# 33$^{rd}$ TOP500 List

ISC'09, Hamburg

# Agenda

- Welcome and Introduction (H.W. Meuer)

- TOP10 and Awards (H.D. Simon)

- Highlights of the 30th TOP500 (E. Strohmaier)

- HPC Power Consumption (J. Shalf)

- Multicore and Manycore and their Impact on HPC (J.J. Dongarra)

- Discussion

TOP 500
SUPERCOMPUTER SITES

# 33rd List: The TOP10

| Rank | Site | Manufacturer | Computer | Country | Cores | Rmax [Tflops] | Power [MW] |
|------|------|--------------|----------|---------|-------|---------------|------------|
| 1 | DOE/NNSA/LANL | IBM | Roadrunner BladeCenter QS22/LS21 | USA | 129,600 | 1,105.0 | 2.48 |
| 2 | Oak Ridge National Laboratory | Cray Inc. | Jaguar Cray XT5 QC 2.3 GHz | USA | 150,152 | 1,059.0 | 6.95 |
| 3 | Forschungszentrum Juelich (FZJ) | IBM | Jugene Blue Gene/P Solution | Germany | 294,912 | 825.50 | 2.26 |
| 4 | NASA/Ames Research Center/NAS | SGI | Pleiades SGI Altix ICE 8200EX | USA | 51,200 | 487.0 | 2.09 |
| 5 | DOE/NNSA/LLNL | IBM | BlueGene/L eServer Blue Gene Solution | USA | 212,992 | 478.2 | 2.32 |
| 6 | University of Tennessee | Cray | Kraken Cray XT5 QC 2.3 GHz | USA | 66,000 | 463.30 | |
| 7 | Argonne National Laboratory | IBM | Intrepid Blue Gene/P Solution | USA | 163,840 | 458.61 | 1.26 |
| 8 | TACC/U. of Texas | Sun | Ranger SunBlade x6420 | USA | 62,976 | 433.2 | 2.0 |
| 9 | DOE/NNSA/LANL | IBM | Dawn Blue Gene/P Solution | USA | 147,456 | 415.70 | 1.13 |
| 10 | Forschungszentrum Juelich (FZJ) | Sun/Bull SA | JUROPA NovaScale /Sun Blade | Germany | 26,304 | 274.80 | 1.54 |

# The TOP500 Project

- Listing the 500 most powerful computers in the world

- Yardstick: Rmax of Linpack
  - Solve Ax=b, dense problem, matrix is random

- Update twice a year:
  - ISC'xy in June in Germany  • SCxy in November in the U.S.

- All information available from the TOP500 web site at: www.top500.org

TOP 500®
SUPERCOMPUTER SITES

# TOP500 Status

- 1st Top500 List in June 1993 at ISC'93 in Mannheim

  -

  -

- 30th Top500 List on November 13, 2007 at SC07 in Reno
- 31st Top500 List on June 18, 2008 at ISC'08  in Dresden
- 32nd Top500 List on November 18, 2008 at SC08 in Austin
- 33rd Top500 List on June 24, 2009 at ISC'09 in Hamburg
- 34th Top500 List on November 18, 2009 at SC09 in Portland


- Acknowledged by HPC-users, manufacturers and media

# 33<sup>rd</sup> List: Highlights

- The Roadrunner system, which broke the petaflop/s barrier one year ago, held on to its No. 1 spot. It still is one of the most energy efficient systems on the TOP500.
- The most powerful system outside the U.S. is an IBM BlueGene/P system at the German Forschungszentrum Juelich (FZJ) at No. 3.  FZJ has a second system in the TOP10, a mixed Bull and Sun system at No. 10.
- Intel dominates the high-end processor market with 79.8 percent of all systems and 87.7 percent of quad-core based systems.
- Intel Core i7 (Nehalem) makes its first appearance in the list with 33 systems
- Quad-core processors are used in 76.6 percent of the systems. Their use accelerates performance growth at all levels.
- Other notable systems are:
  - An IBM BlueGene/P system at the King Abdullah University of Science and Technology (KAUST) in Saudi Arabia at No. 14
  - The Chinese-built Dawning 5000A at the Shanghai Supercomputer Center at No 15. It is the largest system which can be operated with Windows HPC 2008.
- Hewlett-Packard kept a narrow lead in market share by total systems from IBM, but IBM still stays ahead by overall installed performance.
- Cray's XT system series is very popular for big customers 10 systems in the TOP50 (20 percent).

# 33rd List: Notable (New) Systems

- Shaheen, an IBM BlueGene/P system at the King Abdullah University of Science and Technology (KAUST) in Saudi Arabia at No. 14

- The Chinese-built Dawning 5000A at the Shanghai Supercomputer Center at No. 15. It is the largest system which can be operated with Windows HPC 2008.

- The (new) Earth Simulator at No. 22 (SX-9E based - only vector system) – it's back!

- Koi, a Cray internal system using Shanghai six-cores at #128.

- A Grape-DR system at the National Astronomical Observatory (NAO/CfCA)  at # 259
  - target size 2PF/s peak !

# Multi-Core and Many-Core

- Power consumption of chips and systems has increased tremendously, because of 'cheap' exploitation of Moore's Law.
  - Free lunch has ended
  - Stall of frequencies forces increasing concurrency levels, Multi-Cores
  - Optimal core sizes/power are smaller than current 'rich cores', which leads to Many-Cores
- Many-Cores, more (10-100x) but smaller cores:
  - Intel Polaris – 80 cores,
  - Clearspeed CSX600 – 96 cores,
  - nVidia G80 – 128 cores, or
  - CISCO Metro – 188 cores

# Performance Development

# Performance Development

# Replacement Rate

# Vendors / System Share



Legend:
- HP
- IBM
- Cray Inc.
- Dell
- SGI
- Sun
- Appro
- Fujitsu

Pie chart labels:
- HP 46%
- IBM 40%
- Cray 4%
- Dell 3%
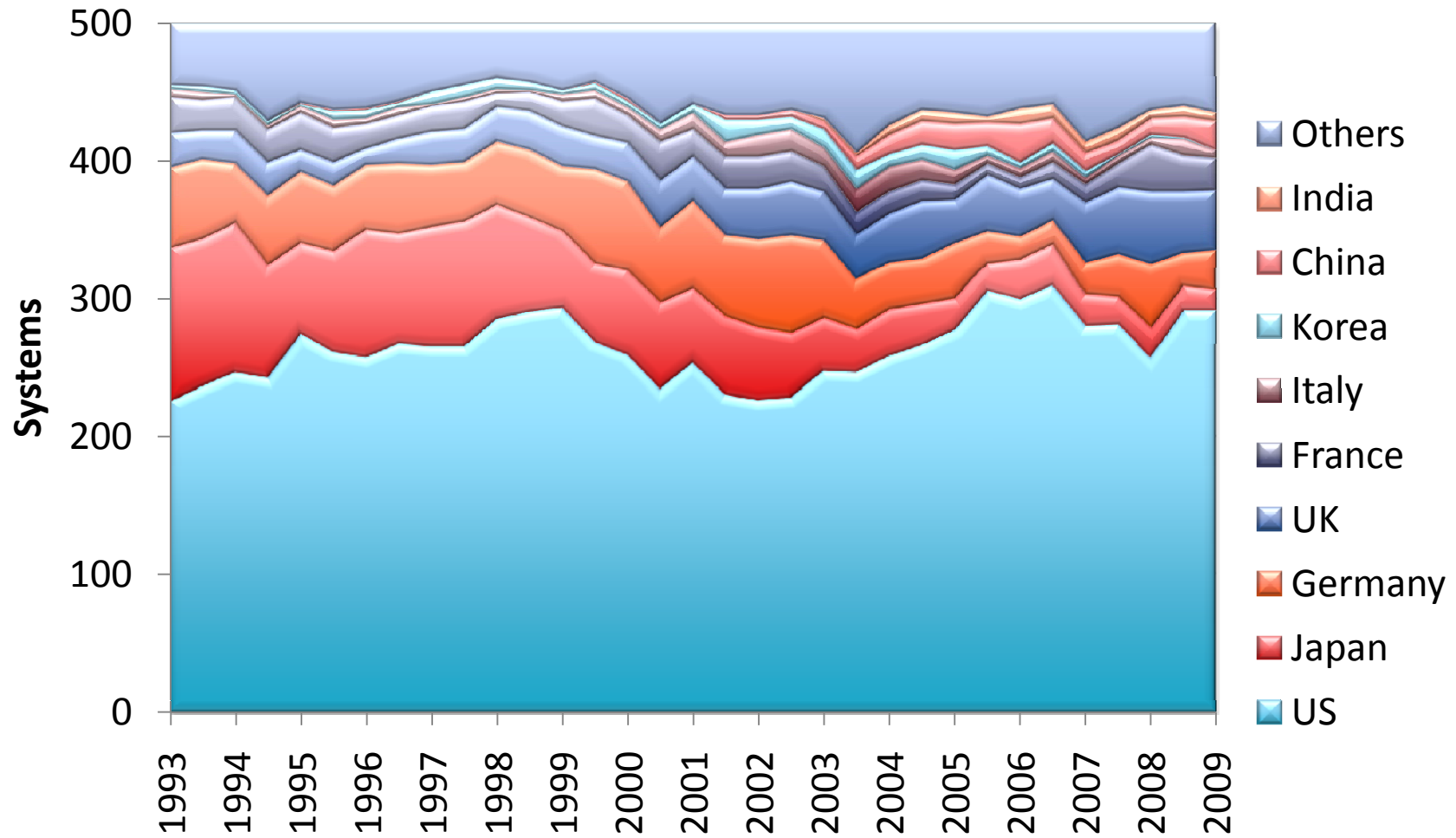- SGI 4%

