

UT 1. Introduction to Computer Architecture

Tema 1.1 Concept of Computer Architecture

J. Duato, J. Flich, P. López, V. Lorente,
A. Pérez, S. Petit, J.C. Ruiz, S. Sáez, J. Sahuquillo

Department of Computer Engineering
Universitat Politècnica de València



Contents

- 1 Concept of Computer Architecture
- 2 Computer Requirements
- 3 Technology, power consumption and cost
- 4 Evolution of processor performance
- 5 Types of computers
- 6 Máster en Ingeniería de Computadores y Redes

Bibliography



John L. Hennessy and David A. Patterson.

Computer Architecture, Sixth Edition: A Quantitative Approach.

Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 6 edition, 2018.

Contents

- 1 Concept of Computer Architecture
- 2 Computer Requirements
- 3 Technology, power consumption and cost
- 4 Evolution of processor performance
- 5 Types of computers
- 6 Máster en Ingeniería de Computadores y Redes

Defines the *hardware* of a computer attending to its requirements, as a design trade-off between performance, available technology and cost. It encompasses with the following levels:

Instruction Set Architecture (ISA).

It includes everything that must be known by assembler programmers: description of instructions and their use, logical organization of memory, etc.

Processor organization

It describes the logical elements enabling the execution of instructions: registers, decoders, ALU operators, memory interfaces, etc.

Implementation

It defines the computer implementation: VLSI design, cooling, power supply, packacking technologies, connections, etc.

- Along the classic period (until the 70's), each computer architecture level was designed in isolation by a different specialist.
- Today, computer architecture engineering is considered a transversal discipline, encompassing all the aforementioned three levels. It focuses on the design of programmable machines that must execute a (foreseeable or not) set of programs in a correct and efficient way.
- The three levels are dependent: design decisions at one level may affect the others.

The computer engineer task

- consider the expected requirements
 - identify existing technology, energy and cost limitations
 - provide the best possible design
- quantify performance, cost and other features, compare and select

Contents

- 1 Concept of Computer Architecture
- 2 Computer Requirements**
- 3 Technology, power consumption and cost
- 4 Evolution of processor performance
- 5 Types of computers
- 6 Máster en Ingeniería de Computadores y Redes

In order to design a computer, an architect must consider

- The type of required computer [▶ go](#)
- The degree of compatibility with other existing computers

Source code	More flexible design, Need of new compilers
Binary	Existing <i>ISA</i> (lack of flexibility) No new software required

- Operating system requirements

Address space	Limits applications size
Memory management	Paging, segmentation, ...
Protection	
Process Management	Multitasking support

(cont.)

■ Market standards

Floating point	IEEE 754
E/S	SATA, SCSI, PCI Express...
Operating systems	Linux, Windows, Mac OSX ...
Networks	Ethernet, InfiniBand ...
Programming languages	C, C++, Java, FORTRAN, ...

Contents

- 1 Concept of Computer Architecture
- 2 Computer Requirements
- 3 Technology, power consumption and cost
- 4 Evolution of processor performance
- 5 Types of computers
- 6 Máster en Ingeniería de Computadores y Redes

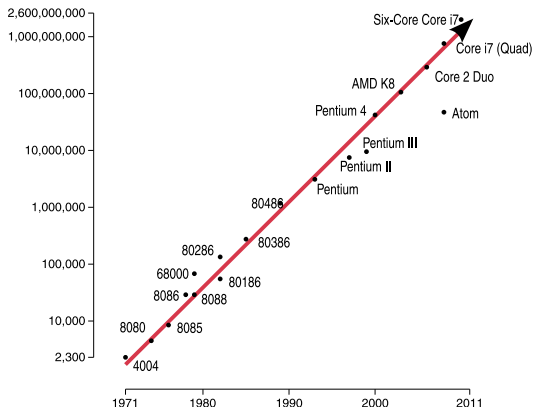
Designers must take into account:

