# Parallel programming
# Introduction



PRINCETON UNIVERSITY

# Parallel Computing

- Moore's law

- Parallel processing

- Hardware evolution

- Parallel computing languages

- Open Multi Processing (OpenMP)

- Message Passing Interface (MPI)

- Graphical Processing Units (GPU)

# Moore's law

**Gordon Moore**: founder of Intel.

**Statement of the law:**
    The number of transistors that can be placed on an integrated circuit at a reasonable cost doubles every two years.

**Hardware evolution**:
- Processor frequency reached a plateau around 3GHz since 2002-2004 but Moore's law is still true
- More cores per chip (many cores architectures, GPU...)

**Energy consumption and cost:**
- Dissipated electric power scales at clock frequency to the cube.
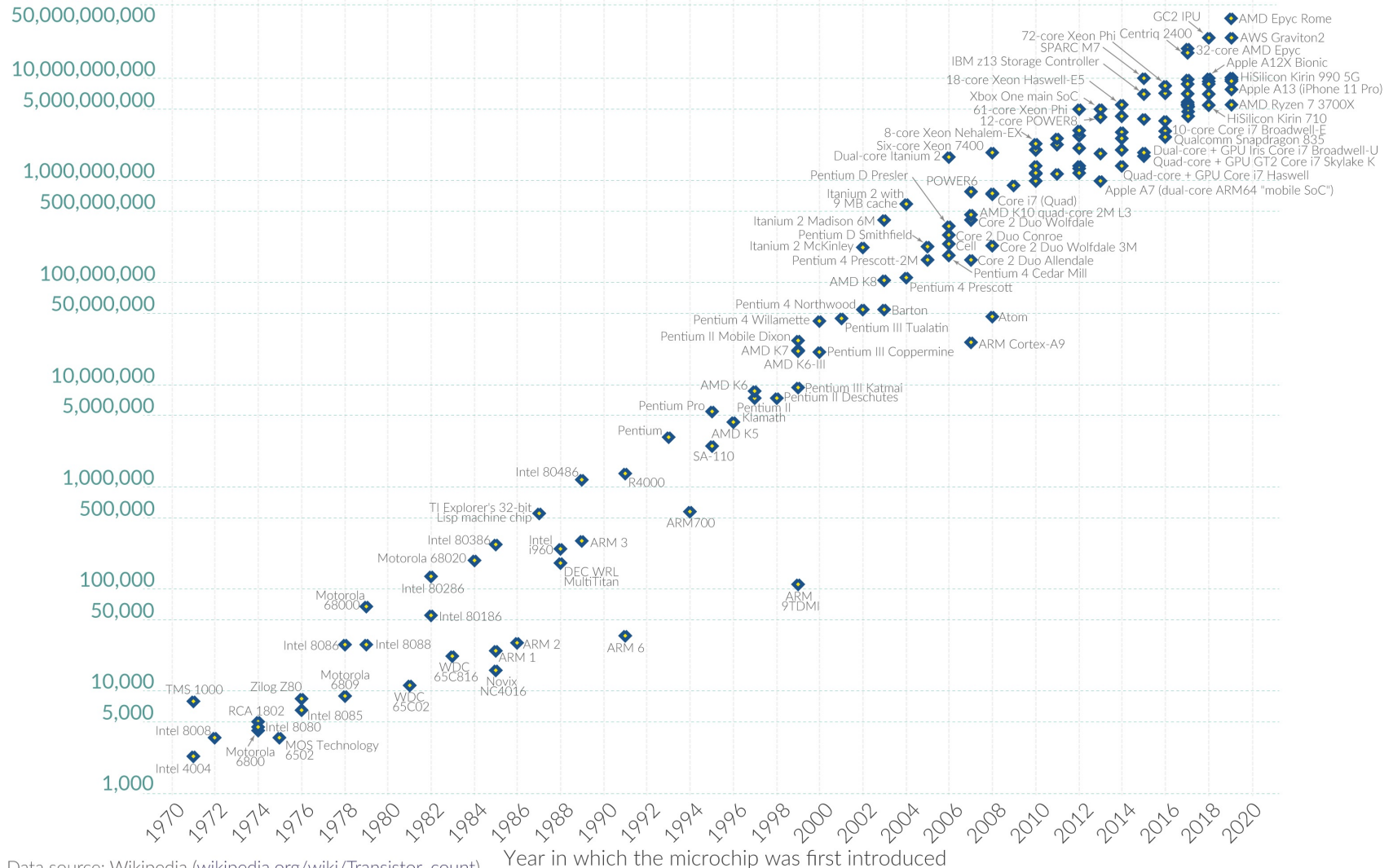- Dissipated power by square cm is limited by cooling

# Moore's law

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

Our World in Data

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.