# Vendors

# Vendors (TOP50)
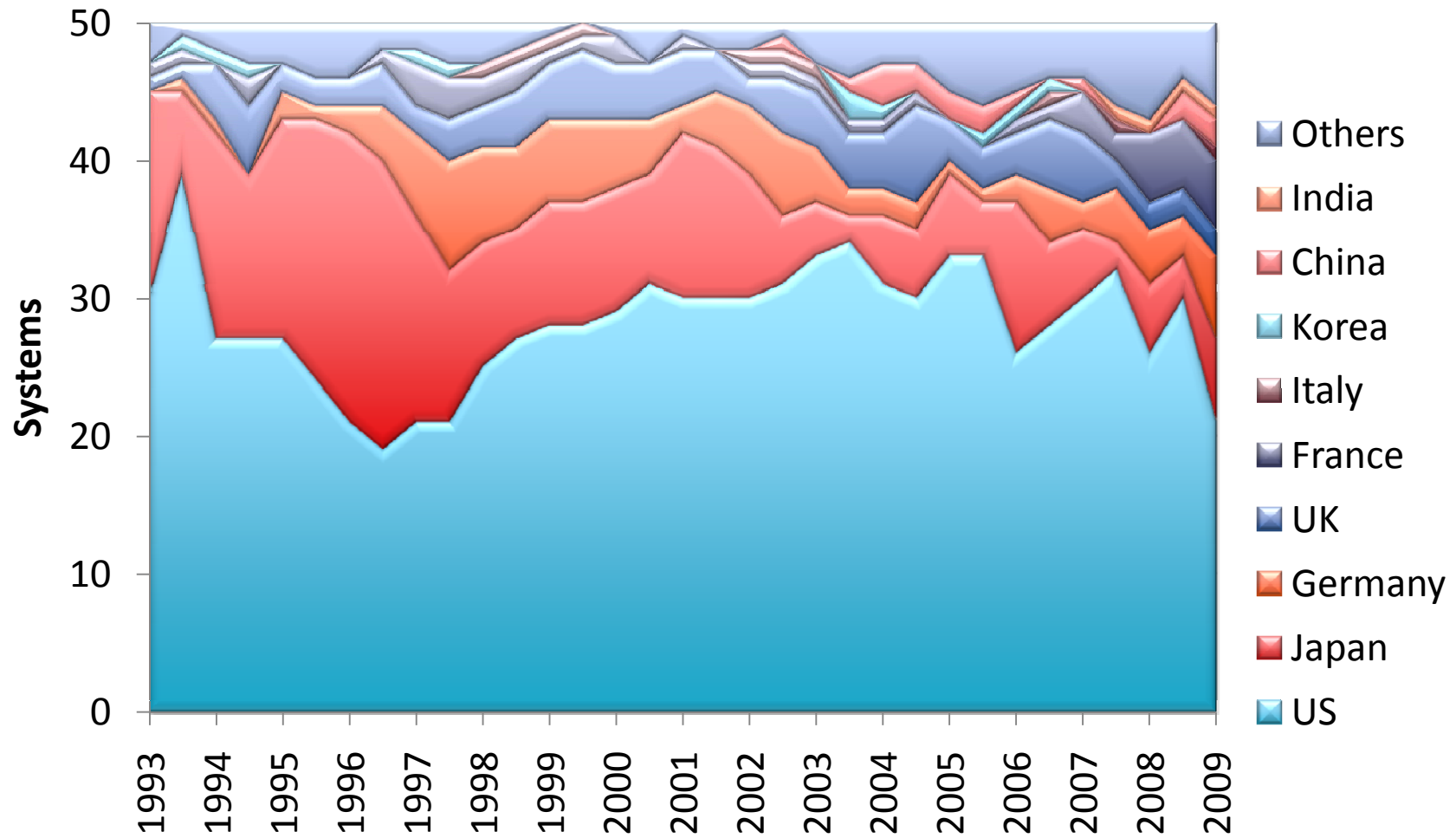
# Customer Segments

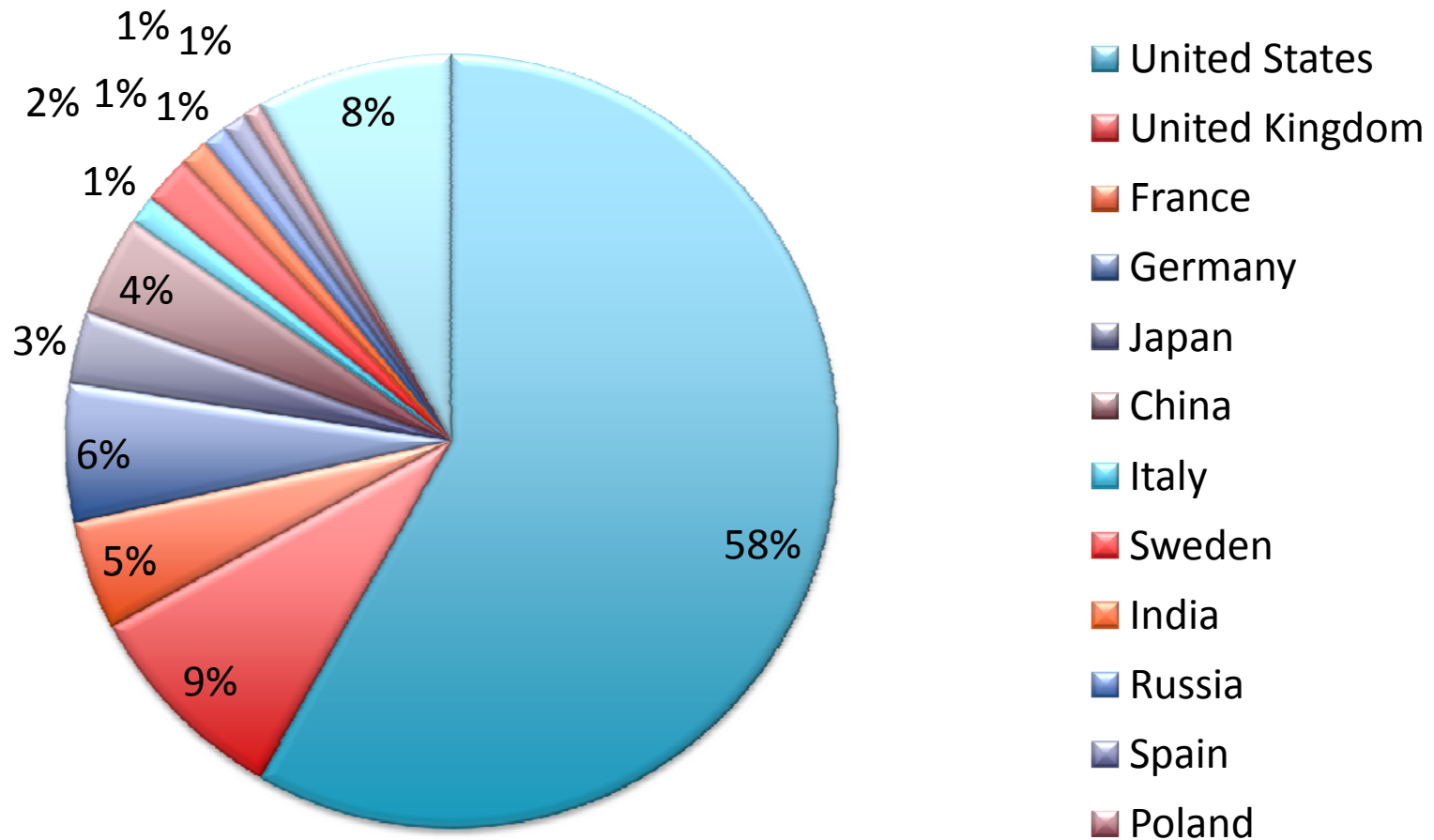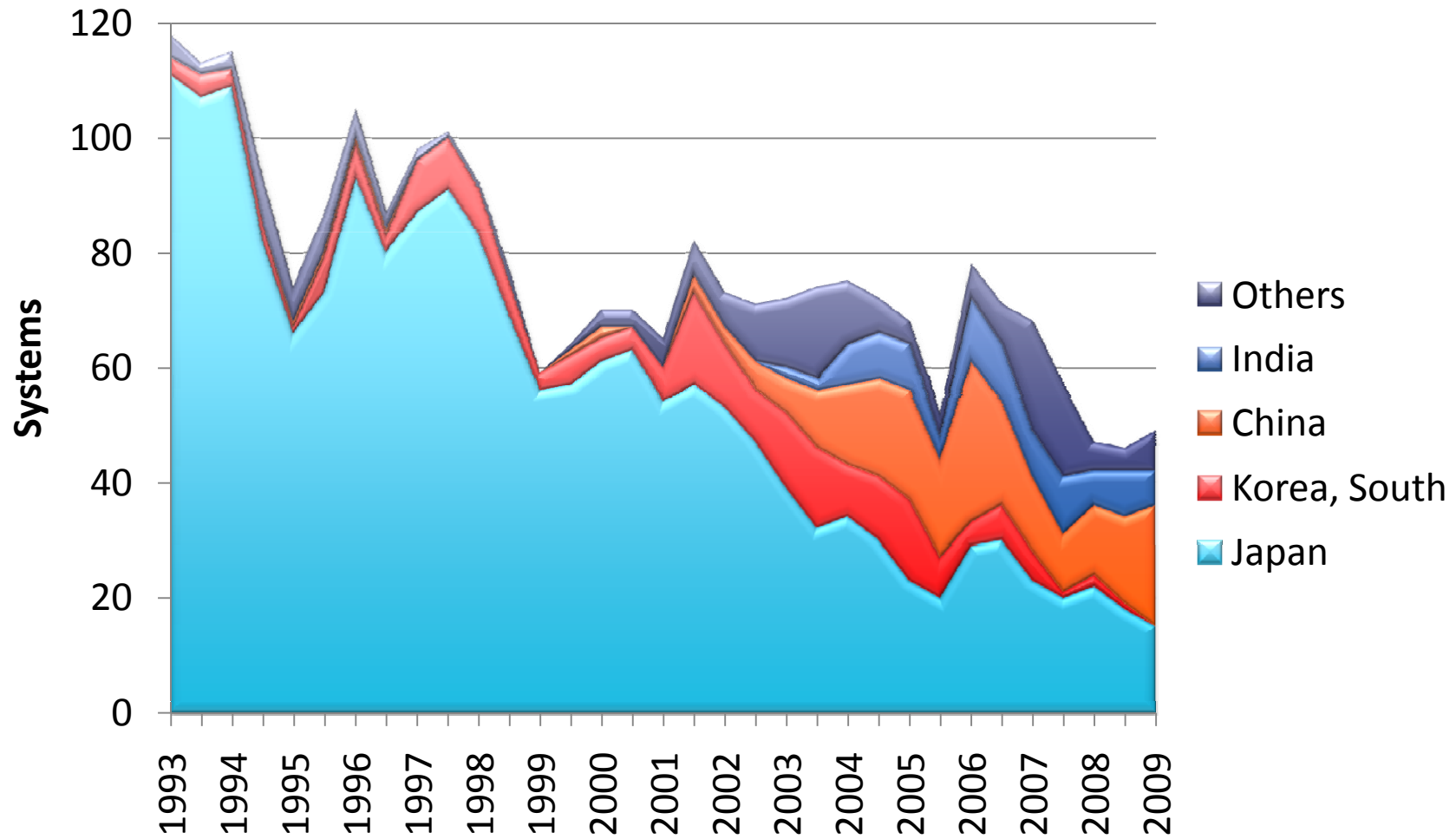# Customer Segments

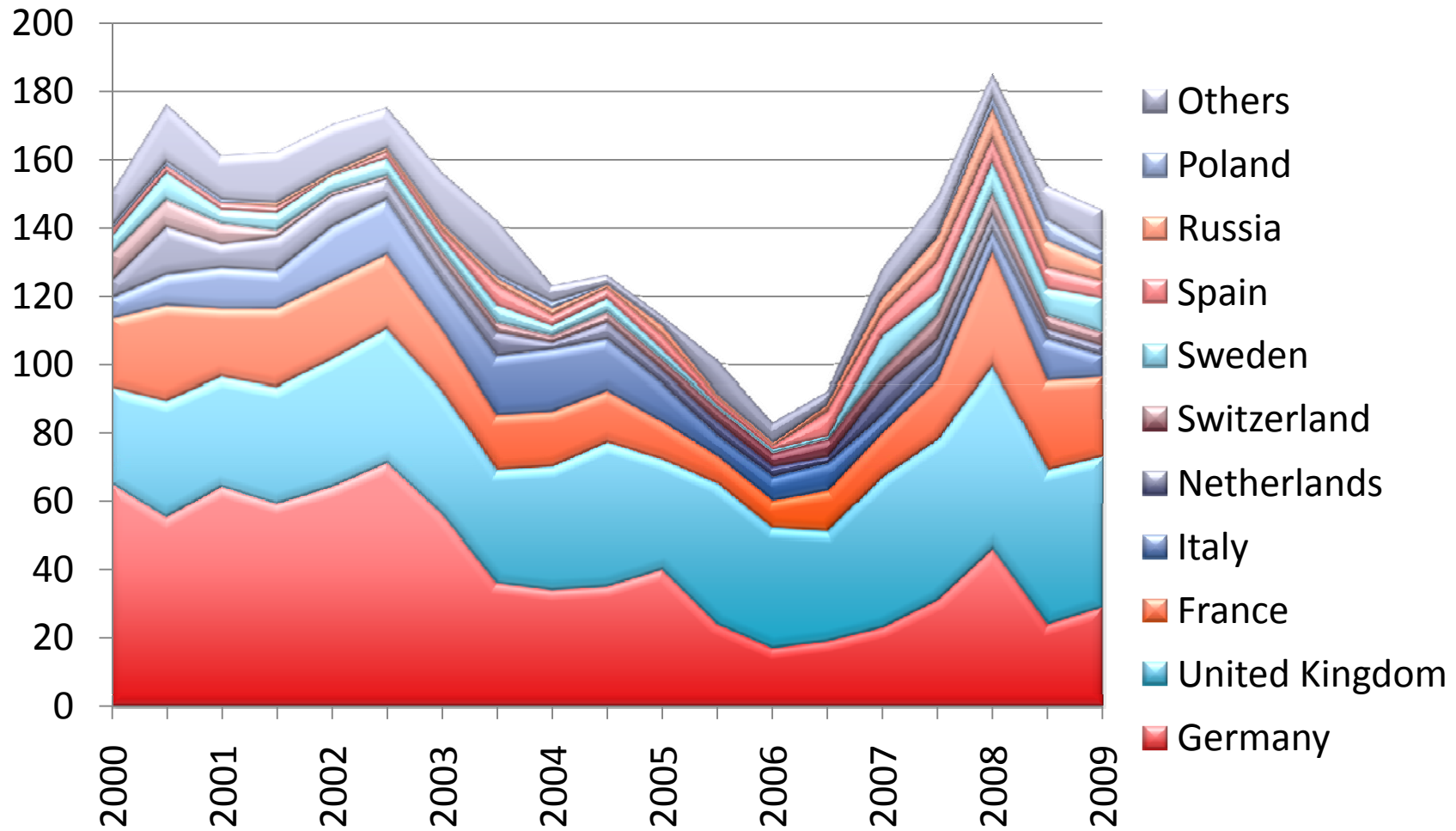# Customer Segments (TOP50)
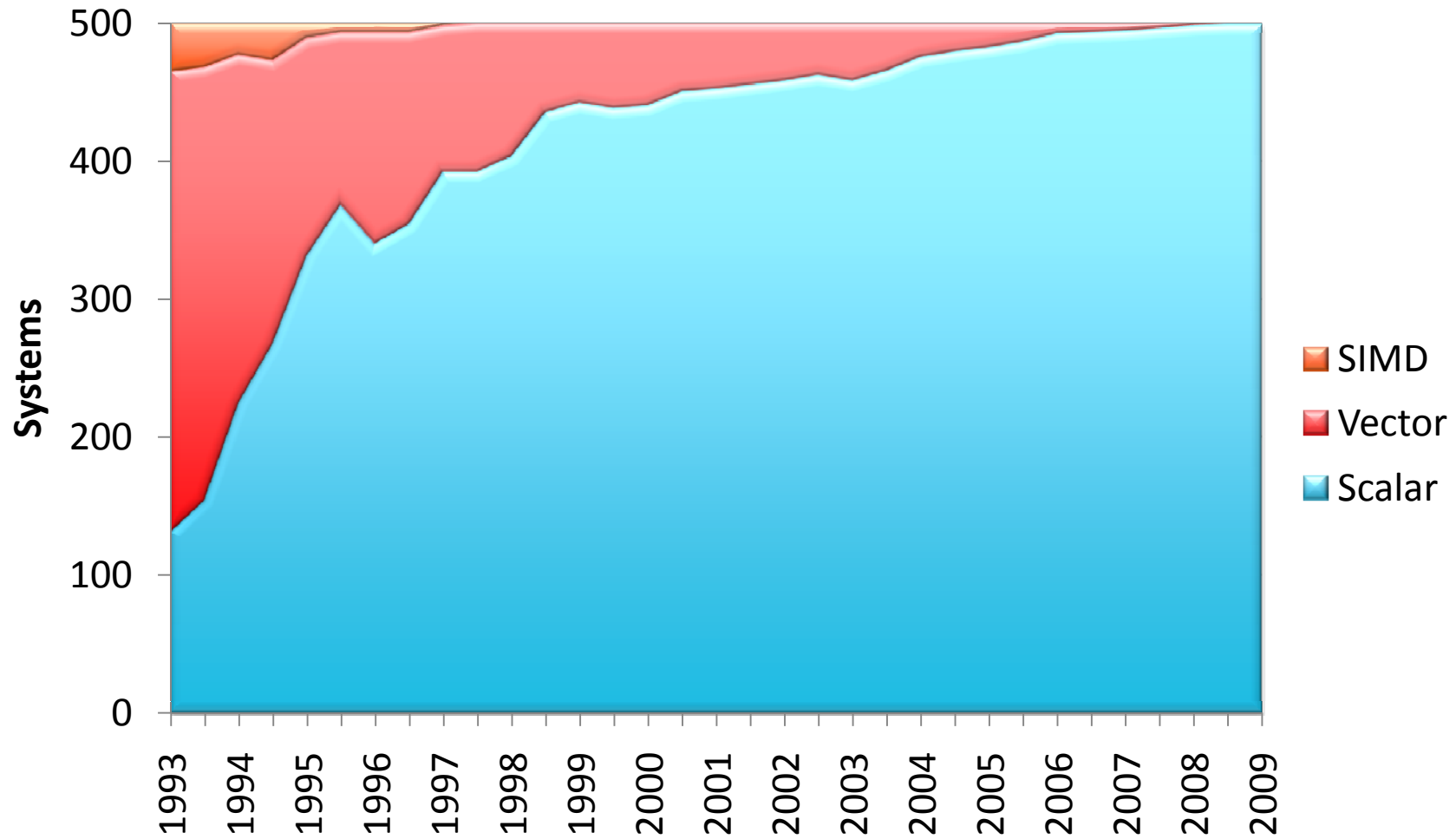
# Continents