- Available technology
- Power and energy limitations
- Cost

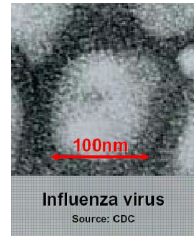
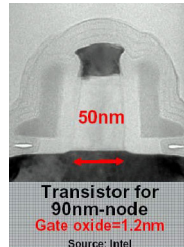
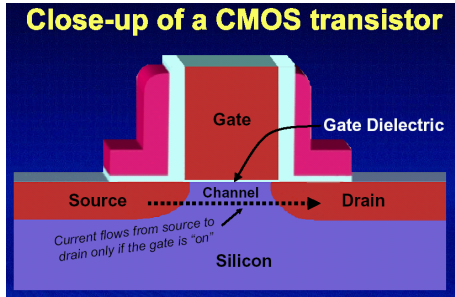
Chip Technology: Moore's Law

The number of transistors per chip increases 40–55% per year (2X every 18 months / 2 years). Two main causes:

- Density of transistors increases by 35% per year
- Die size increases by 10–20% every year



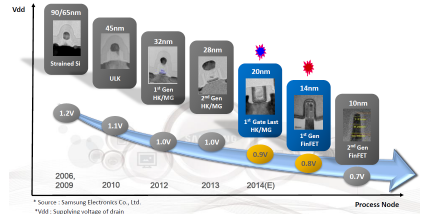
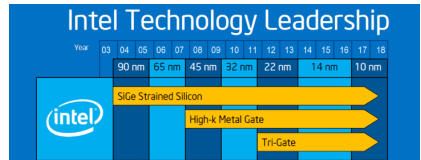
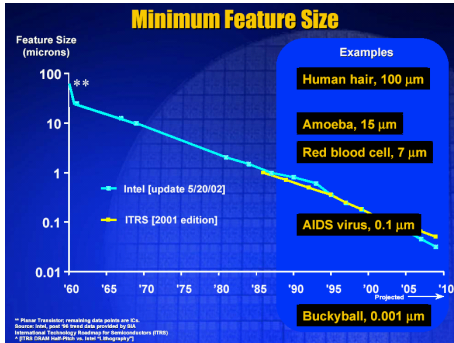
Feature size (transistor size)



- The number of transistors grows quadratically with the reduction of their feature size
- Transistor speed increases linearly with the reduction of their feature size

Feature size (transistor size) (cont.)

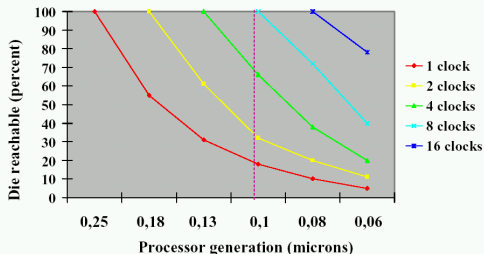
It is measured in *nm* ($10^{-9}m$).



Rise of propagation delay in interconnections

Propagation delay relates to wire resistance and capacitance ($R \cdot C$).

- As feature size is reduced, wire section (width and height) also does, thus increasing R
- Although capacitance of conductor surface decreases, capacitance coupling between wires (also named Interwinding capacitance) increases, thus augmenting C



→ The fraction of the chip reachable in a single cycle is reduced as the feature size does: 100% with $0.25\mu\text{m}$ to 5% with $0.06\mu\text{m}$.

Power consumption and heat dissipation

Power per transistor = $Pw_{\text{dynamic}} + Pw_{\text{static}}$

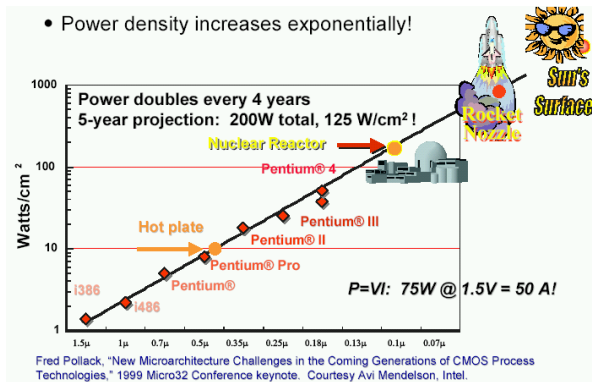
- $Pw_{\text{dynamic}} = \frac{1}{2} \cdot C \cdot V^2 \cdot f \rightarrow$ proportional to frequency (power switching)

