# An Overview of High Performance Computing

**Jack Dongarra**
**University of Tennessee**
**and**
**Oak Ridge National Laboratory**

---

# Overview

♦ **Look at fastest computers**
  ➢ **From the Top500**
♦ **Some of the changes that face us**
  ➢ **Hardware**
  ➢ **Software**
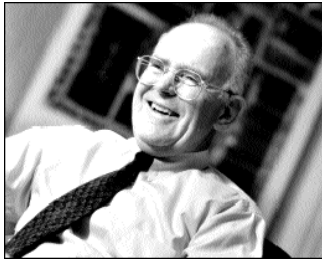  ➢ **Algorithms**

## Technology Trends: Microprocessor Chip Capacity

**Gordon Moore**
**(co-founder of Intel)**
Electronics Magazine, 1965
Number of devices/chip doubles
every 18 months

**2X transistors/Chip Every**
**1.5 years "Moore's Law"**

---

The experts look ahead

## Cramming more components onto integrated circuits

With unit cost falling as the number of components per circuit rises, by 1975 economics may dictate squeezing as many as 65,000 components on a single silicon chip

By Gordon E. Moore

Director, Research and Development Laboratories, Fairchild Semiconductor division of Fairchild Camera and Instrument Corp.

The future of integrated electronics is the future of electronics itself. The advantages of integration will bring about a proliferation of electronics, pushing this science into many new areas.

Integrated circuits will lead to such wonders as home computers—or at least terminals connected to a central computer—automatic controls for automobiles, and personal portable communications equipment. The electronic wristwatch needs only a display to be feasible today.

But the biggest potential lies in the production of large systems. In telephone communications, integrated circuits in digital filters will separate channels on multiplex equipment. Integrated circuits will also switch telephone circuits and perform data processing.

Computers will be more powerful, and will be organized in completely different ways. For example, memories built of integrated electronics may be distributed throughout the

The author

Dr. Gordon E. Moore is one of the new breed of electronic engineers, schooled in the physical sciences rather than in electronics. He earned a B.S. degree in chemistry from the University of California and a Ph.D. degree in physical chemistry from the California Institute of Technology. He was one of the founders of Fairchild Semiconductor and has been director of the research and development laboratories since 1959.

machine instead of being concentrated in a central unit. In addition, the improved reliability made possible by integrated circuits will allow the construction of larger processing units. Machines similar to those in existence today will be built at lower costs and with faster turn-around.

Present and future

By integrated electronics, I mean all the various technologies which are referred to as microelectronics today as well as any additional ones that result in electronics functions supplied to the user as irreducible units. These technologies were first investigated in the late 1950's. The object was to miniaturize electronics equipment to include increasingly complex electronic functions in limited space with minimum weight. Several approaches evolved, including microassembly techniques for individual components, thin-film structures and semiconductor integrated circuits.

Each approach evolved rapidly and converged so that each borrowed techniques from another. Many researchers believe the way of the future to be a combination of the various approaches.

The advocates of semiconductor integrated circuitry are already using the improved characteristics of thin-film resistors by applying such films directly to an active semiconductor substrate. Those advocating a technology based upon films are developing sophisticated techniques for the attachment of active semiconductor devices to the passive film arrays.

Both approaches have worked well and are being used in equipment today.

Electronics, Volume 38, Number 8, April 19, 1965

---

# TOP 500 superCOMPUTER

**H. Meuer, H. Simon, E. Strohmaier, & JD**

- Listing of the 500 most powerful
  Computers in the World
- Yardstick: Rmax from LINPACK MPP
  $Ax=b,$ *dense problem*



TPP performance

Rate

Size

- Updated twice a year
  SC'xy in the States in November
  Meeting in Germany in June

All data available from **www.top500.org**

4

# Performance Development



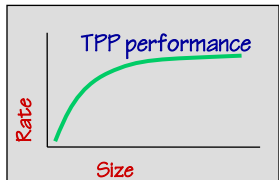| | |
|---|---|
| 1 Pflop/s | 2.3 PF/s |
| 100 Tflop/s | 280.6 TF/s |
| 10 Tflop/s | SUM |
| 1 Tflop/s | 1.167 TF/s   N=1 |
| 100 Gflop/s | 59.7 GF/s |
| 10 Gflop/s | |
| 1 Gflop/s | 0.4 GF/s |
| 100 Mflop/s | |

Labels in chart: IBM BlueGene/L, NEC Earth Simulator, IBM ASCI White LLNL, 1.646 TF/s, Intel ASCI Red Sandia, Fujitsu 'NWT' NAL, N=500, My Laptop

Years: 1993, 1994, 1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005

*OO*                                                                 5

---

# Architecture/Systems Continuum

**Tightly Coupled**

- ♦ **Custom processor with custom interconnect**
  - ➢ **Cray X1**
  - ➢ **NEC SX-8**
  - ➢ **IBM Regatta**
  - ➢ **IBM Blue Gene/L**
- ♦ **Commodity processor with custom interconnect**
  - ➢ **SGI Altix**
    - ➢ **Intel Itanium 2**
  - ➢ **Cray XT3, XD1**
    - ➢ **AMD Opteron**
- ♦ **Commodity processor with commodity interconnect**
  - ➢ **Clusters**
    - ➢ **Pentium, Itanium, Opteron, Alpha**
    - ➢ **GigE, Infiniband, Myrinet, Quadrics**
  - ➢ **NEC TX7**
  - ➢ **IBM eServer**
  - ➢ **Dawning**

**Loosely Coupled**



Chart labels: Custom, Hybrid, Commod

Dates: Jun-93, Dec-93, Jun-94, Dec-94, Jun-95, Dec-95, Jun-96, Dec-96, Jun-97, Dec-97, Jun-98, Dec-98, Jun-99, Dec-99, Jun-00, Dec-00, Jun-01, Dec-01, Jun-02, Dec-02, Jun-03, Dec-03, Jun-04

*OO*                                                                 6

3

# Commodity Processors

- **Intel Pentium Nocona**
  - **3.6 GHz, peak = 7.2 Gflop/s**
  - **Linpack 100  = 1.8 Gflop/s**
  - **Linpack 1000 = 4.2 Gflop/s**

- **Intel Itanium 2**
  - **1.6 GHz, peak = 6.4 Gflop/s**
  - **Linpack 100  = 1.7 Gflop/s**
  - **Linpack 1000 = 5.7 Gflop/s**

- **AMD Opteron**
  - **2.6 GHz, peak = 5.2 Gflop/s**
  - **Linpack 100  = 1.6 Gflop/s**
  - **Linpack 1000 = 3.9 Gflop/s**