# Countries

# Countries (TOP50)

# Countries / System Share



Legend:
- United States
- United Kingdom
- France
- Germany
- Japan
- China
- Italy
- Sweden
- India
- Russia
- Spain
- Poland

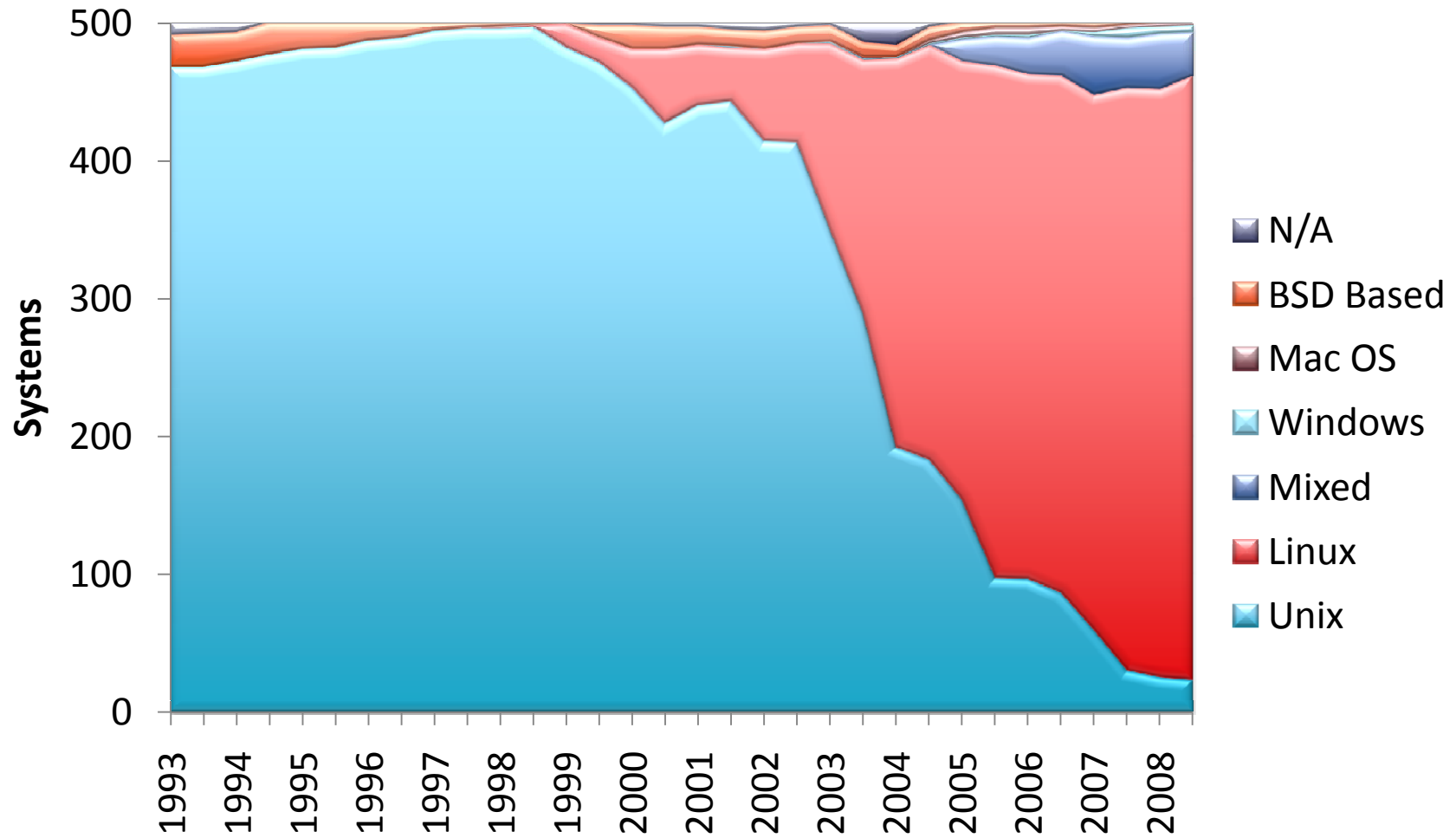Pie chart values: 58%, 9%, 5%, 6%, 3%, 4%, 1%, 1%, 1%, 2%, 1%, 1%, 8%

# Asian Countries

# European Countries

# Processor Architecture / Systems
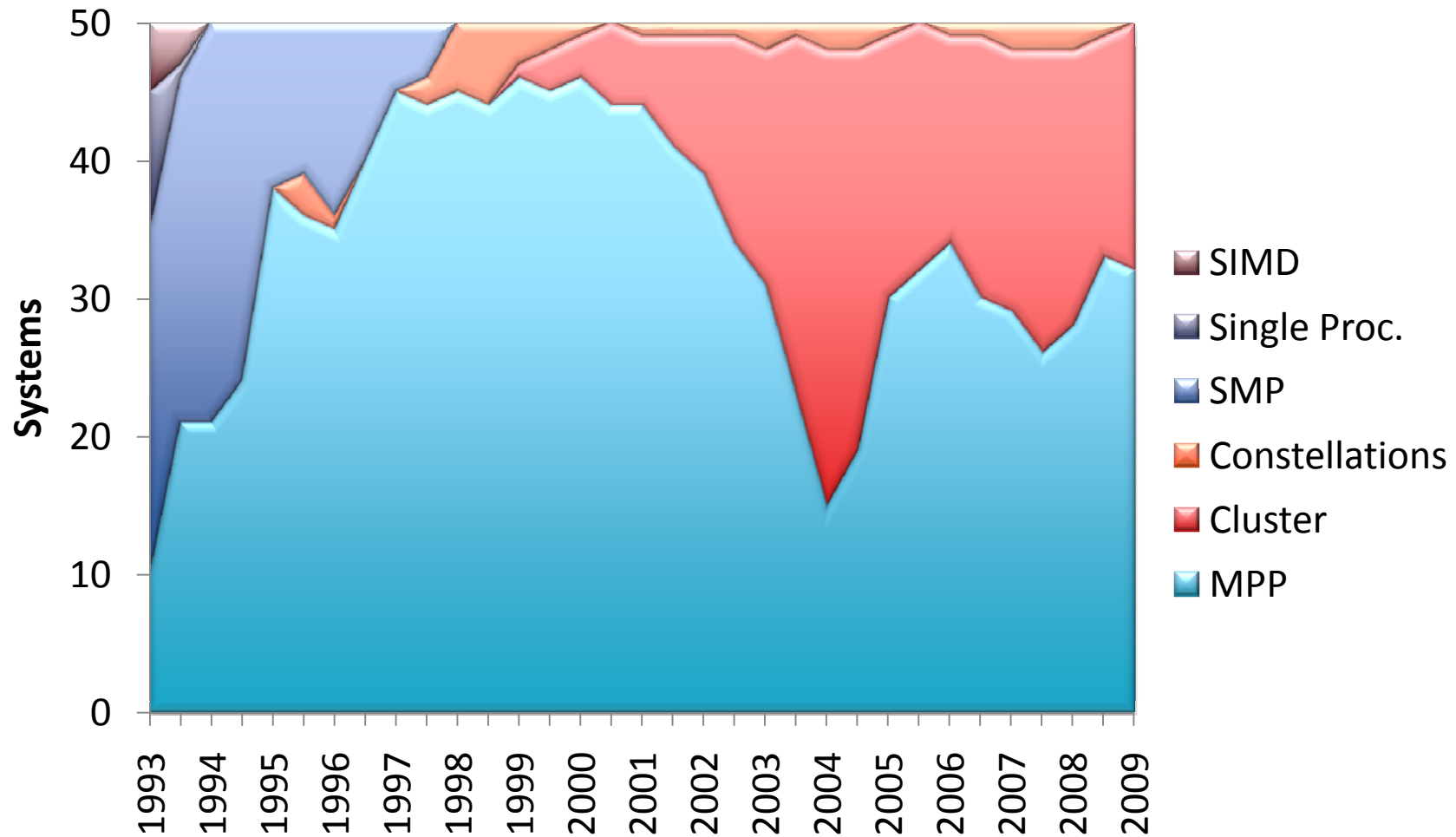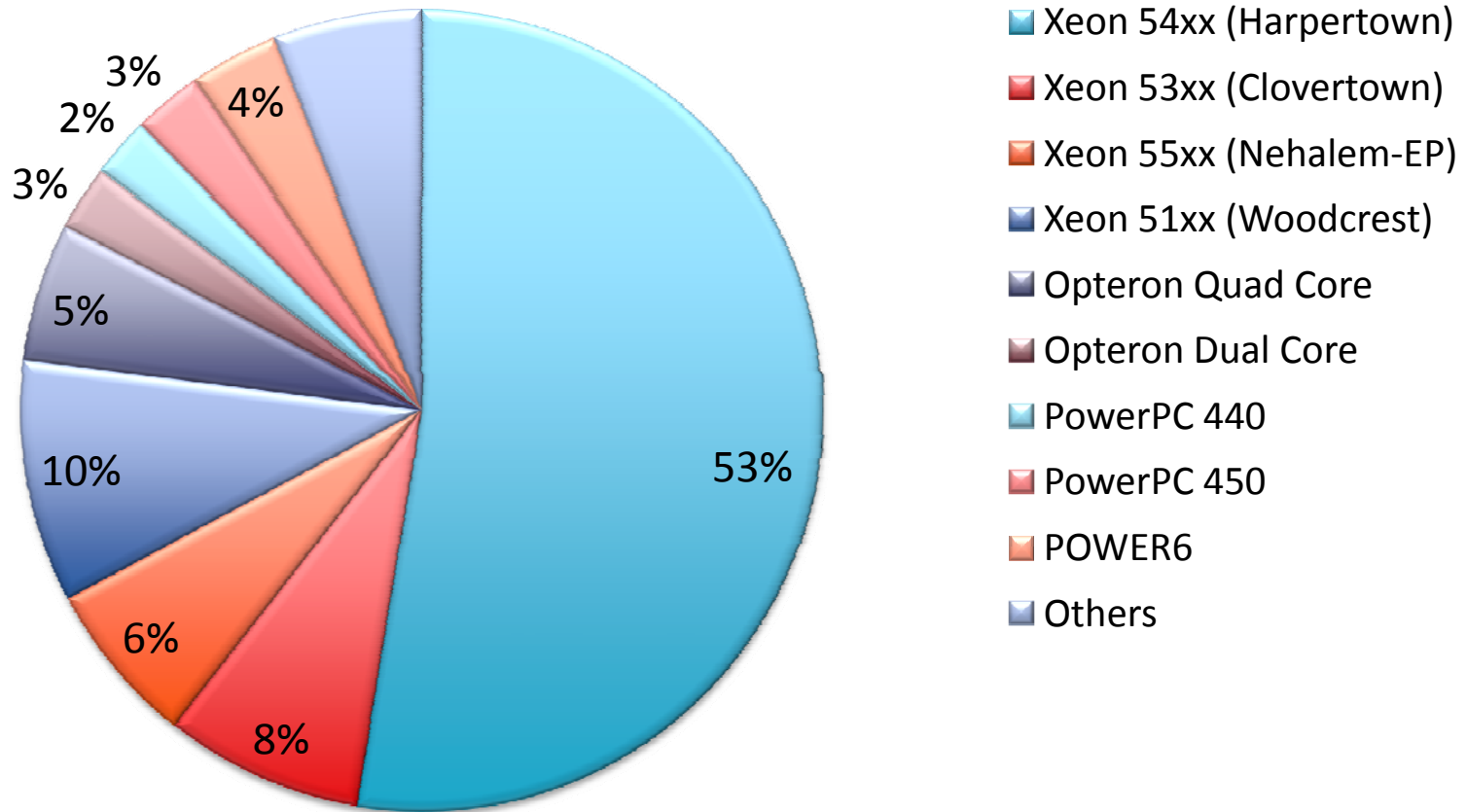
# Operating Systems



