

# National Energy Research Scientific Computing Center (NERSC)



## The Divergence Problem

Horst D. Simon

Director, NERSC Center Division, LBNL

November 19, 2002

# Outline

---

- ✍ Introducing NERSC-3 E
- ✍ The Divergence Problem
- ✍ What NERSC is doing about it



# Combined NERSC-3 Characteristics

---

- ✍ The combined NERSC-3/4 system (NERSC-3Base and NERSC-3Enhanced) will have
  - 416 16 way Power 3+ nodes with each CPU at 1.5 Gflop/s
    - ✍ 380 for computation
  - 6,656 CPUs – 6,080 for computation
  - Total Peak Performance of 10 Teraflop/s
  - Total Aggregate Memory is 7.8 TB
  - Total GPFS disk will be 44 TB
    - ✍ Local system disk is an additional 15 TB
  - Combined SSP-2 measure is 1.238 Tflop/s
  - NERSC-3E be in production by the end of Q1/CY03
    - ✍ Nodes will arrive in the first two weeks of November



# Comparison with Other Systems

---

	<b>NERSC-3 E</b>	<b>ASCI White</b>	<b>ES</b>	<b>Cheetah (ORNL)</b>	<b>PNNL Mid 2003</b>
<b>Nodes</b>	<b>416</b>	<b>512</b>	<b>640</b>	<b>27</b>	<b>700</b>
<b>CPUs</b>	<b>6,656</b>	<b>8,192</b>	<b>5,120</b>	<b>864</b>	<b>1400</b>
<b>Peak(Tflops)</b>	<b>10</b>	<b>12</b>	<b>40</b>	<b>4.5</b>	<b>9.6(8.3)</b>
<b>Memory (TB)</b>	<b>7.8</b>	<b>4</b>	<b>10</b>	<b>1</b>	<b>1.8</b>
<b>Disk(TB)</b>	<b>60</b>	<b>150</b>	<b>700</b>	<b>9</b>	<b>53</b>
<b>SSP(Gflop/s)</b>	<b>1,238</b>	<b>1,652</b>		<b>179</b>	

PNNL system available in Q3 CY2003

SSP = sustained systems performance (NERSC applications benchmark)



# Outline

---

- ✍ Introducing NERSC-3 E
- ✍ **The Divergence Problem**
- ✍ What NERSC is doing about it



# Signposts of Change in HPC

---

In early 2002 there were several signposts, which signal a fundamental change in HPC in the US:

- ✍ Installation and very impressive early performance results of the Earth Simulator System (April 2002)
- ✍ Lack of progress in computer architecture research evident at Petaflops Workshop (WIMPS, Feb. 2002)
- ✍ Poor or non-existing benchmarks on SSP for the NERSC workload (March 2002)

This is happening against the backdrop of:

- ✍ increasing lack of interest in HPC by some US vendors (Sun and SGI),
- ✍ further consolidation and reduction of the number of vendors (Compaq + HP merger)
- ✍ reduced profitability and reduced technology investments (dot com bust)