McKinley microprocessor

OO

7

---

# Architectures / Systems



- SIMD
- Single Proc.
- Cluster (360)
- Constellations
- SMP
- MPP

Cluster: Commodity processors & Commodity interconnect

Constellation: # of procs/node $\geqslant$ nodes in the system

OO

8

# Interconnects / Systems



- Others
- Cray Interconnect
- SP Switch
- Crossbar
- Quadrics
- Infiniband
- Myrinet (101)
- Gigabit Ethernet (249)
- N/A

*OO*

9

# Processor Types



- SIMD
- Sparc
- Vector
- MIPS
- Alpha
- HP
- AMD
- IBM Power
- Intel

*OO*

10

5

# Processors Used in Each of the 500 Systems

91% = 66% Intel
15% IBM
11% AMD

Hitachi SR8000 0%
Sun Sparc 1%
NEC 1%
HP Alpha 1%
Cray 2%
HP PA-RISC 3%
Intel IA-64 9%
AMD x86_64 11%
IBM Power 15%
Intel EM64T 16%
Intel IA-32 41%

# 26th List: The TOP10

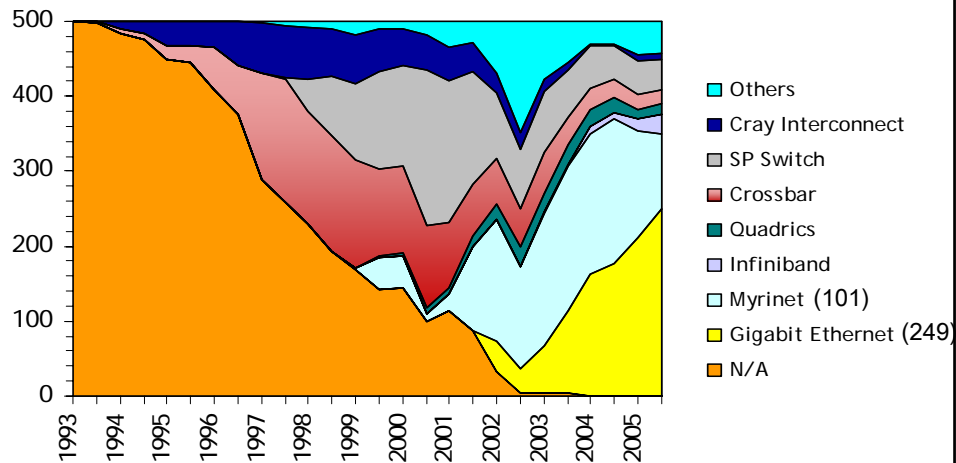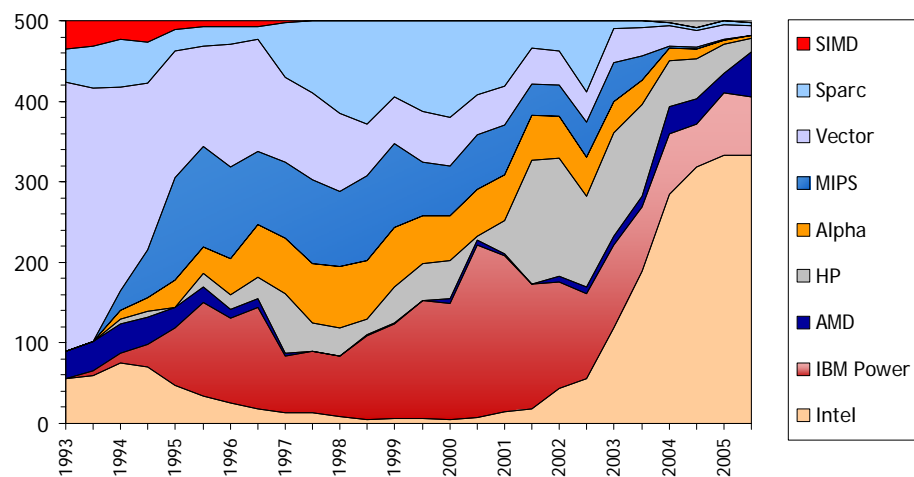| | Manufacturer | Computer | Rmax [TF/s] | Installation Site | Country | Year | #Proc |
|---|---|---|---|---|---|---|---|
| 1 | IBM | BlueGene/L eServer Blue Gene | 280.6 | DOE/NNSA/LLNL | USA | 2005 | 131072 |
| 2 | IBM | BGW eServer Blue Gene | 91.29 | IBM Thomas Watson | USA | 2005 | 40960 |
| 3 | IBM | ASC Purple Power5 p575 | 63.39 | DOE/NNSA/LLNL | USA | 2005 | 10240 |
| 4 | SGI | Columbia Altix, Itanium/Infiniband | 51.87 | NASA Ames | USA | 2004 | 10160 |
| 5 | Dell | Thunderbird Pentium/Infiniband | 38.27 | Sandia | USA | 2005 | 8000 |
| 6 | Cray | Red Storm Cray XT3 AMD | 36.19 | Sandia | USA | 2005 | 10880 |
| 7 | NEC | Earth-Simulator SX-5 | 35.86 | Earth Simulator Center | Japan | 2002 | 5120 |
| 8 | IBM | MareNostrum PPC 970/Myrinet | 27.91 | Barcelona Supercomputer Center | Spain | 2005 | 4800 |
| 9 | IBM | eServer Blue Gene | 27.45 | ASTRON University Groningen | Netherlands | 2005 | 12288 |
| 10 | Cray | Jaguar Cray XT3 AMD | 20.53 | Oak Ridge National Lab | USA | 2005 | 5200 |

# Customer Segments / Performance



ICL UT

Customer Segments / Performance

100% — Government
90% — Classified
80% — Vendor
70% — Academic
60% —
50% — Industry
40% —
30% —
20% —
10% — Research
0% —
1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005

OO                                                                    13

# Countries / Performance



ICL UT

Countries / Performance

100%
90%
80%
70%
60%
50%
40%
30%
20%
10%
0%
1993 1994 1995 1996 1997 1998 1999 2000 2001 2002 2003 2004 2005

- Others
- France    1.0%
- China     2.6%
- Germany 3.0%
- UK        5.4%
- Japan     6.0%
- US        68%

OO                                                                    14

7

# Asian Countries / Systems



Others 17
India 4
Korea 7
China 17
Japan 21

*OO*

15

# Concurrency Levels of the Top500



64k-128k
32k-64k
16k-32k
8k-16k
4k-8k
2049-4096
1025-2048
513-1024
257-512
129-256
65-128
33-64
17-32
9-16
5-8
3-4
2
1

*OO*

16

# Concurrency Levels of the Top500



Chart: # processors (y-axis, logarithmic scale from 1 to 1,000,000) vs. date (x-axis, Jun-93 to Jun-05)

