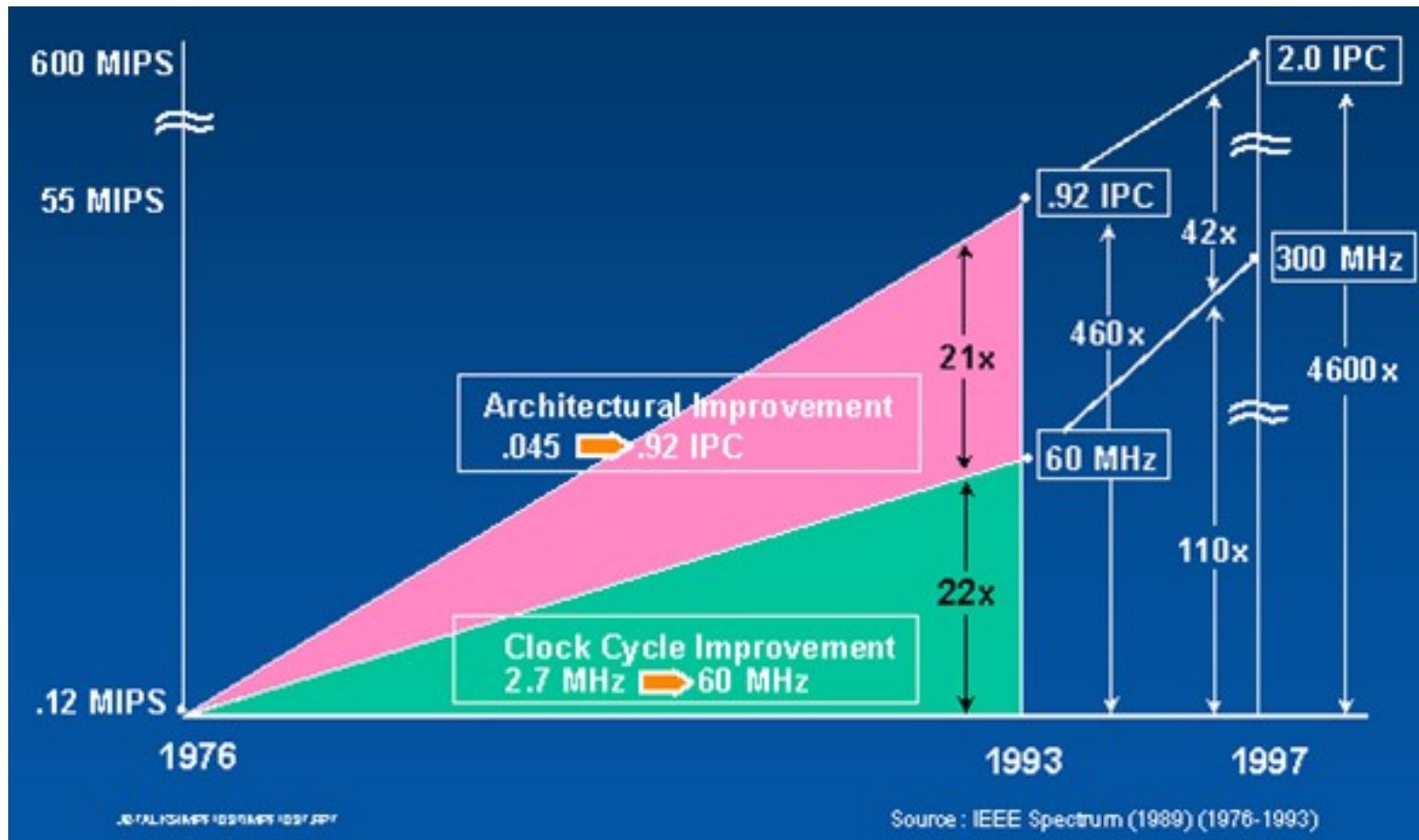