Solution: reduce supply voltage, but a minimum is required for each frequency

- During last 24 years, supply voltage has been reduced from 12 to 1.1 V. This reduces Pw_{dynamic} by a factor of $\frac{12^2}{1.1^2} = 119X$
- But there is a minimum voltage value \rightarrow remaining margin for reduction is low (from 1.1 to 0.7 V): $\frac{1.1^2}{0.7^2} = 2.5X$
- $Pw_{\text{static}} = I_{\text{leakage}} \cdot V$
 - As feature size decreases, $\uparrow I_{\text{leakage}}$
 - As the number of transistors increases, the Pw_{static} also does.

Power consumption and heat dissipation (cont.)

The increase in the number of transistors and frequency prevails over the reduction in voltage supply and capacitance → from 0.01 W in first microprocessors to around 130 W in an Itanium2 processor.

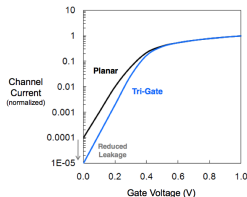
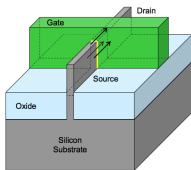


Power consumption and heat dissipation (cont.)

Implications:

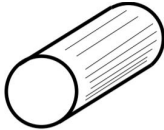
- Microprocessor power (current) distribution. Modern microprocessors integrate hundreds of power supply pins.
- Heat dissipation and cooling
- Development of new materials to reduce the leakage current. Dielectric improvement (*high-k* or *Tri-Gate* from Intel).

Tri-Gate Transistor:

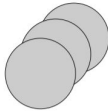


Cost

■ Cost of an integrated circuit.



BARRA DE SILICIO



OBLEA



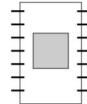
DADOS



DADO
(defectos)



DADO



CIRCUITO INTEGRADO

Final cost $\approx f$ (die surface⁴)

Die surface = f (design complexity)

Cost (cont.)

- Factors decreasing the cost of components:

Learning curve: The cost of a component decreases over time, since throughput increases (faulty component production rate decreases)

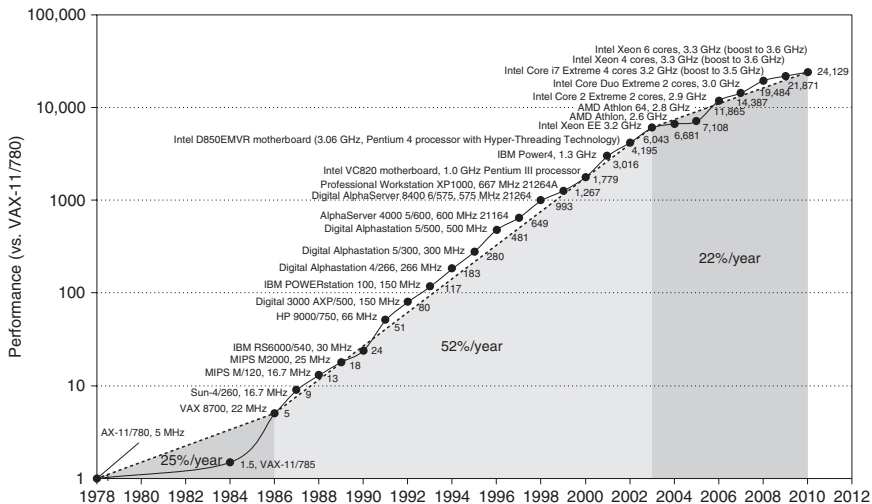
Sales volume: Doubling the sales volume decreases costs by 10%

- Design costs

- Factory costs are inversely proportional to *feature sizes*.
- Size of a design team is inversely proportional to *feature size*: from 3 people in the Intel 4004 design team to more than 300 in the ones designing modern processors.

