



Алгоритмического и программного Проблемы при переходе на Exascale

Jack Dongarra

University of Tennessee
Oak Ridge National Laboratory
University of Manchester



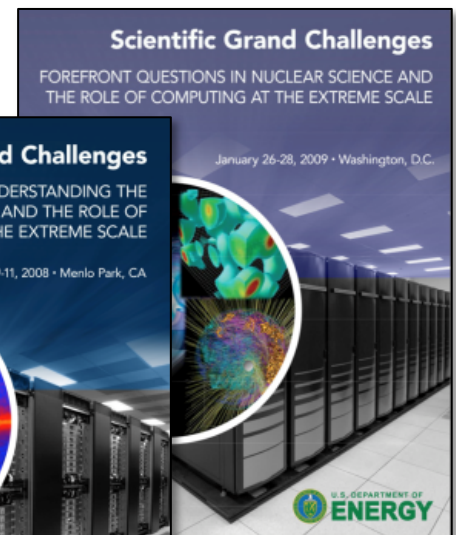
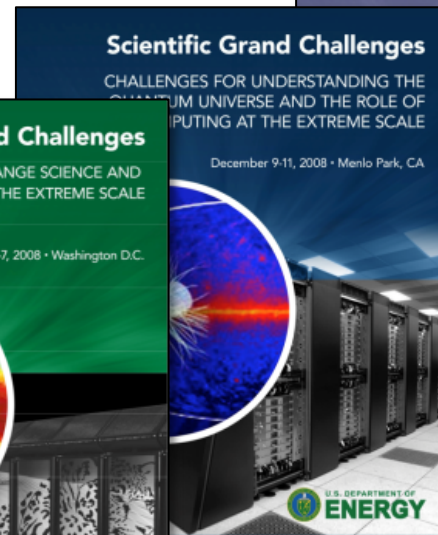
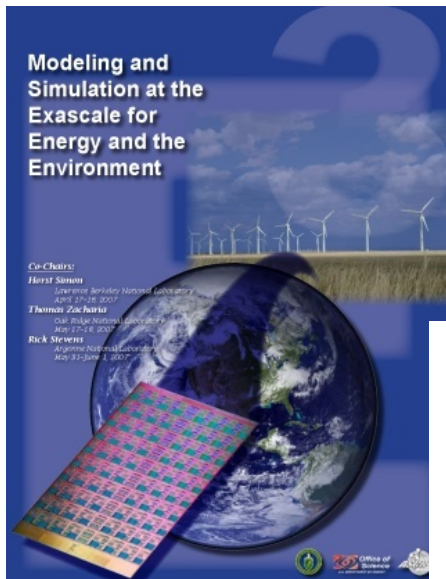
Algorithmic and Software Challenges when Moving Towards Exascale

Jack Dongarra

University of Tennessee
Oak Ridge National Laboratory
University of Manchester

Key Message

- Exascale has been discussed in numerous workshops, conferences, planning meetings for about five years.
- Exascale projects have been started in the US and many other countries and regions.
- Progress has been made, but key challenges to exascale remain.





State of Supercomputing in 2013

- Pflops computing fully established with 26 machines.
- Three technology “swim lanes” or architecture possibilities are thriving.
- Interest in supercomputing is now worldwide, and growing in many new markets (over 50% of Top500 computers are in industry).
- Exascale projects exist in many countries and regions.
- Rapid growth predicted by IDC for the next three years.



June 2013: The TOP10

Rank	Site	Computer	Country	Cores	Rmax [Pflops]	% of Peak	Power [MW]	MFlops /Watt
1	National University of Defense Technology	Tianhe-2 NUDT, Xeon 12C 2.2GHz + Intel Xeon Phi (57c) + Custom	China	3,120,000	33.9	62	18	1902
2	DOE / OS Oak Ridge Nat Lab	Titan, Cray XK7 (16C) + Nvidia Kepler GPU (14c) + Custom	USA	560,640	17.6	65	8.3	2143
3	DOE / NNSA L Livermore Nat Lab	Sequoia, BlueGene/Q (16c) + custom	USA	1,572,864	17.1	85	7.9	2177
4	RIKEN Advanced Inst for Comp Sci	K computer Fujitsu SPARC64 VIIIfx (8c) + Custom	Japan	705,024	10.5	93	13	830
5	DOE / OS Argonne Nat Lab	Mira, BlueGene/Q (16c) + Custom	USA	786,432	8.58	85	3.9	2177
6	Texas Advanced Computing Center	Stampede, Dell Intel (8c) + Intel Xeon Phi (61c) + IB	USA	204,900	5.16	61	4.5	1146
7	Forschungszentrum Juelich (FZJ)	JuQUEEN, BlueGene/Q, Power BQC 16C 1.6GHz+Custom	Germany	458,752	5.01	85	2.3	2177
8	DOE / NNSA L Livermore Nat Lab	Vulcan, BlueGene/Q, Power BQC 16C 1.6GHz+Custom	USA	393,216	4.29	85	2.0	2177
9	Leibniz Rechenzentrum	SuperMUC, Intel (8c) + IB	Germany	147,456	2.90	91*	3.4	846
10	Nat. SuperComputer Center in Tianjin	Tianhe-1A, NUDT Intel (6c) + Nvidia Fermi GPU (14c) + Custom	China	186,368	2.57	55	4.0	635
500	Web Company	HP Cluster	USA	17,904	.096	50		