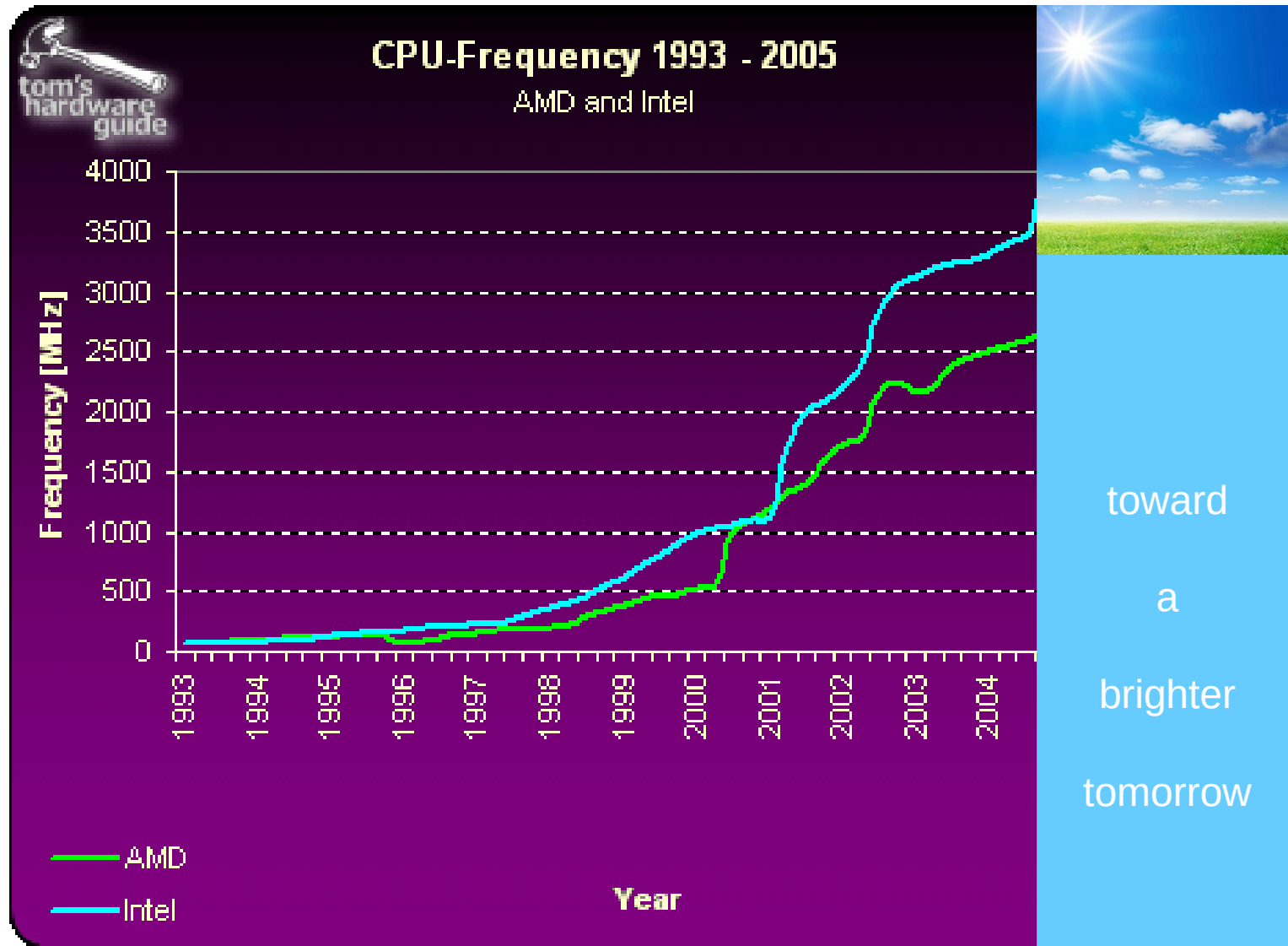# CUDA Programming