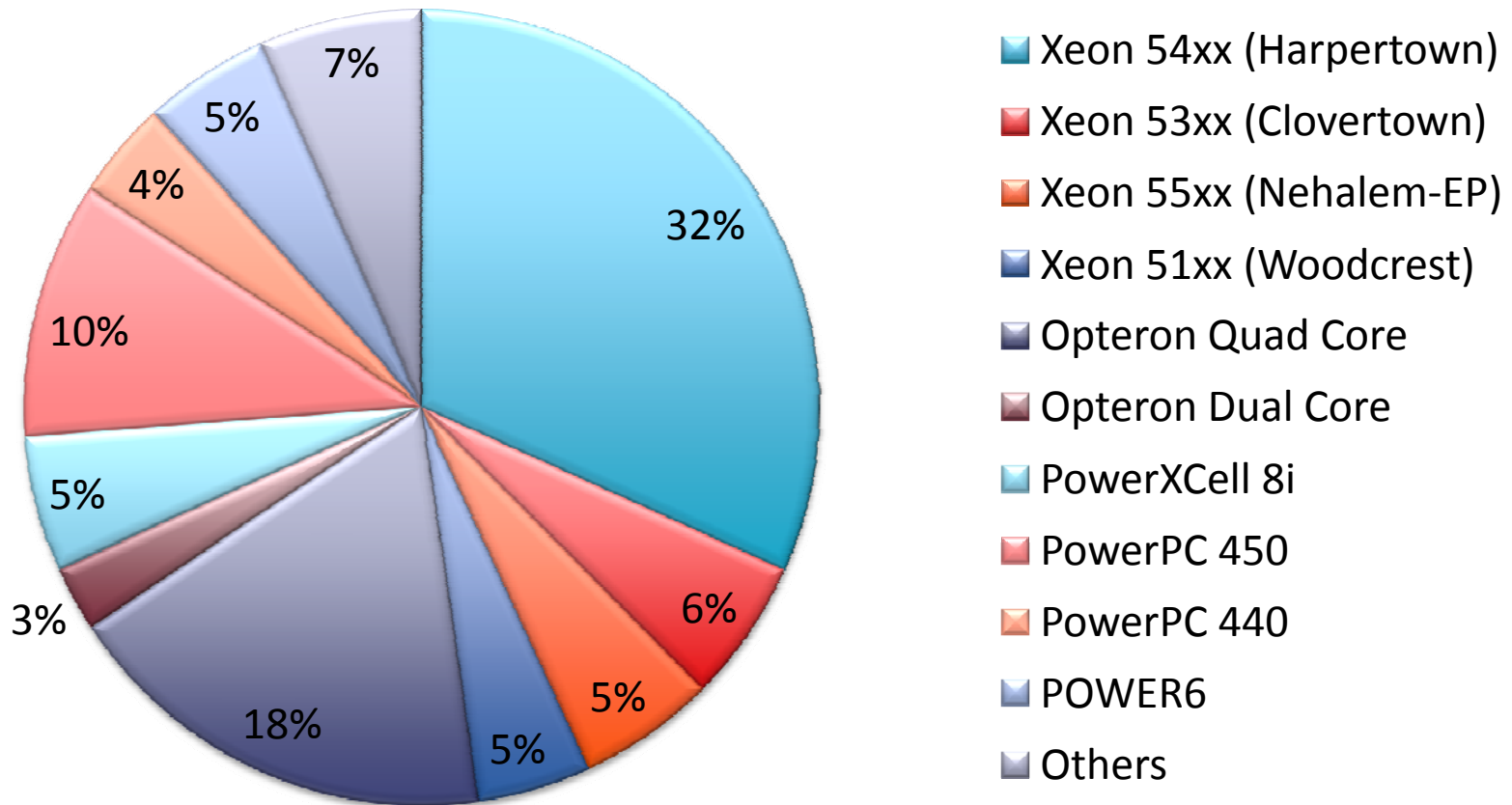Legend: N/A, BSD Based, Mac OS, Windows, Mixed, Linux, Unix
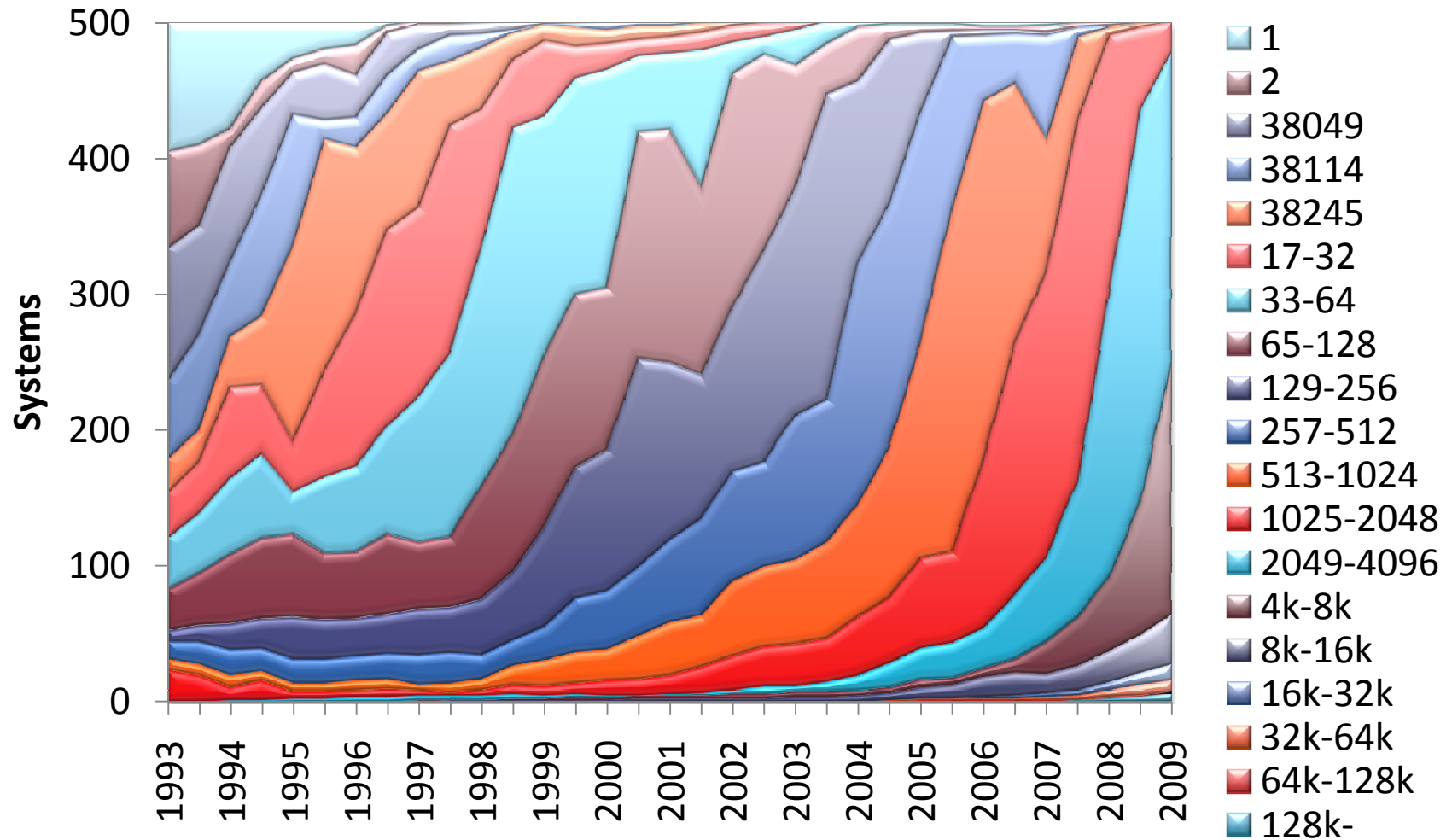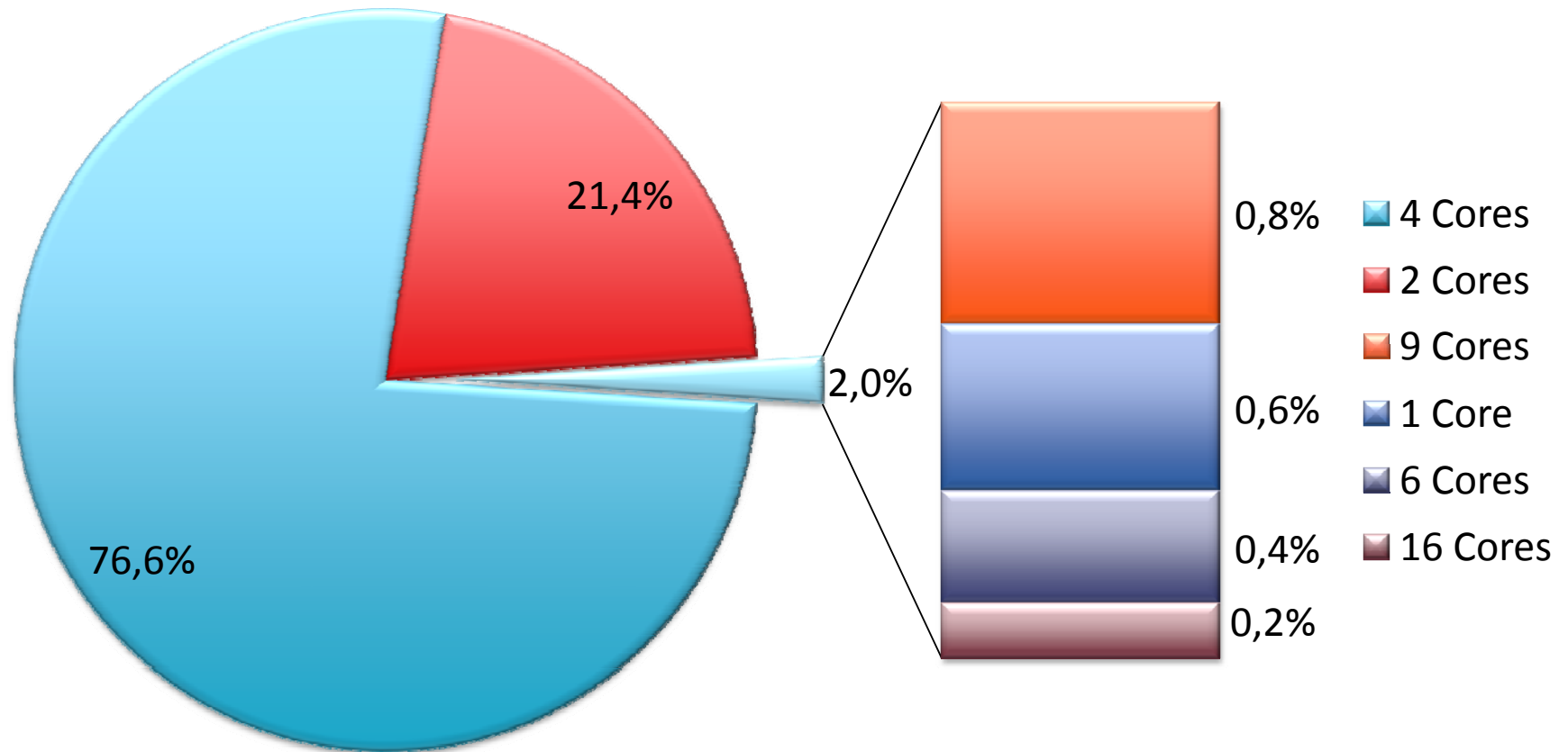
# Architectures