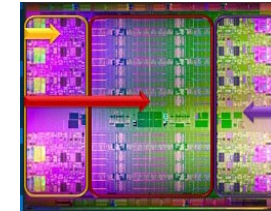


Russian Systems on Top500

Rank	Name	Computer	Site	Manufacturer	Total Cores	Accelerator Cores	Rmax (Tflop/s)	% of Peak	Accelerator/Co-Processor
31	Lomonosov	T-Platforms T-Blade2/1.1, Xeon X5570/X5670/E5630 2.93/2.53 GHz, Nvidia 2070 GPU, PowerXCell 8i Inf QDR	Moscow State University - Research Computing Center	T-Platforms	78,660	29,820	901	53	NVIDIA 2070
72	MVS-10P	RSC Tornado, Xeon E5-2690 8C 2.900GHz, Ind FDR, Intel Xeon Phi SE10X	Joint Supercomputer Center	RSC Group	28,704	25,376	375	72	Intel Xeon Phi
217		Cluster Platform 3000 BL460c Gen8, Xeon E5-2660 8C 2.200GHz, GEnet	IT Services Provider	Hewlett-Packard	18,032	0	160	51	
249	RSC Tornado SUSU	RSC Tornado, Xeon X5680 6C 3.330GHz, Inf QDR, Intel Xeon Phi SE10X	South Ural State University	RSC Group	14,016	11,712	146	62	Intel Xeon Phi
355	MVS-100K	Cluster Platform 3000 BL460c/BL 2x220/SL390, Xeon E5450/5365/X5675 4C 3.000GHz, InF DDR, NVIDIA 2090	Joint Supercomputer Center	Hewlett-Packard	13,004	2432	119	53	NVIDIA 2090
428	Uran	ClusterPlatform SL390s/SL270s, Xeon X5675 6C 3.060GHz, InF QDR, NVIDIA 2090	IMM UrOAN	Hewlett-Packard	5900	5312	105	46	NVIDIA 2090
464		Cluster Platform 3000 BL 2x220, Xeon E5450 4C 3.000GHz, Inf QDR	Kurchatov Institute Moscow	Hewlett-Packard	10,304	0	101	82	
471	SKIF Aurora	SKIF Aurora Platform - Intel Xeon X5680, Infiniband QDR	South Ural State University	RSC Group	8832	0	100	86	

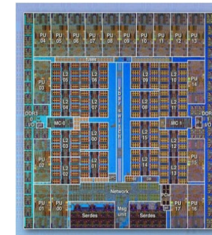
Three Possible Architecture Paths

- **Multicore:** Maintain complex cores, and replicate (x86, SPARC, Power7)
[#4 and 9]



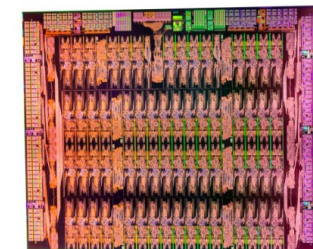
Intel Xeon E7
(10 cores)

- **Manycore/Embedded:** Use many simpler, low power cores from embedded (BlueGene, future ARM)
[#3, 5, 7, and 8]



IBM BlueGene/Q
(16 +2 cores)

- **GPU/Coprocessor/Accelerator:** Use highly specialized processors from graphics market space (NVidia Fermi, Intel Xeon Phi, AMD)
[# 1, 2, 6, and 10]










Intel Xeon Phi
(60 cores)