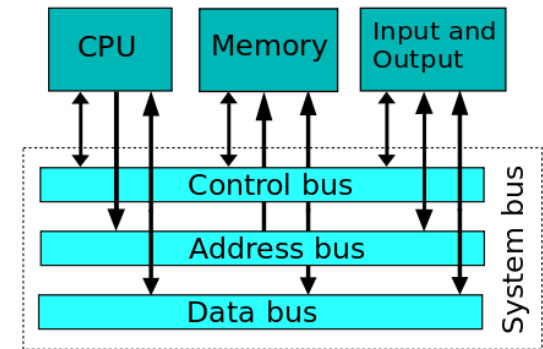Data source: Wikipedia (wikipedia.org/wiki/Transistor_count)
OurWorldinData.org – Research and data to make progress against the world's largest problems.
Licensed under CC-BY by the authors Hannah Ritchie and Max Roser.

4

# The Memory Wall

**Von Neumann Architecture**

Programming instructions (stored-program) and data share the memory unit. They are loaded to the Central Processing Unit (CPU) for execution. Results are loaded back to the memory unit.
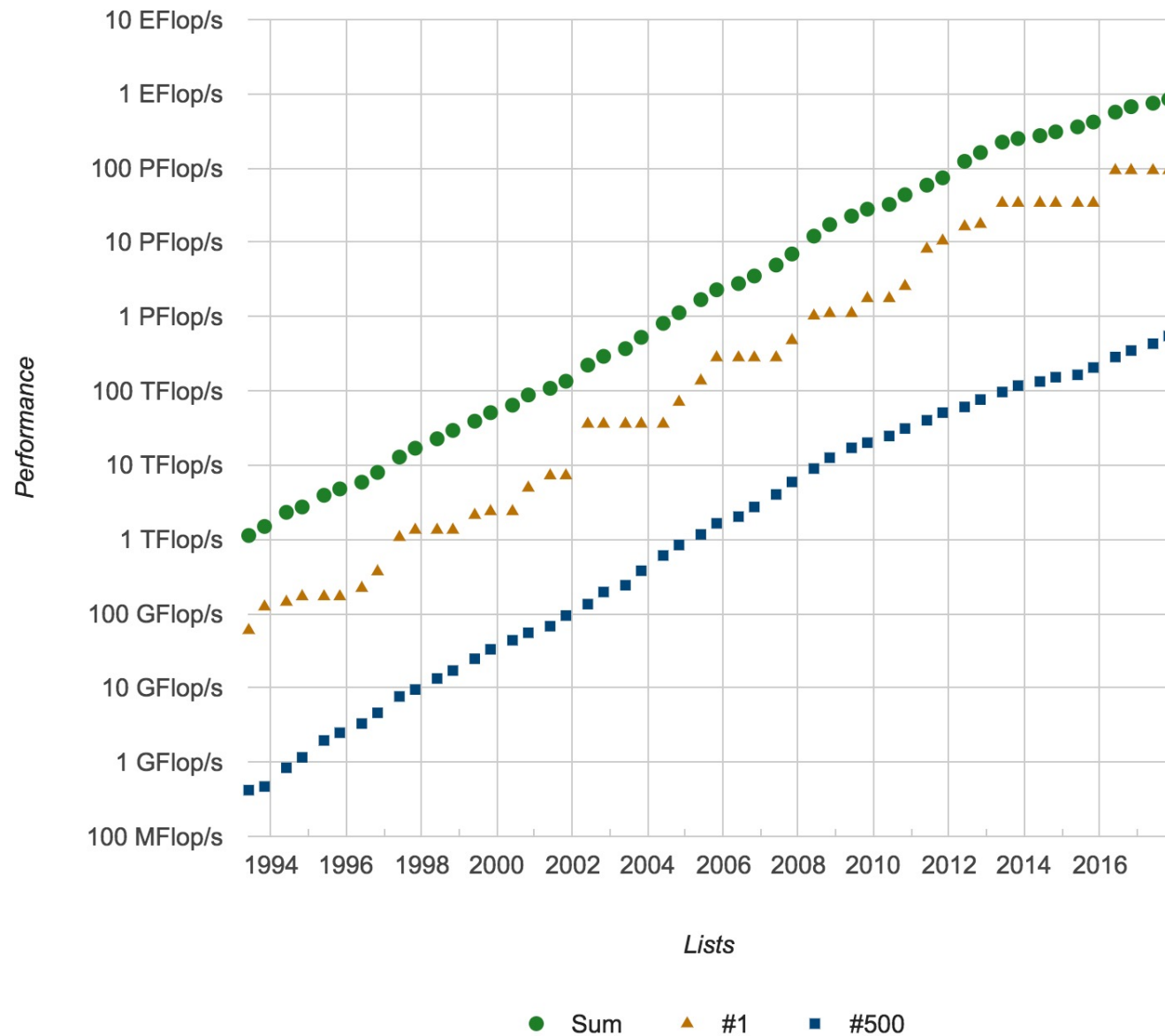


