_{\text{enhanced}}) + \frac{\text{Fraction}_{\text{enhanced}}}{\text{Speedup}_{\text{enhanced}}}}$$



Principles of Computer Design

- The Processor Performance Equation

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

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

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

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

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

Principles of Computer Design

- Different instruction types having different CPIs

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

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

Principles of Computer Design

- Different instruction types having different CPIs

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

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

Fallacies and Pitfalls

- 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

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

Where Are We Headed?

Modern trends:

- Clock speed improvements are slowing (power constraints)
- Difficult to further optimize a single core for performance
- Multi-cores: each new processor generation will accommodate more cores
- Need better programming models and efficient execution for multi-threaded applications
- Need better memory hierarchies
- Need greater energy efficiency
- Reduced data movement
- Emergence of new metrics: security, reliability
- Emergence of new workloads: ML, graphs, genomics



UTM
UNIVERSITI TEKNOLOGI MALAYSIA

*Innovating Solutions
Menginovasi Penyelesaian*