ICL

P
e
t
a
f
l
o
p
s

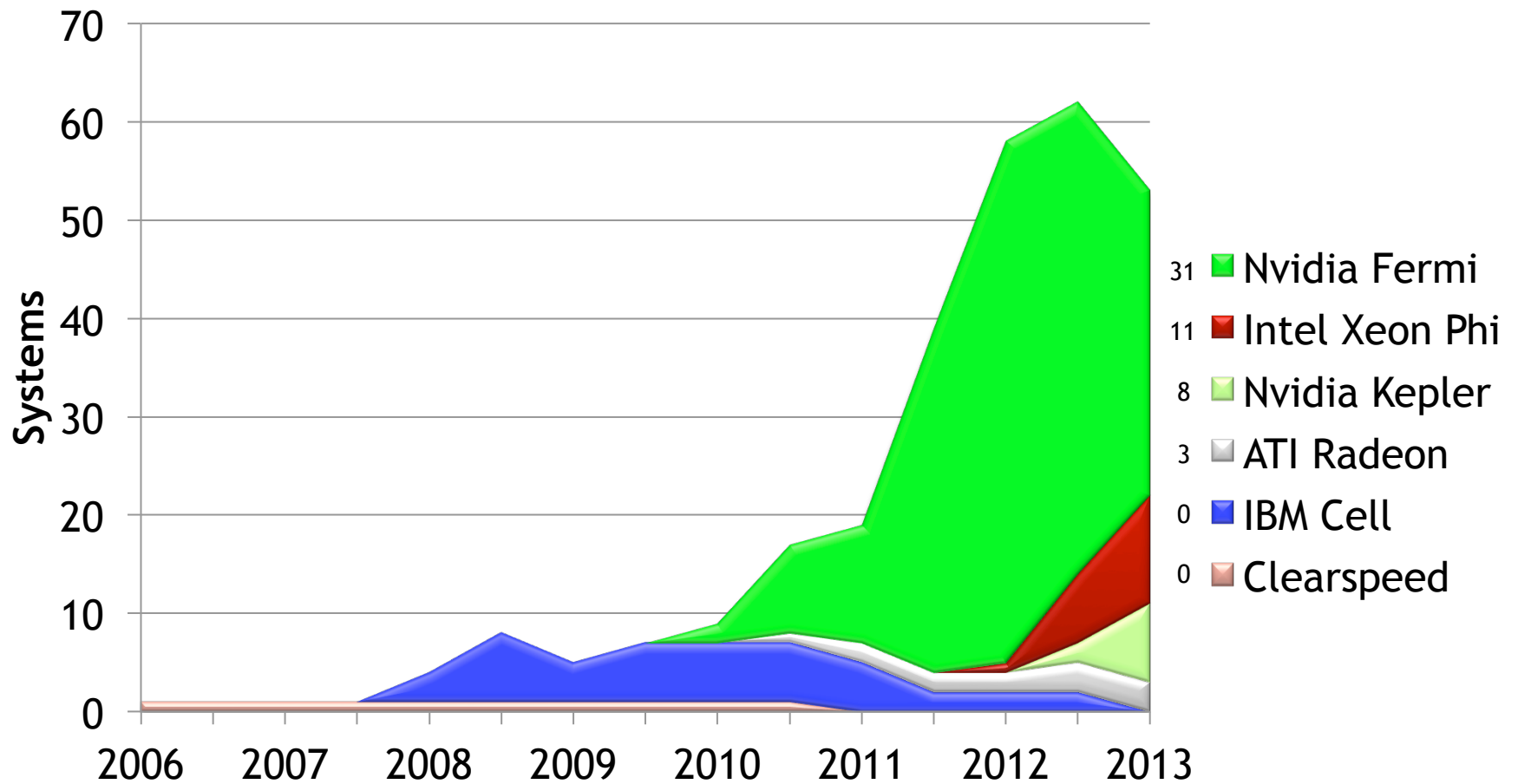
C
l
u
b

Name	Rmax Linpack# Pflops	Country	11  4  2  3  2  3  1 
Tianhe-2 (MilkyWay-2)	33.9	China	NUDT: Hybrid Intel/Intel/Custom
Titan	17.6	US	Cray: Hybrid AMD/Nvidia/Custom
Sequoia	17.2	US	IBM: BG-Q/Custom
K Computer	10.5	Japan	Fujitsu: Sparc/Custom
Mira	8.59	US	IBM: BG-Q/Custom
Stampede	5.17	US	Dell: Hybrid/Intel/Intel/IB
JUQUEEN	5.01	Germany	IBM: BG-Q/Custom
Vulcan	4.29	US	IBM: BG-Q/Custom
SuperMUC	2.90	Germany	IBM: Intel/IB
Tianhe-1A	2.57	China	NUDT: Hybrid Intel/Nvidia/Custom
Pangea	2.10	France	Bull: Intel/IB
Fermi	1.79	Italy	IBM: BG-Q/Custom
DARPA Trial Subset	1.52	US	IBM: Intel/IB
Spirit	1.42	US	SGI: Intel/IB
Curie thin nodes	1.36	France	Bull: Intel/IB
Nebulae	1.27	China	Dawning: Hybrid Intel/Nvidia/IB
Yellowstone	1.26	US	IBM: BG-Q/Custom
Blue Joule	1.25	UK	IBM: BG-Q/Custom
Pleiades	1.24	US	SGI Intel/IB
Helios	1.24	Japan	Bull: Intel/IB
TSUBAME 2.0	1.19	Japan	NEC/HP: Hybrid Intel/Nvidia/IB
Cielo	1.11	US	Cray: AMD/Custom
DiRAC	1.07	K	IBM: BG-Q/Custom
Hopper	1.05	US	Cray: AMD/Custom
Tera-100	1.05	France	Bull: Intel/IB
Oakleaf-FX	1.04	Japan	Fujitsu: Sparc/Custom

6 Hybrid Architectures
8 IBM BG/Q
15 Custom X
11 Infiniband X
9 Look like "clusters"

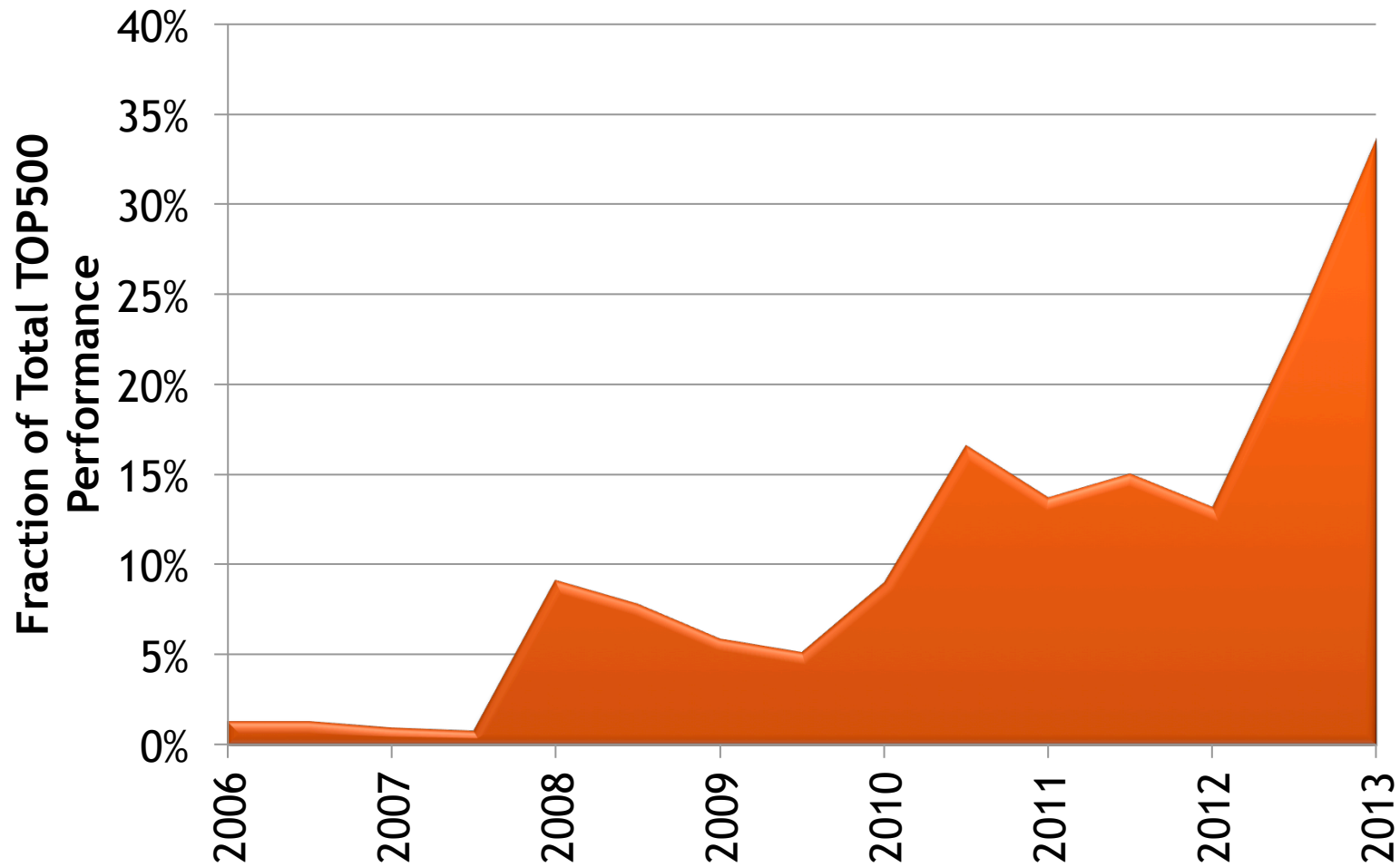


Hybrid/Accelerators (53 Systems)