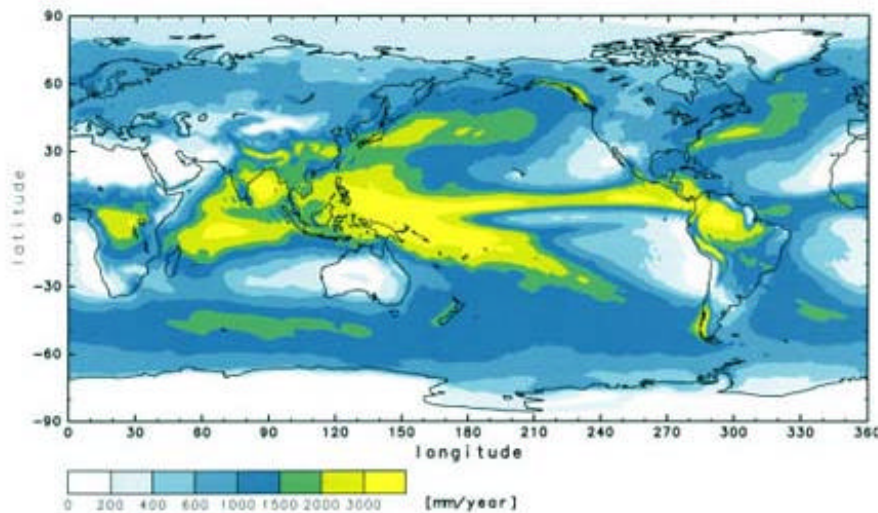
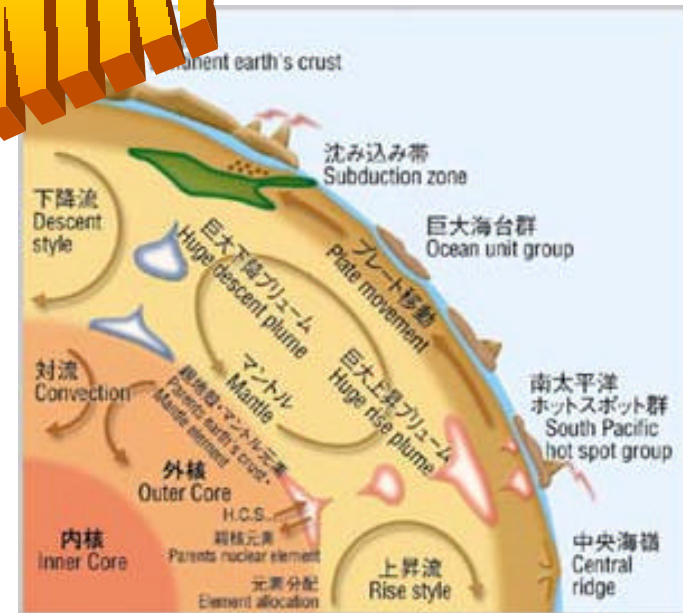
**We are in the middle of a fundamental change of the basic premises of the HPC market in the U.S.**



# The Earth Simulator in Japan

**COMPUTENIK!**

- Linpack benchmark  
TF/s = 87% of 4000
- Completed Apr 2002
- Driven by climate and  
earthquake simulation
- Built by NEC



<u>Understanding and Prediction of Global Climate Change</u>	<u>Understanding of Plate Tectonics</u>
Occurrence prediction of meteorological disaster	Understanding of long-range crustal movements
Occurrence prediction of El Niño	Understanding of mechanism of seismicity
Understanding of effect of global warming	Understanding of migration of underground water and materials transfer in strata
Establishment of simulation technology with 1km resolution	

<http://www.es.jamstec.go.jp/esrdc/eng/menu.html>

# Catalyst for fundamental change in U.S. science policy or call for a small course correction?

---

- ✍ The important event is not a single machine but the commitment of the Japanese government to invest in science-driven computing.
- ✍ U.S. computer industry is driven by commercial applications -- not focused on scientific computing.
- ✍ The Earth Simulator is a direct investment in scientific computing, giving Japanese scientific communities a material advantage and making them more attractive as international collaborators.
- ✍ **The Earth Simulator is not a special purpose machine:** All U.S. scientific computing communities are potentially now at a handicap of 10 to 100 in delivered computing capability.





# Perspective

---

- ✍ Peak performance does not reveal the real impact of the Earth Simulator.
- ✍ Japanese scientific policy is to build strategic partnerships in climate, nanoscience and fusion, moving to dominate simulation and modeling in many disciplines – not just climate modeling.
- ✍ To optimize architectures for scientific computing, it is necessary to establish the feedback between scientific applications and computer design over multiple generations of machines.
- ✍ The Japanese Earth Simulator project implemented one cycle of that feedback, and made dramatic progress.