- Maximum: 131K
- Average: 1462
- Minimum: 50

*OO* — 17

# IBM BlueGene/L
## 131,072 Processors (#1-64K and #2-40K)



1.6 MWatts (1600 homes)
43,000 ops/s/person

System
(64 racks, 64x32x32)
131,072 procs

Rack
(32 Node boards, 8x8x16)
2048 processors

Node Board
(32 chips, 4x4x2)
16 Compute Cards
64 processors

Compute Card
(2 chips, 2x1x1)
4 processors

Chip
(2 processors)

180/360 TF/s
32 TB DDR

2.9/5.7 TF/s
0.5 TB DDR

Full system total of
131,072 processors

90/180 GF/s
16 GB DDR

2.8/5.6 GF/s
4 MB (cache)

5.6/11.2 GF/s
1 GB DDR

The compute node ASICs include all networking and processor functionality.
Each compute ASIC includes two 32-bit superscalar PowerPC 440 embedded
*OO* cores (note that L1 cache coherence is not maintained between these cores).
(13K sec about 3.6 hours; n=1.8M

**"Fastest Computer"**
**BG/L 700 MHz 131K proc**
**64 racks**
**Peak:      367 Tflop/s**
**Linpack:   281 Tflop/s**
**77% of peak**

18

9

# Performance Projection



Chart axis labels (left, top to bottom): 1 Eflop/s, 100 Pflop/s, 10 Pflop/s, 1 Pflop/s, 100 Tflop/s, 10 Tflop/s, 1 Tflop/s, 100 Gflop/s, 10 Gflop/s, 1 Gflop/s, 100 Mflop/s

X-axis: 1993 1995 1997 1999 2001 2003 2005 2007 2009 2011 2013 2015
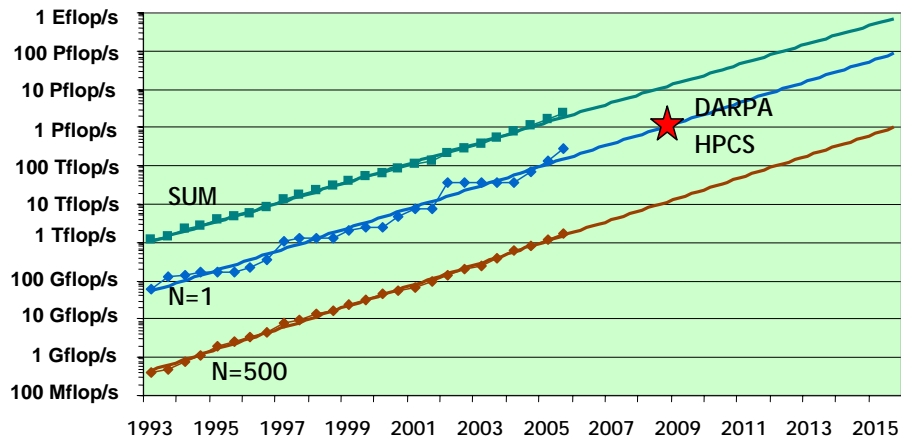
Labels on chart: DARPA HPCS, SUM, N=1, N=500

---

# A PetaFlop Computer by the End of the Decade

♦ **10 Companies working on a building a Petaflop system by the end of the decade.**
  - **Cray**
  - **IBM**
  - **Sun**  } HPCS
  - **Dawning**
  - **Galactic**  } Chinese Companies
  - **Lenovo**
  - **Hitachi**
  - **NEC**  } Japanese "Life Simulator" (10 Pflop/s)
  - **Fujitsu**
  - **Bull**

# Flops per Gross Domestic Product

## Based on the November 2005 Top500

Chart values (approximate):
- Isreal: 191
- USA: 126
- New Zealand: 91
- Switzerland: 89
- Russia: 75
- UK: 54
- Netherlands: 53
- Australian: 49
- South Korea: 44
- Saudi Arabia: 41
- China: 32
- Japan: 29
- Ireland: 26
- Spain: 26
- Germany: 24
- Taiwan: 22
- Canada: 19
- Brazil: 17

---

# KFlop/s per Capita (Flops/Pop)

## Based on the November 2005 Top500

Hint: Peter Jackson had something to do with this

← WETA Digital (Lord of the Rings)

Chart values (approximate):
- United States: 5350
- Switzerland: 4600
- Israel: 3800
- New Zealand: 2450
- United Kingdom: 2100
- Netherlands: 2050
- Australia: 1750
- Japan: 1050
- Germany: 850
- Spain: 720
- Canada: 650
- Korea, South: 620
- Sweden: 550
- Saudia Arabia: 430
- France: 380
- Taiwan: 300
- Italy: 250
- Brazil: 50
- Mexico: 50
- China: 30
- Russia: 30
- India: 30

OO

Has nothing to do with the 47.2 million sheep in NZ

# Fuel Efficiency: GFlops/Watt

GFlops/Watt (y-axis): 0, 0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9

Systems (x-axis):
BlueGene/L, Blue Gene, ASC Purple p5 1.9 GHz, Columbia - SGI Altix 1.5 GHz, Thunderbird - Pentium 3.6 GHz, Red Storm Cray XT3, 2.0 GHz, Earth-Simulator, MareNostrum PPC 970, 2.2 GHz, Blue Gene, Jaguar - Cray XT3, 2.4 GHz, Thunder - Intel Itanium2 1.4GHz, Blue Gene, Blue Gene, Cray XT3, 2.6 GHz, Apple XServe, 2.0 GHz, Cray X1E (4GB), Cray X1E (2GB), ASCI Q - Alpha 1.25 GHz, IBM p5 575 1.9 GHz, System X 2.3 GHz Apple XServe/, SGI Altix 3700 1.6 GHz
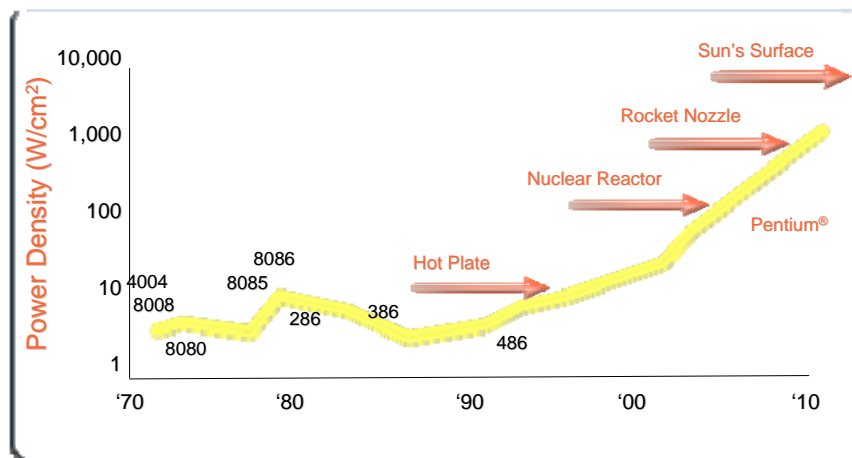
*OO*