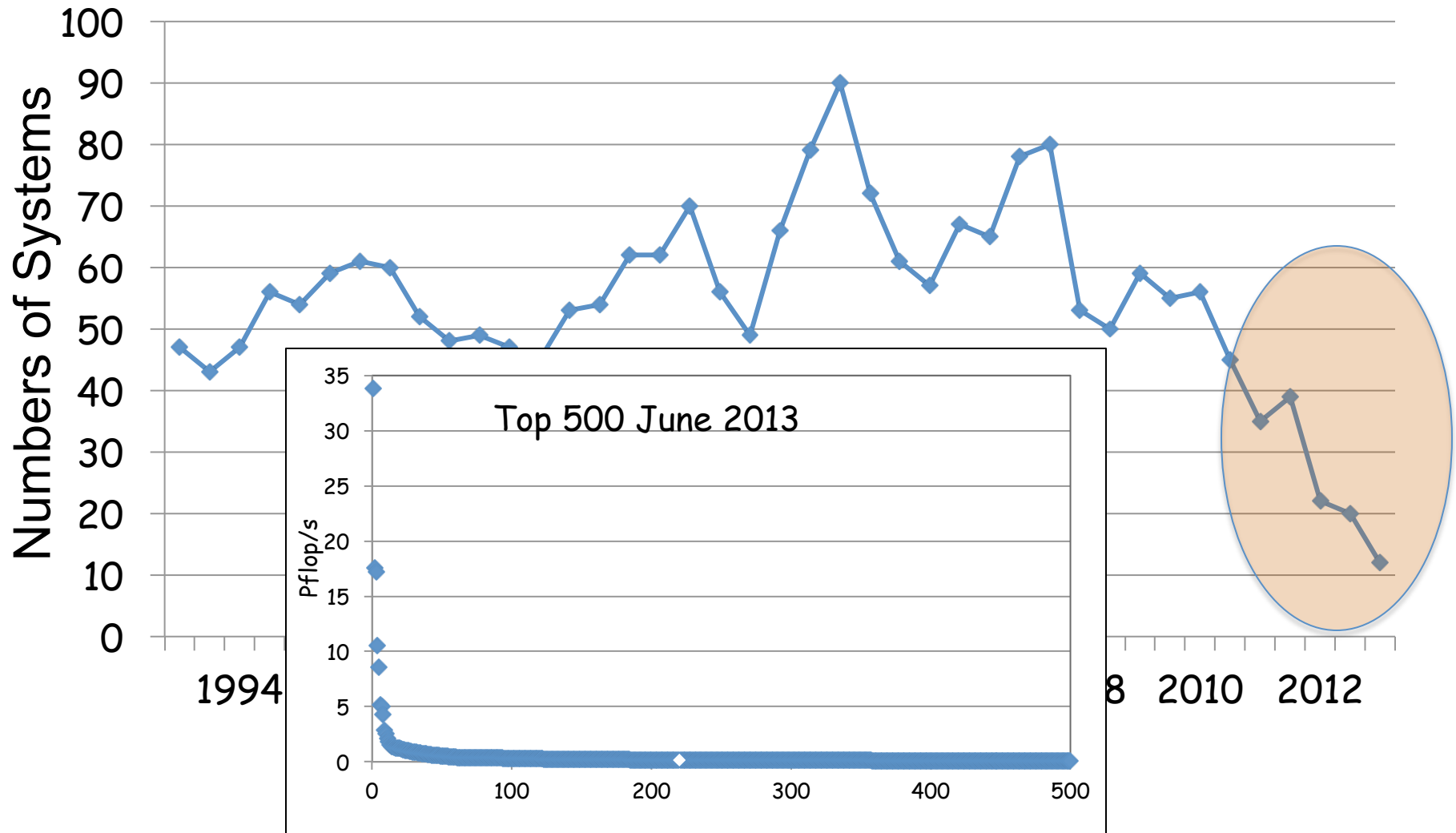




Top500 Performance Share of Accelerators



For the Top 500: Rank at which Half of Total Performance is Accumulated



Commodity plus Accelerator Today

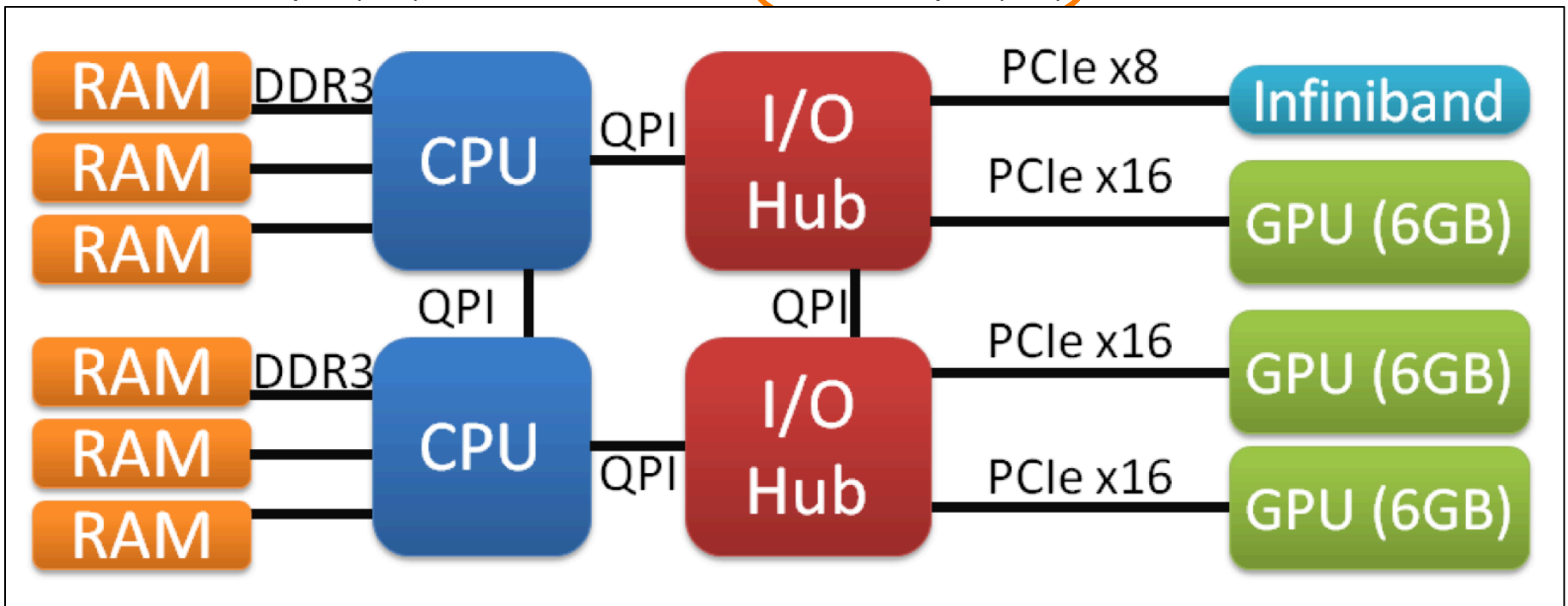
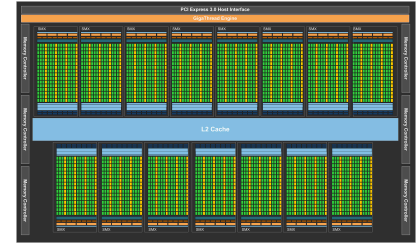
Commodity

Intel Xeon
 8 cores
 3 GHz
 8*4 ops/cycle
 96 Gflop/s (DP)

Accelerator (GPU)

Nvidia K20X "Kepler"
 2688 "Cuda cores"
 .732 GHz
 2688*2/3 ops/cycle
 1.31 Tflop/s (DP)

192 Cuda cores/SMX
 2688 "Cuda cores"



Interconnect
 PCI-X 16 lane
 64 Gb/s (8 GB/s)
 1 GW/s



We Have Seen This Before

- Floating Point Systems FPS-164/MAX Supercomputer (1976)
- Intel Math Co-processor (1980)
- Weitek Math Co-processor (1981)

1976

THREE HUNDRED FORTY ONE MILLION FLOATING POINT OPERATIONS PER SECOND. THE FPS-164/MAX.