**The Von Neumann Bottleneck:**
- Memory bandwidth is not increasing as quickly as processor computing power
- Memory latency is decreasing very slowly
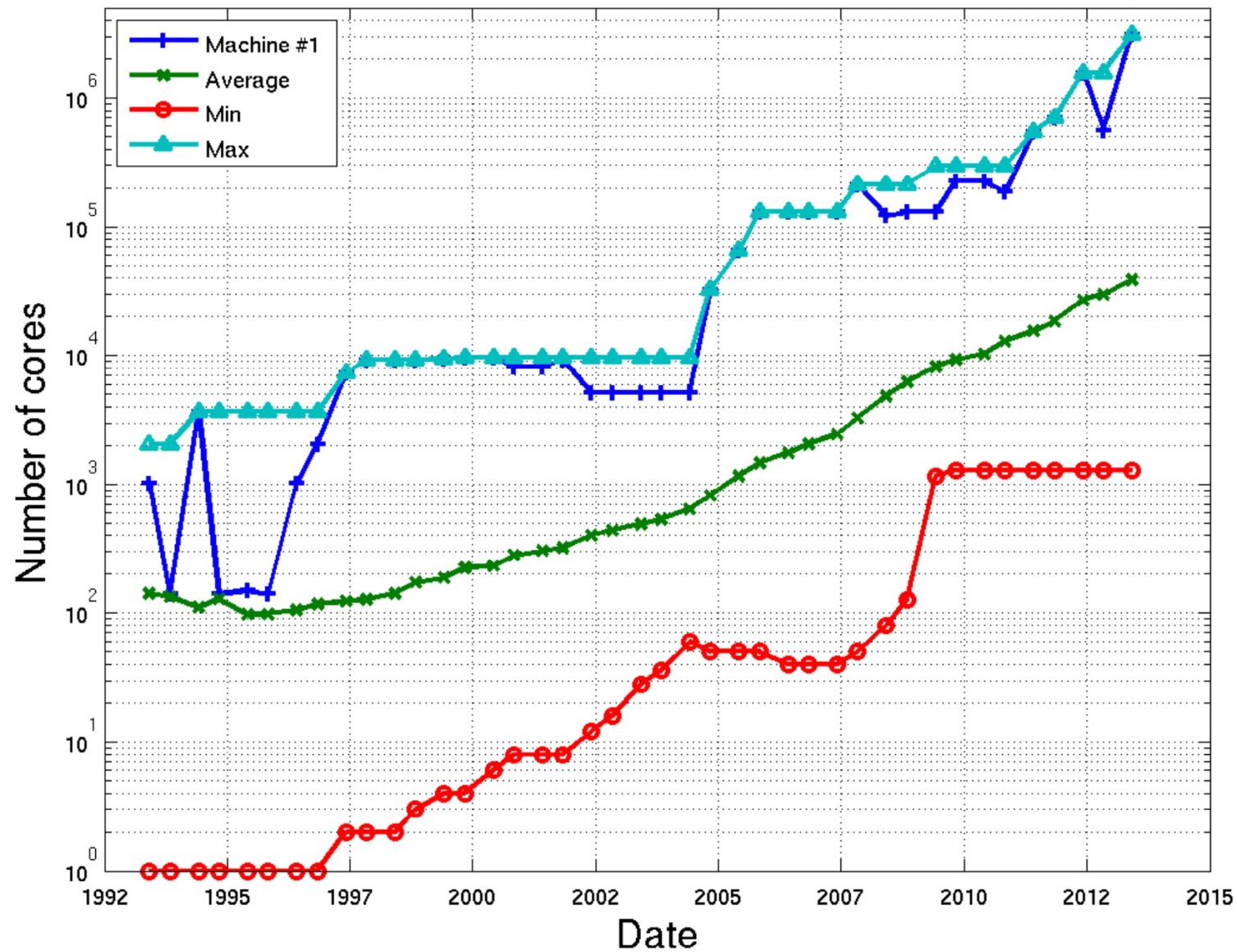- Number of computing cores per memory units is increasing