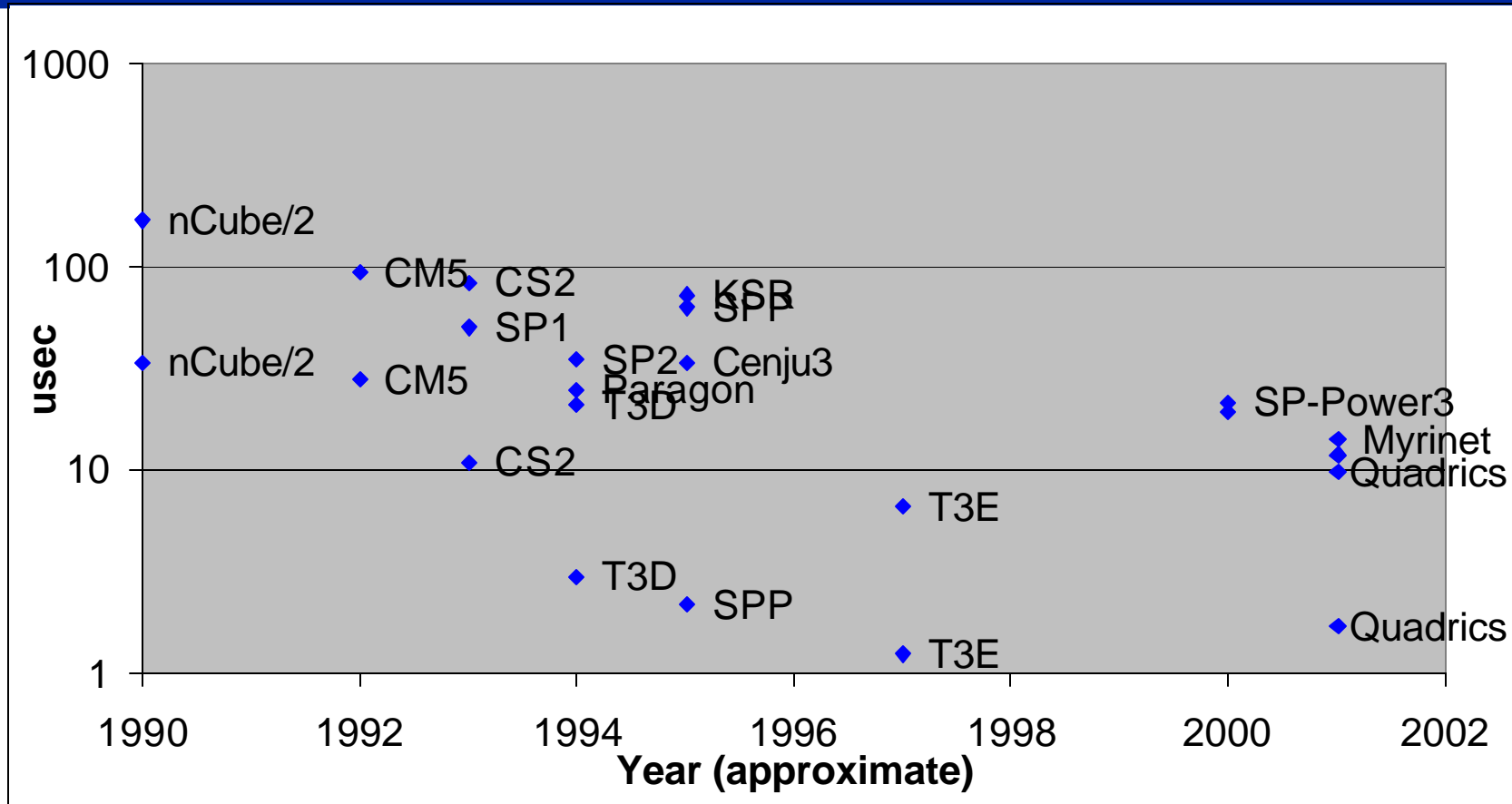
# Basic Research Issues/Observations

---

- ✍ Only a handful of supercomputing relevant computer architecture projects at US universities; versus of the order of 50 in 1992
- ✍ Lack of interest in supporting supercomputing relevant basic research
  - parallel language and tools research has been almost abandoned
  - focus on grid middleware and tools
- ✍ WIMPS2002 = Petaflops 1997
  - no significant progress in five years



