

# ICNPG 2023

## Clase 0

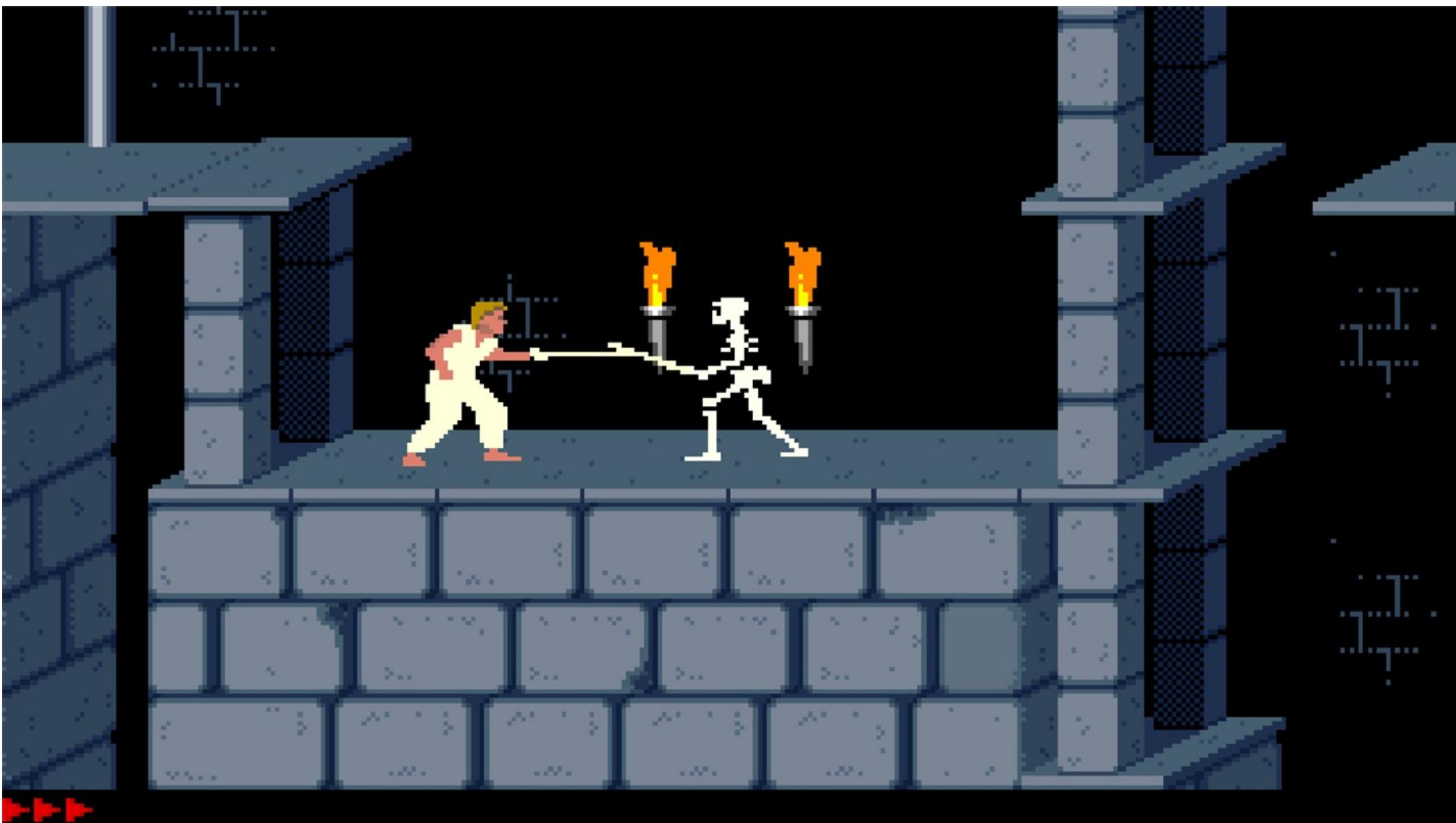


# Introducción al Cálculo Numérico en Procesadores Gráficos

- 16 clases x 4hs = 64 hs. Martes y Miércoles de 9 a 13. Aula Egipcios.
- Laboratorio de computación paralela en GPUs, orientada a aplicaciones científicas y tecnológicas.
- Ejercicios ilustrativos en clase y para hacer en casa (para discutir la clase siguiente).
- Evaluación: Participación en clase, desarrollo y exposición/entrega de un miniproyecto usando GPGPU. **El código tiene que correr en el cluster**
- Obligatorio aprender a usar el cluster de GPUs.
- Consultas fuera de horario: google-chat.
- Todo el material en el google-classroom IB.

# El principe de Persia (1990)

El desarrollo de placas gráficas comenzó con el desarrollo de los videojuegos. Se buscaba disminuir el trabajo en la CPU



# El principe de Persia (2010)



Videos de Yt:

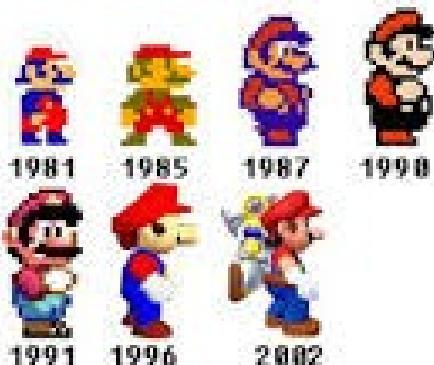
## **The Evolution Of Prince of Persia Games**

<https://youtu.be/yS6DEuuFJ6k>

## **Evolution of Video Game Graphics 1958-2020**

<https://youtu.be/IsPPWWlV-T8>

# The Evolution of Mario



\*Dates based on US release dates.

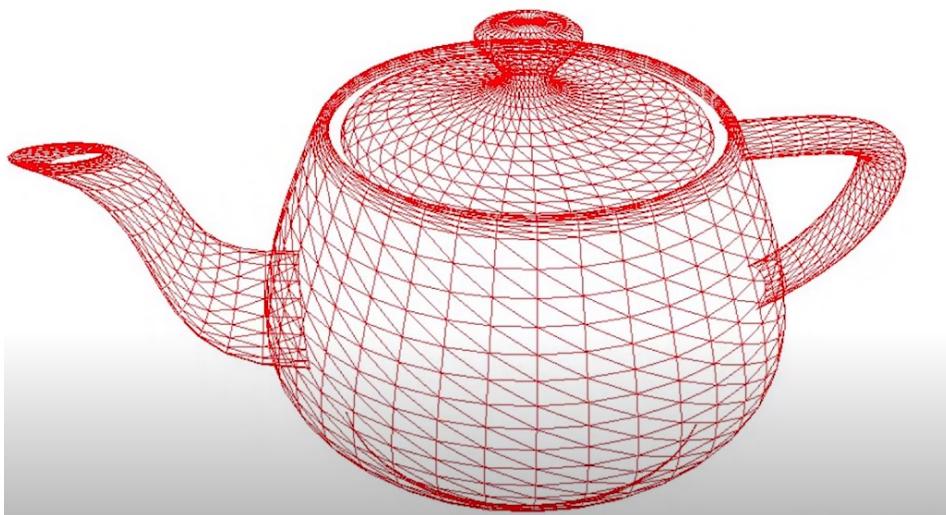
GamingLife.com



# Juegos (2018)



Las GPU aparecieron como un  
Acelerador gráfico 3D



Modelo 3D basado en vértices → Pixels en la pantalla

Para mover, rotar, dar textura, iluminación, etc. es necesario realizar  
*;Un montón de operaciones!*  
que son altamente paralelizables

## ARTIFICIAL INTELLIGENCE

Early artificial intelligence stirs excitement.



1950's    1960's    1970's    1980's    1990's    2000's    2010's

### MACHINE LEARNING

Machine learning begins to flourish.



### DEEP LEARNING

Deep learning breakthroughs drive AI boom.