**Consequences and solutions:**
- The CPU wastes cycles while waiting for data.
- Introduction of cache memory (L1, L2, L3)
- Parallel access to memory (vector architectures, AVX)

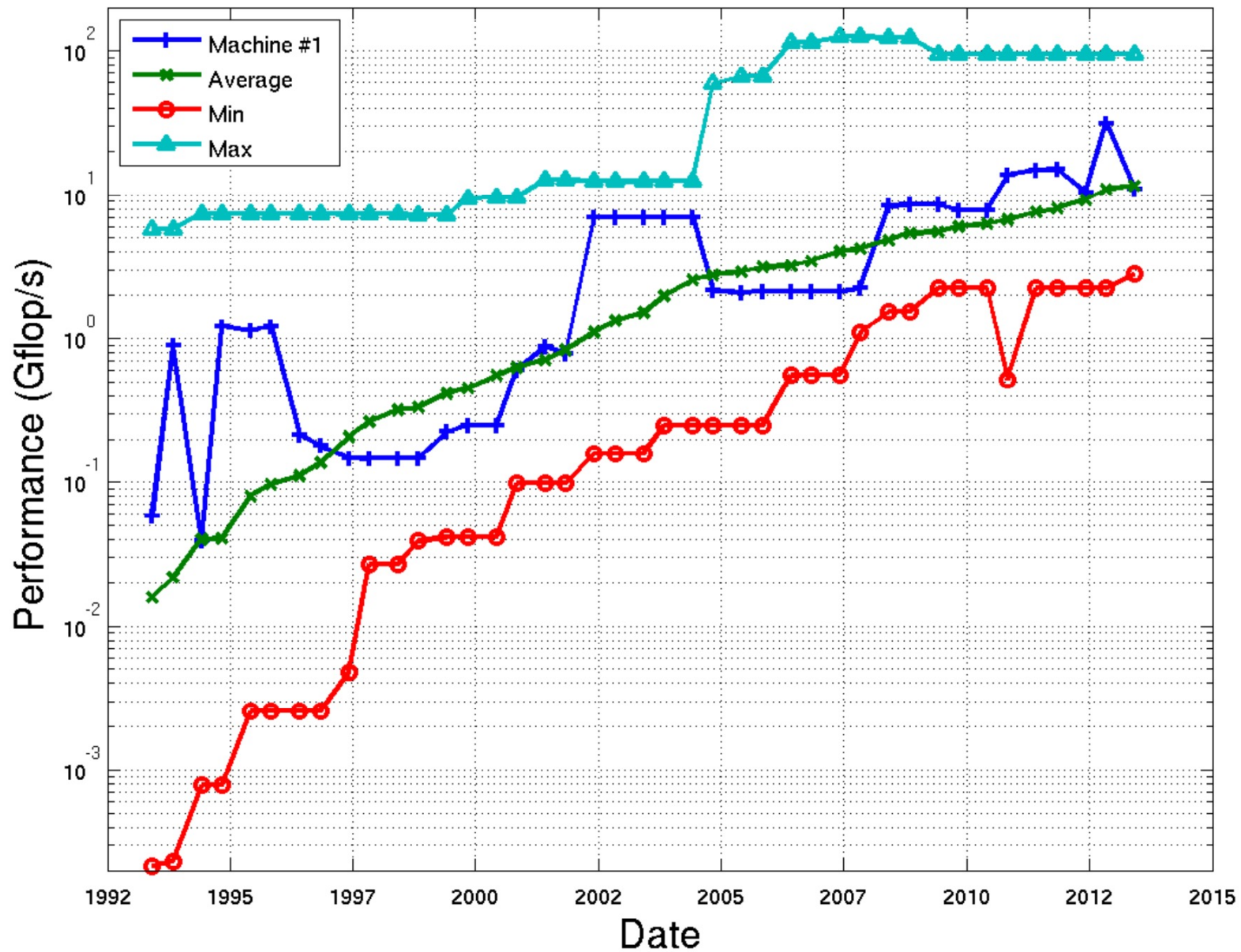# Performance in the Top 500

# Number of cores in the Top 500

# Performance per core in the Top 500

# Parallel Processing

**Gene Amdahl (1967):** Amdahl's Law

**Statement of the law:**
The theoretical maximum speedup obtained by parallelizing a code ideally, for a given problem with a fixed size:

$$\text{Speedup}\,(N) = \frac{T_s}{T_p\,(N)} = \frac{1}{\alpha + \frac{1-\alpha}{N}}$$

$T_s$ : execution time of the serial code
$T_p$ : execution time of the parallel code