Rapid scientific and engineering problems increasingly call for super-computing machines and higher technicality, which leads to calculations on very large numbers. The small size, cost of a super-computer with the speed and accuracy of a super-computer has become a reality.

Now, there's the FPS-164/MAX -- a special purpose, modular supercomputer that matches the size of CRAY CYBER and offers an extremely cost-effective alternative -- at a fraction of the cost.

The FPS-164/MAX is fast.
800 point operations a second FPS-164/MAX can handle 3.5 million floating point operations per second, depending on configuration, according to 70% of all 800 point operations available in the world. But the FPS-164/MAX gives you all the speed and accuracy you need to solve those really complicated engineering problems.

The FPS-164/MAX configuration is able to incorporate up to 24 vector channels at one time, allowing a fully configured FPS-164/MAX to factor a 1,000 by 1,000 matrix in about 2 seconds, complete two 10,000 by 10,000 matrices in two days.

The FPS-164/MAX is powerful.
A parallel pipelined processor designed to run 160,000 or high speed, the FPS-164/MAX has all the vector capability of our original FPS-164. We've just added a lot more power with super-special processing units which amplify the vector processing capability of the original FPS-164 by up to 10 times.

The FPS-164/MAX is cost-effective.
Its modular design, cost-effective operation and design, built from open, microprocessor modules, is any application requiring the handling of large matrices. The FPS-164/MAX offers unparalleled cost efficiency. In fact, it's the only supercomputer that can be built or broken down supercomputer cost less than \$1 million.

Whether you're looking to upgrade your existing FPS-164 -- or something for a completely new system -- you need the supercomputer performance for one million dollars or less.

What's more, the FPS-164/MAX is designed for the immediate installation of Floating Point Systems. We'll have a service office available to you. All major diagnostic capabilities, and a standard of product quality and reliability second to none, you can be sure the FPS-164/MAX will be up, running, and ready to solve your problem solving needs.

For complete information and application, call toll free 1-800-567-8143.

Floating Point Systems, Inc.
P.O. Box 20480
Palo Alto, CA 94303
(415) 324-3533
U.S. MAIL TO: FLOATING POINT SYSTEMS

1980

The Intel® Math CoProcessor™ is for crunching numbers faster.

intel
Personal Computer Enhancement

There's one for every machine.

80387™ Family, for 80386™ based machines.

80287™ Family, for 80286™ based machines.

8087™ Family, for 8086™ and 8088™ based machines.

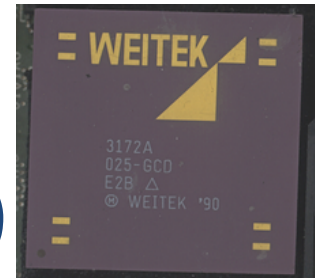
It's FAST!
The Intel Math CoProcessor dramatically speeds up the number crunching that's part of the work you do every day: budgeting, statistical analysis, financial analysis, CAD and other engineering analysis. In fact, the Math CoProcessor is supported by more than 100 commonly used software packages including Lotus 1-2-3, dBase IV, AutoCAD, and most language and statistical packages.

It's EASY!
Intel makes a variety of math coprocessors. Every PC has a built-in socket. Just plug it in and go.

It's SAFE!
Made by Intel, the same people who designed your PC's microprocessor, each and every Math CoProcessor is backed by an industry leading the way warranty and full free technical support. You are assured the highest degree of quality, compatibility, reliability and support for your investment.

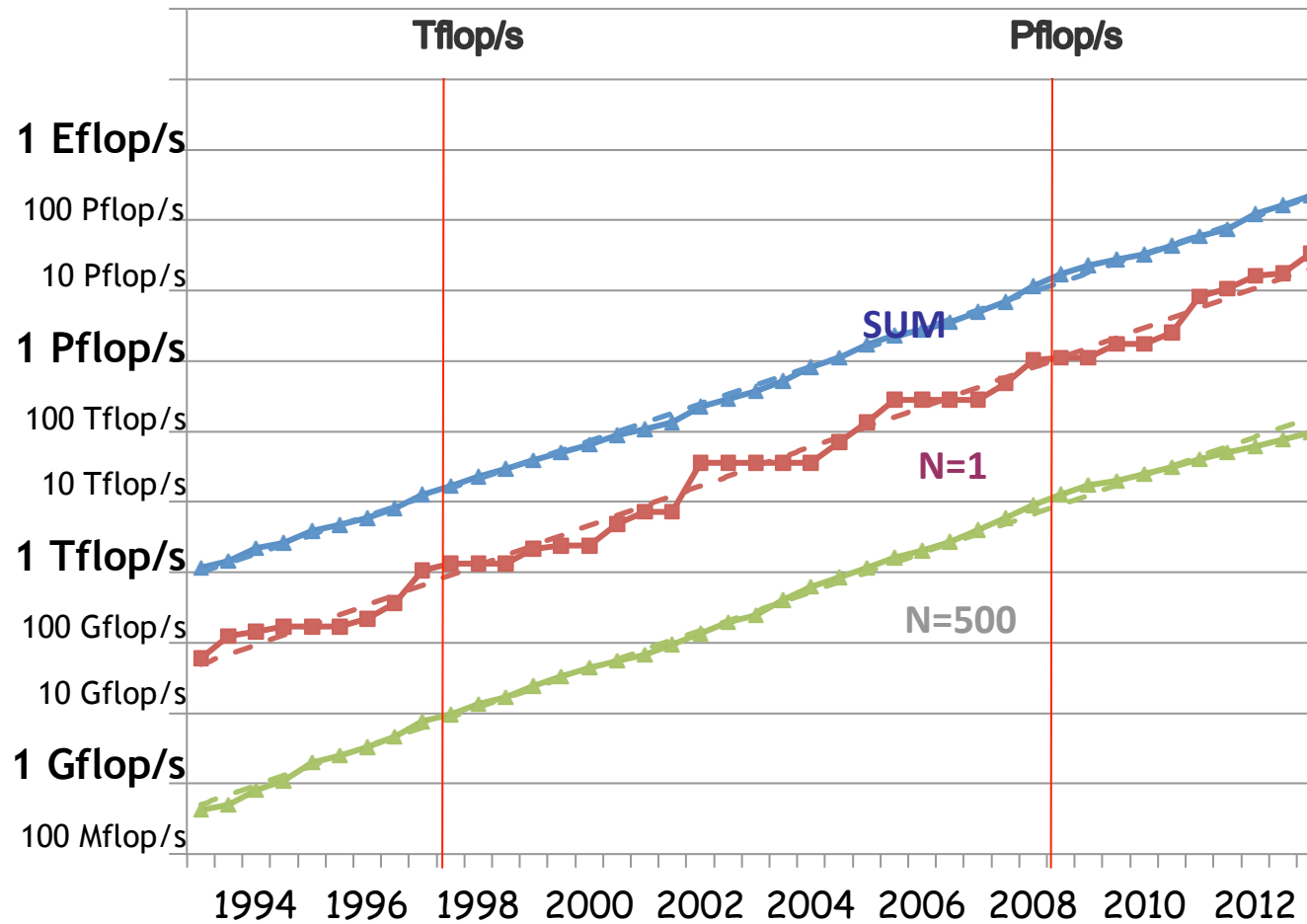
For more information, or technical support call:
(800) 538-3173 in the U.S. and Canada
(510) 652-7154 for International