# End to End Latency Over Time

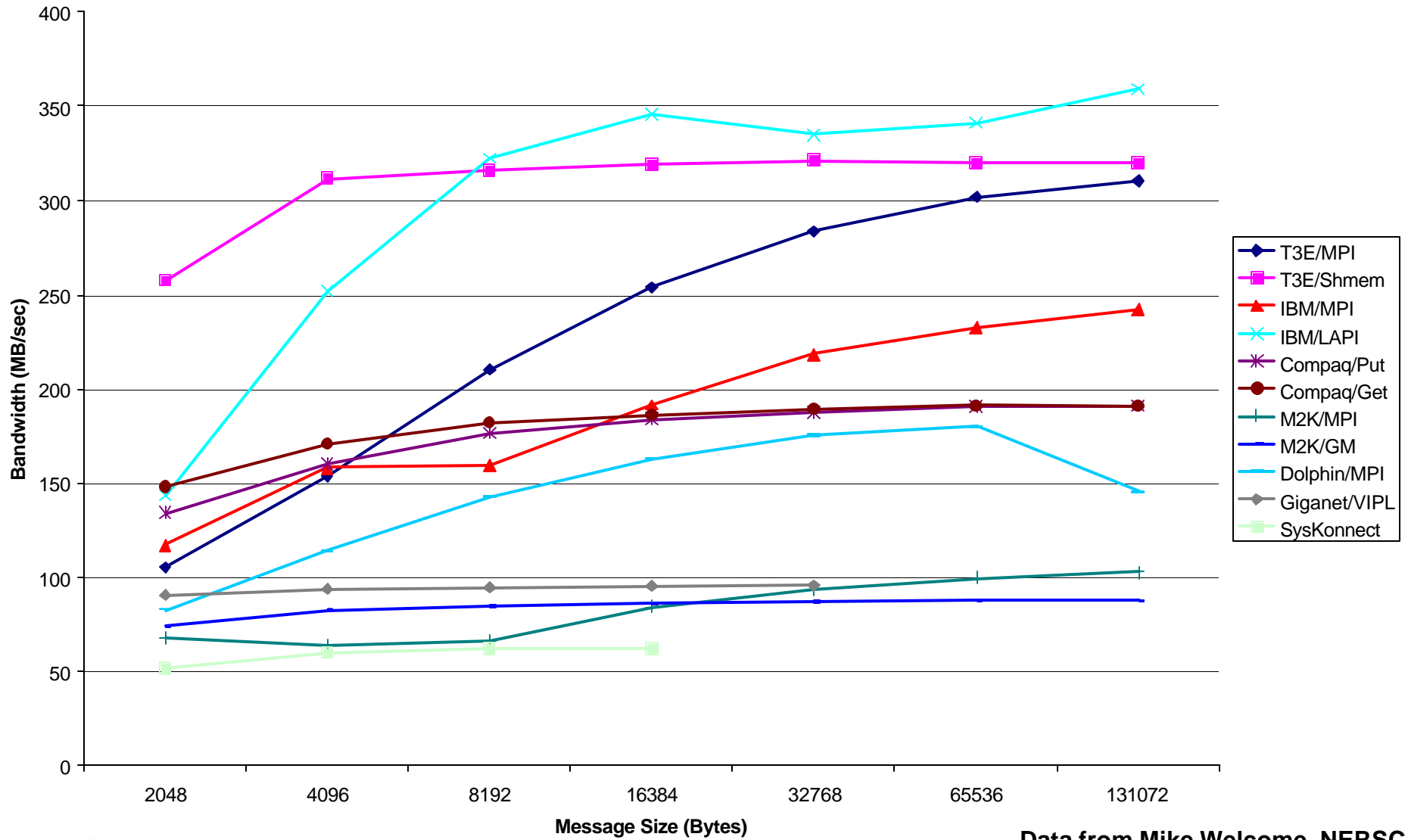


- ✍ Latency has not improved significantly
  - T3E (shmem) was lowest point
  - Federation in 2003 will not reach that level – 7 years later!