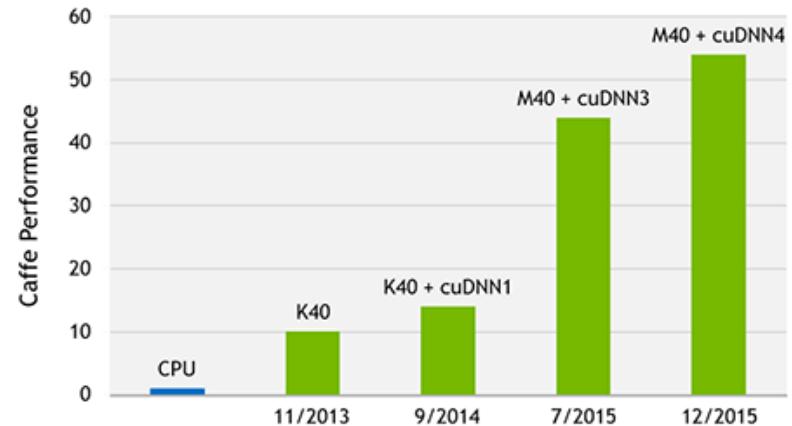


Since an early flush of optimism in the 1950s, smaller subsets of artificial intelligence – first machine learning, then deep learning, a subset of machine learning – have created ever larger disruptions.



Además, lo que vemos hoy en Deep Learning no sería posible si no tuviéramos las GPUs

### 50X BOOST IN DEEP LEARNING IN 3 YEARS



AlexNet training throughput based on 20 iterations,  
CPU: 1x E5-2680v3 12 Core 2.5GHz. 128GB System Memory, Ubuntu 14.04

## Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification

Kaiming He

Xiangyu Zhang

Shaoqing Ren

Jian Sun

Microsoft Research

(2015)

Review Article | [Published: 23 March 2022](#)

## The transformational role of GPU computing and deep learning in drug discovery

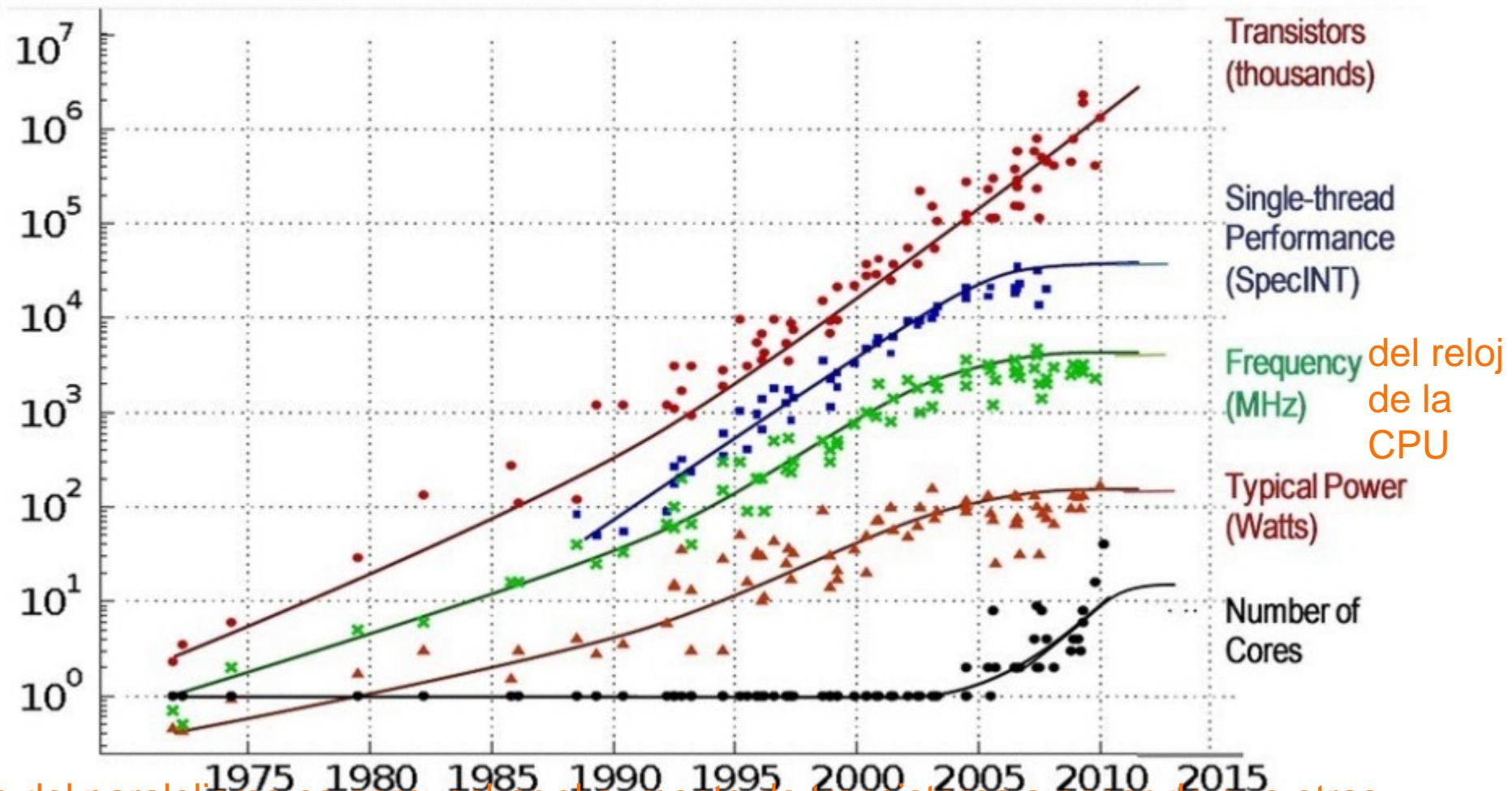
[Mohit Pandey](#), [Michael Fernandez](#), [Francesco Gentile](#), [Olexandr Isayev](#), [Alexander Tropsha](#), [Abraham C. Stern](#)✉ & [Artem Cherkasov](#)✉

[Nature Machine Intelligence](#) **4**, 211–221 (2022) | [Cite this article](#)

Se acabó: Mi viejo programa secuencial de CPU ya no correrá más rápido cuando cambie la computadora...

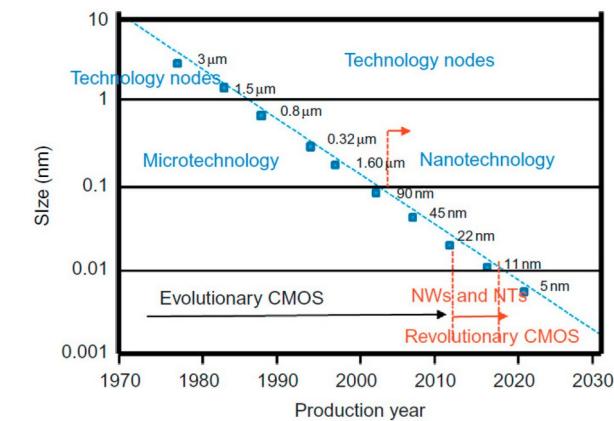
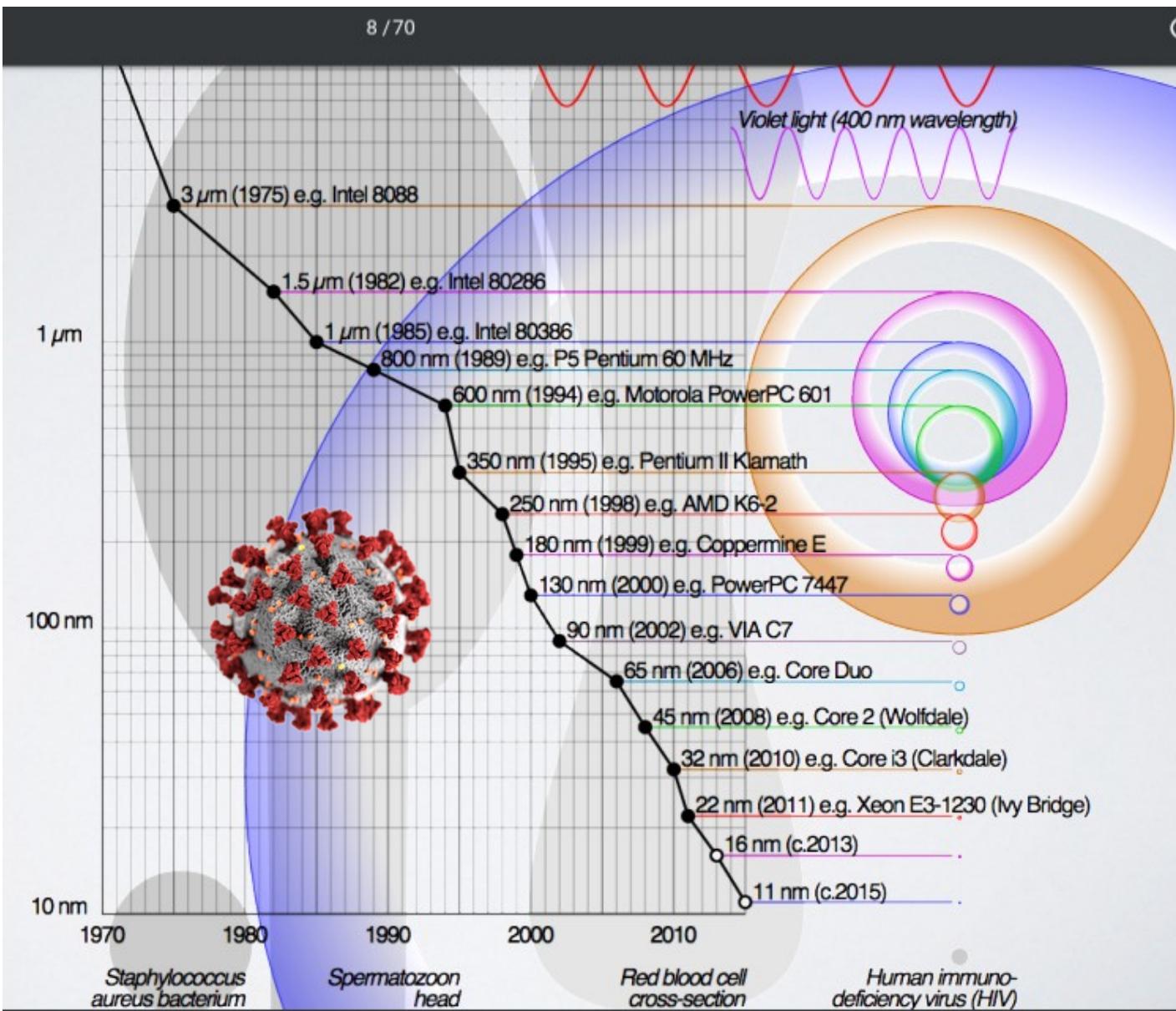
porque la frecuencia en el tiempo satura y por razones físicas no es posible aumentarla