intel
Personal Computer Enhancement



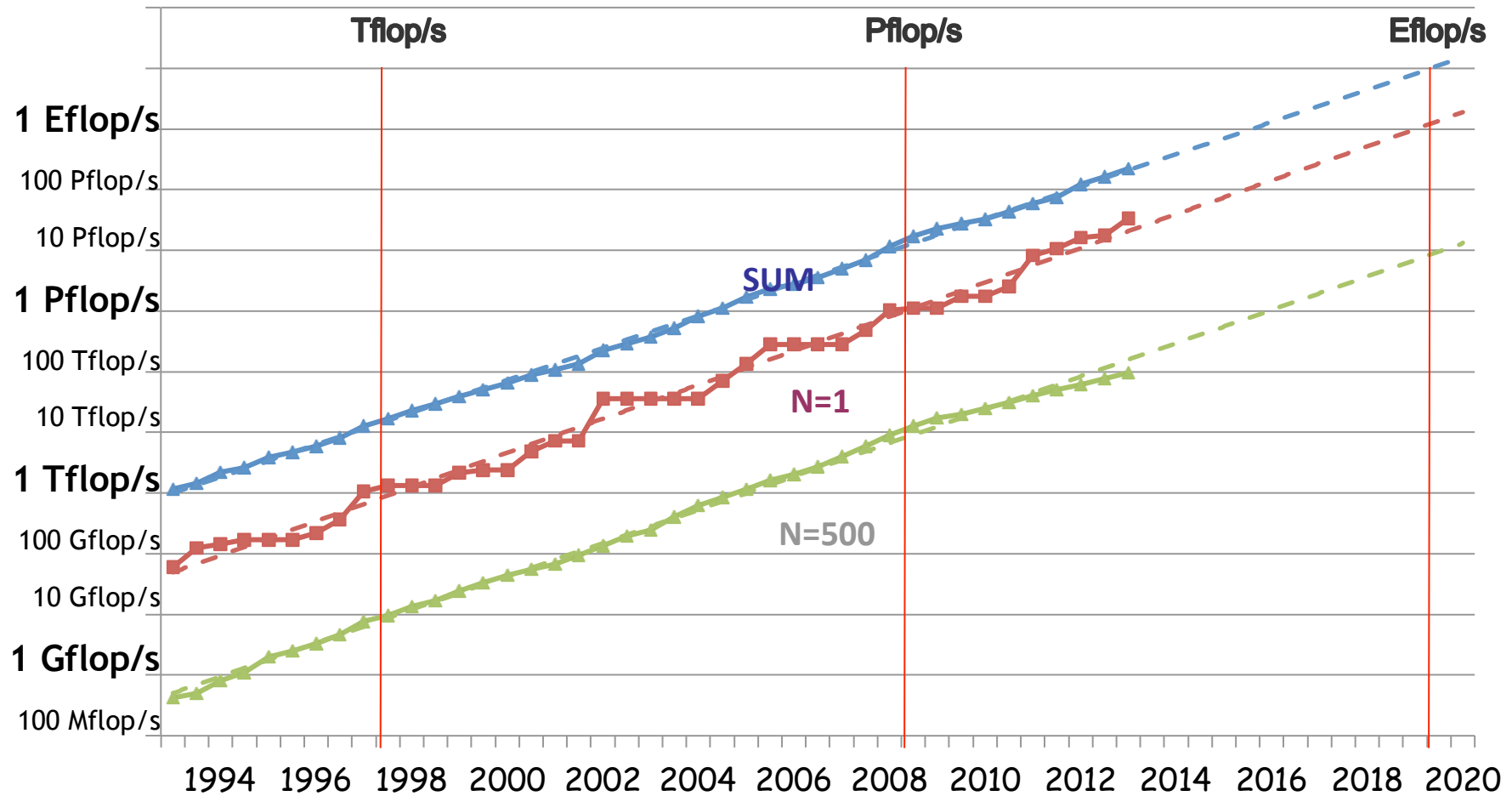


TOP500 Editions (41 so far, 20 years)





TOP500 Editions (53 edition, 26 years)





Today's #1 System

Systems	2013 Tianhe-2
System peak	55 Pflop/s
Power	18 MW (3 Gflops/W)
System memory	1.4 PB (1.024 PB CPU + .384 PB CoP)
Node performance	3.43 TF/s (.4 CPU +3 CoP)
Node concurrency	24 cores CPU + 171 cores CoP
Node Interconnect BW	6.36 GB/s
System size (nodes)	16,000
Total concurrency	3.12 M 12.48M threads (4/core)
MTTF	Few / day



Exascale System Architecture with a cap of \$200M and 20MW

Systems	2013 Tianhe-2
System peak	55 Pflop/s
Power	18 MW (3 Gflops/W)
System memory	1.4 PB (1.024 PB CPU + .384 PB CoP)
Node performance	3.43 TF/s (.4 CPU +3 CoP)
Node concurrency	24 cores CPU + 171 cores CoP
Node Interconnect BW	6.36 GB/s
System size (nodes)	16,000
Total concurrency	3.12 M 12.48M threads (4/core)
MTTF	Few / day