# Architectures (TOP50)

# Processors / Systems

# Processors / Performance
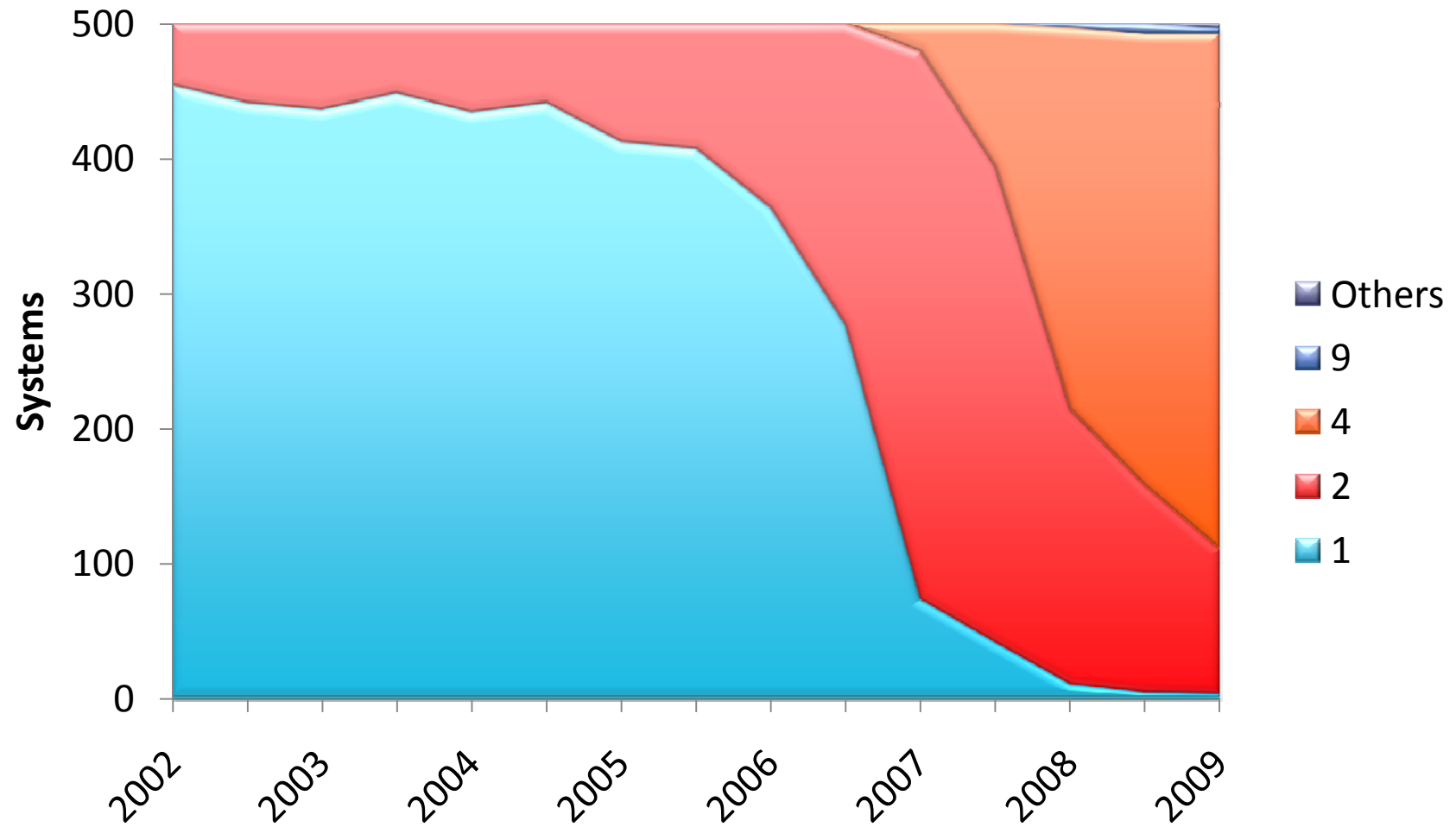


Pie chart legend:
- Xeon 54xx (Harpertown)
- Xeon 53xx (Clovertown)
- Xeon 55xx (Nehalem-EP)
- Xeon 51xx (Woodcrest)
- Opteron Quad Core
- Opteron Dual Core
- PowerXCell 8i
- PowerPC 450
- PowerPC 440
- POWER6
- Others

Chart values: 32%, 6%, 5%, 5%, 18%, 3%, 5%, 10%, 4%, 5%, 7%

# Core Count

# Cores per Socket



- 4 Cores — 76,6%
- 2 Cores — 21,4%
- 9 Cores — 0,8%
- 1 Core — 0,6%
- 6 Cores — 0,4%
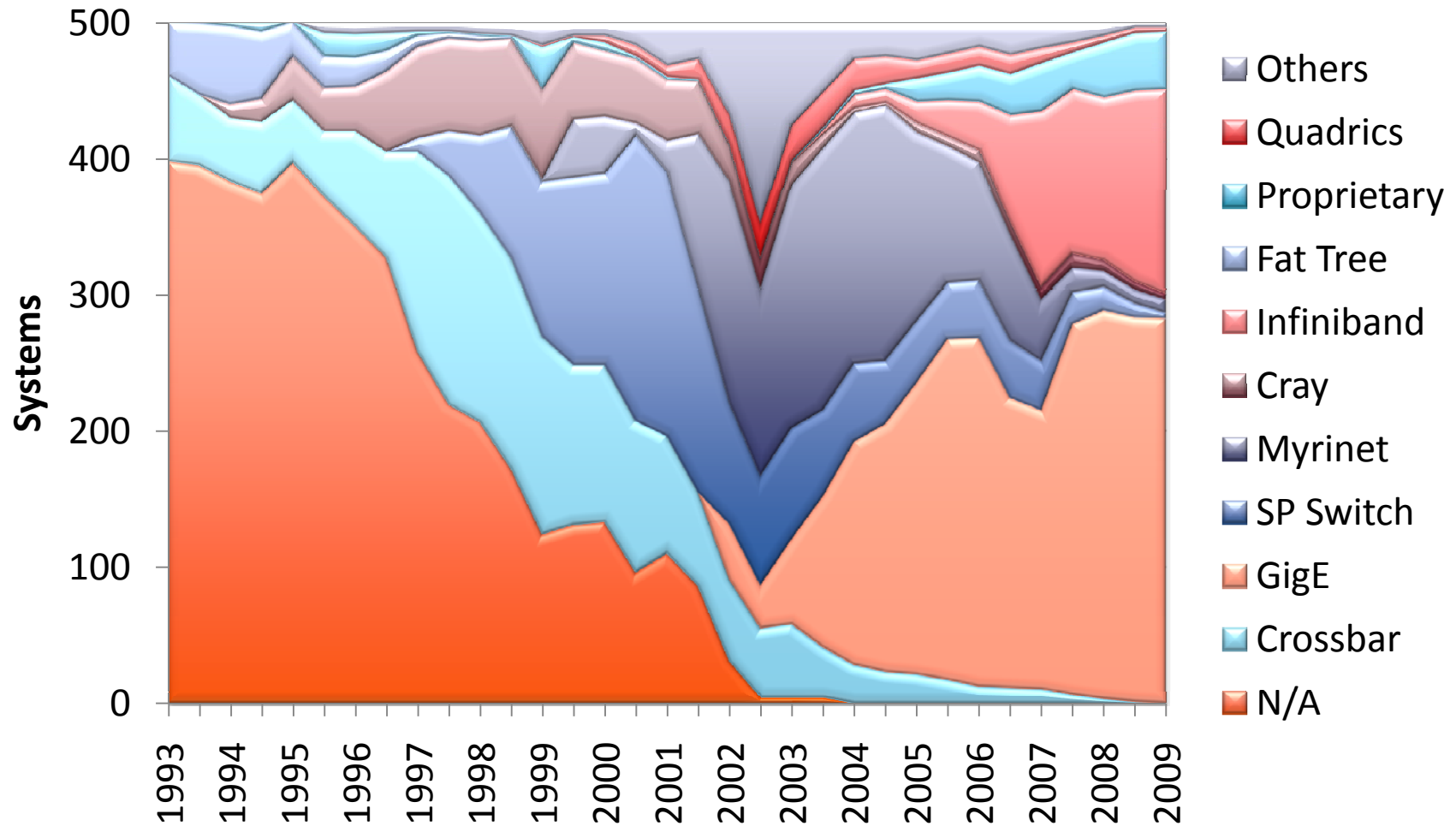- 16 Cores — 0,2%