Data from Kathy Yelick, UCB and NERSC



# Bandwidth Chart

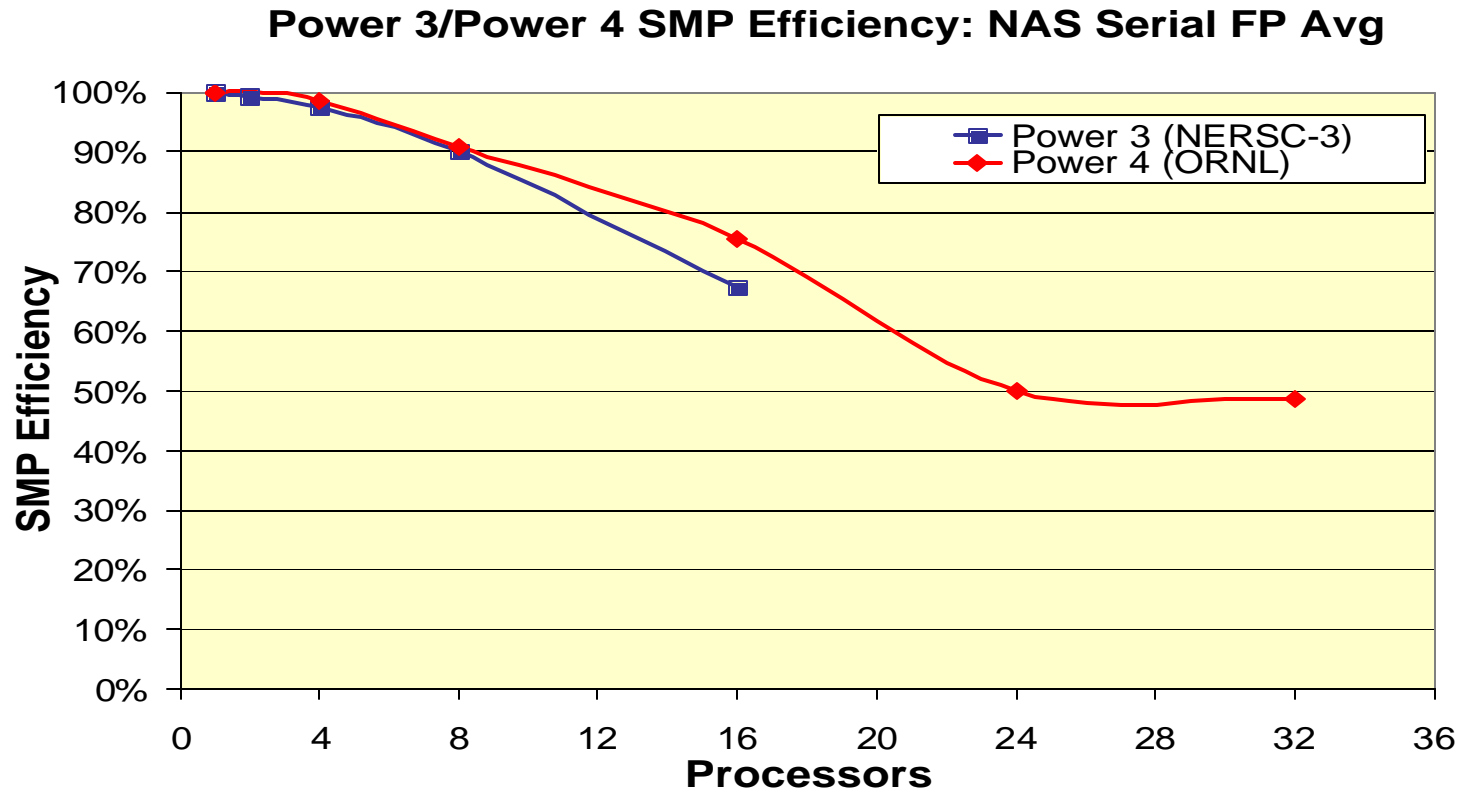


Data from Mike Welcome, NERSC



# Power 4 does not perform as well as expected

Power 4 memory bandwidth does not support 32 CPUs, and Power 4 Memory Latency is only 29% long than Power 3.



# Power 4 in the NERSC Applications Benchmark

---

- ✍ The NERSC – 3 base system delivers  
**618 Gflop/s** on NERSC SSP
- ✍ We measured **179 Gflop/s** on the 4.5 Tflop/s peak Power 4 system at ORNL
- ✍ Assume a Power 4 system with same base cost as NERSC-3:
  - Available to NERSC users only in mid to late 2004
  - Only a 7% performance improvement 3 years after NERSC-3
- ✍ The performance of Power 4 is a clear indication of the **DIVERGENCE PROBLEM**
- ✍ Power 4 was not designed for scientific applications



# We have pursued the logical extreme of the “commodity parts” path.



This

Low cost path



Became



Clusters of Symmetric Multiprocessors:  
Ensembles of Data Servers + Fast Switch

- ✍ The commodity building block was the microprocessor but is now the entire server (SMP).
- ✍ Communications and memory bandwidth are not scaling with processor power.
- ✍ We have arrived at near football-field size computers consuming megawatts of electricity.