Top 20 systems

23

Based on processor power rating only (3,>100,>800)

---

# Today's CPU Architecture:
## Heat becoming an unmanageable problem

Power Density (W/cm$^2$) (y-axis): 1, 10, 100, 1,000, 10,000

Labels: Sun's Surface, Rocket Nozzle, Nuclear Reactor, Hot Plate, Pentium®

Data points: 4004, 8008, 8080, 8085, 8086, 286, 386, 486

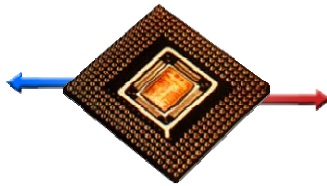X-axis: '70, '80, '90, '00, '10

Intel Developer Forum, Spring 2004 - Pat Gelsinger
(Pentium at 90 W)

Cube relationship between the cycle time and power

24

12

# Increasing CPU Performance:
## A Delicate Balancing Act

**Lower Voltage**

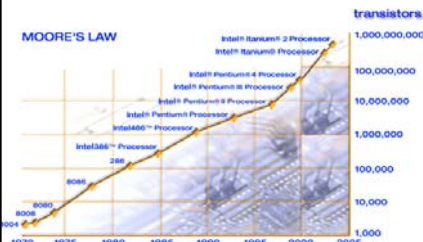**Increase Clock Rate & Transistor Density**

We have seen increasing number of gates on a chip and increasing clock speed.

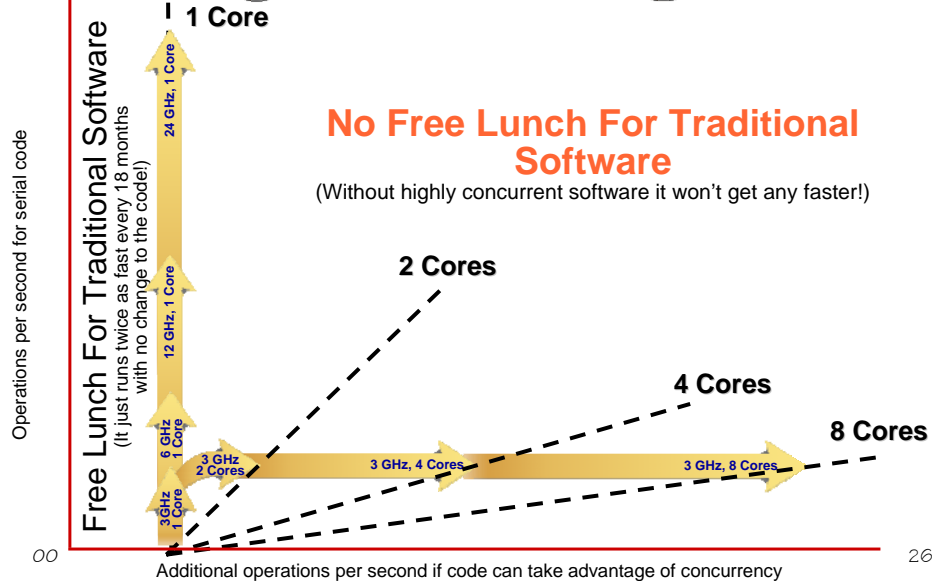Heat becoming an unmanageable problem, Intel Processors > 100 Watts

We will not see the dramatic increases in clock speeds in the future.

However, the number of gates on a chip will continue to increase.

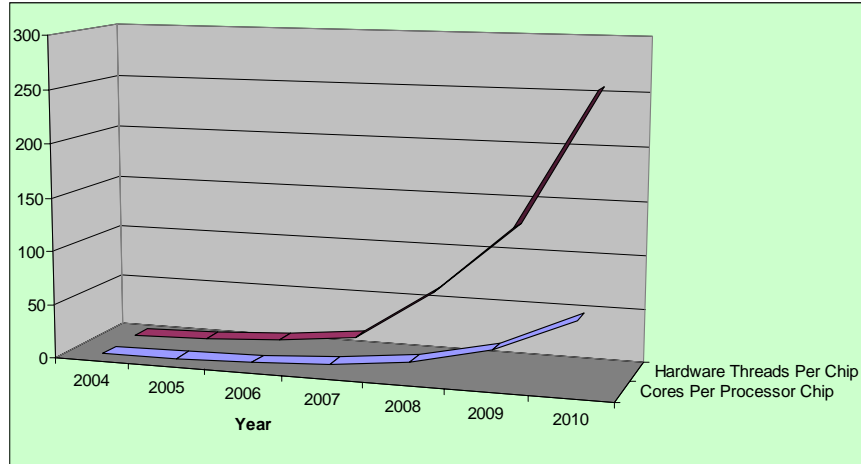Intel Yonah will double the processing power on a per watt basis.

MOORE'S LAW

transistors

Intel® Itanium® 2 Processor
Intel® Itanium® Processor
1,000,000,000

Intel® Pentium® 4 Processor
Intel® Pentium® III Processor
100,000,000

Intel® Pentium® II Processor
Intel® Pentium® Processor
Intel486™ Processor
10,000,000

Intel386™ Processor
286
1,000,000

8086
100,000

8080
8008
4004
10,000

1,000

1970  1975  1980  1985  1990  1995  2000  2005

25

---

# Change Is Coming

**Operations per second for serial code**

**Free Lunch For Traditional Software**
(It just runs twice as fast every 18 months with no change to the code!)

**1 Core**

24 GHz, 1 Core

12 GHz, 1 Core

6 GHz, 1 Core

3 GHz, 1 Core

**No Free Lunch For Traditional Software**
(Without highly concurrent software it won't get any faster!)

**2 Cores**

**4 Cores**

**8 Cores**

3 GHz, 2 Cores

3 GHz, 4 Cores

3 GHz, 8 Cores

Additional operations per second if code can take advantage of concurrency

26

13

# CPU Desktop Trends 2004-2010

- ♦ **Relative processing power will continue to double every 18 months**
- ♦ **256 logical processors per chip in late 2010**



Legend:
- Hardware Threads Per Chip
- Cores Per Processor Chip

Y-axis: 0, 50, 100, 150, 200, 250, 300
X-axis (Year): 2004, 2005, 2006, 2007, 2008, 2009, 2010

---

# Commodity Processor Trends

Bandwidth/Latency is the Critical Issue, not FLOPS

Got Bandwidth?

| | Annual increase | Typical value in 2005 |
|---|---|---|
| Single-chip floating-point performance | 59% | 4 GFLOP/s |
| Front-side bus bandwidth | 23% | 1 GWord/s = 0.25 word/flop |
| DRAM latency | (5.5%) | 70 ns = 280 FP ops = 70 loads |

Source: *Getting Up to Speed: The Future of Supercomputing*, National Research Council, 222 pages, 2004, National Academies Press, Washington DC, ISBN 0-309-09502-6.