TOP 500®
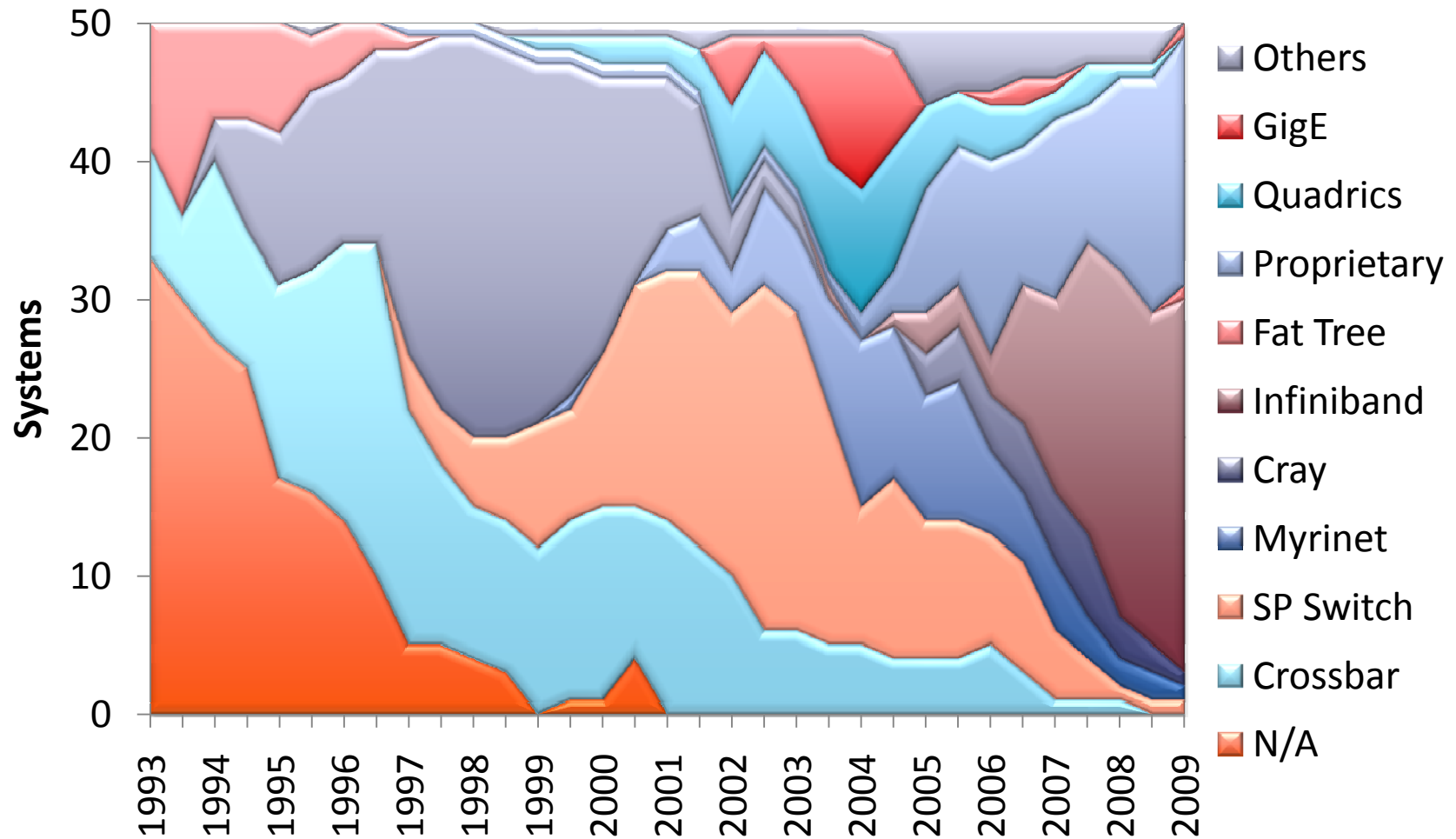SUPERCOMPUTER SITES

# Cores per Socket
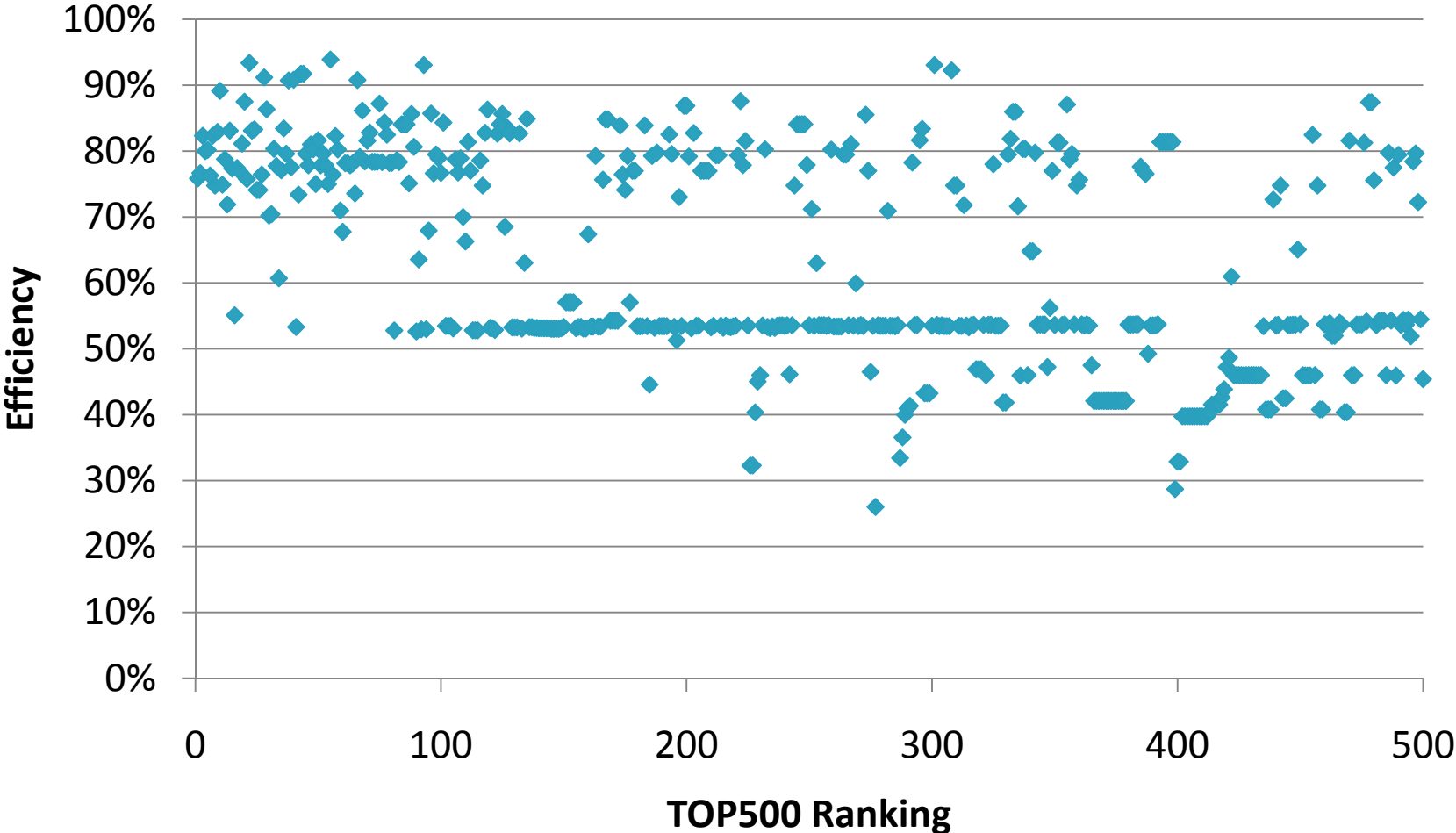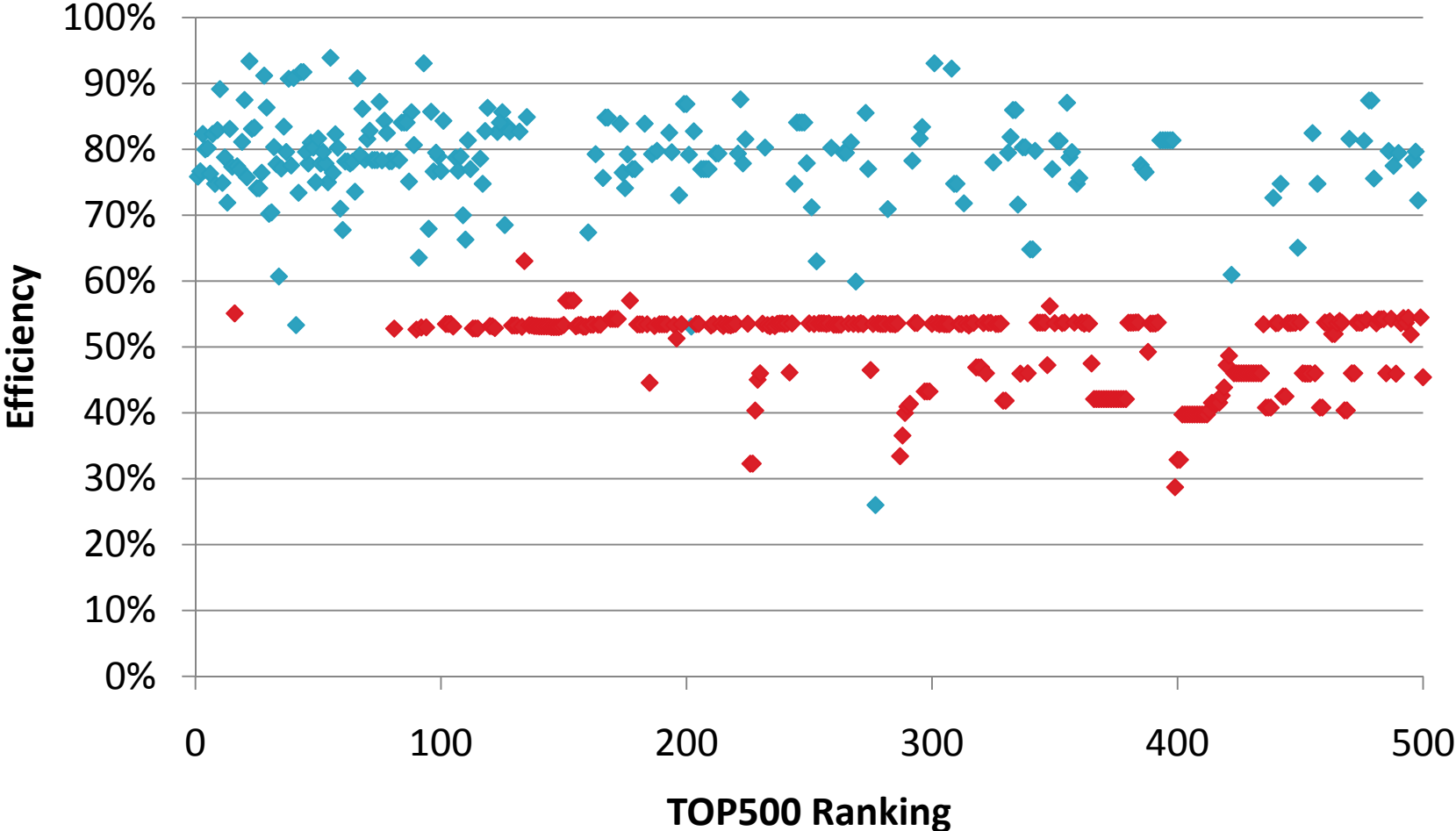
# Cluster Interconnects

# Interconnect Family

# Interconnect Family (TOP50)

# Linpack Efficiency

# Linpack Efficiency

# Bell's Law (1972)

*"**Bell's Law of Computer Class formation** was discovered about 1972. It states that technology advances in semiconductors, storage, user interface and networking advance every decade enable a new, usually lower priced computing platform to form. Once formed, each class is maintained as a quite independent industry structure. This explains mainframes, minicomputers, workstations and Personal computers, the web, emerging web services, palm and mobile devices, and ubiquitous interconnected networks. We can expect home and body area networks to follow this path."*

From Gordon Bell (2007), http://research.microsoft.com/~GBell/Pubs.htm

# HPC Computer Classes and Bell's Law

- Bell's Law states, that:

- Important classes of computer architectures come in cycles of about 10 years.

- It takes about a decade for each phase
  - Early research
  - Early adoption and maturation
  - Prime usage
  - Phase out past its prime

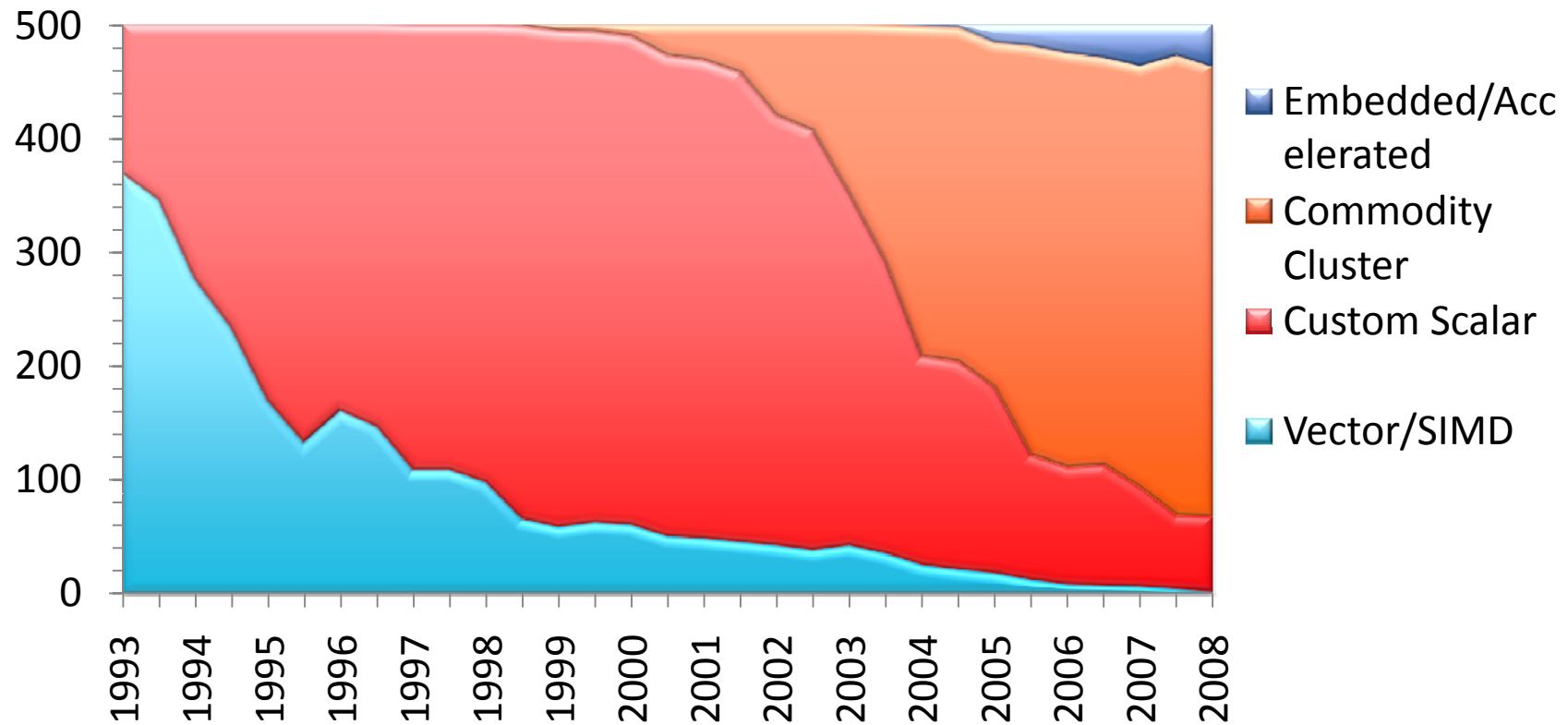- Can we use Bell's Law to classify computer architectures in the TOP500?

# HPC Computer Classes and Bell's Law

- Gordon Bell (1972): 10 year cycles for computer classes
- Computer classes in HPC based on the TOP500:
- Data Parallel Systems:
  - Vector (Cray Y-MP and X1, NEC SX, …)
  - SIMD (CM-2, …)
- Custom Scalar Systems:
  - MPP (Cray T3E and XT3, IBM SP, …)
  - Scalar SMPs and Constellations (Cluster of big SMPs)
- Commodity Cluster: NOW, PC cluster, Blades, …
- Power-Efficient Systems (BG/L as first example of low-power / embedded systems = potential new class ?)
  - Tsubame with Clearspeed, Roadrunner with Cell, ?
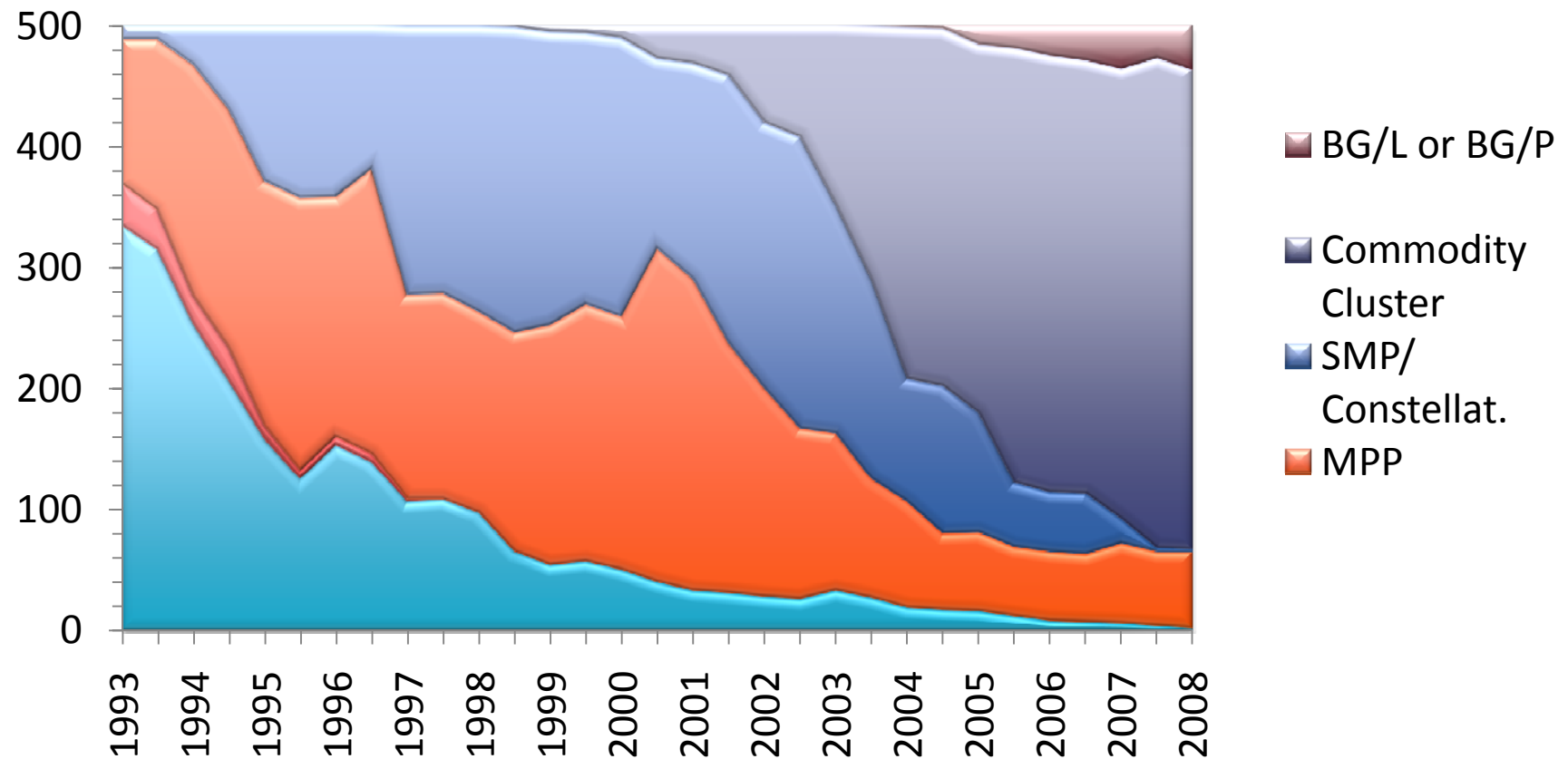  - This is a fundamental rapture and new paradigm for computing at all levels!