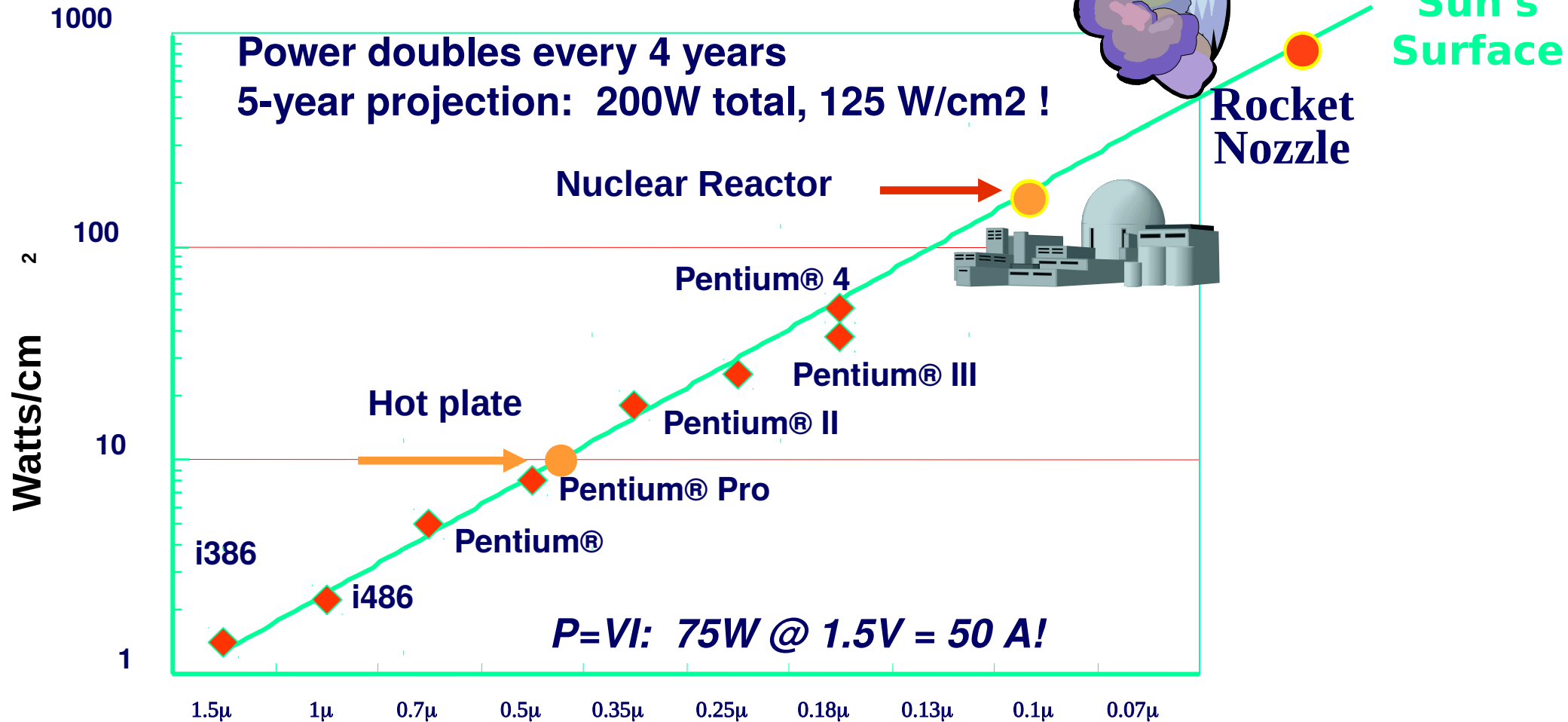# The Good Old Days for Software

**Source: J. Birnbaum**



- Single-processor performance experienced dramatic improvements from **clock**, and **architectural** improvement (Pipelining, Instruction-Level-Parallelism).
- Applications experienced **automatic** performance improvement.
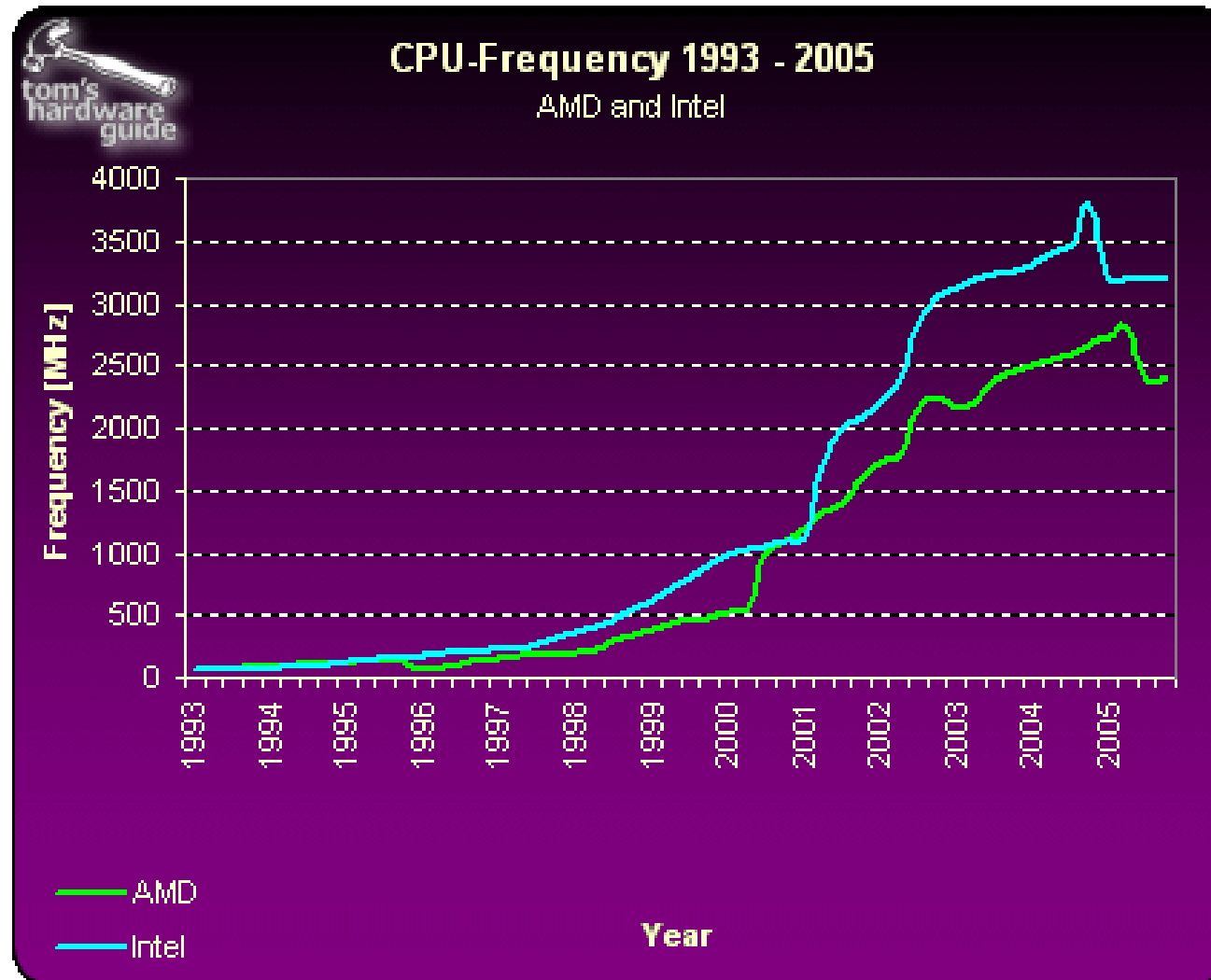
# Hitting the Power Wall



CPU-Frequency 1993 - 2005
AMD and Intel

toward a brighter tomorrow

http://img.tomshardware.com/us/2005/11/21/the_mother_of_all_cpu_charts_2005/cpu_frequency.gif

# Hitting the Power Wall

**Watts/cm²**

1000

**Power doubles every 4 years**
**5-year projection:  200W total, 125 W/cm2 !**

100

**Nuclear Reactor**

**Pentium® 4**

**Pentium® III**

**Hot plate**

**Pentium® II**

10

**Pentium® Pro**

**i386**

**Pentium®**

**i486**

*P=VI:  75W @ 1.5V = 50 A!*

1

**Rocket Nozzle**

**Sun's Surface**

| 1.5μ | 1μ | 0.7μ | 0.5μ | 0.35μ | 0.25μ | 0.18μ | 0.13μ | 0.1μ | 0.07μ |

"New Microarchitecture Challenges in the Coming Generations of CMOS Process Technologies" – Fred Pollack, Intel Corp. Micro32 conference key note - 1999. Courtesy Avi Mendelson, Intel.

4

# Hitting the Power Wall



http://img.tomshardware.com/us/2005/11/21/the_mother_of_all_cpu_charts_2005/cpu_frequency.gif

