

# MC732 - Computer Organization and Design

Pedro Henrique Limeira da Cruz

August 7, 2023

# Contents

<b>1</b>	<b>Computer Abstractions and Technology</b>	<b>3</b>
1.1	Abstraction . . . . .	3
1.2	Memory . . . . .	3
1.3	Networks . . . . .	3
1.4	Technology Trends . . . . .	3
1.5	Performance . . . . .	4
1.5.1	Indicators . . . . .	4
1.5.2	Relative Performance . . . . .	4
1.5.3	Measuring Execution Time . . . . .	4
1.6	Power Consumption . . . . .	5
1.7	Multiprocessors . . . . .	6
1.8	SPEC CPU Benchmark . . . . .	6
1.9	Optimization Pitfall: Amdahl's Law . . . . .	6
<b>2</b>	<b>Instruction Set Architecture - ISA</b>	<b>6</b>
2.1	Instruction Set . . . . .	6
2.2	Arithmetic Operations . . . . .	6
2.3	Memory Operands . . . . .	7

# 1 Computer Abstractions and Technology

## 1.1 Abstraction

Abstraction helps us deal with complexity, by hiding lower-level implementation details. An example of that is the ISA (instruction set architecture), which is the abstraction from the hardware to the software, or like an API (common on OS, web, etc).

## 1.2 Memory

Memory, which is a place to store data, can be divided into two different large types:

1. *Volatile memory*: Needs power to maintain stored data, like RAM.
2. *Non-volatile memory*: Persistent data storage, like magnetic disk, DVD, etc.

Although it might seem like better to always have non-volatile memory, it comes with a change in both paradigm (it may impact safety, configuration, etc) and pricing (it is cheaper to have large quantities of non-volatile memory, or at least it was that way).

## 1.3 Networks

In today's day and age, networking has become an integral part on the world of computing. It is responsible for the communication between devices, resource sharing, non local access, etc.

Can be divided into three general areas:

1. Local Area Network (LAN): Ethernet
2. Wide Area Network (WAN): The internet
3. Wireless network: WiFi, Bluetooth, ...

## 1.4 Technology Trends

In general, since the invention of computers, the evolution of the processing power of the CPUs have basically grown exponentially. The speed of the memories, however, have grown at a much slower rate (almost linear in comparison). That results in what is called *the memory wall*, which is the name given to the discrepancy between the speeds of the processors and memories.

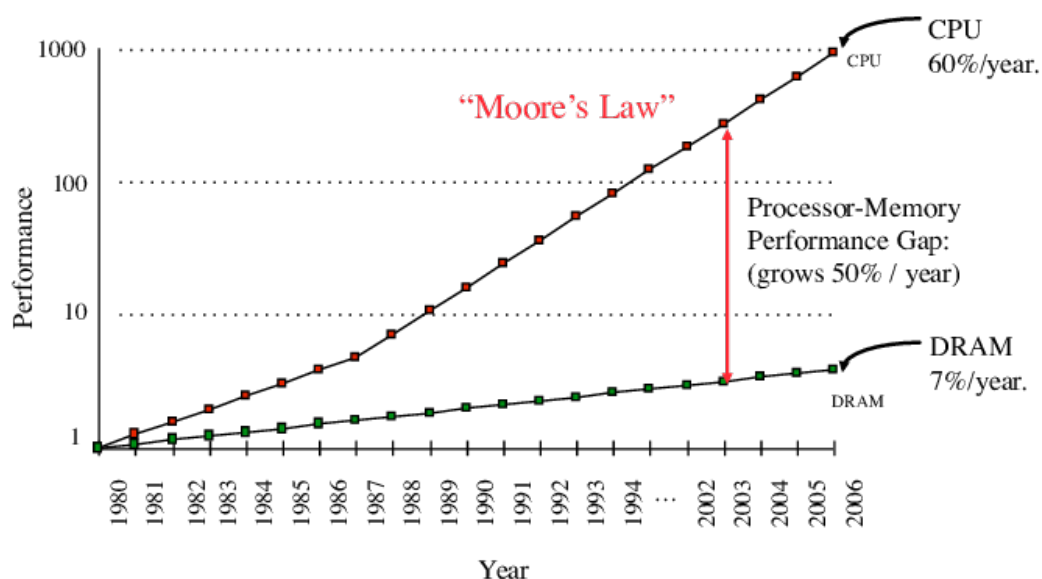


Figure 1: Performance of Memory vs Processor over the years

## 1.5 Performance

### 1.5.1 Indicators

Although when we talk about performance, we always associate it with speed, that might not give the whole picture. A clear example of that is the bellow chart, describing different types of performance metrics that might be used to characterize an airplane, which gives a general idea that there are different ways of measuring how a system behaves, and that each might be interesting to a different stakeholder.