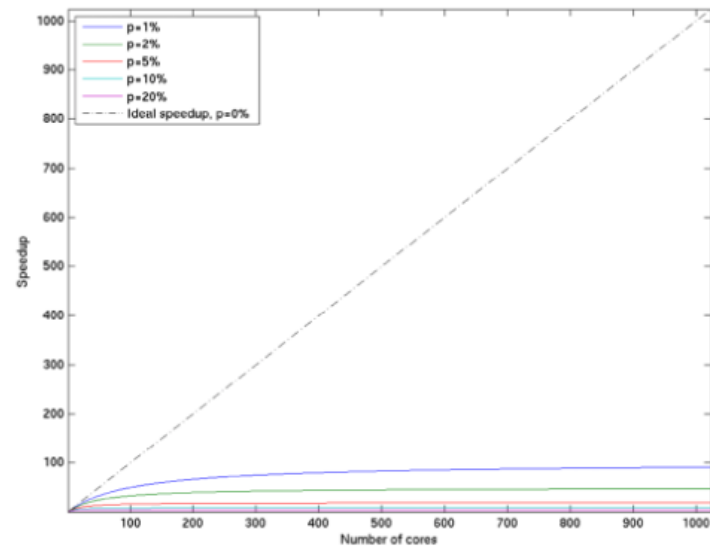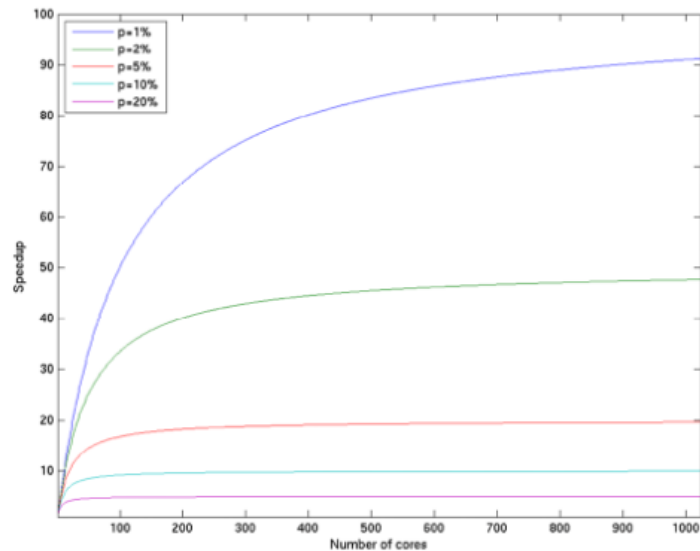$\alpha$ : fraction of the code that is not parallel
$N$ : number of processors

$$\text{Speedup}\,(N) \to \frac{1}{\alpha} \text{ for } N \to +\infty$$

# Strong Scaling

## Theoretical Maximum Speedup

| Cores | $\alpha$ (%) | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | 0 | 0.01 | 0.1 | 1 | 2 | 5 | 10 | 25 | 50 |
| 10 | 10 | 9.99 | 9.91 | 9.17 | 8.47 | 6.90 | 5.26 | 3.08 | 1.82 |
| 100 | 100 | 99.0 | 91.0 | 50.2 | 33.6 | 16.8 | 9.17 | 3.88 | 1.98 |
| 1000 | 1000 | 909 | 500 | 91 | 47.7 | 19.6 | 9.91 | 3.99 | 1.998 |
| 10000 | 10000 | 5000 | 909 | 99.0 | 49.8 | 19.96 | 9.99 | 3.99 | 2 |
| 100000 | 100000 | 9091 | 990 | 99.9 | 49.9 | 19.99 | 10 | 4 | 2 |
| $\infty$ | $\infty$ | 10000 | 1000 | 100 | 50 | 20 | 10 | 4 | 2 |