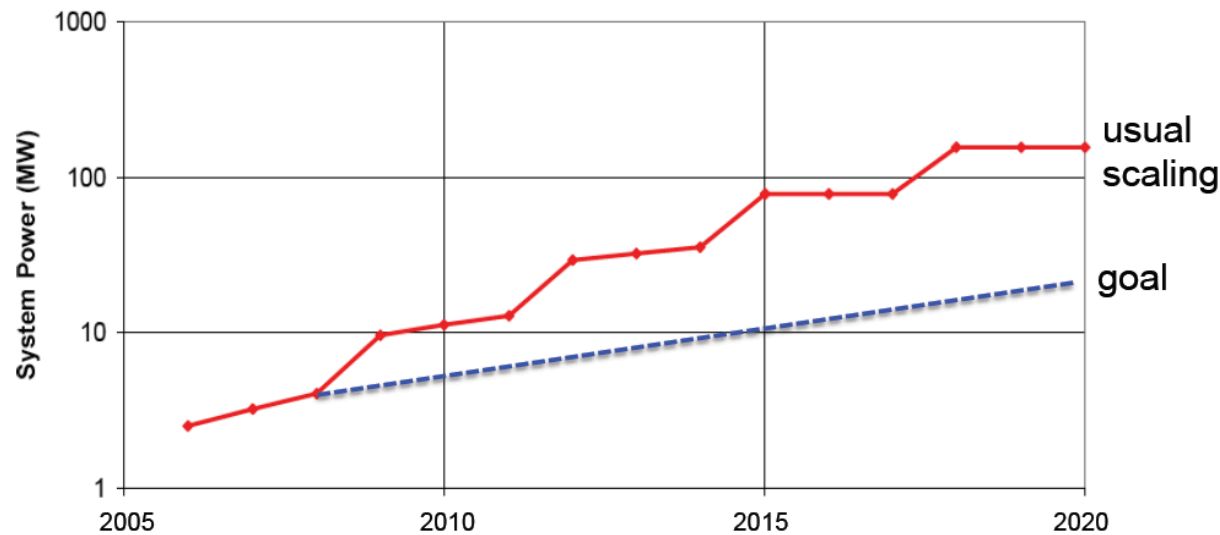
Exascale System Architecture with a cap of \$200M and 20MW

Systems	2013 Tianhe-2	2020-2022	Difference Today & Exa
System peak	55 Pflop/s	1 Eflop/s	~20x
Power	18 MW (3 Gflops/W)	~20 MW (50 Gflops/W)	O(1) ~15x
System memory	1.4 PB (1.024 PB CPU + .384 PB CoP)	32 - 64 PB	~50x
Node performance	3.43 TF/s (.4 CPU +3 CoP)	1.2 or 15TF/s	O(1)
Node concurrency	24 cores CPU + 171 cores CoP	O(1k) or 10k	~5x - ~50x
Node Interconnect BW	6.36 GB/s	200-400GB/s	~40x
System size (nodes)	16,000	O(100,000) or O(1M)	~6x - ~60x
Total concurrency	3.12 M 12.48M threads (4/core)	O(billion)	~100x
MTTF	Few / day	Many / day	O(?)

Energy Cost Challenge

At ~\$1M per MW energy costs are substantial

- 10 Pflop/s in 2011 uses ~10 MWs
- 1 Eflop/s in 2020 > 100 MWs



- DOE Target: 1 Eflop/s around 2020-2022 at 20 MWs

The High Cost of Data Movement

- Flop/s or percentage of peak flop/s become much less relevant

Approximate power costs (in picoJoules)

	2011
DP FMADD flop	100 pJ
DP DRAM read	4800 pJ
Local Interconnect	7500 pJ
Cross System	9000 pJ

Source: John Shalf, LBNL

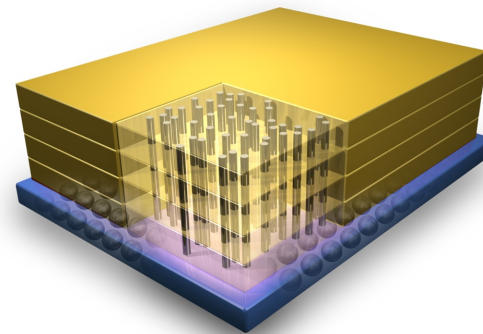
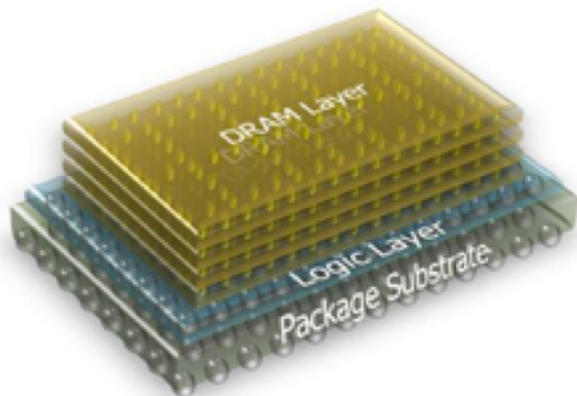
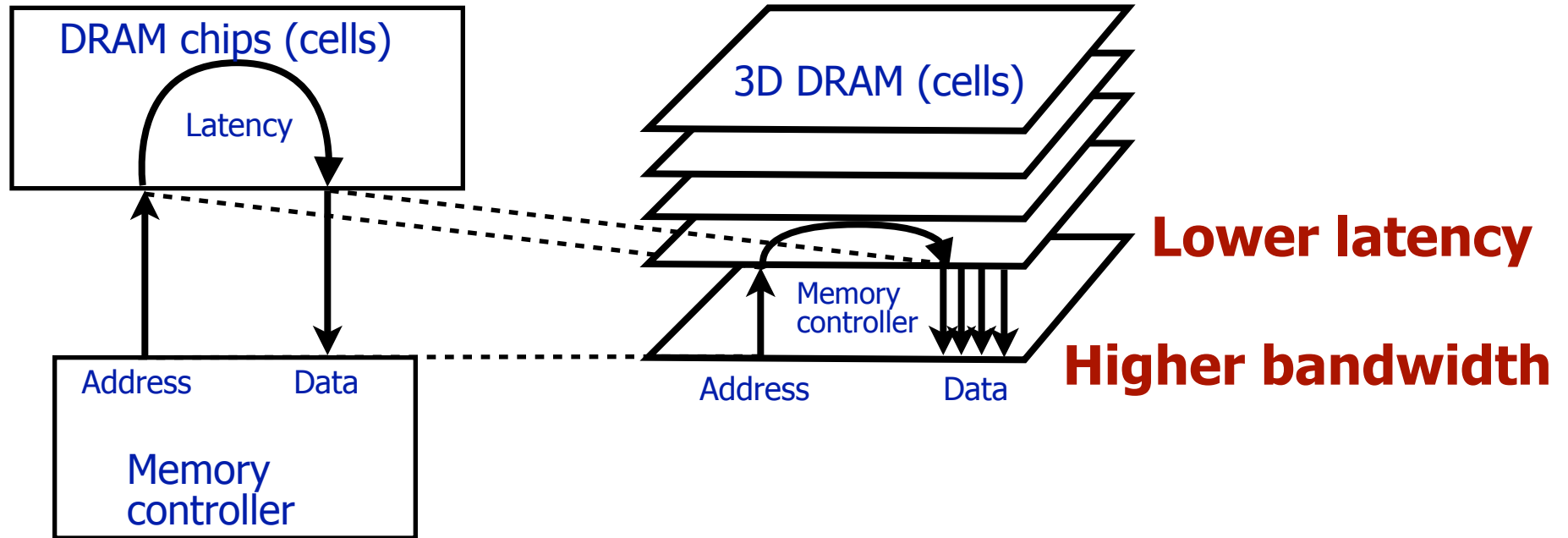
- Algorithms & Software: minimize data movement; perform more work per unit data movement.