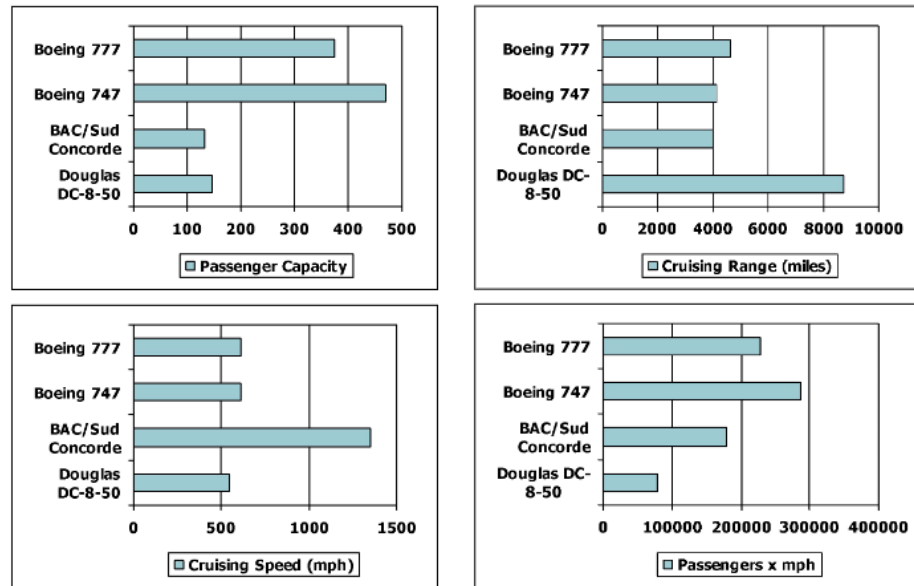


Figure 2: Different Performance Metrics for Airplanes

When it comes down to processor/computer systems, the two main indicators are:

1. Response Time: How long it takes to do a task, which might be important for the processor that runs the flight controller of an airplane.
2. Throughput: Total work done per unit time, which might be important for a server being requested hundreds of times a second.

### 1.5.2 Relative Performance

When choosing the right hardware/system for your problem, it is important to be able to have metrics to compare the two. Taking into account *Execution Time*, we might have:

- *Performance*: Defined as  $1/Execution\ Time$
- *Speed Up/Down*: We say "X is  $n$  times faster the Y", where  $n = Execution\ Time_y / Execution\ time_x$ <sup>1</sup>

### 1.5.3 Measuring Execution Time

We just saw that it is important to be able to measure relative performance, but for that we need to be able to measure the time required to run a determined application or bench mark. To do so, there are a couple of ways, each with their draw backs:

- Elapsed Time: Total response time, including all external factors, that might not be related to the task at hand (e.g I/O). Used a lot when referring to user experience.
- CPU Time: Time spent processing a given job, without considering other job's share, used heavily when analyzing scientific computing.

<sup>1</sup>It is important to notice that the order inverts, to say there was a *speed up*, the lowest value should be in the denominator

To check elapsed time, we can simply use an external timer, or the one on the OS. For CPU time, on the other hand, it is calculated based on the number of clock cycles and its period, shown below (where CLK Rate refers to the Clocks frequency, and CLK Cycles refers to the number of cycles required for the application, which some times it's referred to as *CLK Count*):

$$CPU\ Time = CPU\ CLK\ Cycles \times CLK\ Cycle\ Time \quad (1)$$

$$= \frac{CPU\ CLK\ Cycle}{CLK\ Rate} \quad (2)$$

There is, however, a **direct relationship between CLK frequency and cycle count**, where, if the frequency gets too high, the cycle count must also increase, for the time per cycle is not enough for the tasks to be finished.

It is also important to note that the execution time is strictly intertwined with the environment which the code was generated (compiler, architecture, ...).

Another way of calculating the CPU time is by using the *Cycles per instruction*, which is an average number of clock cycles needed for instructions. The CPU time is given by:

$$CPU\ Time = \frac{Instruction\ Count \times CPI}{CLK\ Rate} \quad (3)$$

If different instruction classes exist, however, the CPI becomes a weighted average:

$$CPI = \frac{CLK\ Cycles}{Instruction\ Count} = \sum_{i=1} \left( CPI_i \times \underbrace{\frac{Instruction\ Count_i}{Instruction\ Count}}_{Relative\ Frequency} \right) \quad (4)$$

**CPU TIME EXAMPLE** → A computer A has a 2GHz clock, and a CPU Time of 10s. The target is to develop a computer B, that aims to have a 6s CPU time and a faster clock, which results in a clock count of  $1.2 \times clock\ cycles$ . How fast must computer B's clock be?