# Bell's Law



**Computer Classes / Systems**

Legend:
- Embedded/Accelerated
- Commodity Cluster
- Custom Scalar
- Vector/SIMD

# Bell's Law



Computer Classes - refined / Systems

# Bell's Law

## HPC Computer Classes

| Class | Early Adoption starts: | Prime Use starts: | Past Prime Usage starts: |
|---|---|---|---|
| Data Parallel Systems | Mid 70's | Mid 80's | Mid 90's |
| Custom Scalar Systems | Mid 80's | Mid 90's | Mid 2000's |
| Commodity Cluster | Mid 90's | Mid 2000's | Mid 2010's ??? |
| BG/L or BG/P | Mid 2000's | Mid 2010's ??? | Mid 2020's ??? |

# Power Consumption Data

- Have been planning this for years.
- Started in June 2008
- Independent from the Green500, but we try to learn from each other.
- Collect power consumption for:
  - Linpack as workload
  - Including all essential parts of a system (processor, memory, …)
  - Excluding features related to machine room (Most disk, UPS, …)
  - Full system or part large enough to include all shared components (fans, power supplies, …)
- Analyze these data carefully!

# Power Ranking and How Not to do it!

- To rank ob
  - Weight c
  - Rmax (T
    - A 'lar
- The ratio c
  - (weight/
  - Perform
- One can-n
  - Density
  - A piece of lead is not heavier or larger than one piece of wood.
- Linpack (sub-linear) / Power (linear)
will always sort smaller systems before larger ones!

# Power Ranking and How Not to do it!

- To rank objects by "size" one needs extensive properties:
  - Weight or Volume
  - Rmax (TOP500)
    - A 'larger' system should have a larger Rmax.
- The ratio of 2 extensive properties is an intensive one:
  - (weight/volume = density)
  - Performance / Power Consumption = Power_efficiency
- One can-not 'rank' objects with densities BY SIZE:
  - Density does not tell anything about size of an object
  - A piece of lead is not heavier or larger than one piece of wood.
- Linpack (sub-linear) / Power (linear)
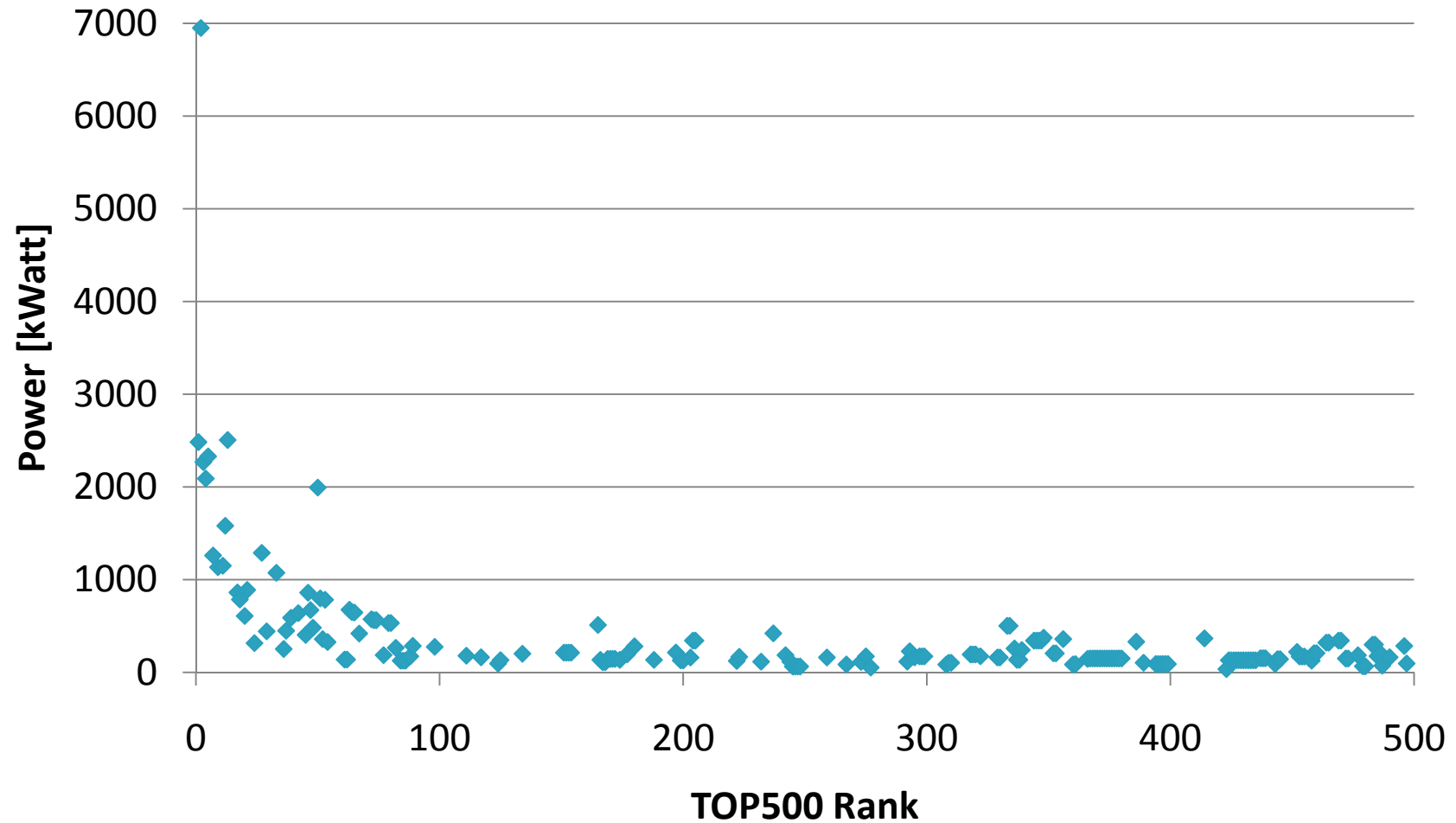  will always sort smaller systems before larger ones!
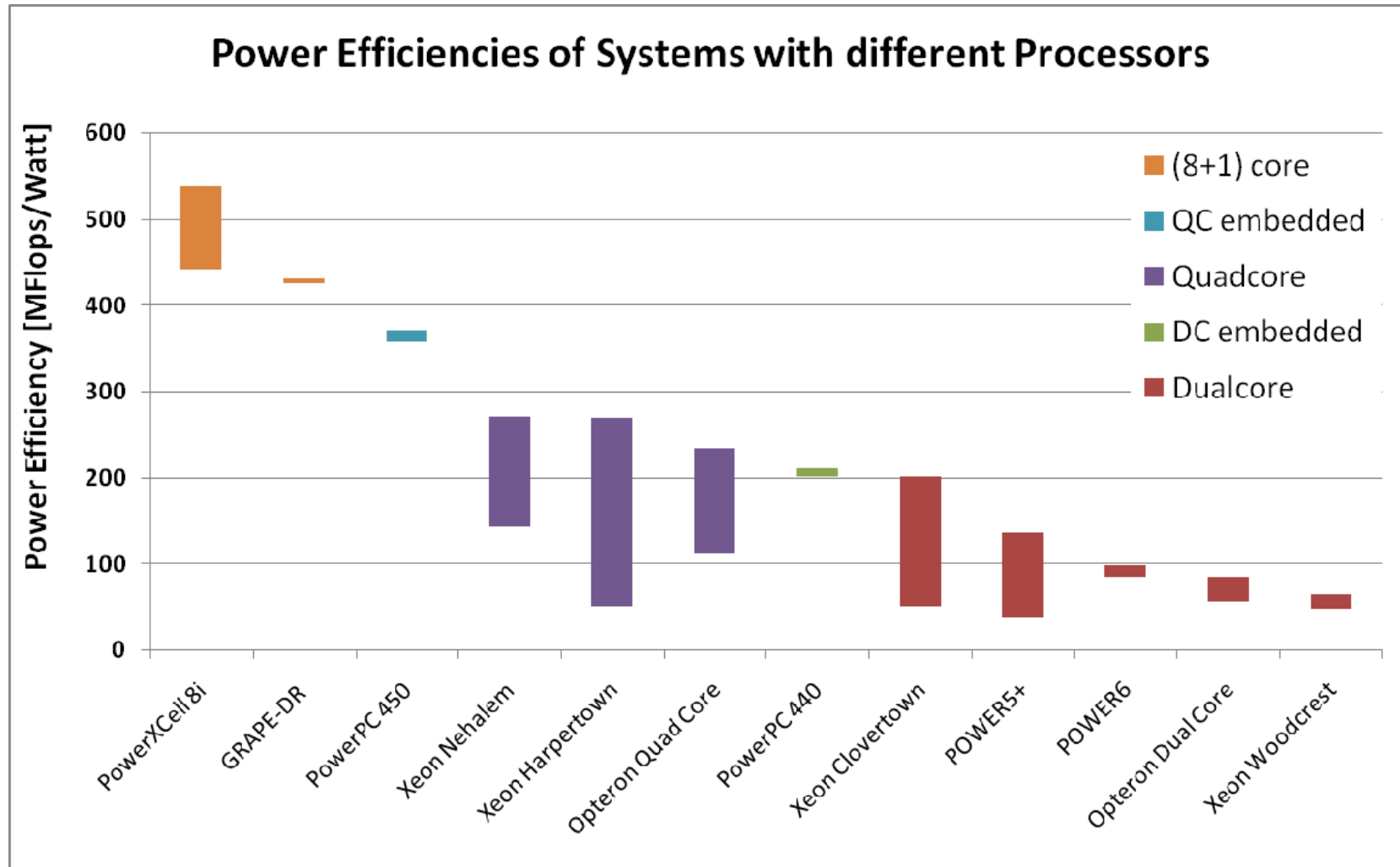
# Multi-Core and Many-Core