**2004 – Intel cancels Tejas and Jayhawk due to**
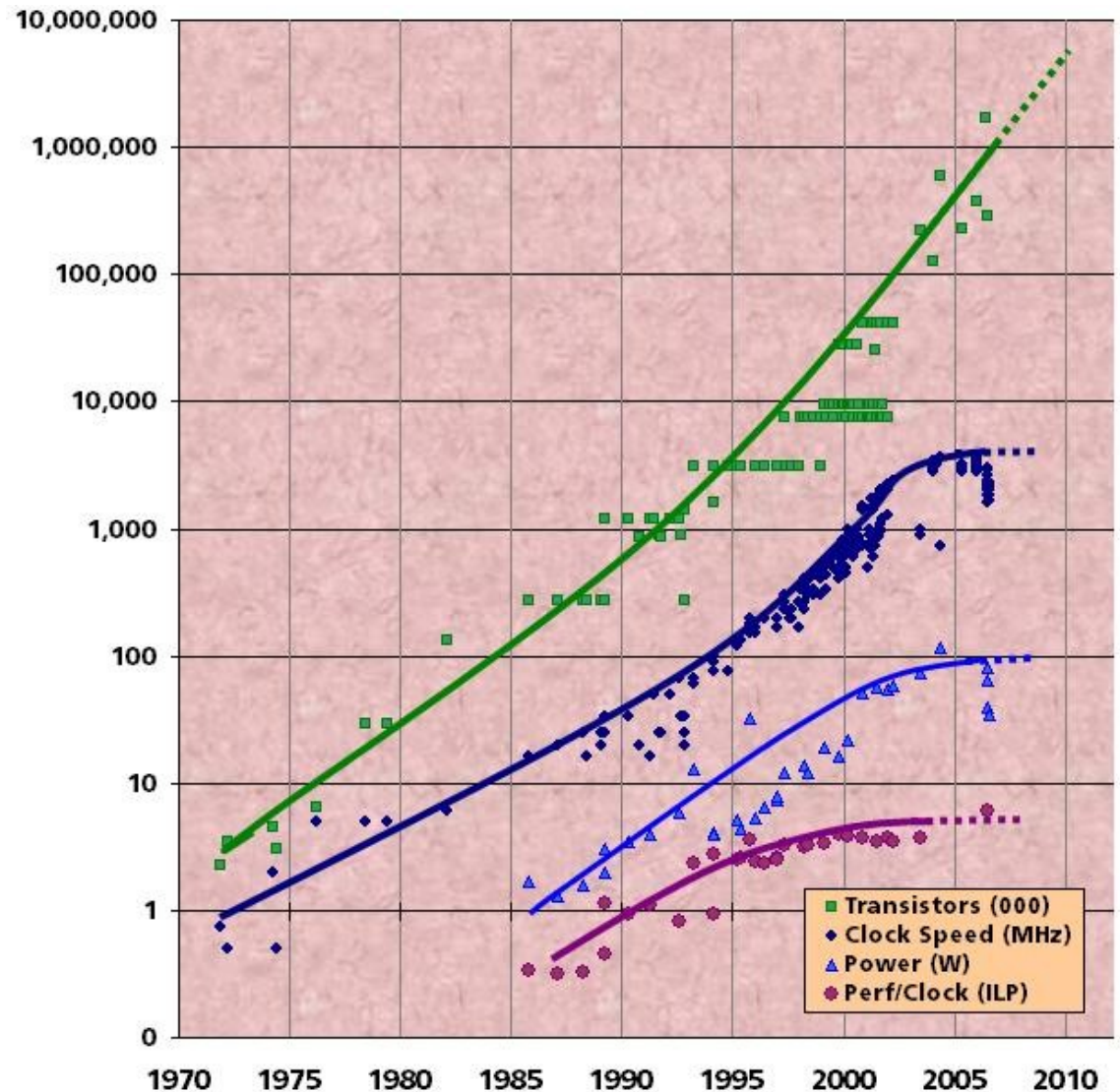*heat problems due to the extreme power consumption of the core ...*