- When he asks about "fast" he actually means the clock rate, which is give by:

$$CLK\ Rate_B = \frac{CLK\ Cycle_B}{CPU\ Time_B} = \frac{1.2 \times CLK\ Cycles_A}{6s}$$

- As we can see, we need to calculate the clock cycles of A, which is described by:

$$\begin{aligned} CLK\ Cycles_A &= CPU\ Time_A \times CLK\ Rate_A \\ &= 10s \times 2GHz = 20 \times 10^9 \end{aligned}$$

- Therefore we have the following clock rate for B:

$$CLK\ Rate_B = 4GHz$$

## 1.6 Power Consumption

All modern computers use the CMOS IC technology, which consumes the following amount of power:

$$Power = Capacitive\ load \times Voltage^2 \times Frequency \quad (5)$$

The most important thing to keep in mind is that the power is proportional to the square of the voltage used in the processor. That voltage has seen a drastic drop in modern computing, reaching a physical limit, where we cannot go any lower, but because of the incredible high clock speeds, the power generate is also big, to a point where we also cannot remove more heat.

The solution to having more processing power but not worrying about temperature dissipation in one die (one silicon chip), we moved from single core to multi-core architectures.

## 1.7 Multiprocessors

As state previously, a multiprocessor architecture is comprised of multiple processors per chip (as in a quad-core, octa-core CPU), which is very different from a single core architecture. That shift in paradigma creates a series of difficulties, since there is a need for explicit parallel programming, which by itself is hard because:

- It is harder to program for performance in a multi-core architecture
- The need for a load balancing strategy rises
- Optimizing communication and synchronization

## 1.8 SPEC CPU Benchmark

When we talked about performance, we were talking about the performance measured based on a program/algorithm. In order to standardize which tests/algorithms are run to compare the results between CPU architectures, we run a *Benchmark*, which is a specialized program made to emulate a typical workload and measure the performance.

Furthermore, in order to have more statistical relevance, the performance is measured and normalized against a reference machine, and summarized as geometric mean of performance ratios, following the formula below:

$$\sqrt[n]{\prod_i Execution\ time\ ratio_i} \quad (6)$$

## 1.9 Optimization Pitfall: Amdahl's Law

In a lot of cases, when people optimize a certain aspect of the computing, they expect an overall improvement on performance. That, however, is not possible, since the improvement is applicable only to the time spent by such operation.

An example is that, if you improve the multiplication computing, which is responsible for 80s out of the 100s of the whole computation, 5× faster, we won't have a 5× boost in overall performance, but rather the improved final time is given by:

$$T_{improved} = \frac{T_{affected}}{improvement\ factor} + T_{unaffected} \quad (7)$$

# 2 Instruction Set Architecture - ISA

## 2.1 Instruction Set

The instruction set is just the name given to the group/repertoire of instructions that the processor is capable of handling. Therefore, it is important to note that the ISA is strictly related to the computer/architecture.

In the very beginning, computers had very simple instruction sets, but as the processors became more powerful, they became more complex. In recent years, however, the industry went back to simpler instruction sets, since they result in simpler chips and, therefore, lower power consumption, hence the proliferation of **RISC architectures** (Reduced Instruction Set Computers), having RISC-V (open source) and ARM as the big players.

## 2.2 Arithmetic Operations

The base of any ISA, and therefore present in the RISC-V ISA, are arithmetic operations, all which have the following form:

add t0, t1, t2 // t0 gets t1 + t2

This decision (having all arithmetic operations follow the above structure), and only being able to access registers for operands reflects the first and second design principles used, respectively:

**Design Principle 1** → Simplicity favours regularity, which makes implementation simpler and enables higher performance at lower cost.

**Design Principle 2** → Smaller is faster (e.g. main memory has millions of locations, registers are very limited and directly coupled to data bus).

## 2.3 Memory Operands

In the RISC-V ISA, main memory is used for composite data, such as arrays, structures, and dynamic data. As said before, however, apply any arithmetic operation, first you need to load values from memory into registers, and then stores the resulting values from the registers to memory.

Furthermore:

- Memory is byte addressed, but most of the times you align the values to the word (32 bits  $\therefore$  4 bytes)
- RISC-V is Little Endian, which means that the least-significant byte at least address of a word
- RISC-V does not require words to be aligned.