## 35 Years of Microprocessor Trend Data



El objetivo del paralelismo es aprovechar el aumento de transistores a pesar de que otras características están saturando

# Gracias a la miniaturización (nanociencia/nanotecnología)



2023: ~5 nm

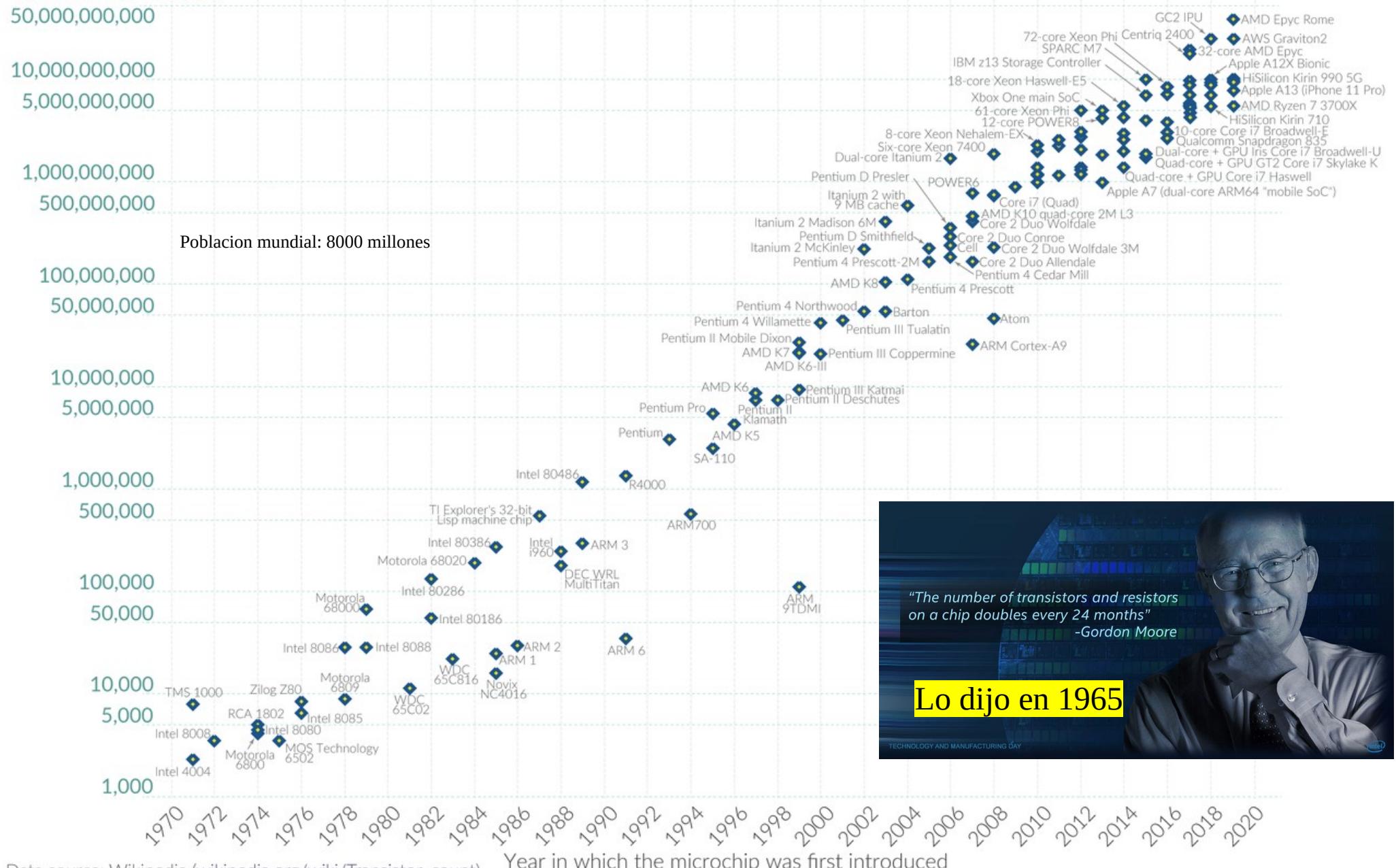
No se puede ver con luz visible



# Moore's Law: The number of transistors on microchips doubles every two years

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.

## Transistor count



Data source: Wikipedia ([wikipedia.org/w/index.php?title=Transistor\\_count&oldid=983010000](https://en.wikipedia.org/w/index.php?title=Transistor_count&oldid=983010000))