- Power consumption of chips and systems has increased tremendously, because of 'cheap' exploitation of Moore's Law.
  - Free lunch has ended
  - Stall of frequencies forces increasing concurrency levels, Multi-Cores
  - Optimal core sizes/power are smaller than current 'rich cores', which leads to Many-Cores
- Many-Cores, more (10-100x) but smaller cores:
  - Intel Polaris – 80 cores,
  - Clearspeed CSX600 – 96 cores,
  - nVidia G80 – 128 cores, or
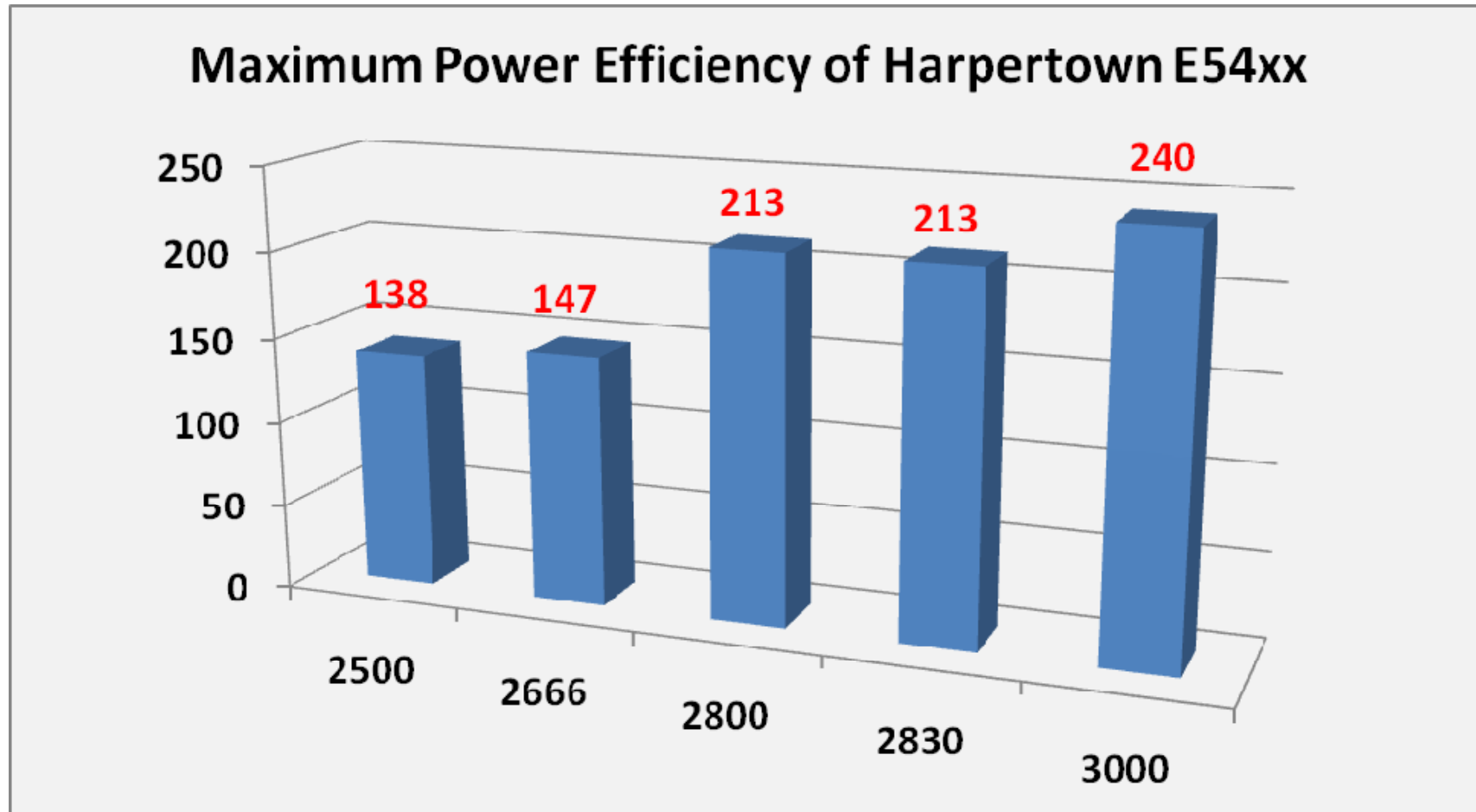  - CISCO Metro – 188 cores
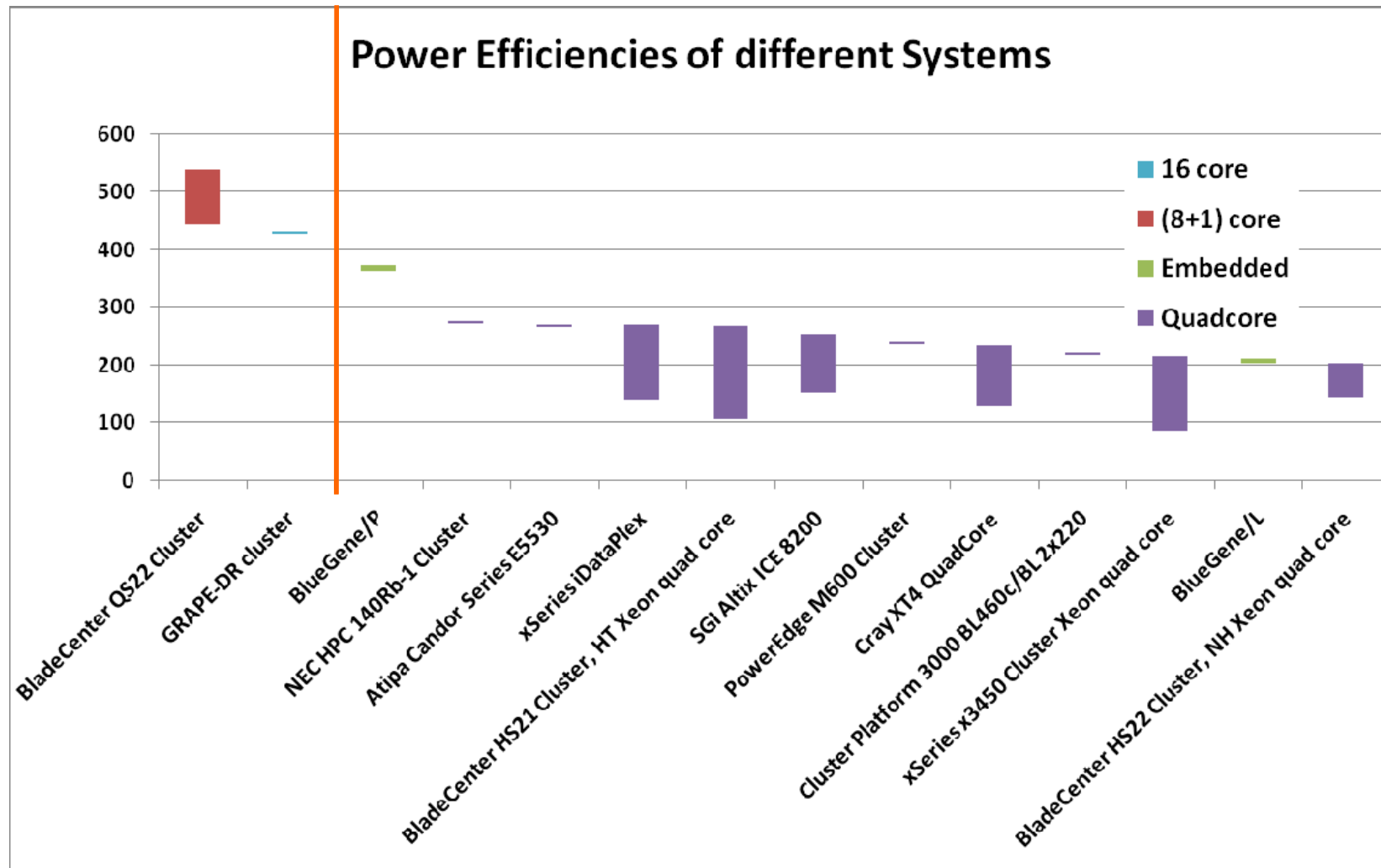
# Absolute Power Levels
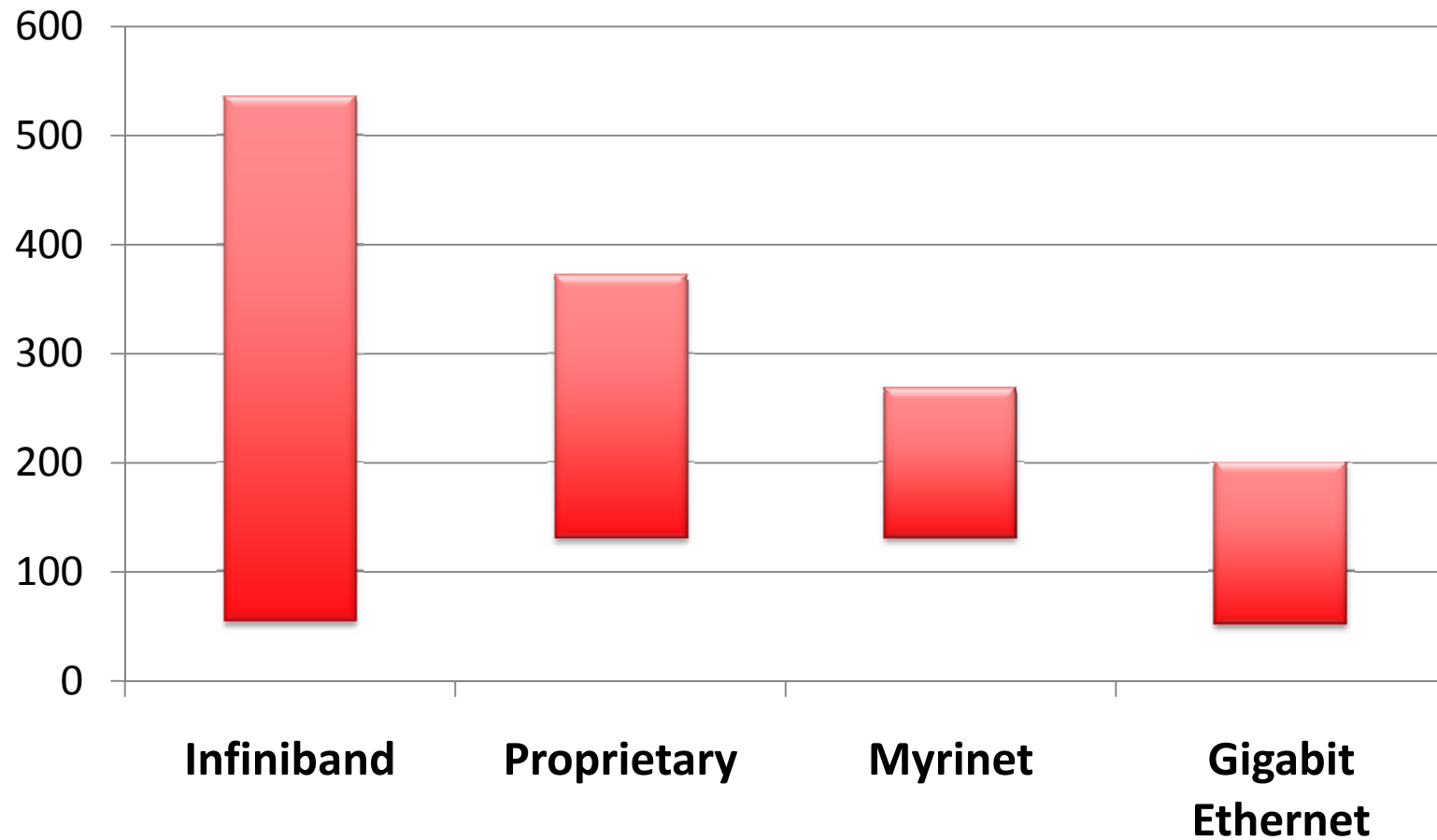
# Power Efficiency related to Processors



Power Efficiencies of Systems with different Processors

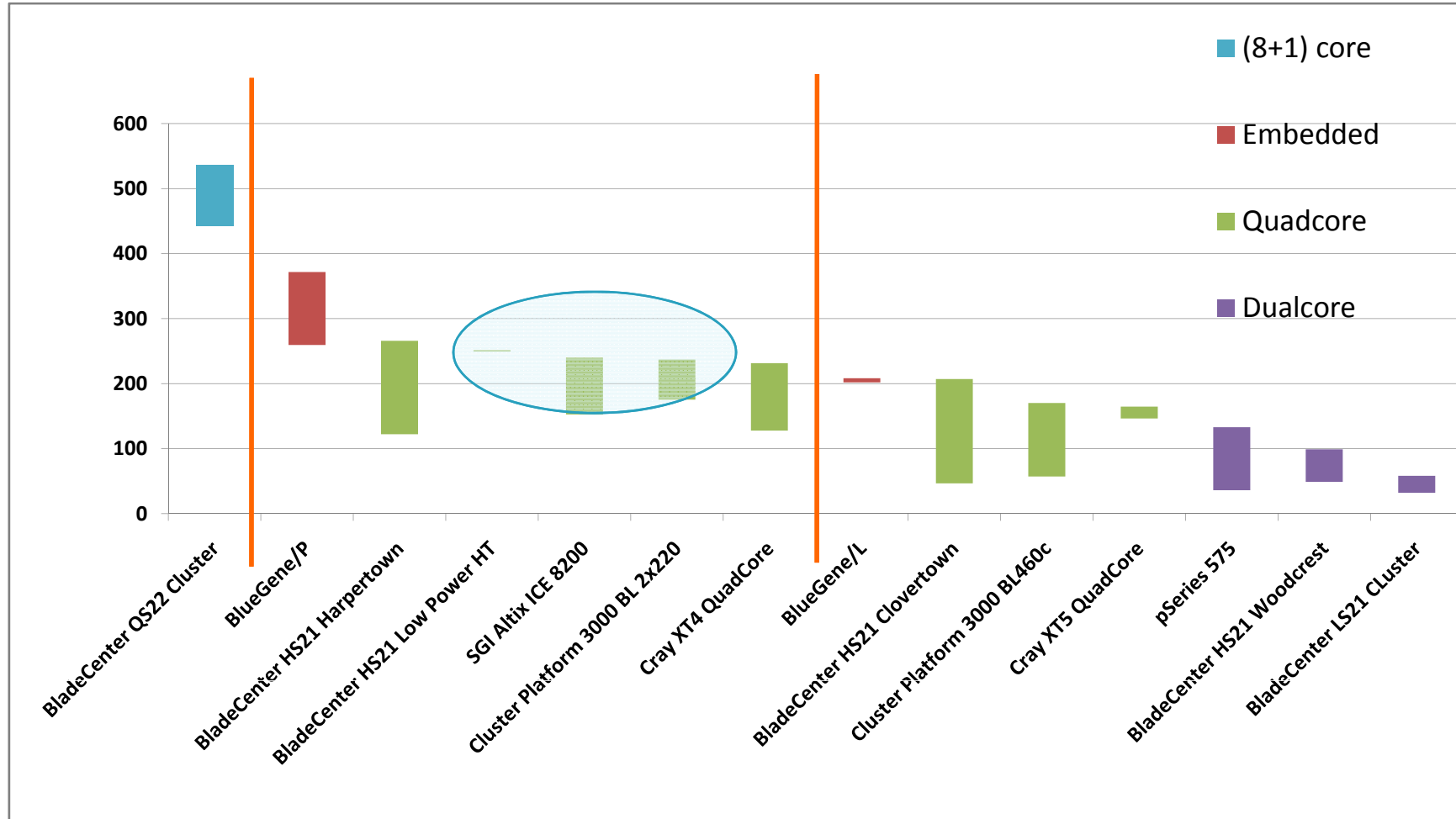# Frequencies and Power Efficiency

# Power Efficiencies of different Systems

# Power Efficiency related to Interconnects

# Power Efficiencies of different Systems

# Power Efficiencies of different Systems