Contents

- 1 Concept of Computer Architecture
- 2 Computer Requirements
- 3 Technology, power consumption and cost
- 4 Evolution of processor performance**
- 5 Types of computers
- 6 Máster en Ingeniería de Computadores y Redes



- Period I** (1978–1986). Performance grows at a rate of 25% per year, mainly due to technology enhancements.
- Period II** (1986–2003). Performance grows at a rate of 52% per year, due to technology enhancements and **architectural improvements**: RISC (*Reduced Instruction Set Computer*) architecture, pipelining, instruction level parallelism (ILP), caches, etc.
- Period III** (2003–2011). Performance grows at a rate of 23% per year. Limits in ILP, Memory latency, power consumption, and heat dissipation → performance improvement can only be achieved through parallelism.
- Period IV** (2011–2015). Performance grows at a rate of 12% per year due to technological and parallelism limitations (Amdahl's law).

Period V (2015–2018). Negligible performance growth (6.5% per year)

- End of Moore's law
- Technology does not improve energy efficiency
- Reaching the limits of parallelism

Current situation:

- Performance and energy efficiency through specialization → use of domain-specific accelerators, such as Google Tensor Processing Units (TPU) multiplying per 80 the performance/watt wrt a CPU in the context of neural networks-based inference.
- Research and Development of new technologies, as in:
<https://bit.ly/2lU7DN9>.

Architectural improvements

Some examples:

RISC architecture: Simple instructions that execute very fast. Simple hardware.

Pipelining: Splitting the instruction cycle into phases that are executed concurrently.

ILP exploitation: Execution of instructions out of program order.

Parallelism:

DLP *Data-Level Paralellism*. Performing an operation on multiple data.

TLP *Thread-Level Paralelism*. Performing several tasks in parallel.

Many of these improvements are only possible when the number of available transistors is high enough.

Some architectural ideas are motivated by the technology available

Some examples:

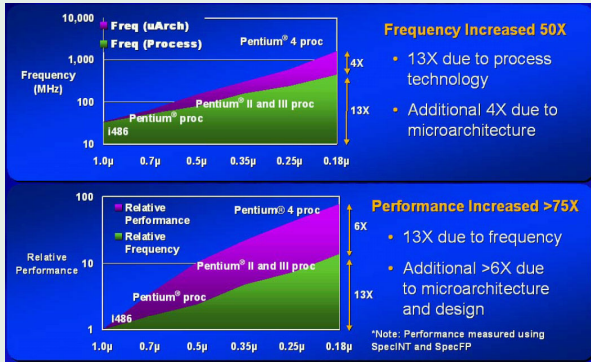
Caches The speed difference between the processor and the memory motivated the introduction of the *cache* memory between the processor and memory. The continuous growth of this difference has motivated the use of various levels (L1, L2, L3) of *cache* in the designs.

Stages to propagate signals The increasing wire delay motivated the inclusion in the Pentium 4 of two pipeline stages intended only to propagate signals through the wires.

Chip multiprocessors A sophisticated processor that works at very high frequency consumes a lot of power and dissipates a lot of heat. Instead, we can put several simpler processors that operate at lower frequency and lower voltage.

The relevance of architectural enhancements

Example: Intel processors



→ The architectural improvements delivered a speedup of 7X over the one that would have been achieved with only technological improvements.

Contents

- 1 Concept of Computer Architecture
- 2 Computer Requirements
- 3 Technology, power consumption and cost
- 4 Evolution of processor performance
- 5 Types of computers**
- 6 Máster en Ingeniería de Computadores y Redes

Types of computers

Computer market evolves as available technology and society consumption habits do.

Computer industry adapts its offer to this evolution by offering different types of devices.

For instance, during the 70's, as first microprocessors were under construction, two types of computers were commercially available:

Mainframes: very big and expensive computers only affordable to big corporations.

Minicomputers: medium-size computers, largely used at universities.



Types of computers (cont.)

Since 2012, thanks to advances in microprocessor design, communications and human–machine interfaces, the application domain of computers has become wider:

- Processors and memories have become small enough to be integrated in the design of mass-market electronic devices, some of them portable.
- A computer can integrate a large number of processors with memory that collaborate to improve the overall system performance.