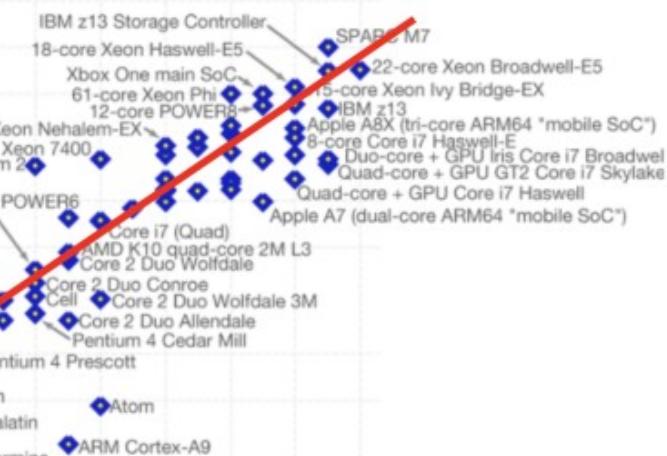
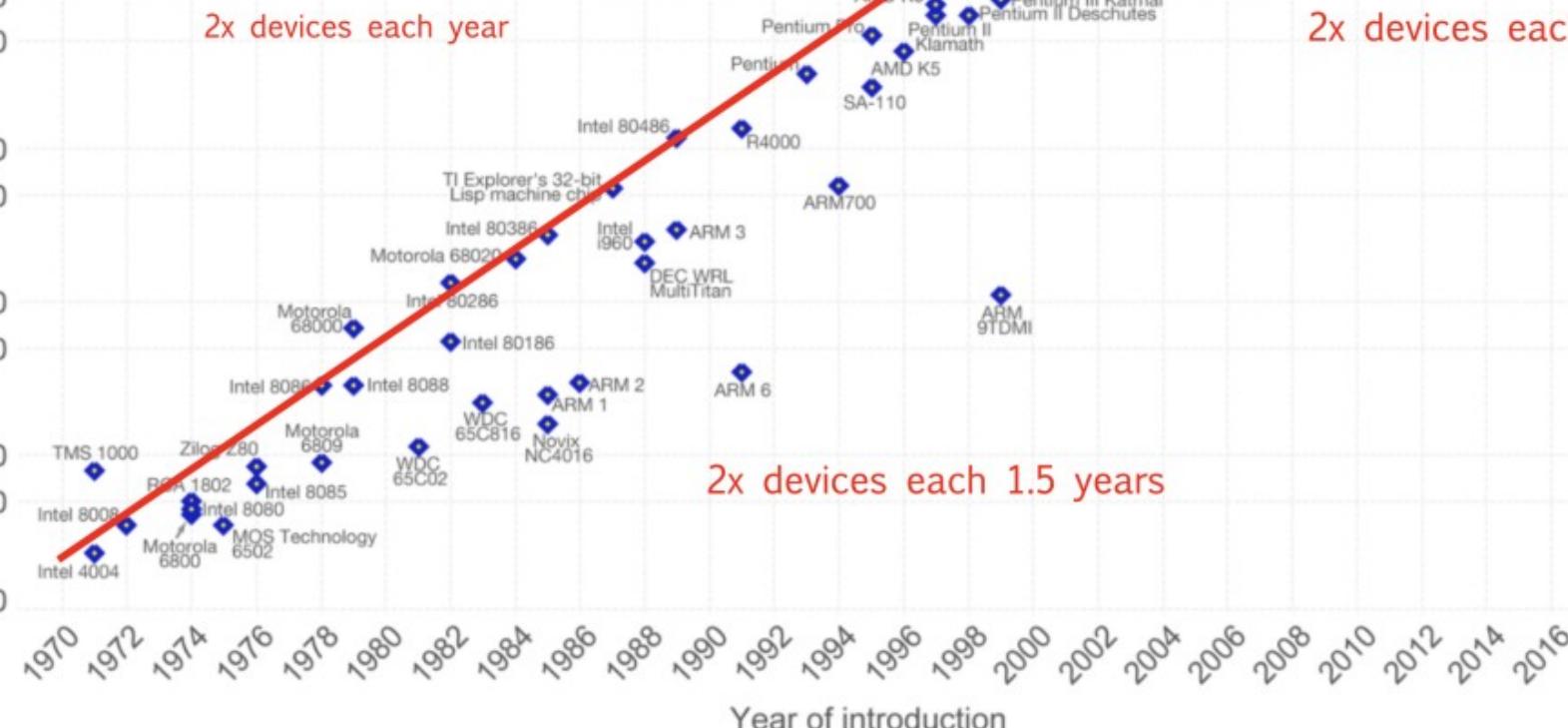
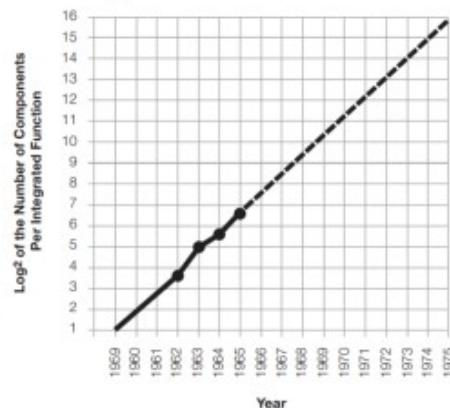
OurWorldInData.org – Research and data to make progress against the world's largest problems.

Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.

Transistor count

20,000,000,000  
10,000,000,000  
5,000,000,000  
1,000,000,000  
500,000,000  
100,000,000  
50,000,000  
10,000,000  
5,000,000

## Original Moore's chart



2x devices each 2 years

2x devices each 1.5 years

# ‘Moore’s Law’s dead,’ Nvidia CEO Jensen Huang says in justifying gaming-card price hike

Last Updated: Sept. 22, 2022 at 7:43 a.m. ET

First Published: Sept. 21, 2022 at 6:16 p.m. ET

**MORE ABOUT MOORE —**

## Intel: “Moore’s law is not dead” as Arc A770 GPU is priced at \$329

Expected performance somewhere near Nvidia’s RTX 3060 Ti—at least, for DirectX 12.

SAM MACHKOVECH - 9/27/2022, 2:40 PM

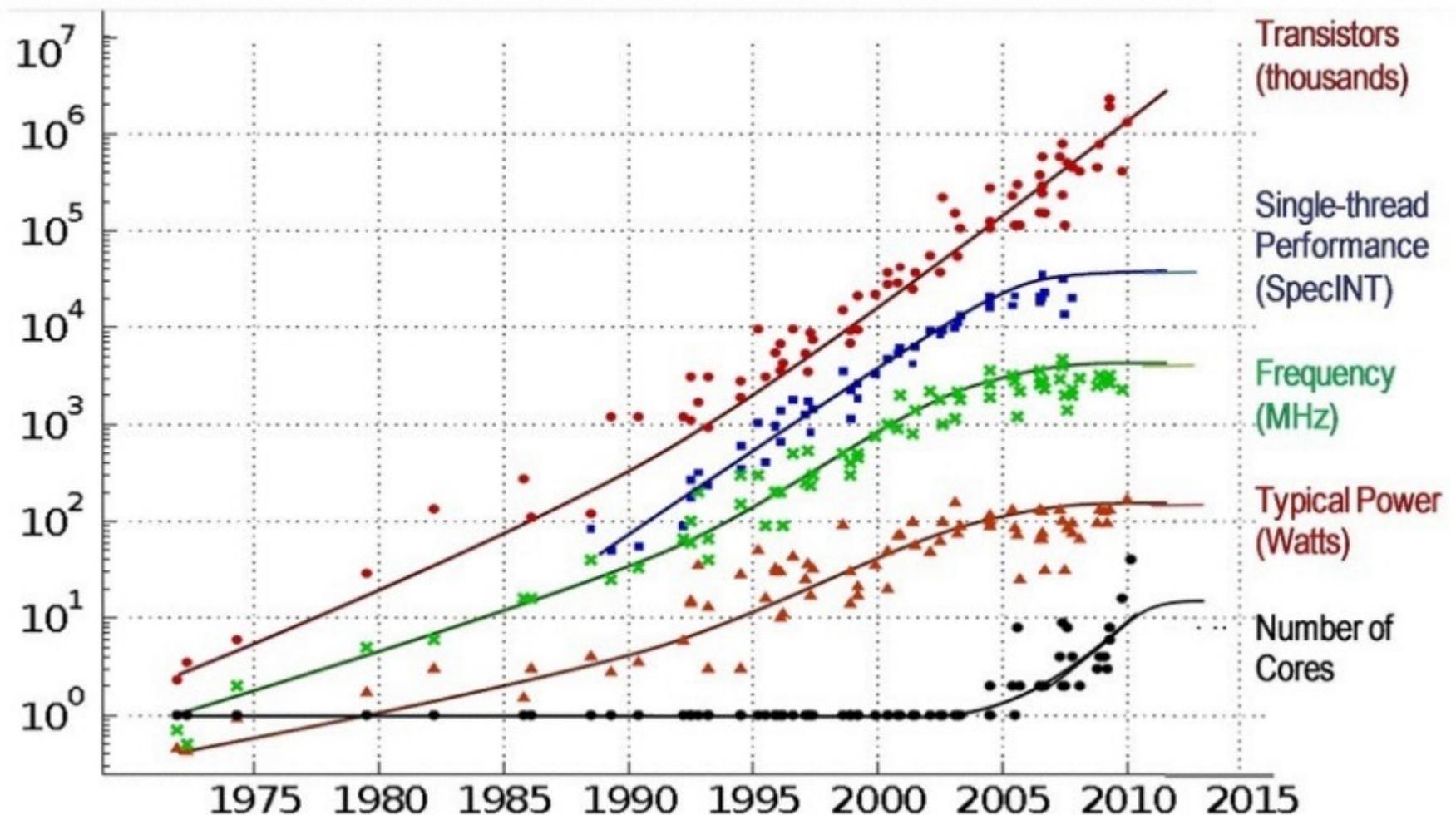
Actualmente las placas más conocidas son Nvidia, Intel y AMD. Con CUDA se puede aprender a usar las Nvidia y para las demás hay que usar otros lenguajes de características similares



(al final del curso vamos a ver brevemente como se programan las gpus de intel!)