# The Divergence Problem

---

- ✍ The requirements of high performance computing for science and engineering, and the requirements of the commercial market are diverging.
- ✍ The commercial cluster of SMP approach is no longer sufficient to provide the highest level of performance
  - Lack of memory bandwidth and latency
  - High interconnect latency
  - Lack of interconnect bandwidth
  - High cost of ownership for large scale systems
- ✍ U.S. computer industry is driven by commercial applications -- not focused on scientific computing.
- ✍ The decision for NERSC-3 E can be seen as a first indication of the divergence problem: Power 4 had a low SSP number





# The State of the American Computer Industry – In Scientific Computing

---

- ✍ The major players that are still active in scientific supercomputing are
  - ✍ IBM
  - ✍ Hewlett Packard
  - ✍ Cray (a small surviving and evolved portion)
  - ✍ Sun
  - ✍ SGI
- ✍ We don't have a building block optimized for scientific computation.
- ✍ The target commercial market is data and web serving, and that market dominates the economics of the computer industry beyond the personal computer.
- ✍ The architectural barriers for scientific computing stem from this situation
  - Memory bandwidth and latency (optimized for databases)
  - Interconnect bandwidth and latency (optimized for transaction processing)
- ✍ If you don't have a viable market for those building blocks, then how do you cause them to be created?



# Gone, But Not Forgotten: Evidence of Enormous Creativity in Computing in the U.S.

- ✍ ACRI
- ✍ Alliant
- ✍ American Supercomputer
- ✍ Ametek
- ✍ Applied Dynamics
- ✍ Astronautics
- ✍ BBN
- ✍ CDC
- ✍ Convex
- ✍ Cray Computer
- ✍ Cray Research
- ✍ Culler-Harris
- ✍ Culler Scientific
- ✍ Cydrome
- ✍ Dana/Ardent/Stellar/Stardent
- ✍ Denelcor
- ✍ Elexsi
- ✍ ETA Systems
- ✍ Evans and Sutherland Computer
- ✍ Floating Point Systems
- ✍ Galaxy YH-1
- ✍ Goodyear Aerospace MPP
- ✍ Gould NPL
- ✍ Guiltech
- ✍ Intel Scientific Computers
- ✍ International Parallel Machines
- ✍ Kendall Square Research
- ✍ Key Computer Laboratories
- ✍ MasPar
- ✍ Meiko
- ✍ Multiflow
- ✍ Myrias
- ✍ Numerix
- ✍ Prisma
- ✍ Tera
- ✍ Thinking Machines
- ✍ Saxpy
- ✍ Scientific Computer Systems (SCS)
- ✍ Soviet Supercomputers
- ✍ Supertek
- ✍ Supercomputer Systems
- ✍ Suprenum
- ✍ Vitesse Electronics

## But this is not 1990

---

- ✍ Starting a number of new small companies seeded by federal research investment (as DARPA did in the HPCCI) is probably not feasible .
- ✍ There is now a much larger commercial market, and the industry dynamics are different.
- ✍ The Earth Simulator “event” has motivated IBM and others to better address the needs of the scientific community.
- ✍ There is still a significant scientific market for high performance computing outside of supercomputer centers.
- ✍ For this new environment, we need a new, sustainable strategy for the future of scientific computing.



# Outline

---

- ✍ Introducing NERSC-3 E
- ✍ The Divergence Problem
- ✍ What NERSC is doing about it



# Need for a Sustainable Effort

---

- ✍ Without a sustained effort, scientific communities cannot invest their efforts and resources to adapt their computing strategy to new classes of hardware.
- ✍ Parallel computing itself required a decade to find scalable algorithms to make it useful, and the process is still continuing.
- ✍ The U.S. policy should not be to create one machine just to show we can do it, but should be a long-term program that ensures preeminence in scientific computing.
- ✍ The most powerful of these systems need to be available to the open, scientific community (in addition to any special communities)



# Why Does Cost Matter?

---

## **If this is so important, why does cost matter?**

- ✍ If effective scientific supercomputing is only available at high cost, it will have impact on only a small part of the scientific community.
- ✍ So, need to leverage the resources of mainstream IT companies like IBM, HP and Intel as well as any special architecture companies like Cray.
- ✍ And the national science policy should motivate them to participate durably.