# Fault Tolerance: Motivation

- **Trends in HPC:**
  - High end systems with thousand of processors

- **Increased probability of a node failure**
  - Most systems nowadays are robust

- **MPI widely accepted in scientific computing**
  - Process faults not tolerated in MPI model

**Mismatch between hardware and (non fault-tolerant) programming paradigm of MPI.**

# Reliability of Leading-Edge HPC Systems

| System | CPUs | Reliability |
|--------|------|-------------|
| LANL ASCI Q | 8,192 | MTBI: 6.5 hours. Leading outage sources: storage, CPU, memory. |
| LLNL ASCI White | 8,192 | MTBF: 5.0 hours ('01) and 40 hours ('03). Leading outage sources: storage, CPU, 3rd-party HW. |
| Pittsburgh Lemieux | 3,016 | MTBI: 9.7 hours. |

MTBI: mean time between interrupts = wall clock hours / # downtime periods
MTBF: mean time between failures (measured)

- **100K processor systems**
  - are here
  - we have fundamental challenges in dealing with machines of this size
  - ... and little in the way of programming support

# Future Challenge: Developing the Ecosystem for HPC

From the NRC Report on "The Future of Supercomputing":

♦ Hardware, software, algorithms, tools, networks, institutions, applications, and people who solve supercomputing applications can be thought of collectively as an ecosystem

♦ Research investment in HPC should be informed by the ecosystem point of view - progress must come on a broad front of interrelated technologies, rather than in the form of individual breakthroughs.



A supercomputer ecosystem is a continuum of computing platforms, system software, algorithms, tools, networks, and the people who know how to exploit them to solve computational science applications.

---

# Real Crisis With HPC Is With The Software

♦ **Our ability to configure a hardware system capable of 1 PetaFlop ($10^{15}$ ops/s) is without question just a matter of time and \$\$.**

♦ **A supercomputer application and software are usually much more long-lived than a hardware**
  ➢ **Hardware life typically five years at most…. Apps 20-30 years**
  ➢ **Fortran and C are the main programming models (still!!)**

♦ **The REAL CHALLENGE is Software**
  ➢ **Programming hasn't changed since the 70's**
  ➢ **HUGE manpower investment**
      ➢ **MPI… is that all there is?**
  ➢ **Often requires HERO programming**
  ➢ **Investments in the entire software stack is required (OS, libs, etc.)**

♦ **Software is a major cost component of modern technologies.**
  ➢ **The tradition in HPC system procurement is to assume that the software is free… SOFTWARE COSTS (over and over)**

# Summary of Current Unmet Needs

- **Performance / Portability**
- **Fault tolerance**
- **Memory bandwidth/Latency**
- **Adaptability: Some degree of autonomy to self optimize, test, or monitor.**
  - **Able to change mode of operation: static or dynamic**
- **Better programming models**
  - **Global shared address space**
  - **Visible locality**
- **Maybe coming soon (incremental, yet offering real benefits):**
  - **Global Address Space (GAS) languages: UPC, Co-Array Fortran, Titanium, Chapel)**
    - **"Minor" extensions to existing languages**
    - **More convenient than MPI**
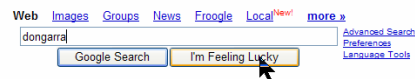    - **Have performance transparency via explicit remote memory references**

*OO* 33

---

# Collaborators / Support

- **Top500 Team**
  - **Erich Strohmaier, NERSC**
  - **Hans Meuer, Mannheim**
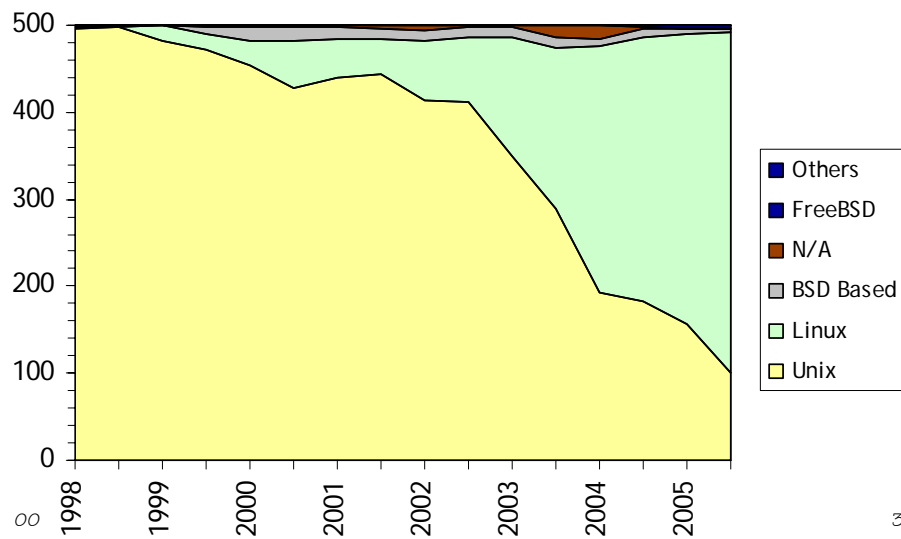  - **Horst Simon, NERSC**