Currently, most relevant computers are: Mobile devices, Embedded Systems, Personal computers, Servers, Clusters, Supercomputers

Personal Mobile devices

Include (but not limited to): smartphones, tablets, PDAs, etc.

Main features:

- Very limited energy consumption
 - Battery powered
 - No forced cooling
- *Responsive* and *predictable* design
 - Need of multimedia GUIs
 - Video frames and audio block must be processed on time (*real-time* requirements).
- Reduced main memory capacity
 - Optimized code (small footprint)
- Secondary memory of *flash* type

Embedded systems

Include: automotive and aerospace embedded systems, appliances, (air, marine, space) navigation, etc. . . .

Mobile devices are a concrete case of embedded system.

Differences wrt mobile devices:

- Wide spectrum of design and performance requirements
- Specific software developed by manufacturers (no third-party code)

Personal computers

Examples: Desktop, Laptops, Netbooks.

Main features:

- Well-balanced numerical and graphical computing power.
- Optimal relation between performance and cost
- Wide spectrum of configurations

Servers

Computer providing services (data, email, printing, . . .) in a network.
Use to be part of the computer infrastructure of a corporation.

Main features:

- Availability is critical
- Scalable design
- Design for *responsiveness*
- Throughput is a key attribute

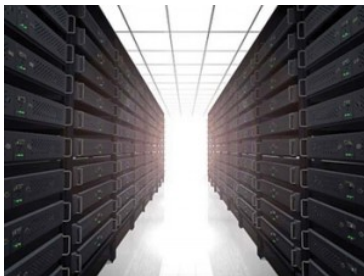


Cluster

It is a collection of computers, each one with its own operating system, interconnected through a network. Is externally seen as a single computer.

Main Features:

- Significant investment in power supply and cooling
- Scalability: The system can be easily adapted to changes (typically peaks) in workloads



Large Scale Clusters

Provide support to the big internet services: social networks, search engines, e-commerce and file sharing platforms, etc.

Main features:

- Need of a huge Internet bandwidth and secondary storage
- Good balance between performance and cost.
- Dependability-oriented design: Availability is critical, but low-cost components are used → need of redundancy

Supercomputer

Machines designed for high performance, despite their cost.

Main features:

- Execution of large scientific distributed applications with limited user interaction (weather forecasting, protein folding, ...)
- Very high floating point arithmetic throughput
- Cluster-based supercomputers are becoming more and more common



Contents

- 1 Concept of Computer Architecture
- 2 Computer Requirements
- 3 Technology, power consumption and cost
- 4 Evolution of processor performance
- 5 Types of computers
- 6 Máster en Ingeniería de Computadores y Redes

¿Interesado en Aprender Más?

- Título oficial, 1 año de duración (60 ECTS).
 - Más información en <http://mic.disca.upv.es>.
- Asignaturas relacionadas:
 - Arquitectura y Tecnología de los Procesadores Multinúcleo
 - Claves de las prestaciones de los procesadores multinúcleo.
 - Diseño de los procesadores multinúcleo actuales.
 - Conceptos avanzados sobre caches compartidas.
 - Memoria principal con controladores de memoria compartidos.
 - Redes en Chip
 - Diseño avanzado y construcción de redes en chip.
 - Prestaciones de las redes en chip.
 - Encaminamiento, reconfiguración, topologías y conmutación.
 - Tecnologías futuras y ejemplos actuales de redes en chip.

¿Interesado en Aprender Más? (cont.)

- Arquitectura de Redes de Altas Prestaciones
 - Diseño y construcción de redes de altas prestaciones.
 - Técnicas de control de flujo, tráfico, tolerancia a fallos, etc.
 - Reducción del consumo en redes de altas prestaciones.
 - Diseño de routers.
- Configuración, Administración y Utilización de Clusters
 - Diseño y configuración de clusters.
 - Sistemas de almacenamiento.
 - Equilibrado de carga.
 - Clusters de alta disponibilidad.