Mi nuevo programa paralelo si correrá más rápido cuando cambie mi GPU

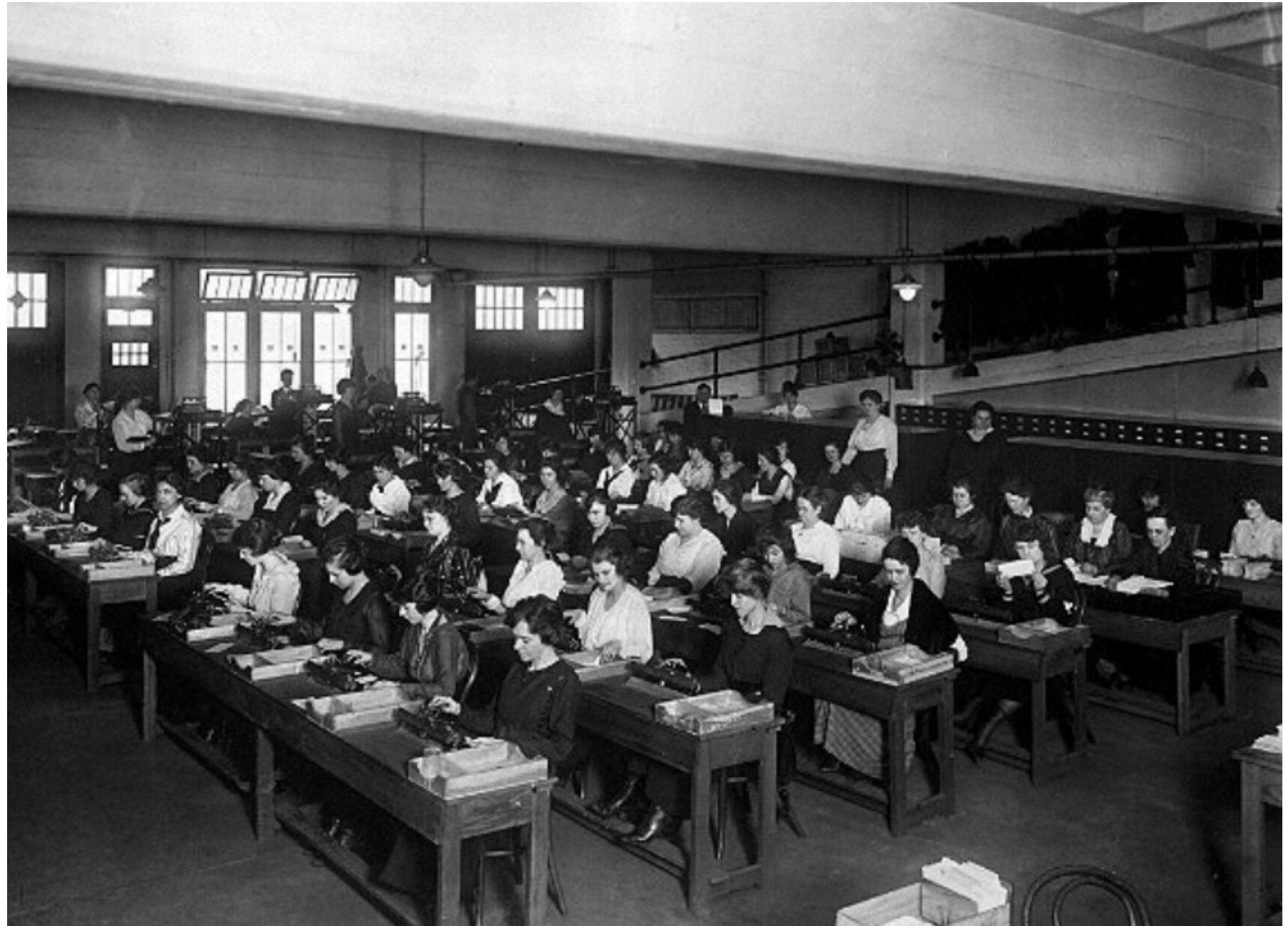
## 35 Years of Microprocessor Trend Data



# Performance → paralelismo



Women at work tabulating during World War II



<https://www.atomicheritage.org/history/human-computers-los-alamos>

# Realismo → paralelismo

- For example, imagine modeling these serially:



Galaxy Formation



Planetary Movements



Climate Change



Rush Hour Traffic

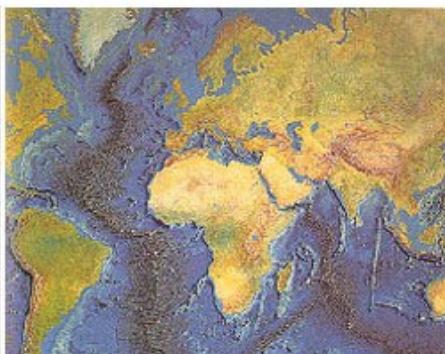


Plate Tectonics



Weather



Auto Assembly

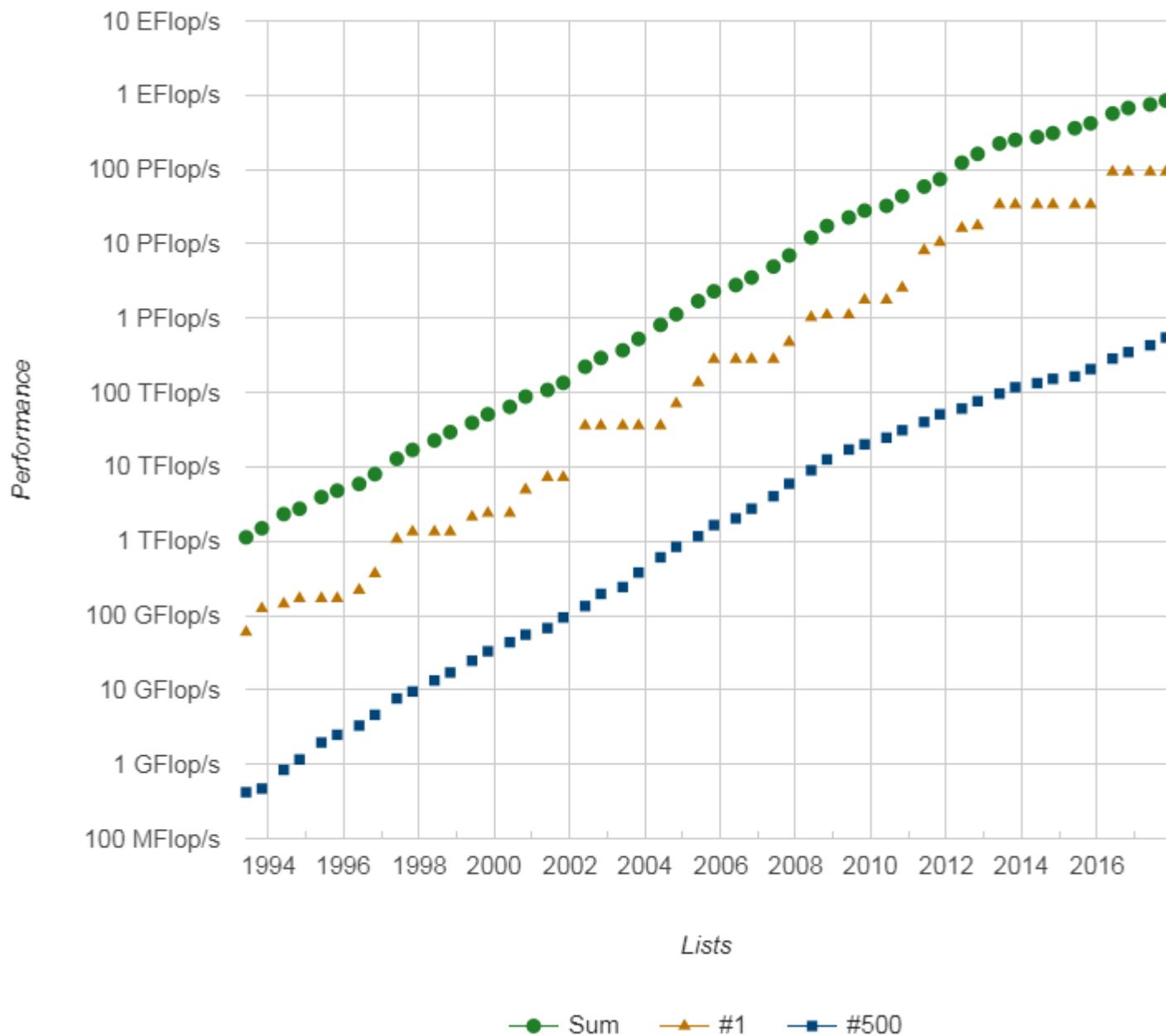


Jet Construction

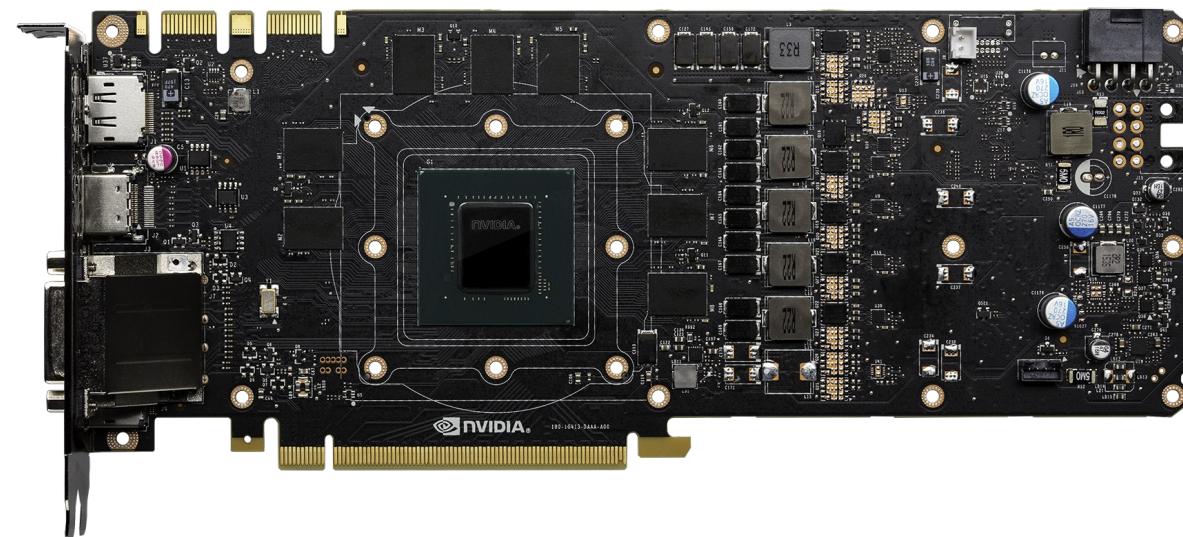
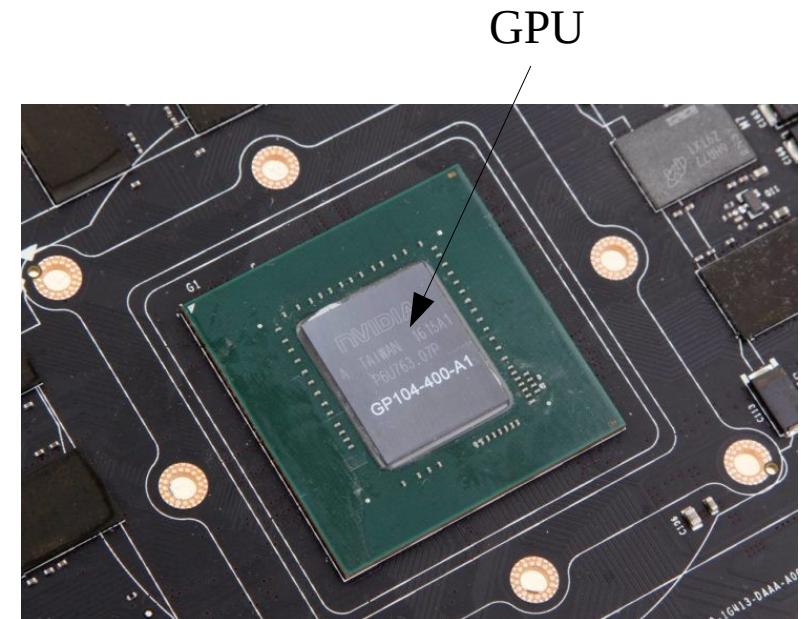


Drive-thru Lunch

## Performance Development



# GPGPU: “computadoras paralelas”



# Tipos de GPUs

- **Dedicadas:** aplicaciones graficas, juegos. Alta performance. Precios accesibles...
- **Integradas:** están integradas con la CPU. Menos poderosas que las dedicadas, pero más baratas.
- **Móviles:** procesadores gráficos para laptops, tablets, celulares. Menos poderosas que una dedicada pero optimizadas para consumir poco.
- **En la nube:** son gpus virtualizadas que pueden ser accedidas en forma remota. Se usa mucho en machine learning, etc.
- **Workstation:** para aplicaciones profesionales, modelado 3D, edición de video, computación científica. Alta performance, confiables, y permiten cosas avanzadas como multi-GPU y disponer de mucha memoria, etc. Caras.

# Marcas

- **NVIDIA**: es la líder, para consumo y profesional. Línea GeForce para juegos. Quadro y Tesla para aplicaciones profesionales (modelado 3D, ciencia, inteligencia artificial, etc).
- **AMD**: Consumo y profesional. Línea Radeon para juegos, Radeon Pro, Instinct para aplicaciones profesionales.
- **Intel**: aunque es conocida principalmente por sus CPUs, hacen GPUs integradas con sus CPUs. Línea Iris and Iris Pro, y Arc, se usan en laptops and desktops. Ponte Vecchio para HPC (**¡vamos a tener un cluster con estas últimas en Argentina pronto!**).
- **Qualcomm**: GPUs para “mobile devices”, smartphones y tablets. Línea Adreno es popular para Android.
- **Imagination Technologies**: mobile devices, automotive y IoT (Internet of Things). Smartphones y tablets.
- **ARM**: mobile devices, automotive y IoT. Línea Mali usada en smartphones y tablets.