# Weak Scaling

**John Gustafson and Edwin Barsis (1988):** Gustafson's Law

**Statement of the law**:
The theoretical maximum speedup obtained by parallelizing a code ideally *for a problem of constant size per core*:

$$\mathrm{Speedup}\left(N\right) = \frac{T_s\left(N\right)}{T_p\left(N\right)} = \alpha + (1 - \alpha)N$$

Assuming that the execution time of the non-parallel part of the code remains constant, one has:

$$T_s\left(N\right) = \alpha T_0 + N(1 - \alpha)T_0$$
$$T_p\left(N\right) = \alpha T_0 + (1 - \alpha)T_0 = T_0$$

This law is more optimistic than Amdahl's law:

$$\mathrm{Speedup}\left(N\right) \to (1 - \alpha)N \text{ for } N \to +\infty$$

In both laws, everything is determined by $\alpha$, the non-parallel fraction of the code. Reducing $\alpha$ is called "parallel code optimization".

# Hardware Evolution

**Technical trends**:
- Computing power is doubling every year
- Massively parallel and many-cores architectures are dominant
- Since 2010, Graphical Processing Units become a key player
- Increasing hardware complexity (hybrid system, multiple cache layers)
- Memory per core is plateauing  or decreasing
- Performance per core is plateauing
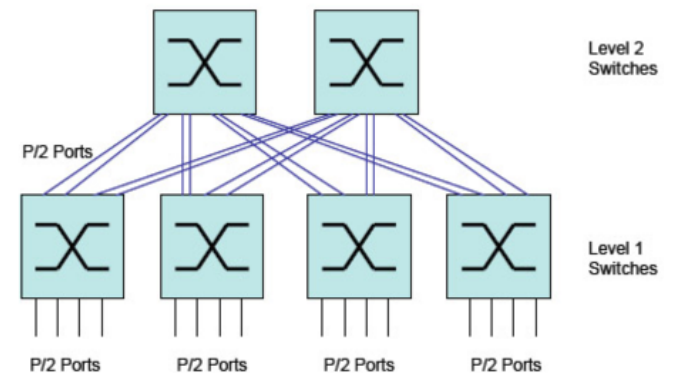- Disk Input/Output bandwidth is increasing very slowly



Cray XC50 with GPU

**Consequences**:
- It is necessary to exploit many (relatively slow) cores
- The memory per core is constant or even decreasing (memory limited)
- Raw performance of individual core not increasing anymore
- Higher level of parallelism and I/O bottlenecks
- More complex architecture (cache levels, accelerators, latency…)
- Multi-disciplinary approach, concept of "co-design"

# Communication Network

Connecting millions of processors requires a high -performance network
Most existing systems use "Infiniband" or variations.
Need a network switch to interconnect all the nodes
Need high quality fiber cables to connect each node to the switch



Fat-tree network: non-blocking network topology invented by Charles Clos (1953)
Depending on how many switches you can afford, you might choose blocking configurations.
http://www.mellanox.com/clusterconfig/
http://clusterdesign.org/fat-trees/

Key network parameters (switches and cables) are the **latency** and the **bandwidth**
**Infiniband network:** latency 0.5-1 microsec, bandwidth 5-10 GB/sec

# Parallel Programming Languages

**Evolution of programming methods:**

- MPI is still the dominant programming technique

- Hybrid OpenMP/MPI approach most effective on supercomputers

- GPU programming develops quickly
    - CUDA
    - OpenACC and OpenMP

- Message Passing directly within the GPU

- New specific parallel programming languages are developed:
    - Co-array Fortran, PGAS, X10, Chapel…

- New runtime systems to handle task-based parallelism:
    - Charm++, HPX, Kokkos

# Distributed memory versus shared memory

**MPI uses a distributed memory paradigm**
- Data are transferred explicitly between nodes through the network

**OpenMP uses a shared memory paradigm**
- Data are shared implicitly within the node through the Random-Access Memory.