# The Only Option: Use Many Cores

Chip density is increasing by ~2x every 2 years

- Clock speed is not
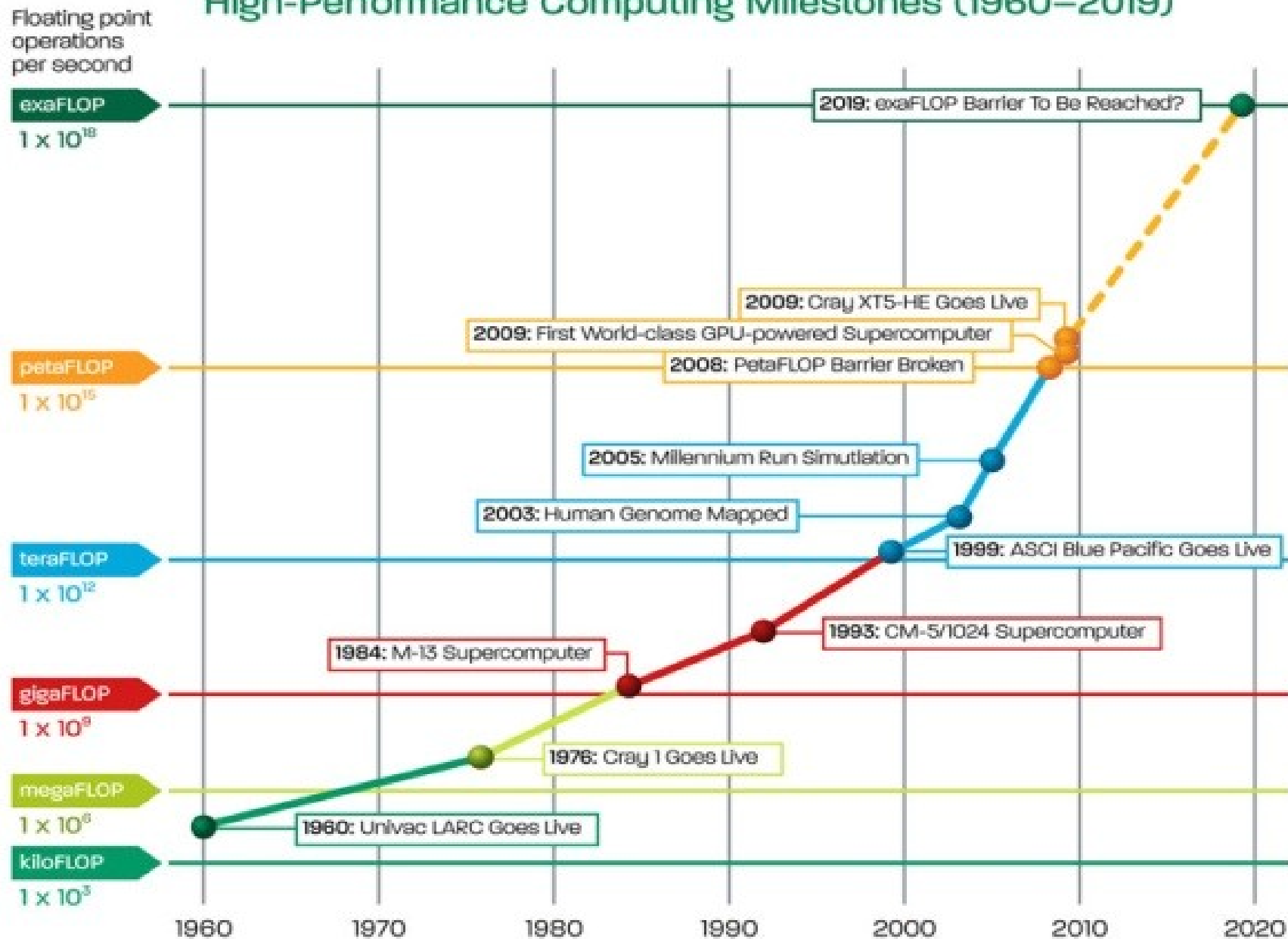
- Number of processor cores may double

There is little or no more hidden parallelism (ILP) to be found

Parallelism must be exposed to and managed by software



Legend:
- Transistors (000)
- Clock Speed (MHz)
- Power (W)
- Perf/Clock (ILP)

High-Performance Computing Milestones (1960–2019)

Floating point operations per second

exaFLOP $1 \times 10^{18}$
petaFLOP $1 \times 10^{15}$
teraFLOP $1 \times 10^{12}$
gigaFLOP $1 \times 10^{9}$
megaFLOP $1 \times 10^{6}$
kiloFLOP $1 \times 10^{3}$

2019: exaFLOP Barrier To Be Reached?
2009: Cray XT5-HE Goes Live
2009: First World-class GPU-powered Supercomputer
2008: PetaFLOP Barrier Broken
2005: Millennium Run Simutlation
2003: Human Genome Mapped
1999: ASCI Blue Pacific Goes Live
1993: CM-5/1024 Supercomputer
1984: M-13 Supercomputer
1976: Cray 1 Goes Live
1960: Univac LARC Goes Live