# Creating a New Class of Computer Architectures for Scientific Computing

---

- ✍ Sustained cooperative development of new computer architectures
- ✍ A focus on sustained performance of scientific applications – not on peak performance!
- ✍ Addressing the key bottlenecks of bandwidth and latency for memory and processor interconnection
- ✍ A strategy to pursue several architectures at multiple sites
- ✍ A new investment in the computer science research and scientific research communities



# A New Architecture Strategy: Beyond Evaluation to Cooperative Development

---

**A proposal to establish feedback between science and computer design lasting for generations of machines**

- ✍ Application teams to drive the design of new architectures
- ✍ Continued, simultaneous evaluation of multiple scientific applications replacing “rules of thumb” for computer designers
  - Example is the Performance Evaluation Research Center (PERC)
- ✍ Leveraging current components and research prototypes into new architectures
- ✍ Continual redesign and testing of prototypes in a vendor partnership to create new scientific computers
- ✍ Addressing the scientific market beyond lab and academic supercomputer centers





# Cooperative Development – NERSC/ANL/IBM Workshop

---



- Held two joint workshops
  - Sept 2002 – defining the Blue Planet architecture
  - Nov. 2002 – IBM gathered input for Power 6
- Developed White Paper "Creating Science-Driven Computer Architecture: A New Path to Scientific Leadership," available at <http://www.nersc.gov/news/blueplanet.html>

# Selection is Based on Scientific Applications

	AMR	Coupled Climate	Astrophysics		Nanoscience	
			MADCAP	Cactus	FLAPW	LSMS
Sensitive to global bisection	<b>X</b>	<b>X</b>	<b>X</b>		<b>X</b>	
Sensitive to processor to memory latency	<b>X</b>	<b>X</b>			<b>X</b>	
Sensitive to network latency	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	
Sensitive to point to point communications	<b>X</b>	<b>X</b>				<b>X</b>
Sensitive to OS interference in frequent barriers				<b>X</b>	<b>X</b>	
Benefits from deep CPU pipelining	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>
Benefits from Large SMP nodes	<b>X</b>					

# “Blue Planet”: Extending IBM Power Technology and “Virtual Vector” Processing

---

Addressing the key barriers to effective scientific computing

- Memory bandwidth and latency
- Interconnect bandwidth and latency
- Programmability for scientific applications

✍ The Strategy is to get back “inside the box” of commercial servers (SMPs)

- Increasing memory and switch bandwidth using commercial parts available over the the next two years