# Mas operaciones por segundo...

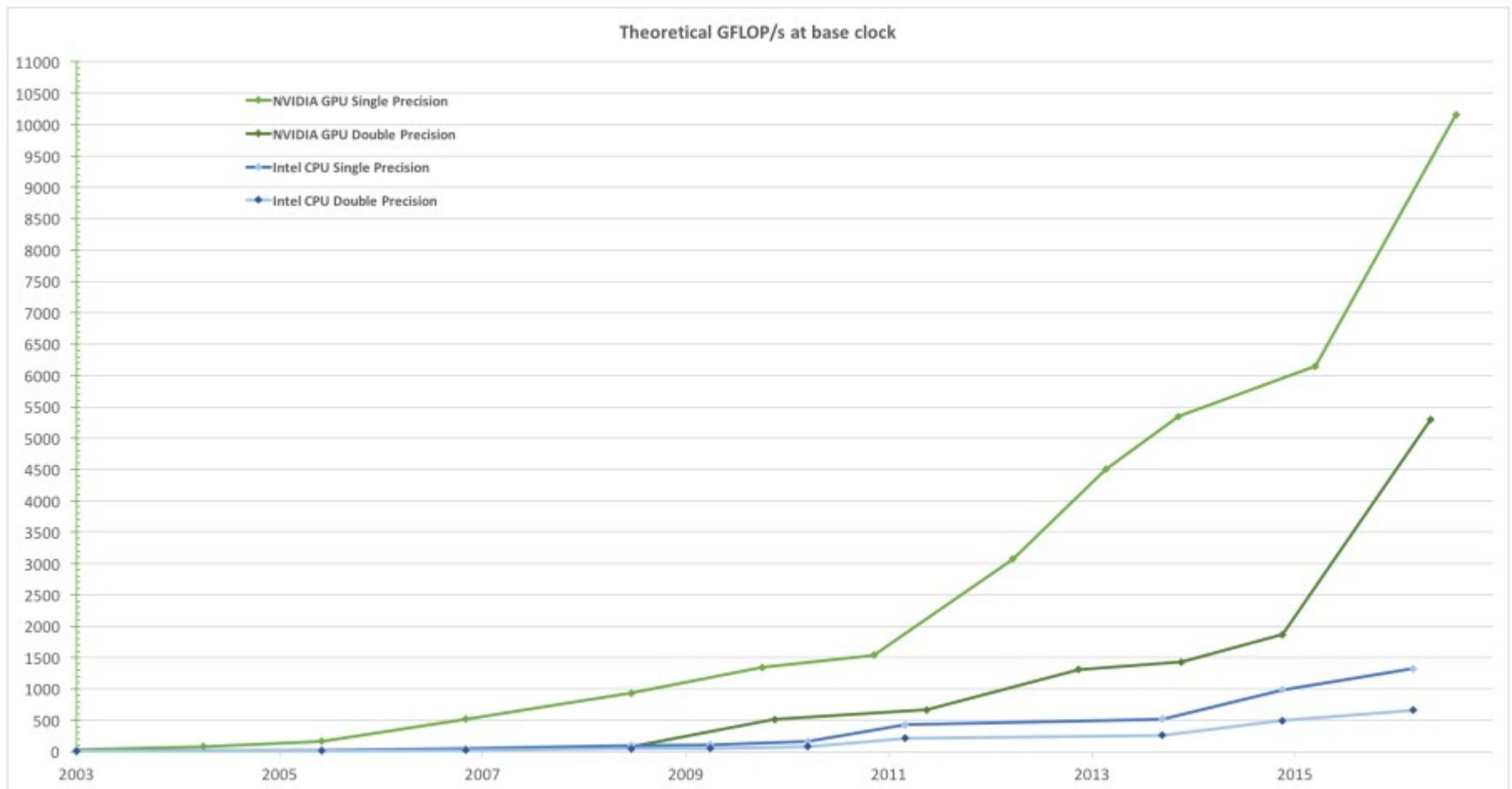


Figure 1 Floating-Point Operations per Second for the CPU and GPU

# Memorias más rápidas...

Aumentó la velocidad de transferencia de datos

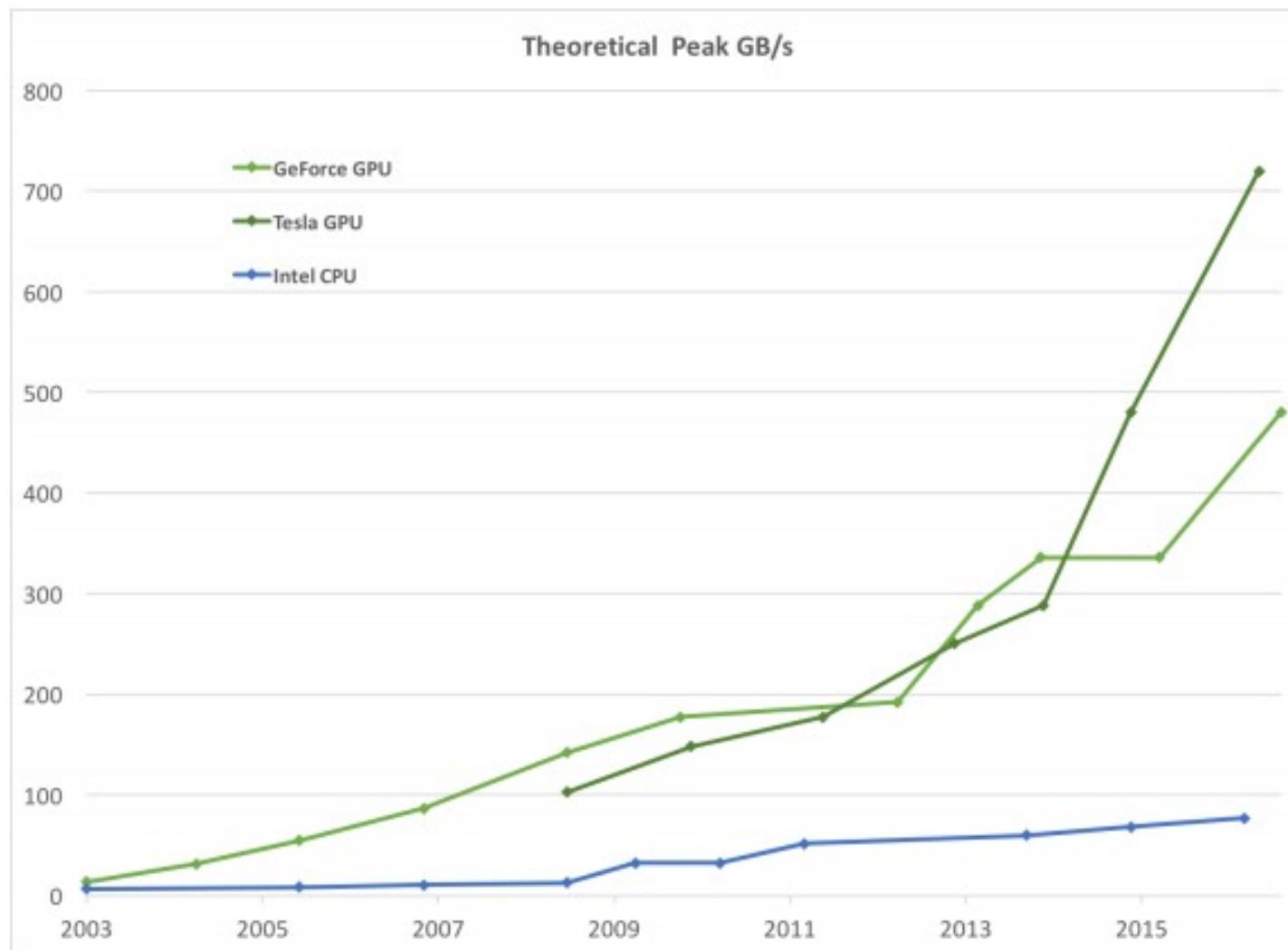


Figure 2 Memory Bandwidth for the CPU and GPU

# ¿Porque sirven como aceleradoras de cálculo?

The reason behind the discrepancy in floating-point capability between the CPU and the GPU is that the GPU is specialized for compute-intensive, highly parallel computation - exactly what graphics rendering is about - and therefore designed such that more transistors are devoted to data processing rather than data caching and flow control, as schematically illustrated by [Figure 3](#).



**Figure 3 The GPU Devotes More Transistors to Data Processing**  
Los núcleos de una GPU son menos poderosos que los de una CPU pero la ventaja es que son muchos