http://www.top500.org/



*OO*

# Next Steps

- ♦ **Software to determine the checkpointing interval and number of checkpoint processors from the machine characteristics.**
  - ➢ Perhaps use historical information.
  - ➢ Monitoring
  - ➢ Migration of task if potential problem
- ♦ **Local checkpoint and restart algorithm.**
  - ➢ Coordination of local checkpoints.
  - ➢ Processors hold backups of neighbors.
- ♦ **Have the checkpoint processes participate in the computation and do data rearrangement when a failure occurs.**
  - ➢ Use p processors for the computation and have k of them hold checkpoint.
- ♦ **Generalize the ideas to provide a library of routines to do the diskless check pointing.**
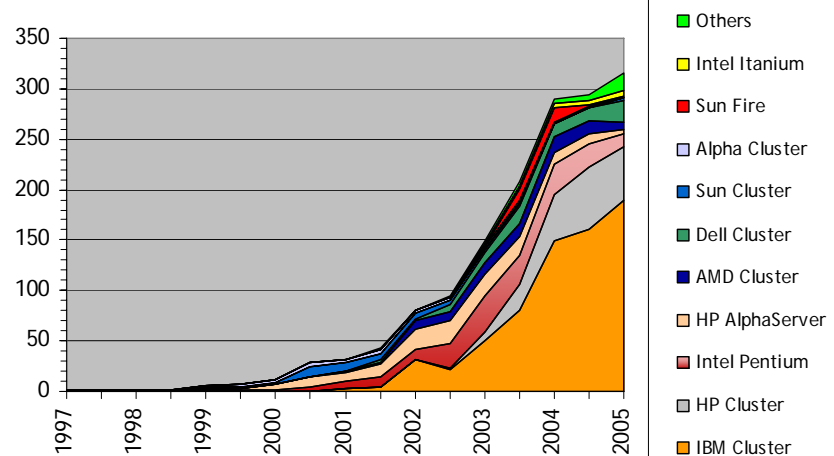- ♦ **Look at "real applications" and investigate "Lossy" algorithms.**
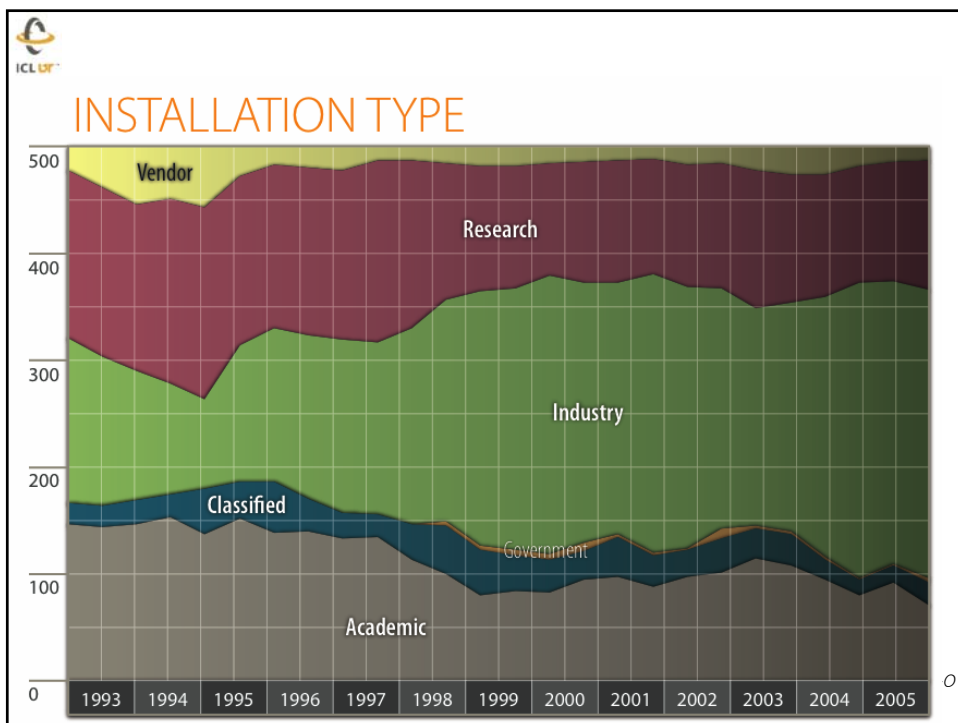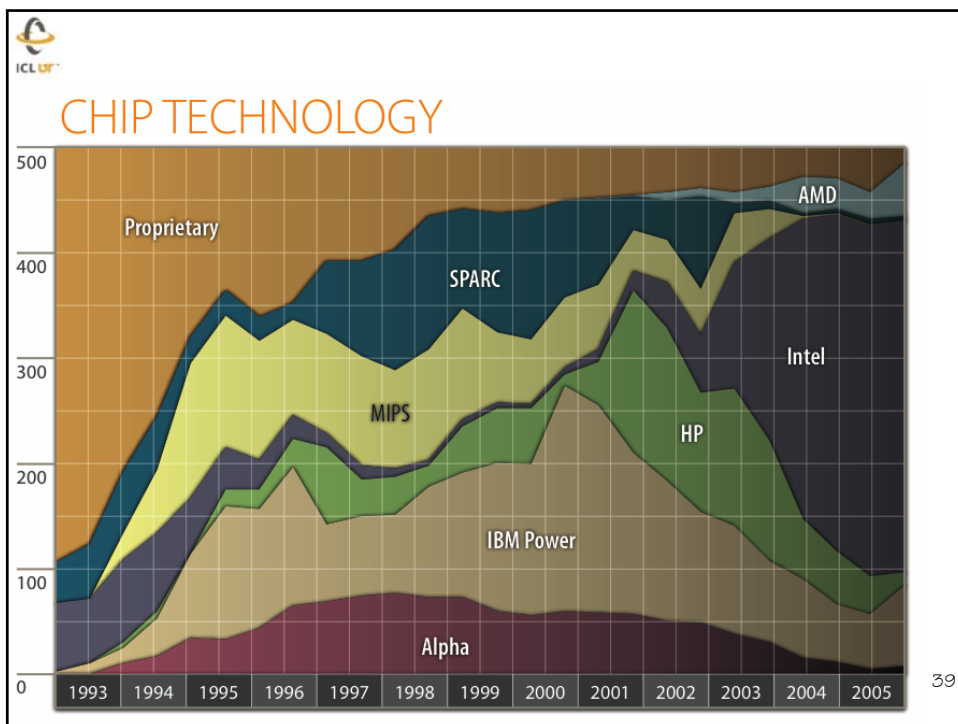
## Operating Systems / Systems



*OO*

## Clusters / Systems



*OO*

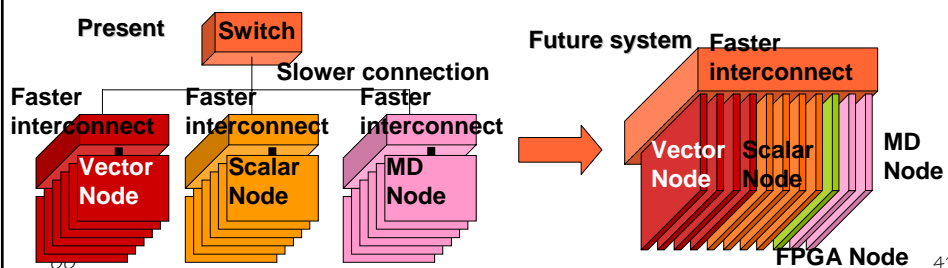# CHIP TECHNOLOGY



# INSTALLATION TYPE

## Japanese:
## Tightly-Coupled Heterogeneous System

- ♦ **Would like to get to 10 PetaFlop/s by 2011**
- ♦ **Scalable, fits any computer center**
  - ➤ **Size, cost, ratio of components**
- ♦ **Easy and low-cost to develop new component**
- ♦ **Scale merit of components**