✍ Exploration of new architectures with the IBM design team

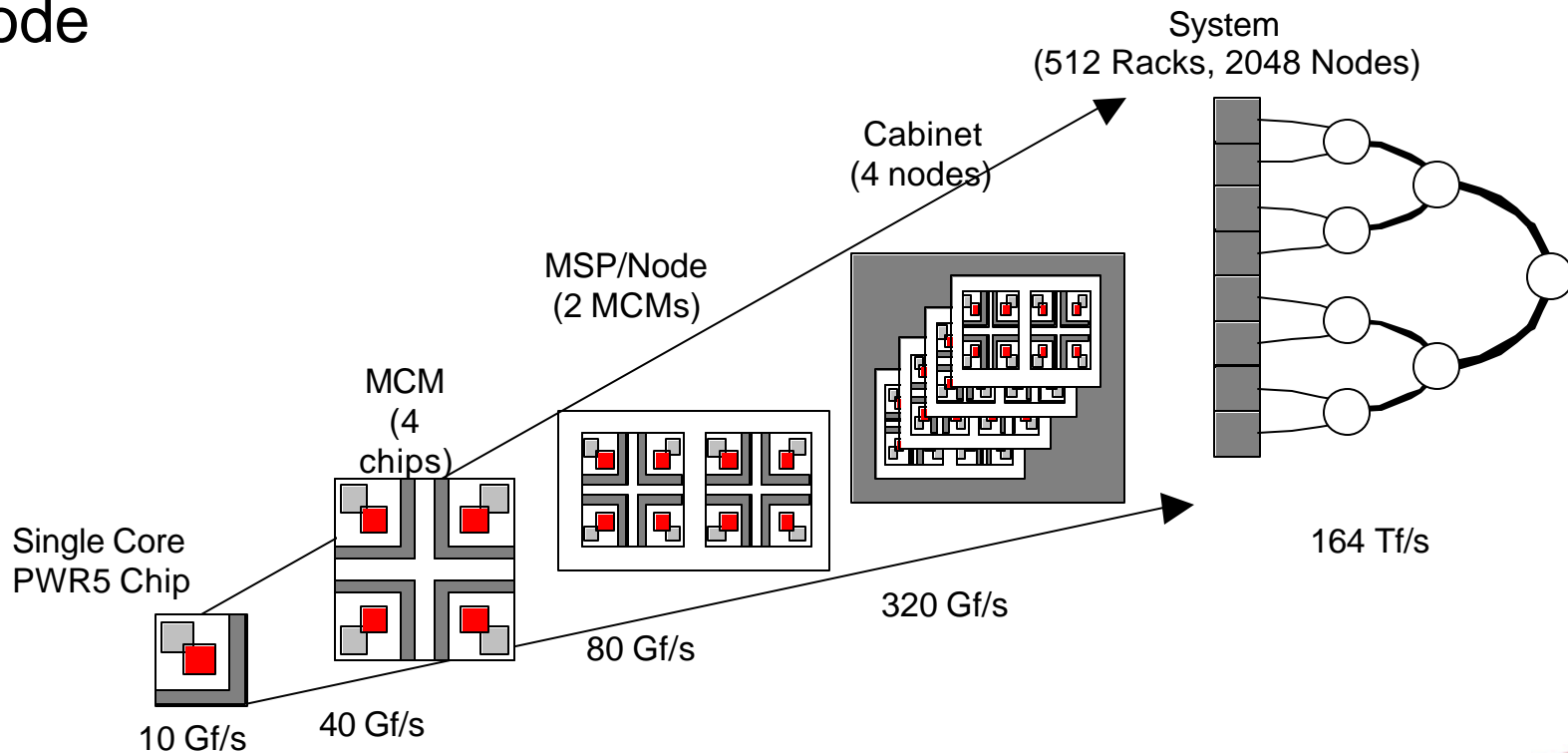
✍ Enabling the vector programming model inside a Power 5 SMP node

✍ Changing the design of subsequent generations of microprocessors



# Blue Planet: A Conceptual View

- ✍ Increasing memory bandwidth – single core chips with dedicated caches for 8 way nodes
- ✍ Increasing switch bandwidth and decreasing latency
- ✍ Enabling “vector” programming model inside each SMP node



# Why this is not Business as Usual for IBM

---

- ✍ Introducing 8 way Power 5+ nodes with single cores early is entirely new packaging
  - For power 4, 8 way nodes came out 18 months after full size SMP (32 CPUs)
  - Each CPU will have its own L1, L2 and L3 cache
  - Each node will have twice the number of memory buses as standard nodes
  - 8 way nodes will run at full clock rate (as opposed to the slower dual core 8 way nodes soon to be introduced).
- ✍ Synchronizing CPUs (“Virtual Vectors) is not in their plan
  - Both hardware and compiler technology involved
- ✍ An additional stage (level) in the Federation switch is not in their plan
  - Increases a factor of 4 in number of links.
- ✍ Decreasing switch latency is not in their plan
  - Requires a radical redesign of their software stack
- ✍ Operating System, Compiler, Library and Scalability Improvements



# Managing Long-Term Architecture Development

---

- ✍ DOE Lab system is ideally suited to manage large-scale, long-term research and development
- ✍ We believe that long-term participation from the universities is critical to the success of this proposed initiative
- ✍ We need to engage architects, scientists, computer scientists in a way that is accountable to one agency
  - And to do that over multiple generations
  - And with multiple vendors
- ✍ These have to be run as closed-loop integrated projects
- ✍ We need to avoid the past failure modes of interagency development



# Conclusion

---

- ✍ We have pursued the logical extreme of the “commodity parts” path.
- ✍ This path was a cost-efficient “free ride” on a Moore’s Law growth curve
- ✍ The divergence problem shows that this free ride is coming to an end.
- ✍ Business as usual will not preserve U.S. leadership in advanced scientific computing
- ✍ New computer architectures optimized for scientific computing are critical to enable 21<sup>st</sup> Century Science
- ✍ The HPC center and user community needs to develop these in a new mode of sustainable partnership with the vendors

**U.S. science requires a strategy to create cost-effective, science-driven computer architectures.**