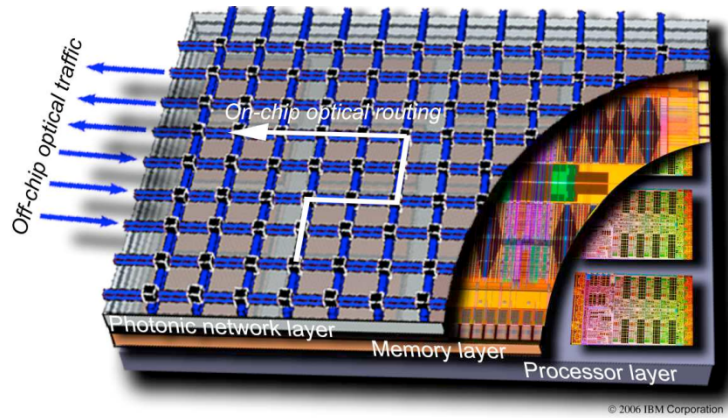
Evolution Over the Last 30 Years

- Initially, commodity PCs where decentralized systems
- As chip manufacturing process shrank to less than a micron, they started to integrate features on-die:
 - 1989: FPU (Intel 80486DX)
 - 1999: SRAM (Intel Pentium III)
 - 2009: GPU (AMD Fusion)
 - 2016: DRAM on chip (3D stacking)

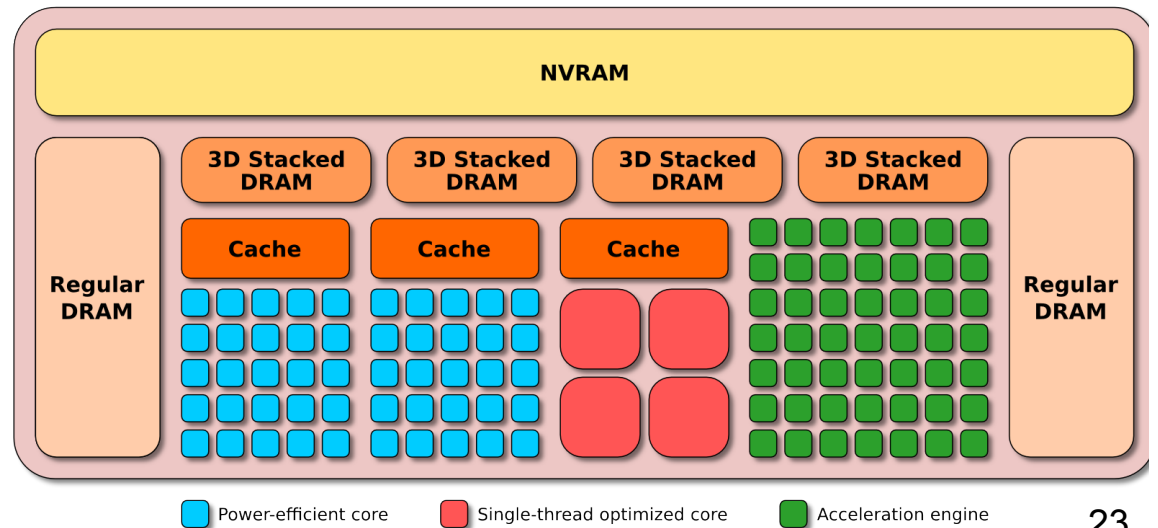


Future Systems May Be Composed of Different Kinds of Cores





Node





Critical Issues at Peta & Exascale for Algorithm and Software Design

- .. **Synchronization-reducing algorithms**
 - **Break Fork-Join model**
- .. **Communication-reducing algorithms**
 - **Use methods which have lower bound on communication**
 - **Cache aware**
- .. **Mixed precision methods**
 - **2x speed of ops and 2x speed for data movement**
- .. **Autotuning**
 - **Today's machines are too complicated, build "smarts" into software to adapt to the hardware**
- .. **Fault resilient algorithms**
 - **Implement algorithms that can recover from failures/bit flips**
- .. **Reproducibility of results**
 - **Today we can't guarantee this. We understand the issues, but some of our "colleagues" have a hard time with this.**



Summary

- .. **Major Challenges are ahead for extreme computing**
 - **Parallelism $O(10^9)$**
 - Programming issues
 - **Hybrid**
 - Peak and HPL may be very misleading
 - No where near close to peak for most apps
 - **Fault Tolerance**
 - Today Sequoia BG/Q node failure rate is 1.25 failures/day
 - **Power**
 - 50 Gflops/w (today at 2 Gflops/w)

- .. **We will need completely new approaches and technologies to reach the Exascale level**



Collaborators / Software / Support

- ◆ **PLASMA**
<http://icl.cs.utk.edu/plasma/>
- ◆ **MAGMA**
<http://icl.cs.utk.edu/magma/>
- ◆ **Quark (RT for Shared Memory)**
<http://icl.cs.utk.edu/quark/>
- ◆ **PaRSEC**(Parallel Runtime Scheduling and Execution Control)
<http://icl.cs.utk.edu/parsec/>



- ◆ Collaborating partners
University of Tennessee, Knoxville
University of California, Berkeley
University of Colorado, Denver

MAGMA



PLASMA