**Present**

**Switch**

**Slower connection**

**Faster interconnect** **Faster interconnect** **Faster interconnect**

**Vector Node** **Scalar Node** **MD Node**

**Future system** **Faster interconnect**

**Vector Node** **Scalar Node** **MD Node**

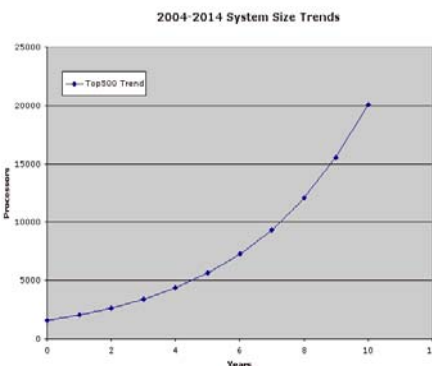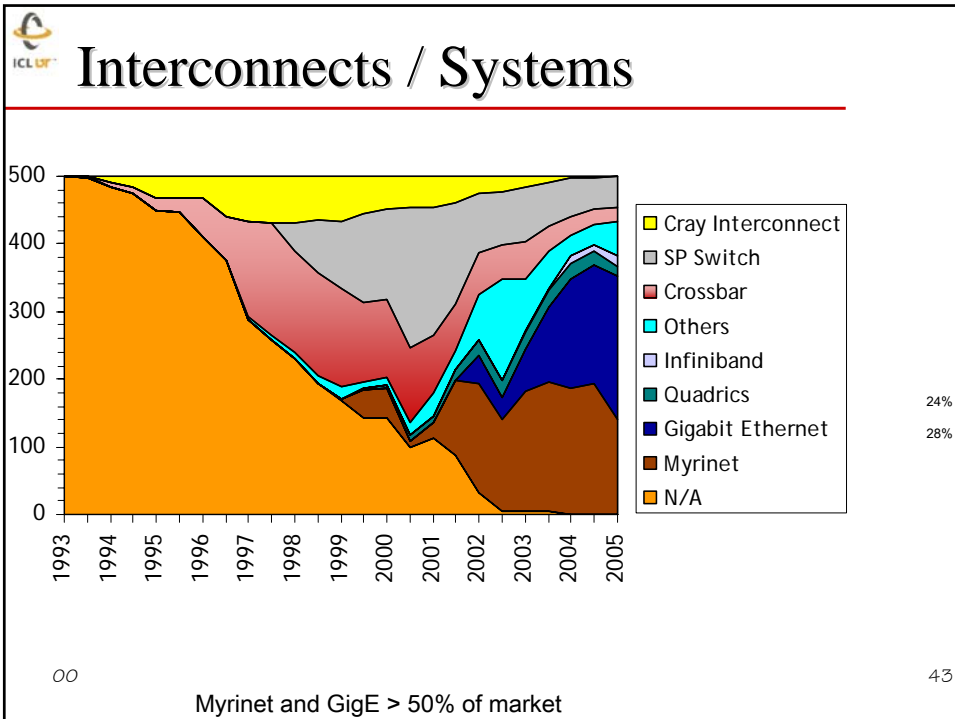**FPGA Node**

41

---

# How Big Is Big?

- ♦ **Every 10X brings new challenges**
  - ➤ **64 processors was once considered large**
    - ➤ **it hasn't been "large" for quite a while**
  - ➤ **1024 processors is today's "medium" size**
  - ➤ **8096 processors is today's "large"**
    - ➤ **we're struggling even here**

- ♦ **100K processor systems**
  - ➤ **are in construction**
  - ➤ **we have fundamental challenges in dealing with machines of this size**
  - ➤ **… and little in the way of programming support**

2004-2014 System Size Trends

Top500 Trend

Processors

Years

# Interconnects / Systems



Legend:
- Cray Interconnect
- SP Switch
- Crossbar
- Others
- Infiniband
- Quadrics — 24%
- Gigabit Ethernet — 28%
- Myrinet
- N/A