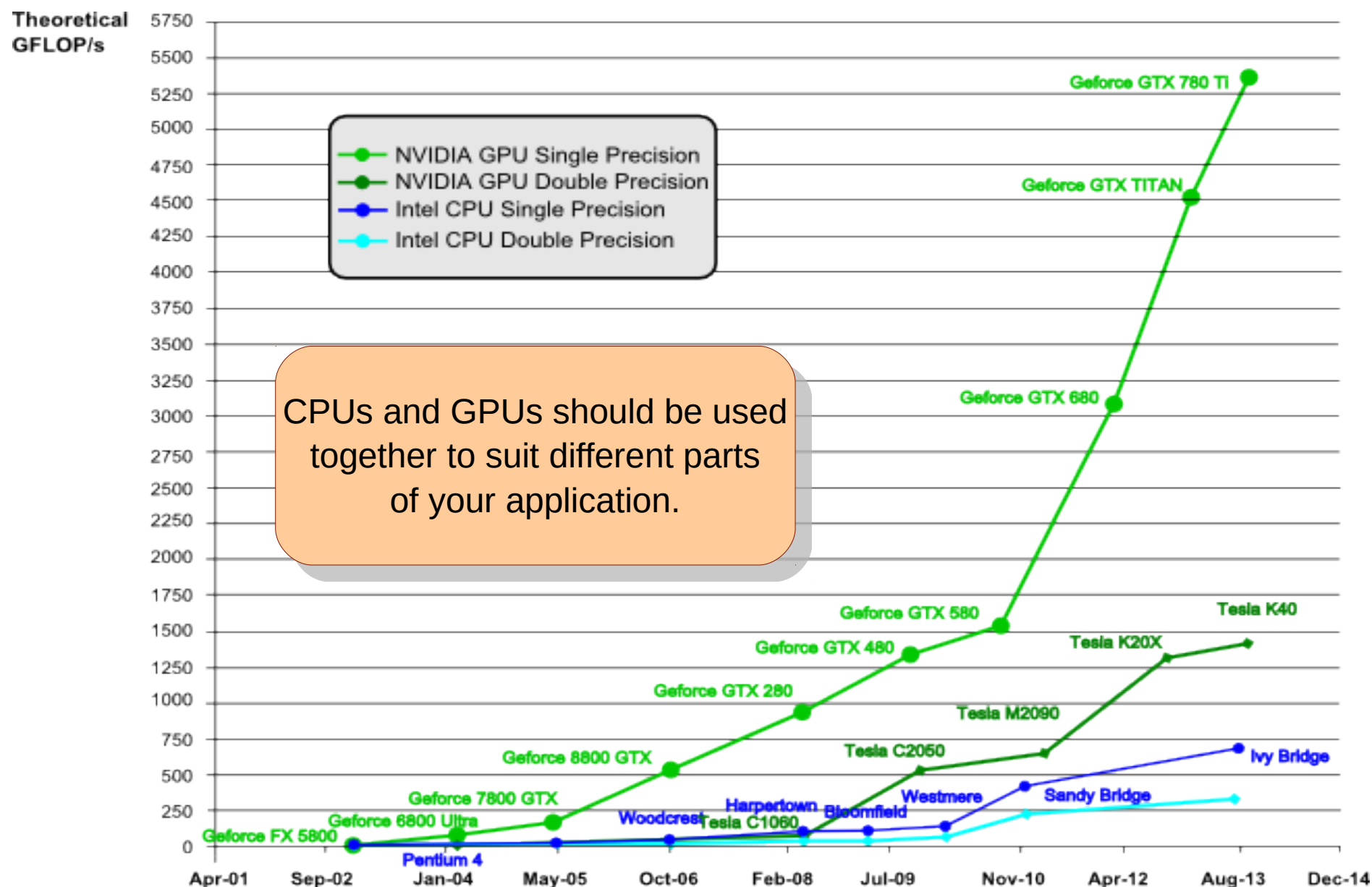1960  1970  1980  1990  2000  2010  2020

# Parallel Platforms

- Shared memory systems (multi-core)

- Distributed systems (cluster)

- Graphics Processing Units (many-core)

- Field-Programmable Gate Arrays (configurable after manufacturing)

- Application-Specific Integrated Circuits

- Heterogeneous Systems

# GPU-CPU Performance Comparison



CPUs and GPUs should be used together to suit different parts of your application.

Source: Thorsten Thormählen

# In this course...

- Basic GPU Programming

  - Computation, Memory, Synchronization, Debugging

- Topics in GPU Programming

  - Unified virtual memory, multi-GPU, peer access

# Logistics

- You need to arrange for your GPU.

  - Your laptop may have one.

  - With gmail account, you get some GPU time on Google colab.

  - You can use the central computing facilities at your institute.