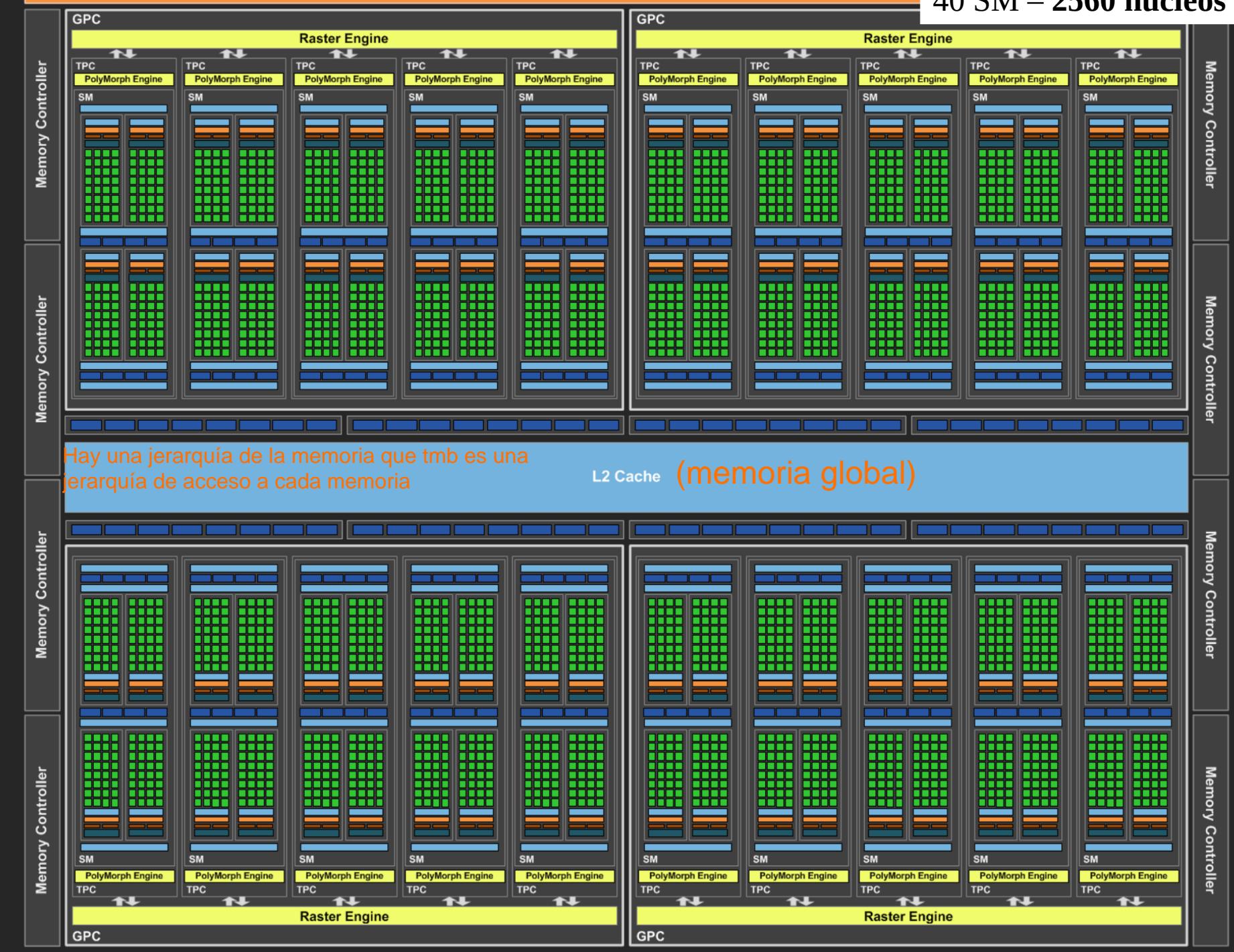
# Esquema de una verdadera GPU

PCI Express 3.0 Host Interface

GigaThread Engine

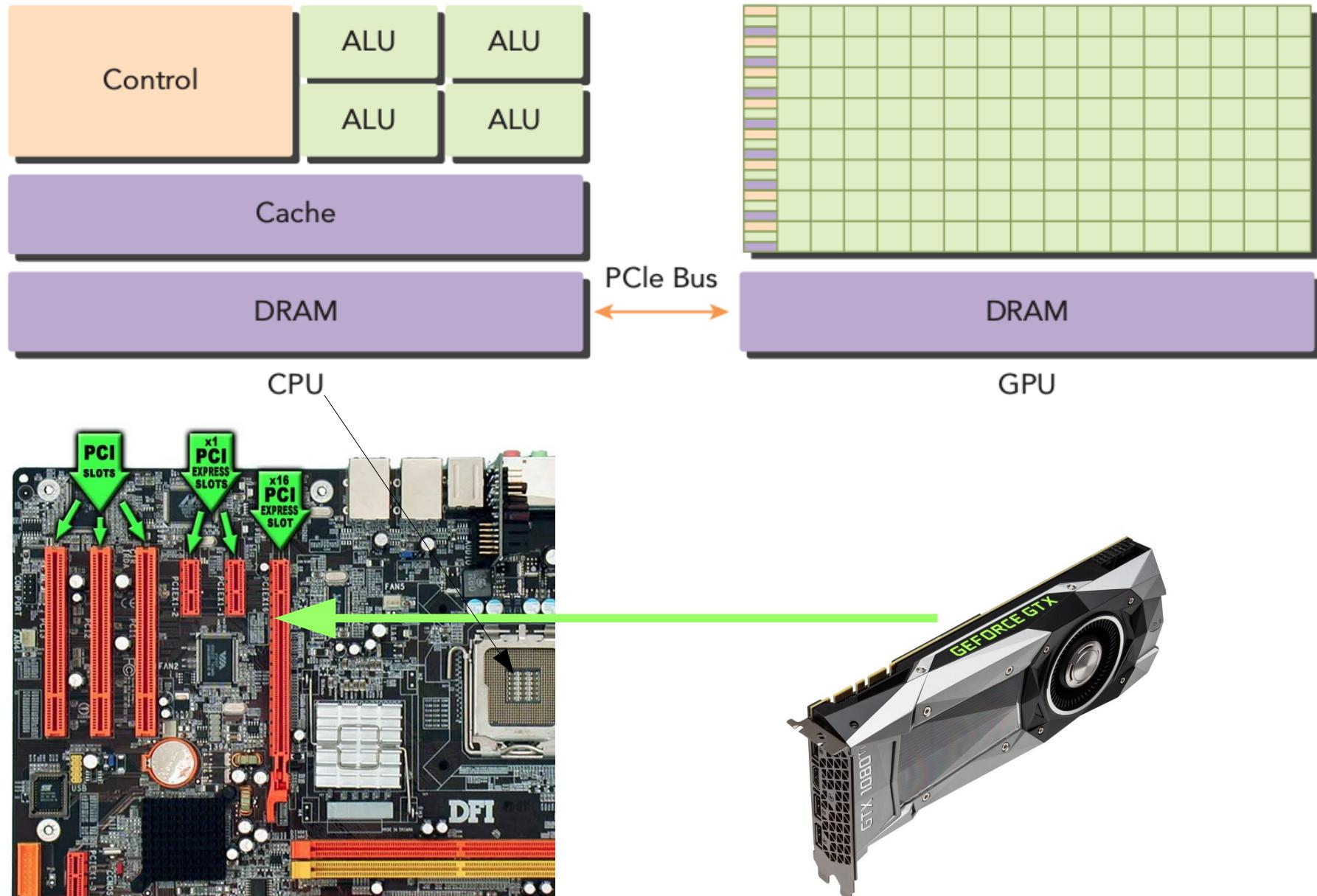
1 SM – 64 núcleos

40 SM – 2560 núcleos



# Arquitectura paralela híbrida

La conexión entre GPU y CPU se da a través de un puerto que es lento. Además, las GPU consumen mucho y es necesario tener una fuente acorde



# Todas se programan igual



Las GPU tamb se pueden colocar en sistemas envevidos

Procesamiento paralelo con GPU sobre un robot o drone



Jetson TK1

127 x 127 mm, 5 Watts, 70 UDS

GPU: NVIDIA Kepler "GK20a" GPU with 192 SM3.2 CUDA cores (upto 326 GFLOPS)

CPU: NVIDIA "4-Plus-1" 2.32GHz ARM quad-core Cortex-A15 CPU with Cortex-A15 battery-saving shadow-core



Jetson nano

70 x 45 mm

90USD



## **Inception Spotlight: New Skydio 2 Drone Powered by NVIDIA Jetson (2019)**

<https://youtu.be/imt2qZ7uw1s>

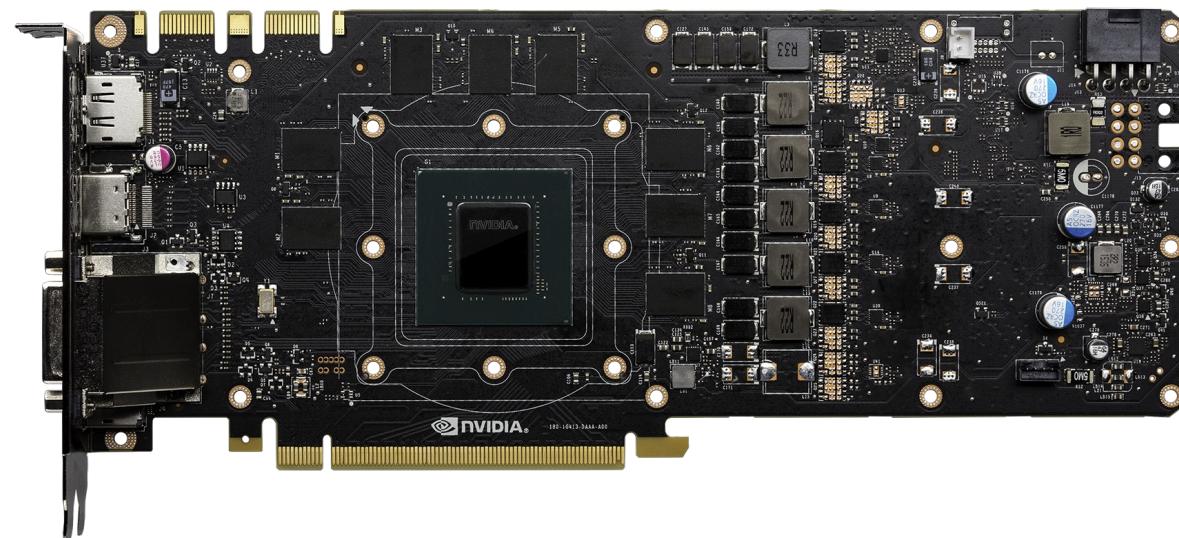


<https://youtu.be/USYlt9t0lZY>

## **Autonomous Drone Flight Over 1 Kilometer Forest Trail (2017)**

**ETC**

# ¿Como se programan estas cosas?



# Complicado hasta que salió CUDA en 2007

(estaba demasiado orientado a gráficos)