*OO*

43

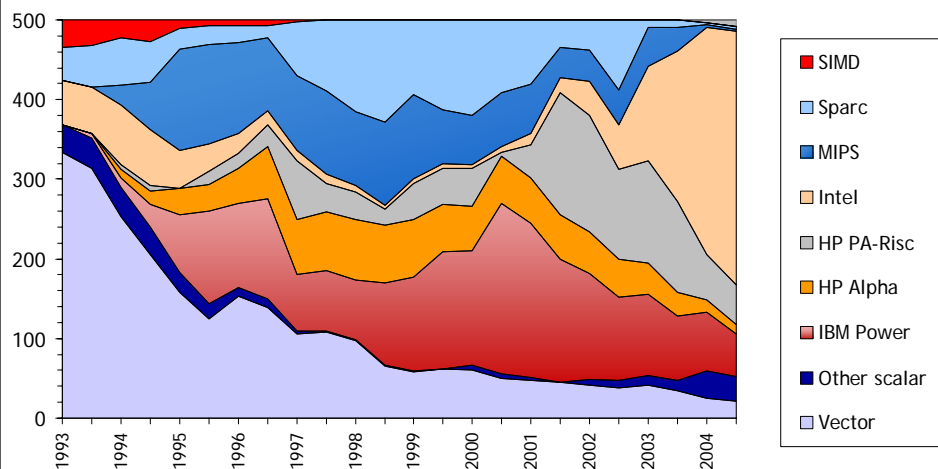Myrinet and GigE > 50% of market

---

# Real Crisis With HPC Is With The Software

♦ **Programming is stuck**
  ➢ **Arguably hasn't changed since the 60's**
♦ **It's time for a change**
  ➢ **Complexity is rising dramatically**
    ➢ **highly parallel and distributed systems**
      ➢ From 10 to 100 to 1000 to 10000 to 100000 of processors!!
    ➢ **multidisciplinary applications**
♦ **A supercomputer application and software are usually much more long-lived than a hardware**
  ➢ **Hardware life typically five years at most.**
  ➢ **Fortran and C are the main programming models**
♦ **Software is a major cost component of modern technologies.**
  ➢ **The tradition in HPC system procurement is to assume that the software is free.**
♦ **We have too few ideas about how to solve this problem.**
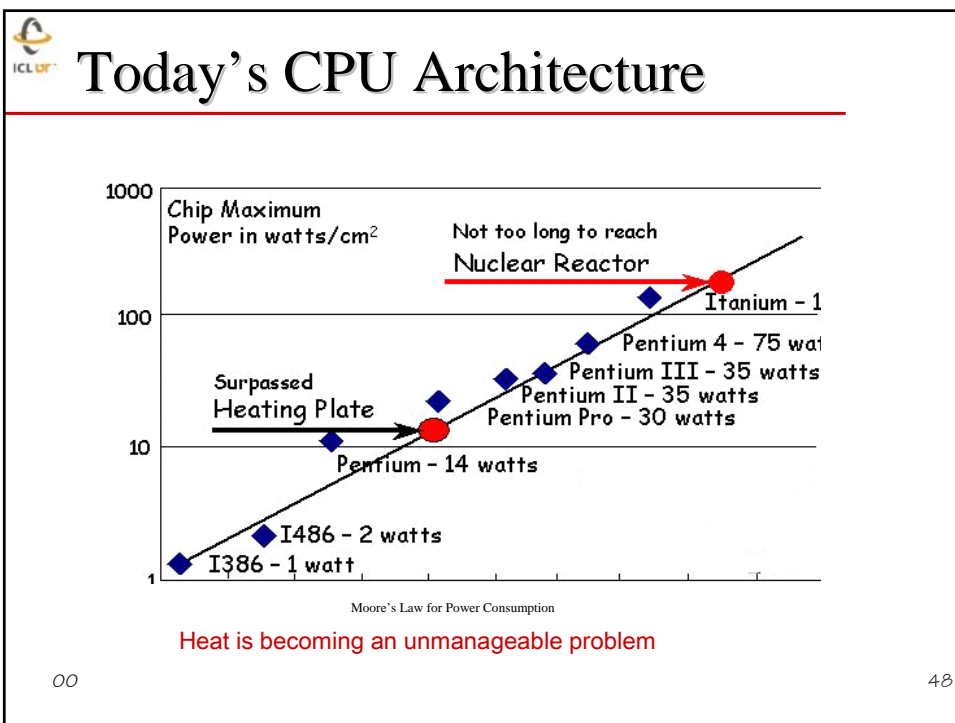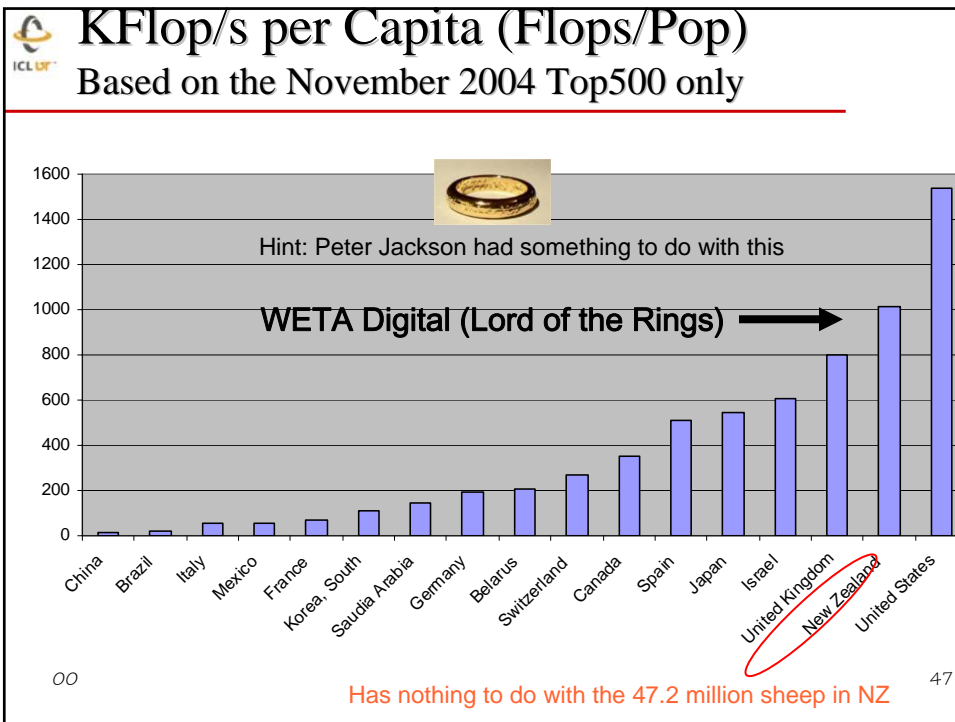
*OO*

44

# Processor Types

---

# Today's Processors

- ♦ **pipelining (superscalar, OOO, VLIW, branch prediction, predication)**
- ♦ **simultaneous multithreading (SMT, Hyper-Threading, multi-core)**
- ♦ **SIMD vector instructions (VIS, MMX/SSE, AltiVec)**
- ♦ **caches and the memory hierarchy**
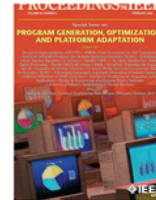- ♦ **Intel added 36 instructions per year to IA-32, or 3 instructions per month!**

## KFlop/s per Capita (Flops/Pop)
### Based on the November 2004 Top500 only

Hint: Peter Jackson had something to do with this

WETA Digital (Lord of the Rings) ⟶

Chart categories (left to right): China, Brazil, Italy, Mexico, France, Korea, South, Saudia Arabia, Germany, Belarus, Switzerland, Canada, Spain, Japan, Israel, United Kingdom, New Zealand, United States

Y-axis: 0, 200, 400, 600, 800, 1000, 1200, 1400, 1600

*OO*

Has nothing to do with the 47.2 million sheep in NZ

---

## Today's CPU Architecture

Chip Maximum
Power in watts/cm$^2$

Not too long to reach
Nuclear Reactor

Itanium – 1

Surpassed
Heating Plate

Pentium 4 – 75 wat
Pentium III – 35 watts
Pentium II – 35 watts
Pentium Pro – 30 watts

Pentium – 14 watts

I486 – 2 watts
I386 – 1 watt

Y-axis: 1, 10, 100, 1000

Moore's Law for Power Consumption

Heat is becoming an unmanageable problem

*OO*

# Self Adapting Numerical Software

♦ **The process of arriving at an efficient solution involves many decisions by an expert.**

- ➢ **Algorithm decisions**
- ➢ **Data decisions**
- ➢ **Management of the computing environment**
- ➢ **Processor specific tuning**

Proceedings of the IEEE, V: 93 #: 2 Feb. 2005 Issue on Program Generation, Optimization, and Platform Adaptation

Complex set of interaction between
Users' applications
Algorithm
Programming language
Compiler
Machine instruction
Hardware
Many layers of translation from the application to the hardware. Changing with each generation of hardware.